

Optimal Complexity in Decentralized Training

Yucheng Lu^{*} and Christopher De Sa[†]

Department of Computer Science, Cornell University

Abstract

Decentralization is a promising method of scaling up parallel machine learning systems. In this paper, we provide a tight lower bound on the iteration complexity for such methods in a stochastic non-convex setting. Our lower bound reveals a theoretical gap in known convergence rates of many existing decentralized training algorithms, such as D-PSGD. We prove by construction this lower bound is tight and achievable. Motivated by our insights, we further propose DeTAG, a practical gossip-style decentralized algorithm that achieves the lower bound with only a logarithm gap. Empirically, we compare DeTAG with other decentralized algorithms on image classification tasks, and we show DeTAG enjoys faster convergence compared to baselines, especially on unshuffled data and in sparse networks.

1 Introduction

Parallelism is a ubiquitous method to accelerate model training [1, 2, 3, 4]. A parallel learning system usually consists of three layers (Table 1): an application to solve, a communication protocol deciding how parallel workers coordinate, and a network topology determining how workers are connected. Traditional design for these layers usually follows a centralized setup: in the application layer, training data is required to be shuffled and shared among parallel workers; while in the protocol and network layers, workers either communicate via a fault-tolerant single central node (e.g. Parameter Server) [5, 6, 7] or a fully-connected topology (e.g. AllReduce) [8, 9]. This centralized design limits the scalability of learning systems in two aspects. First, in many scenarios, such as Federated Learning [10, 11] and Internet of Things (IoT) [12], a shuffled dataset or a complete (bipartite) communication graph is not possible or affordable to obtain. Second, a centralized communication protocol can significantly slow down the training, especially with a low-bandwidth or high-latency network [13, 14, 15].

Table 1: Design choice of centralization and decentralization in different layers of a parallel machine learning system. The protocol specifies how workers communicate. The topology refers to the overlay network that logically connects all the workers.

| Layer | Centralized | Decentralized |
|---------------------|---|---|
| Application | Shuffled Data | Unshuffled Data (Federated Learning) |
| Protocol | AllReduce/AllGather Parameter Server | Gossip |
| Network Topology | Complete- (Bipartite) Graph | Arbitrary Graph |

^{*}Corresponds to: yl2967@cornell.edu

[†]Corresponds to: cdesa@cs.cornell.edu

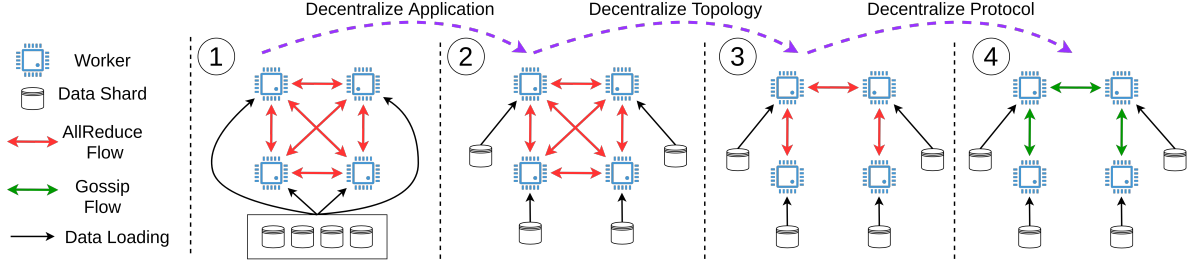


Figure 1: Figure illustrating how decentralization in different layers lead to different learning systems. From left to right: ①: A fully centralized system where workers sample from shared and shuffled data; ②: Based on ①, workers maintain their own data sources, making it decentralized in the application layer; ③: Based on ②, workers are decentralized in the topology layer; ④: A fully decentralized system in all three layers where the workers communicate via Gossip. Our framework and theory are applicable to all kinds of decentralized learning systems.

The rise of decentralization. To mitigate these limitations, decentralization comes to the rescue. Decentralizing the application and network allows workers to learn with unshuffled local datasets [16] and arbitrary topologies [17, 18]. Furthermore, the decentralized protocol, i.e. Gossip, helps to balance load, and has been shown to outperform centralized protocols in many cases [19, 20, 21, 22].

Understanding decentralization with layers. Many decentralized training designs have been proposed, which can lead to confusion as the term “decentralization” is used inconsistently in the literature. Some works use “decentralized” to refer to approaches that can tolerate non-iid or unshuffled datasets [16], while others use it to mean gossip communication [19], and still others use it to mean a sparse topology graph [23]. To eliminate this ambiguity, we formulate Table 1, which summarizes the different “ways” a system can be decentralized. Note that the choices to decentralize different layers are *independent*, e.g., the centralized protocol AllReduce can still be implemented on a decentralized topology like the Ring graph [23].

The theoretical limits of decentralization. Despite the empirical success, the best convergence rates achievable by decentralized training—and how they interact with different notions of decentralization—remains an open question. Previous works often show complexity of a given decentralized algorithm with respect to the number of iterations T or the number of workers n , ignoring other factors including network topologies, function parameters or data distribution. Although a series of decentralized algorithms have been proposed showing theoretical improvements—such as using variance reduction [24], acceleration [17], or matching [25]—we do not know how close they are to an “optimal” rate or whether further improvement is possible.

In light of this, a natural question is: *What is the optimal complexity in decentralized training? Has it been achieved by any algorithm yet?* Previous works have made initial attempts on this question, by analyzing this theoretical limit in a non-stochastic or (strongly) convex setting [17, 26, 27, 28, 29, 30]. These results provide great heuristics but still leave the central question open, since stochastic methods are usually used in practice and many real-world problems of interest are non-convex (e.g. deep learning). In this paper we give the first full answer to this question: our contributions are as follows.

- In Section 4, we prove the first (to our knowledge) tight lower bound for decentralized training in a stochastic non-convex setting. Our results reveal an asymptotic gap between our lower bound and known convergence rates of existing algorithms.
- In Section 5, we prove our lower bound is tight by exhibiting an algorithm called DeFacto that achieves it—albeit while only being decentralized in the sense of the application and network layers.
- In Section 6, we propose DeTAG, a practical algorithm that achieves the lower bound with only a logarithm gap and that is decentralized in all three layers.
- In Section 7, we experimentally evaluate DeTAG on the CIFAR benchmark and show it converges faster compared to decentralized learning baselines.

Table 2: Complexity comparison among different algorithms in the stochastic non-convex setting on arbitrary graphs. The blue text are the results from this paper. Definitions to all the parameters can be found in Section 3. Other algorithms like EXTRA [69] or MSDA [70] are not comparable since they are designed for (strongly) convex problems. Additionally, Liu and Zhang [71] provides alternative complexity bound for algorithms like D-PSGD which improves upon the spectral gap. However, the new bound would compromise the dependency on ϵ , which does not conflict with our comparison here.

| | Source | Protocol | Sample Complexity | Comm. Complexity | Gap to Lower Bound |
|-------------|---------------------|-----------|---|--|---|
| Lower Bound | Theorem 1 | Central | $\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $\Omega\left(\frac{\Delta L D}{\epsilon^2}\right)$ | / |
| | Corollary 1 | Decentral | $\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $\Omega\left(\frac{\Delta L}{\epsilon^2 \sqrt{1-\lambda}}\right)$ | / |
| Upper Bound | DeFacto (Theorem 2) | Central | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\Delta L D}{\epsilon^2}\right)$ | $O(1)$ |
| | DeTAG (Theorem 3) | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\Delta L \log\left(\frac{\varsigma_0 n}{\epsilon \sqrt{\Delta L}}\right)}{\epsilon^2 \sqrt{1-\lambda}}\right)$ | $O\left(\log\left(\frac{\varsigma_0 n}{\epsilon \sqrt{\Delta L}}\right)\right)$ |
| | D-PSGD [19] | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\Delta L n \varsigma}{\epsilon^2 (1-\lambda)^2}\right)$ | $O\left(\frac{n \varsigma}{(1-\lambda)^{\frac{3}{2}}}\right)$ |
| | SGP [72] | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\Delta L n \varsigma}{\epsilon^2 (1-\lambda)^2}\right)$ | $O\left(\frac{n \varsigma}{(1-\lambda)^{\frac{3}{2}}}\right)$ |
| | D ² [24] | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\lambda^2 \Delta L n \varsigma_0}{\epsilon^2 (1-\lambda)^3}\right)$ | $O\left(\frac{\lambda^2 n \varsigma_0}{(1-\lambda)^{\frac{5}{2}}}\right)$ |
| | DSGT [42] | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\lambda^2 \Delta L n \varsigma_0}{\epsilon^2 (1-\lambda)^3}\right)$ | $O\left(\frac{\lambda^2 n \varsigma_0}{(1-\lambda)^{\frac{5}{2}}}\right)$ |
| | GT-DSGD [44] | Decentral | $O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ | $O\left(\frac{\lambda^2 \Delta L n \varsigma_0}{\epsilon^2 (1-\lambda)^3}\right)$ | $O\left(\frac{\lambda^2 n \varsigma_0}{(1-\lambda)^{\frac{5}{2}}}\right)$ |

2 Related Work

Decentralized Training. In the application layer, decentralized training usually denotes federated learning [31]. Research on decentralization in this sense investigates convergence where each worker samples only from a local dataset which is not independent and identically distributed to other workers’ datasets [32, 33, 34, 35]. Another line of research on decentralization focuses on the protocol layer—with average gossip [36, 37], workers communicate by averaging their parameters with neighbors on a graph. D-PSGD [19] is one of the most basic algorithms that scales SGD with this protocol, achieving a linear parallel speed up. Additional works extend D-PSGD to asynchronous and variance-reduced cases [13, 24, 38, 39, 40, 41]. After those, Zhang and You [42], Xin et al. [43, 44] propose adding gradient trackers to D-PSGD. Other works discuss the application of decentralization on specific tasks such as linear models or deep learning [45, 46]. Zhang and You [47] treats the case where only directed communication can be performed. Wang et al. [25] proposes using matching algorithms to optimize the gossip protocol. Multiple works discuss using compression to decrease communication costs in decentralized training [10, 22, 48, 49, 50], and other papers connect decentralized training to other parallel methods and present a unified theory [4, 27, 51]. In some even earlier works like [52, 53], full local gradients on a convex setting is investigated.

Lower Bounds in Stochastic Optimization. Lower bounds are a well studied topic in non-stochastic optimization, especially in convex optimization [54, 55, 56, 57, 58]. In the stochastic setting, Allen-Zhu [59] and Foster et al. [60] discuss the complexity lower bound to find stationary points on convex problems. Other works study the lower bound in a convex, data-parallel setting [61, 62, 63], and Colin et al. [64] extends the result to a model-parallel setting. In the domain of non-convex optimization, Carmon et al. [65, 66] propose a zero-chain model that obtains tight bound for a first order method to obtain stationary points. Zhou and Gu [67] extends this lower bound to a finite sum setting, and Arjevani et al. [68] proposes a probabilistic zero-chain model that obtains tight lower bounds for first-order methods on stochastic and non-convex problems.

3 Setting

In this section, we introduce the notation and assumptions we will use. Throughout the paper, we consider the standard *data-parallel* training setup with n parallel workers. Each worker i stores a copy of the model $\mathbf{x} \in \mathbb{R}^d$ and a local dataset \mathcal{D}_i . The model copy and local dataset define a local loss function (or empirical risk) f_i . The ultimate goal of the parallel workers is to output a target model $\hat{\mathbf{x}}$ that minimizes the average over all the local loss functions, that is,

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left[f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}_{\xi_i \sim \mathcal{D}_i} f_i(\mathbf{x}; \xi_i)}_{f_i(\mathbf{x})} \right]. \quad (1)$$

Here, ξ_i is a data sample from \mathcal{D}_i and is used to compute a stochastic gradient via some oracle, e.g. back-propagation on a mini-batch of samples. The loss functions can (potentially) be non-convex so finding a global minimum is NP-Hard; instead, we expect the workers to output a point $\hat{\mathbf{x}}$ at which $f(\hat{\mathbf{x}})$ has a small gradient magnitude in expectation: $\mathbb{E} \|\nabla f(\hat{\mathbf{x}})\| \leq \epsilon$, for some small ϵ .¹ The assumptions our theoretical analysis requires can be categorized by the layers from Table 1: in each layer, “being decentralized” corresponds to certain assumptions (or lack of assumptions). We now describe these assumptions for each layer separately.

3.1 Application Layer

Application-layer assumptions comprise constraints on the losses f_i from (1) and the gradient oracle via which they are accessed by the learning algorithm, as these are constraints on the learning task itself.

Function class (Δ and L). As is usual in this space, we assume the local loss functions $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ are L -smooth,

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, \quad (2)$$

for some constant $L > 0$, and that the total loss f is range-bounded by Δ in the sense that $f(\mathbf{0}) - \inf_{\mathbf{x}} f(\mathbf{x}) \leq \Delta$. We let the **function class** $\mathcal{F}_{\Delta, L}$ denote the set of all functions that satisfy these conditions (for any dimension $d \in \mathbb{N}^+$).

Oracle class (σ^2). We assume each worker interacts with its local function f_i only via a stochastic gradient oracle \tilde{g}_i , and that when we query this oracle with model \mathbf{x} , it returns an independent unbiased estimator to $\nabla f_i(\mathbf{x})$ based on some random variable z with distribution Z (e.g. the index of a mini-batch randomly chosen for backprop). Formally,

$$\mathbb{E}_{z \sim Z} [\tilde{g}_i(\mathbf{x}, z)] = \nabla f_i(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (3)$$

As per the usual setup, we additionally assume the local estimator has bounded variance: for some constant $\sigma > 0$,

$$\mathbb{E}_{z \sim Z} \|\tilde{g}_i(\mathbf{x}, z) - \nabla f_i(\mathbf{x})\|^2 \leq \sigma^2, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (4)$$

We let \mathcal{O} denote a set of these oracles $\{\tilde{g}_i\}_{i \in [n]}$, and let the **oracle class** \mathcal{O}_{σ^2} denote the class of all such oracle sets that satisfy these two assumptions.

Data shuffling (ς^2 and ς_0^2). At this point, an analysis with a centralized application layer would make the additional assumption that all the f_i are equal and the \tilde{g}_i are identically distributed: this roughly corresponds to the assumption that the data all comes independently from a single centralized source. *We do not make this assumption*, and lacking such an assumption is what makes an analysis decentralized in the application layer. Still, some assumption that bounds the f_i relative to each other somehow is needed: we now discuss two such assumptions used in the literature, from which we use the weaker (and more decentralized) one.

One commonly made assumption [10, 19, 22, 48, 50] in decentralized training is

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \varsigma^2, \quad \forall \mathbf{x} \in \mathbb{R}^d, \quad (5)$$

¹There are many valid stopping criteria. We adopt ϵ -stationary point as the success signal. $\mathbb{E} \|\nabla f(\hat{\mathbf{x}})\|^2 \leq \epsilon^2$ is another commonly used criterion; we adopt the non-squared one following [66]. Other criterions regarding stationary points can be converted to hold in our theory.

for some constant ς , which is said to bound the “outer variance” among workers. This is often unreasonable, as it suggests the local datasets on workers must have close distribution: in practice, ensuring this often requires some sort of shuffling or common centralized data source. We **do not assume** (5) but instead adopt the much weaker assumption

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{0}) - \nabla f(\mathbf{0})\|^2 \leq \varsigma_0^2, \quad (6)$$

for constant $\varsigma_0 > 0$.² This assumption only requires a bound at point $\mathbf{0}$, which is, to the best of our knowledge, the weakest assumption of this type used in the literature [24, 42]. Requiring such a weak assumption allows workers to (potentially) sample from different distributions or vary largely in their loss functions (e.g. in a federated learning environment).

3.2 Protocol Layer

Protocol-layer assumptions comprise constraints on the parallel learning algorithm itself, and especially on the way that the several workers communicate to approach consensus.

Algorithm class (B). We consider algorithms A that divide training into multiple iterations, and between two adjacent iterations, there must be a synchronization process among workers (e.g. a barrier) such that they start each iteration simultaneously.³ Each worker running A has a local copy of the model, and we let $\mathbf{x}_{t,i} \in \mathbb{R}^d$ denote this model on worker i at iteration t . We assume without loss of generality that A initializes each local model at zero: $\mathbf{x}_{0,i} = \mathbf{0}$ for all i . At each iteration, each worker makes *at most* B queries to its gradient oracle \tilde{g}_i , for some constant $B \in \mathbb{N}^+$, and then uses the resulting gradients to update its model. We do not make any explicit rules for output and allow the output of the algorithm $\hat{\mathbf{x}}_t$ at the end of iteration t (the model that A would output if it were stopped at iteration t) to be any linear combination of all the local models, i.e.

$$\hat{\mathbf{x}}_t \in \text{span}(\{\mathbf{x}_{t,j}\}_{j \in [n]}) = \{\sum_{j=1}^n c_j \mathbf{x}_{t,j} \mid c_j \in \mathbb{R}\}. \quad (7)$$

Beyond these basic properties, we further require A to satisfy the following “zero-respecting” property from Carmon et al. [65]. Specifically, if \mathbf{z} is any vector worker i queries its gradient oracle with at iteration t , then for any $k \in [d]$, if $\mathbf{e}_k^\top \mathbf{z} \neq 0$, then there exists a $s \leq t$ and a $j \in [n]$ such that either $j = i$ or j is a neighbor of i in the network connectivity graph G (i.e. $(i, j) \in \{(i, i)\} \cup G$) and $(\mathbf{e}_k^\top \mathbf{x}_{s,j}) \neq 0$. More informally, the worker will not query its gradient oracle with a nonzero value for some weight unless that weight was already nonzero in the model state of the worker or one of its neighbors at some point in the past. Similarly, for any $k \in [d]$, if $(\mathbf{e}_k^\top \mathbf{x}_{t+1,i}) \neq 0$, then either there exists an $s \leq t$ and j such that $(i, j) \in \{(i, i)\} \cup G$ and $(\mathbf{e}_k^\top \mathbf{x}_{s,j}) \neq 0$, or one of the gradient oracle’s outputs \mathbf{v} on worker i at iteration t has $\mathbf{e}_k^\top \mathbf{v} \neq 0$. Informally, a worker’s model will not have a nonzero weight unless either (1) that weight was nonzero on that worker or one of its neighbors at a previous iteration, or (2) the corresponding entry in one of the gradients the worker sampled at that iteration was nonzero.

Intuitively, we are requiring that algorithm A will not modify those coordinates that remain zero in all previous oracle outputs and neighboring models.⁴ This lets A use a wide space of accessible information in communication and allows our class to cover first-order methods including SGD [73], Momentum SGD [74], Adam [75], RMSProp [76], Adagrad [77], and AdaDelta [78]. We let **algorithm class** \mathcal{A}_B denote the set of all algorithms A that satisfy these assumptions.

So far our assumptions in this layer cover both centralized and decentralized protocols. Decentralized protocols, however, must satisfy the additional assumption that they communicate via *gossip* (see Section 2) [36, 37]. A single step of gossip protocol can be expressed as

$$\mathbf{z}_{t,i} \leftarrow \sum_{j \in \mathcal{N}_i} \mathbf{y}_{t,j} \mathbf{W}_{ji}, \quad \forall i \in [n] \quad (8)$$

²As we only use ς_0 for upper bounds, not lower bounds, we do not define a “class” that depends on this parameter.

³We consider synchronous algorithms only here for simplicity of presentation; further discussion of extension to asynchronous algorithms is included in the supplementary material.

⁴On the other hand, it is possible to even drop the zero-respecting requirement and extend A to all the deterministic (not in the sense of sampling but the actual executions) algorithms. At a cost, we would need the function class to follow an “orthogonal invariant” property, and the model dimension needs to be large enough. We leave this discussion to the appendix.

for some constant doubly stochastic matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ called the *communication matrix* and \mathbf{y} and \mathbf{z} are the input and output of the gossip communication step, respectively. The essence of a single Gossip step is to take weighted average over the neighborhood specified by a fixed matrix. To simplify later discussion, we further define the **gossip matrix class** \mathcal{W}_n as the set of all matrices $\mathbf{W} \in \mathbb{R}^{n \times n}$, where \mathbf{W} is doubly stochastic and $\mathbf{W}_{ij} \neq 0$ only if $(i, j) \in G$. We call every $\mathbf{W} \in \mathcal{W}_n$ a gossip matrix and we use $\lambda = \max\{|\lambda_2|, |\lambda_n|\} \in [0, 1]$ to denote its general second-largest eigenvalue, where λ_i denotes the i -th largest eigenvalue of \mathbf{W} . We let **gossip algorithm class** $\mathcal{A}_{B, \mathbf{W}}$ denote the set of all algorithms $A \in \mathcal{A}_B$ that only communicate via gossip using a single matrix $\mathbf{W} \in \mathcal{W}_n$. It trivially holds that $\mathcal{A}_{B, \mathbf{W}} \subset \mathcal{A}_B$.

3.3 Topology Layer

Topology-layer assumptions comprise constraints on how workers are connected topologically. We let the **graph class** $\mathcal{G}_{n, D}$ denote the class of graphs G connecting n workers (vertices) with diameter D , where diameter of a graph measures the maximum distance between two arbitrary vertices (so $1 \leq D \leq n - 1$). A centralized analysis here typically will also require that G be either complete or complete-bipartite (with parameter servers and workers as the two parts): lacking this requirement and allowing arbitrary graphs is what makes an analysis decentralized in the topology layer.

3.4 Complexity Measures

Now that we have defined the classes we are interested in, we can use them to define the complexity measures we will bound in our theoretical results. Given a loss function $f \in \mathcal{F}_{\Delta, L}$, a set of underlying oracles $O \in \mathcal{O}_{\sigma^2}$, a graph $G \in \mathcal{G}_{n, D}$, and an algorithm $A \in \mathcal{A}_B$, let $\hat{\mathbf{x}}_t^{A, f, O, G}$ denote the output of algorithm A at the end of iteration t under this setting. Then the *iteration complexity* of A solving f under O and G is defined as

$$T_\epsilon(A, f, O, G) = \min \left\{ t \in \mathbb{N} \mid \mathbb{E} \left\| \nabla f(\hat{\mathbf{x}}_t^{A, f, O, G}) \right\| \leq \epsilon \right\},$$

that is, the least number of iterations required by A to find a ϵ -stationary-in-expectation point of f .

4 Lower Bound

Given the setup in Section 3, we can now present and discuss our lower bound on the iteration complexity. Note that in the formulation of protocol layer, the algorithm class \mathcal{A}_B only specifies the information available for each worker, and thus \mathcal{A}_B covers both centralization and decentralization in the protocol layer. Here, we show our lower bound in two parts: first a general bound where an arbitrary protocol that follows \mathcal{A}_B is allowed, and then a corollary bound for the case where only decentralized protocol is allowed.

4.1 Lower Bound for Arbitrary Protocol

We start from the general bound. We expect this lower bound to show given arbitrary setting (functions, oracles and graph), the smallest iteration complexity we could obtain from \mathcal{A}_B , i.e.

$$\inf_{A \in \mathcal{A}_B} \sup_{f \in \mathcal{F}_{\Delta, L}} \sup_{O \in \mathcal{O}_{\sigma^2}} \sup_{G \in \mathcal{G}_{n, D}} T_\epsilon(A, f, O, G), \quad (9)$$

it suffices to construct a hard instance containing a loss function $\hat{f} \in \mathcal{F}_{\Delta, L}$, a graph $\hat{G} \in \mathcal{G}_{n, D}$ and a set of oracles $\hat{O} \in \mathcal{O}_{\sigma^2}$ and obtain a valid lower bound on $\inf_{A \in \mathcal{A}_B} T_\epsilon(A, \hat{f}, \hat{O}, \hat{G})$ since Equation (9) is always lower bounded by $\inf_{A \in \mathcal{A}_B} T_\epsilon(A, \hat{f}, \hat{O}, \hat{G})$.

For the construction, we follow the idea of probabilistic zero-chain model [65, 66, 67, 68], which is a special loss function where adjacent coordinates are closely dependent on each other like a “chain.” Our main idea is to use this function as f and split this chain onto different workers. Then the workers must conduct a sufficient number of optimization steps and rounds of communication to make progress.⁵ From this, we obtain the following lower bound.

⁵For brevity, we leave details in the supplementary material.

Theorem 1 For function class $\mathcal{F}_{\Delta,L}$, oracle class \mathcal{O}_{σ^2} and graph class $\mathcal{G}_{n,D}$ defined with any $\Delta > 0$, $L > 0$, $n \in \mathbb{N}^+$, $D \in \{1, 2, \dots, n-1\}$, $\sigma > 0$, and $B \in \mathbb{N}^+$, there exists $f \in \mathcal{F}_{\Delta,L}$, $O \in \mathcal{O}_{\sigma^2}$, and $G \in \mathcal{G}_{n,D}$, such that no matter what $A \in \mathcal{A}_B$ is used, $T_\epsilon(A, f, O, G)$ will always be lower bounded by

$$\Omega \left(\frac{\Delta L \sigma^2}{n B \epsilon^4} + \frac{\Delta L D}{\epsilon^2} \right). \quad (10)$$

Dependency on the parameters. The bound in Theorem 1 consists of a sample complexity term, which is the dominant one for small ϵ , and a communication complexity term. We can see the increase of query budget B will only reduce the sample complexity. On the other hand, as the diameter D of a graph will generally increase as the number of vertices n increases, we can observe a trade-off between two terms when the system scales up: when more workers join the system, the communication complexity will gradually become the dominant term.

Consistency with the literature. Theorem 1 is tightly aligned with the state-of-the-art bounds in many settings. With $n = B = D = 1$, we recover the tight bound for sequential stochastic non-convex optimization $\Theta(\Delta L \sigma^2 \epsilon^{-4})$ as shown in Arjevani et al. [68]. With $\sigma = 0, D = 1$, we recover the tight bound for sequential non-stochastic non-convex optimization $\Theta(\Delta L \epsilon^{-2})$ as shown in Carmon et al. [66]. With $B = 1, D = 1$, we recover the tight bound for centralized training $\Theta(\Delta L \sigma^2 (n \epsilon^4)^{-1})$ given in Li et al. [6].

Improvement upon previous results. Previous works like Seaman et al. [17], Scaman et al. [26] provide similar lower bounds in a convex setting which relates to the diameter. However, these results treat D as a fixed value, i.e., $D = n - 1$, and thus makes the bound to be only tight on linear graph. By comparison, Theorem 1 allows D to be chosen independently to n .

4.2 Lower Bound for Decentralized Protocol

The bound in Theorem 1 holds for both centralized and decentralized protocols. A natural question is: *How would the lower bound adapt if the protocol is restricted to be decentralized?* i.e., the quantity of

$$\inf_{A \in \mathcal{A}_{B, \mathbf{W}}} \sup_{f \in \mathcal{F}_{\Delta,L}} \sup_{O \in \mathcal{O}_{\sigma^2}} \sup_{G \in \mathcal{G}_{n,D}} T_\epsilon(A, f, O, G),$$

we can extend the lower bound to Gossip in the following corollary.

Corollary 1 For every $\Delta > 0$, $L > 0$, $n \in \{2, 3, 4, \dots\}$, $\lambda \in [0, \cos(\pi/n)]$, $\sigma > 0$, and $B \in \mathbb{N}^+$, there exists a loss function $f \in \mathcal{F}_{\Delta,L}$, a set of underlying oracles $O \in \mathcal{O}_{\sigma^2}$, a gossip matrix $\mathbf{W} \in \mathcal{W}_n$ with second largest eigenvalue being λ , and a graph $G \in \mathcal{G}_{n,D}$, such that no matter what $A \in \mathcal{A}_{B, \mathbf{W}}$ is used, $T_\epsilon(A, f, O, G)$ will always be lower bounded by

$$\Omega \left(\frac{\Delta L \sigma^2}{n B \epsilon^4} + \frac{\Delta L}{\epsilon^2 \sqrt{1 - \lambda}} \right). \quad (11)$$

Gap in the existing algorithms. Comparing this lower bound with many state-of-the-art decentralized algorithms (Table 2), we can see they match on the sample complexity but leave a gap on the communication complexity. In many cases, the spectral gap significantly depends on the number of workers n and thus can be arbitrarily large. For example, when the graph G is a cycle graph or a linear graph, the gap of those baselines can increase by up to $O(n^6)$ [79, 80]!

5 DeFacto: Optimal Complexity in Theory

In the previous section we show the existing algorithms have a gap compared to the lower bound. This gap could indicate the algorithms are suboptimal, but it could also be explained by our lower bound being loose. In this section we address this issue by proposing DeFacto, an example algorithm showing the lower bound is achievable, which verifies the tightness of our lower bound—showing that (10) would hold with equality and $\Theta(\cdot)$, not just $\Omega(\cdot)$.

We start with the following insight on the theoretical gap: the goal of communication is to let all the workers obtain information from neighbors. Ideally, the workers would, at each iteration, perform (8) with $\mathbf{W}^* = \mathbf{1}_n \mathbf{1}_n^\top / n$, where $\mathbf{1}_n$ is the n -dimensional all-one vector. We call this matrix the *Average Consensus*

Algorithm 1 Decentralized Stochastic Gradient Descent with Factorized Consensus Matrices (DeFacto) on worker i

Input: initialized model $\mathbf{x}_{0,i}$, a copy of model $\tilde{\mathbf{x}}_{0,i} \leftarrow \mathbf{x}_{0,i}$, gradient buffer $\mathbf{g} = \mathbf{0}$, step size α , a sequence of communication matrices $\{\mathbf{W}_r\}_{1 \leq r \leq R}$ of size R , number of iterations T , neighbor list \mathcal{N}_i

```

1: for  $t = 0, 1, \dots, T - 1$  do
2:    $k \leftarrow \lfloor t/2R \rfloor$ .
3:    $r \leftarrow t \bmod 2R$ .
4:   if  $0 \leq r < R$  then
5:     Spend all  $B$  oracle budgets to compute stochastic gradient  $\tilde{\mathbf{g}}$  at point  $\mathbf{x}_{k,i}$  and accumulate it to gradient
     buffer:  $\mathbf{g} \leftarrow \mathbf{g} + \tilde{\mathbf{g}}$ .
6:   else
7:     Update model copy with the  $r$ -th matrix in  $\{\mathbf{W}_r\}_{1 \leq r \leq R}$ :

$$\tilde{\mathbf{x}}_{t+1,i} \leftarrow \sum_{j \in \mathcal{N}_i \cup \{i\}} \tilde{\mathbf{x}}_{t,j} [\mathbf{W}_r]_{ji} \tag{12}$$

8:   end if
9:   if  $r = 2R - 1$  then
10:    Update Model:  $\mathbf{x}_{t+1,i} \leftarrow \tilde{\mathbf{x}}_{t+1,i} - \alpha \frac{\mathbf{g}}{R}$ .
11:    Reinitialize gradient buffer:  $\mathbf{g} \leftarrow \mathbf{0}$ .
12:    Copy the current model:  $\tilde{\mathbf{x}}_{t+1,i} \leftarrow \mathbf{x}_{t+1,i}$ .
13:   end if
14: end for
15: return  $\hat{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{T,i}$ 

```

matrix. The Average Consensus is statistically equivalent to centralized communication (All-Reduce operation). However, due to the graph constraints, we can not use this \mathbf{W}^* unless workers are fully connected; instead, a general method is to repeatedly apply a sequence communication matrices in consecutive iterations and let workers achieve or approach the Average Consensus. Previous work uses Gossip matrix \mathbf{W} and expect $\prod_{r=1}^R \mathbf{W} \approx \mathbf{1}_n \mathbf{1}_n^\top / n$ for some R . This R is known to be proportional to the mixing time of the Markov Chain \mathbf{W} defines [4, 22], which is related to the inverse of its spectral gap [81]. This limits convergence depending on the spectrum of the \mathbf{W} chosen. The natural question to ask here is: can we do better? What are the limits of how fast we can reach average consensus on a connectivity graph G ? This question is answered by the following lemma.

Lemma 1 For any $G \in \mathcal{G}_{n,D}$, let \mathcal{W}_G denote the set of $n \times n$ matrices such that for all $\mathbf{W} \in \mathcal{W}_G$, $\mathbf{W}_{ij} = 0$ if edge (i, j) does not appear in G . There exists a sequence of R matrices $\{\mathbf{W}_r\}_{r \in [R]}$ that belongs to \mathcal{W}_G such that $R \in \{D, D + 1, \dots, 2D\}$ and

$$\mathbf{W}_{R-1} \mathbf{W}_{R-2} \cdots \mathbf{W}_0 = \frac{\mathbf{1}_n \mathbf{1}_n^\top}{n} = \mathbf{W}^*.$$

Lemma 1 is a classic result in the literature of graph theory. The formal proof and detailed methods to identify these matrices can be found in many previous works [82, 83, 84]. Here we treat this as a black box procedure.⁶

Lemma 1 shows that we can achieve the exact average consensus by factorizing the matrix $\mathbf{1}_n \mathbf{1}_n^\top / n$, and we can obtain the factors from a preprocessing step. From here, the path to obtain an optimal rate becomes clear: starting from $t = 0$, workers first spend R iterations only computing stochastic gradients and then another R iterations to reach consensus communicating via factors from Lemma 1; they then repeat this process until a stationary point is found. We call this algorithm DeFacto (Algorithm 1).

DeFacto is statistically equivalent to centralized SGD operating $T/2R$ iterations with a mini-batch size of BR . It can be easily verified that DeFacto holds membership in \mathcal{A}_B . A straightforward analysis gives the convergence rate of DeFacto shown in the following Theorem.

Theorem 2 Let A_1 denote Algorithm 1. For $\mathcal{F}_{\Delta,L}$, \mathcal{O}_{σ^2} and $\mathcal{G}_{n,D}$ defined with any $\Delta > 0$, $L > 0$, $n \in \mathbb{N}^+$, $D \in \{1, 2, \dots, n - 1\}$, $\sigma > 0$, and $B \in \mathbb{N}^+$, the convergence rate of A_1 running on any loss function $f \in \mathcal{F}_{\Delta,L}$,

⁶We cover specific algorithms and details in the supplementary.

Algorithm 2 Decentralized Stochastic Gradient Tracking with By-Phase Accelerated Gossip (DeTAG) on worker i

Input: initialized model $\mathbf{x}_{0,i}$, a copy of model $\tilde{\mathbf{x}}_{0,i} \leftarrow \mathbf{x}_{0,i}$, gradient tracker $\mathbf{y}_{0,i}$, gradient buffer $\mathbf{g}_{(0)} = \mathbf{g}_{(-1)} = \mathbf{0}$, step size α , a gossip matrix \mathbf{W} , number of iterations T , neighbor list \mathcal{N}_i

```

1: for  $t = 0, 1, \dots, T - 1$  do
2:    $k \leftarrow \lfloor t/R \rfloor$ .
3:    $r \leftarrow t \bmod R$ .
4:   Perform the  $r$ -th step in Accelerate Gossip:
       
$$\tilde{\mathbf{x}}_{t+1,i} \leftarrow AG(\tilde{\mathbf{x}}_{t,i}, \mathbf{W}, \mathcal{N}_i, i) \tag{13}$$

       
$$\mathbf{y}_{t+1,i} \leftarrow AG(\mathbf{y}_{t,i}, \mathbf{W}, \mathcal{N}_i, i) \tag{14}$$

5:   Spend all  $B$  oracle budgets to compute stochastic gradient  $\tilde{\mathbf{g}}$  at point  $\mathbf{x}_{k,i}$  and accumulate it to gradient buffer:
        $\mathbf{g}_{(k)} \leftarrow \mathbf{g}_{(k)} + \tilde{\mathbf{g}}$ .
6:   if  $r = R - 1$  then
7:     Update gradient tracker and model:
       
$$\mathbf{x}_{t+1,i} \leftarrow \tilde{\mathbf{x}}_{t+1,i} - \alpha \mathbf{y}_i \tag{15}$$

       
$$\mathbf{y}_{t+1,i} \leftarrow \mathbf{y}_{t+1,i} + \mathbf{g}_{(k)} - \mathbf{g}_{(k-1)} \tag{16}$$

8:     Reinitialize gradient buffer:  $\mathbf{g}_{(k-1)} \leftarrow \mathbf{g}_{(k)}$  and then  $\mathbf{g}_{(k)} \leftarrow \mathbf{0}$ .
9:     Copy the current model:  $\tilde{\mathbf{x}}_{t+1,i} \leftarrow \mathbf{x}_{t+1,i}$ .
10:  end if
11: end for
12: return  $\hat{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{T,i}$ 
```

Algorithm 3 Accelerated Gossip (AG) with R steps

Input: $\mathbf{z}_{0,i}$, \mathbf{W} , \mathcal{N}_i , i

```

1:  $\mathbf{z}_{-1,i} \leftarrow \mathbf{z}_{0,i}$ 
2:  $\eta \leftarrow \frac{1 - \sqrt{1 - \lambda^2}}{1 + \sqrt{1 - \lambda^2}}$ 
3: for  $r = 0, 1, 2, \dots, R - 1$  do
4:    $\mathbf{z}_{r+1,i} \leftarrow (1 + \eta) \sum_{j \in \mathcal{N}_i \cup \{i\}} \mathbf{z}_{r,j} \mathbf{W}_{ji} - \eta \mathbf{z}_{r-1,i}$ 
5: end for
6: return  $\mathbf{z}_{R,i}$ 
```

any graph $G \in \mathcal{G}_{n,D}$, and any oracles $O \in \mathcal{O}_{\sigma^2}$ is bounded by

$$T_\epsilon(A_1, f, O, G) \leq O \left(\frac{\Delta L \sigma^2}{n B \epsilon^4} + \frac{\Delta L D}{\epsilon^2} \right). \tag{17}$$

Comparing Theorem 1 and Theorem 2, DeFacto achieves the optimal rate asymptotically. *This shows that our lower bound in Theorem 1 is tight.*

Despite its optimality, the design of DeFacto is unsatisfying in three aspects: (1) It compromises the throughput⁷ by a factor of two because in each iteration, a worker either communicates with neighbors or computes gradients but not both. This fails to overlap communication and computation and creates extra idle time for the workers. (2) It needs to iterate over all the factor matrices before it can query the gradient oracle at subsequent parameters. When diameter D increases, the total time to finish such round will increase proportionally. (3) DeFacto works with decentralized data and arbitrary graph, achieving decentralization in both application and topology layers. However, the matrices used in Lemma 1 are not Gossip matrices as defined in \mathcal{W}_n , and thus it fails to be decentralized in the protocol-layer sense.

6 DeTAG: Optimal Complexity in Practice

To address the limitations of DeFacto, a natural idea is to replace all the factor matrices in Lemma 1 with a gossip matrix \mathbf{W} . The new algorithm after this mild modification is statistically equivalent to a D-PSGD

⁷The number of stochastic gradients computed per iteration.

variant: every R iterations, it updates the model the same as one iteration in D-PSGD with a mini-batch size of BR and communicate with a matrix \mathbf{W}' whose second largest eigenvalue $\lambda' = \lambda^R$, with T/R iterations in total. However, even with arbitrarily large R , the communication complexity in this “updated D-PSGD” is still $O(\Delta L n \varsigma \epsilon^{-2})$ (Table 2), leaving an $O(n\varsigma)$ gap compared to our lower bound.

To close this gap, we adopt two additional techniques:⁸ one is a gradient tracker \mathbf{y} that is used as reference capturing gradient difference in the neighborhood; the other is using acceleration in gossip as specified in Algorithm 3. Modifying DeFacto results in Algorithm 2, which we call DeTAG. DeTAG works as follows: it divides the total number of iterations T into several phases where each phase contains R iterations. In each iteration, the communication process calls Accelerated Gossip to update a model replica $\tilde{\mathbf{x}}$ and the gradient tracker (line 4) while the computation process constantly computes gradients at the same point (line 5). At the end of each phase, model \mathbf{x} , its replica $\tilde{\mathbf{x}}$ and gradient tracker \mathbf{y} are updated in line 7-10 and then DeTAG steps into the next phase. Aside from the two additional techniques, the main difference between DeTAG and DeFacto is that the communication matrix in DeTAG is a fixed gossip matrix \mathbf{W} , which allows DeTAG to benefit from decentralization in the protocol layer as well as to adopt arbitrary $R \geq 1$ in practice (allowing R to be tuned independently of G).

Improvement on design compared to baselines. Comparing with other baselines in Table 2, the design of DeTAG improves in the sense that (1) It removes the dependency on the outer variance ς . (2) It drops the requirement⁹ on the gossip matrix assumed in Tang et al. [24]. (3) The baseline DSGT [42] and GT-DSGD [44] can be seen as special cases of taking $R = 1$ and $\eta = 0$ in DeTAG. That implies in practice, a well tuned DeTAG can never perform worse than the baseline DSGT or GT-DSGD.

The convergence rate of DeTAG is given in the following theorem.

Theorem 3 *Let A_2 denote Algorithm 2. For $\mathcal{F}_{\Delta,L}$, \mathcal{O}_{σ^2} and $\mathcal{G}_{n,D}$ defined with any $\Delta > 0$, $L > 0$, $n \in \mathbb{N}^+$, $\lambda \in [0, 1)$, $\sigma > 0$, and $B \in \mathbb{N}^+$, under the assumption of Equation (6), if we set the phase length R to be*

$$R = \frac{\max\left(\frac{1}{2} \log(n), \frac{1}{2} \log\left(\frac{\varsigma_0^2 T}{\Delta L}\right)\right)}{\sqrt{1 - \lambda}},$$

the convergence rate of A_2 running on any loss function $f \in \mathcal{F}_{\Delta,L}$, any graph $G \in \mathcal{G}_{n,D}$, and any oracles $O \in \mathcal{O}_{\sigma^2}$ is bounded by

$$T_\epsilon(A_2, f, O, G) \leq O\left(\frac{\Delta L \sigma^2}{n B \epsilon^4} + \frac{\Delta L \log\left(n + \frac{\varsigma_0 n}{\epsilon \sqrt{\Delta L}}\right)}{\epsilon^2 \sqrt{1 - \lambda}}\right).$$

Comparing Theorem 1 and Theorem 3, DeTAG achieves the optimal complexity with only a logarithm gap.

Improvement on complexity. Revisiting Table 2, we can see the main improvement of DeTAG’s complexity is in the two terms on communication complexity: (1) DeTAG only depends on the outer variance term ς_0 inside a log, and (2) It reduces the dependency on the spectral gap $1 - \lambda$ to the lower bound of square root, as shown in Corollary 1.

Understanding the phase length R . In DeTAG, the phase length R is a tunable parameter. Theorem 3 provides a suggested value for R . Intuitively, the value of R captures the level of consensus of workers should reach before they step into the next phase. Theoretically, we observe R is closely correlated to the mixing time of \mathbf{W} : if we do not use acceleration in Gossip, then R will become $\tilde{O}\left(\frac{1}{1 - \lambda}\right)$, which is exactly the upper bound on the mixing time of the Markov Chain \mathbf{W} defines [81].

7 Experiments

In this section we empirically compare the performance among different algorithms. All the models and training scripts in this section are implemented in PyTorch and run on an Ubuntu 16.04 LTS cluster using a SLURM workload manager running CUDA 9.2, configured with 8 NVIDIA GTX 2080Ti GPUs. We launch

⁸Note that neither of these techniques is our original design, and we do not take credit for them. Our main contribution here is to prove their combination leads to optimal complexity.

⁹Tang et al. [24] requires the gossip matrix to be symmetric and its smallest eigenvalue is lower bounded by $-\frac{1}{3}$.

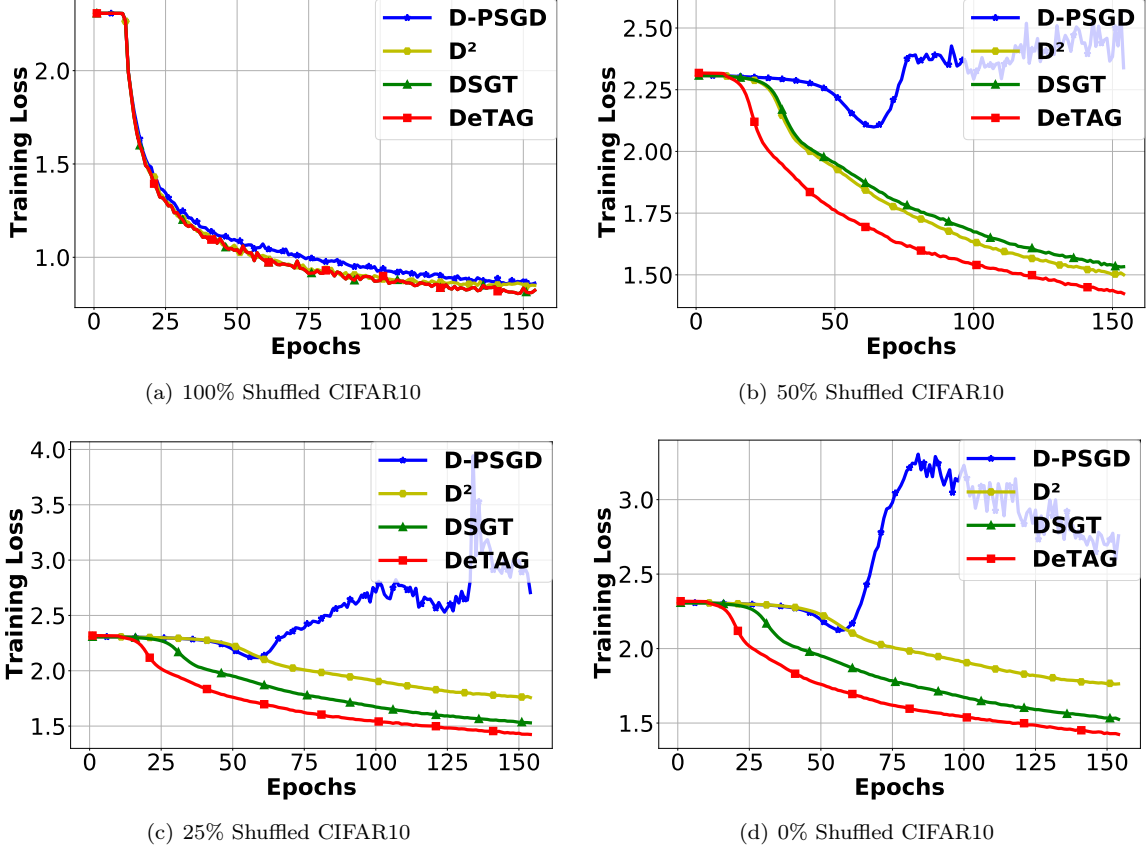


Figure 2: Fine tuned results of training LeNet on CIFAR10 with different shuffling strategies.

one process from the host as one worker and let them use `gloo` as the communication backend. In each experiment, we compare the following algorithms¹⁰: D-PSGD [19], D^2 [24], DSGT [42] and DeTAG. Note that GT-DSGD [44] and DSGT [42] are essentially the same algorithm so we omit the comparison to GT-DSGD. Also note that SGP [72] reduces to D-PSGD for symmetric mixing matrices in undirected graphs. Throughout the experiment we use Ring graph. Hyperparameters can be found in the supplementary material.

Convergence over different outer variance. In the first experiments, we investigate the correlation between convergence speed and the outer variance $\varsigma(\varsigma_0)$. We train LeNet on CIFAR10 using 8 workers, which is a standard benchmark experiment in the decentralized data environment [24, 42]. To create the decentralized data, we first sort all the data points based on its labels, shuffle the first $X\%$ data points and then evenly split to different workers. The X controls the degree of decentralization, we test $X = 0, 25, 50, 100$ and plot the results in Figure 2.

We can see in Figure 2(a) when the dataset is fully shuffled, all the algorithms converge at similar speed while D-PSGD converges a little slower than other variance reduced algorithms. From Figure 2(b) to Figure 2(d) we can see when we shuffle less portion of the dataset, i.e., the dataset becomes more decentralized, D-PSGD fails to converge even with fine-tuned hyperparameter. Meanwhile, among D^2 , DSGT and DeTAG, we can see DeTAG converges the fastest. When dataset becomes more decentralized, DSGT seems to receive more stable performance than D^2 .

Convergence over different spectral gaps. In the second experiments, we proceed to explore the relation between convergence speed and spectral gap $1 - \lambda$ of the gossip matrix \mathbf{W} . We use 16 workers connected

¹⁰Since DeFacto is a only a "motivation" algorithm and in practice we observe it performs bad, we do not include the discussion of that.

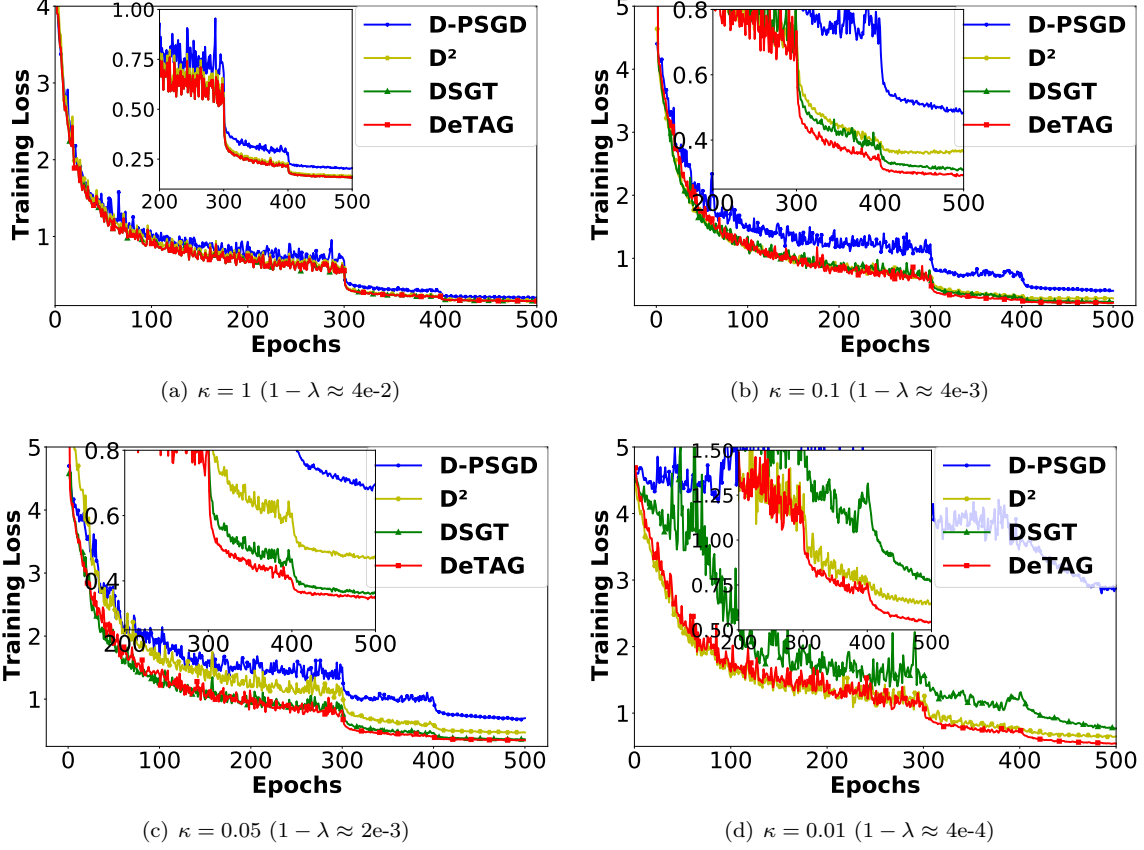


Figure 3: Fine tuned results of training Resnet20 on CIFAR100 with different spectral gaps.

with a Ring graph to train Resnet20 on CIFAR100, and we generate a \mathbf{W}_0 on such graph using Metropolis method. Then we adopt the slack matrix method to modify the spectral gap [4]: $\mathbf{W}_\kappa = \kappa \mathbf{W}_0 + (1 - \kappa) \mathbf{I}$, where κ is a control parameter. We test $\kappa = 1, 0.1, 0.05, 0.01$ and plot the results in Figure 3. We can see with different κ , DeTAG is able to achieve faster convergence compared to baselines. When the network becomes sparse, i.e., κ decreases, DeTAG enjoys more robust convergence.

8 Conclusion

In this paper, we investigate the tight lower bound on the iteration complexity of decentralized training. We propose two algorithms, DeFacto and DeTAG, that achieve the lower bound in terms of different decentralization in a learning system. DeTAG uses Gossip protocol, and is shown to be empirically competitive to many baseline algorithms, such as D-PSGD. In the future, we plan to investigate the variants of the complexity bound with respect to communication that are compressed, asynchronous, etc.

Acknowledgement

This work is supported by NSF IIS-2046760. The authors would like to thank A. Feder Cooper, Jerry Chee, Zheng Li, Ran Xin, Jiaqi Zhang and anonymous reviewers from ICML 2021 for providing valuable feedbacks on earlier versions of this paper.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.
- [2] Dan Alistarh. A brief tutorial on distributed and concurrent machine learning. In *Proceedings of the 2018 ACM Symposium on Principles of Distributed Computing*, pages 487–488, 2018.
- [3] Dan Alistarh, Bapi Chatterjee, and Vyacheslav Kungurtsev. Elastic consistency: A general consistency model for distributed stochastic gradient descent. *arXiv preprint arXiv:2001.05918*, 2020.
- [4] Yucheng Lu, Jack Nash, and Christopher De Sa. Mixml: A unified analysis of weakly consistent parallel learning. *arXiv preprint arXiv:2005.06706*, 2020.
- [5] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, and Bor-Yiing Su. Scaling distributed machine learning with the parameter server. In *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, pages 583–598, 2014.
- [6] Mu Li, David G Andersen, Alexander J Smola, and Kai Yu. Communication efficient distributed machine learning with the parameter server. In *Advances in Neural Information Processing Systems*, pages 19–27, 2014.
- [7] Qirong Ho, James Cipar, Henggang Cui, Seunghak Lee, Jin Kyu Kim, Phillip B Gibbons, Garth A Gibson, Greg Ganger, and Eric P Xing. More effective distributed ml via a stale synchronous parallel parameter server. In *Advances in neural information processing systems*, pages 1223–1231, 2013.
- [8] William Gropp, Rajeev Thakur, and Ewing Lusk. *Using MPI-2: Advanced features of the message passing interface*. MIT press, 1999.
- [9] Pitch Patarasuk and Xin Yuan. Bandwidth optimal all-reduce algorithms for clusters of workstations. *Journal of Parallel and Distributed Computing*, 69(2):117–124, 2009.
- [10] Anastasia Koloskova, Tao Lin, Sebastian U Stich, and Martin Jaggi. Decentralized deep learning with arbitrary communication compression. *arXiv preprint arXiv:1907.09356*, 2019.
- [11] H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.
- [12] Ameeth Kanawaday and Aditya Sane. Machine learning for predictive maintenance of industrial machines using iot sensor data. In *2017 8th IEEE International Conference on Software Engineering and Service Science (ICSESS)*, pages 87–90. IEEE, 2017.
- [13] Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous decentralized parallel stochastic gradient descent. *arXiv preprint arXiv:1710.06952*, 2017.
- [14] Hanlin Tang, Xiangru Lian, Chen Yu, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. *arXiv preprint arXiv:1905.05957*, 2019.
- [15] Chen Yu, Hanlin Tang, Cedric Renggli, Simon Kassing, Ankit Singla, Dan Alistarh, Ce Zhang, and Ji Liu. Distributed learning over unreliable networks. *arXiv preprint arXiv:1810.07766*, 2018.
- [16] Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- [17] Kevin Seaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3027–3036. JMLR. org, 2017.

- [18] Uday Shankar Shanthamallu, Andreas Spanias, Cihan Tepedelenlioglu, and Mike Stanley. A brief survey of machine learning methods and their sensor and iot applications. In *2017 8th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pages 1–8. IEEE, 2017.
- [19] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340, 2017.
- [20] Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. *arXiv preprint arXiv:1905.03817*, 2019.
- [21] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. Dadam: A consensus-based distributed adaptive gradient method for online optimization. *arXiv preprint arXiv:1901.09109*, 2019.
- [22] Yucheng Lu and Christopher De Sa. Moniqua: Modulo quantized communication in decentralized sgd. *arXiv preprint arXiv:2002.11787*, 2020.
- [23] Xinchen Wan, Hong Zhang, Hao Wang, Shuihai Hu, Junxue Zhang, and Kai Chen. Rat-resilient allreduce tree for distributed machine learning. In *4th Asia-Pacific Workshop on Networking*, pages 52–57, 2020.
- [24] Hanlin Tang, Xiangru Lian, Ming Yan, Ce Zhang, and Ji Liu. D2: Decentralized training over decentralized data. *arXiv preprint arXiv:1803.07068*, 2018.
- [25] Jianyu Wang, Anit Kumar Sahu, Zhouyi Yang, Gauri Joshi, and Soumya Kar. Matcha: Speeding up decentralized sgd via matching decomposition sampling. *arXiv preprint arXiv:1905.09435*, 2019.
- [26] Kevin Scaman, Francis Bach, Sébastien Bubeck, Laurent Massoulié, and Yin Tat Lee. Optimal algorithms for non-smooth distributed optimization in networks. In *Advances in Neural Information Processing Systems*, pages 2740–2749, 2018.
- [27] Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian U Stich. A unified theory of decentralized sgd with changing topology and local updates. *arXiv preprint arXiv:2003.10422*, 2020.
- [28] Blake E Woodworth, Jiale Wang, Adam Smith, Brendan McMahan, and Nati Srebro. Graph oracle models, lower bounds, and gaps for parallel stochastic optimization. In *Advances in neural information processing systems*, pages 8496–8506, 2018.
- [29] Darina Dvinskikh and Alexander Gasnikov. Decentralized and parallelized primal and dual accelerated methods for stochastic convex programming problems. *arXiv preprint arXiv:1904.09015*, 2019.
- [30] Haoran Sun and Mingyi Hong. Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms. *IEEE Transactions on Signal processing*, 67(22): 5912–5928, 2019.
- [31] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- [32] Keith Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman, Vladimir Ivanov, Chloe Kiddon, Jakub Konečný, Stefano Mazzocchi, H Brendan McMahan, et al. Towards federated learning at scale: System design. *arXiv preprint arXiv:1902.01046*, 2019.
- [33] Nguyen H Tran, Wei Bao, Albert Zomaya, Minh NH Nguyen, and Choong Seon Hong. Federated learning over wireless networks: Optimization model design and analysis. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, pages 1387–1395. IEEE, 2019.
- [34] Qiang Yang, Yang Liu, Yong Cheng, Yan Kang, Tianjian Chen, and Han Yu. Federated learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 13(3):1–207, 2019.

- [35] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [36] Stephen Boyd, Arpita Ghosh, Balaji Prabhakar, and Devavrat Shah. Gossip algorithms: Design, analysis and applications. In *Proceedings IEEE 24th Annual Joint Conference of the IEEE Computer and Communications Societies.*, volume 3, pages 1653–1664. IEEE, 2005.
- [37] Stephen Boyd, Arpita Ghosh, Balaji Prabhakar, and Devavrat Shah. Randomized gossip algorithms. *IEEE transactions on information theory*, 52(6):2508–2530, 2006.
- [38] Ye Tian, Ying Sun, and Gesualdo Scutari. Achieving linear convergence in distributed asynchronous multiagent optimization. *IEEE Transactions on Automatic Control*, 65(12):5264–5279, 2020.
- [39] Jiaqi Zhang and Keyou You. Asyspa: An exact asynchronous algorithm for convex optimization over digraphs. *IEEE Transactions on Automatic Control*, 65(6):2494–2509, 2019.
- [40] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. Asynchronous accelerated proximal stochastic gradient for strongly convex distributed finite sums. *arXiv preprint arXiv:1901.09865*, 2019.
- [41] Ran Xin, Usman A Khan, and Soumya Kar. A hybrid variance-reduced method for decentralized stochastic non-convex optimization. *arXiv preprint arXiv:2102.06752*, 2021.
- [42] Jiaqi Zhang and Keyou You. Decentralized stochastic gradient tracking for empirical risk minimization. *arXiv preprint arXiv:1909.02712*, 2019.
- [43] Ran Xin, Usman A Khan, and Soumya Kar. Variance-reduced decentralized stochastic optimization with gradient tracking. *arXiv preprint arXiv:1909.11774*, 2019.
- [44] Ran Xin, Usman A Khan, and Soumya Kar. An improved convergence analysis for decentralized online stochastic non-convex optimization. *IEEE Transactions on Signal Processing*, 69:1842–1858, 2021.
- [45] Lie He, An Bian, and Martin Jaggi. Cola: Decentralized linear learning. In *Advances in Neural Information Processing Systems*, pages 4536–4546, 2018.
- [46] Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Michael Rabbat. Stochastic gradient push for distributed deep learning. *arXiv preprint arXiv:1811.10792*, 2018.
- [47] Jiaqi Zhang and Keyou You. Asynchronous decentralized optimization in directed networks. *arXiv preprint arXiv:1901.08215*, 2019.
- [48] Anastasia Koloskova, Sebastian U Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. *arXiv preprint arXiv:1902.00340*, 2019.
- [49] Hanlin Tang, Xiangru Lian, Shuang Qiu, Lei Yuan, Ce Zhang, Tong Zhang, and Ji Liu. Deepsqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. *arXiv preprint arXiv:1907.07346*, 2019.
- [50] Hanlin Tang, Shaoduo Gan, Ce Zhang, Tong Zhang, and Ji Liu. Communication compression for decentralized training. In *Advances in Neural Information Processing Systems*, pages 7652–7662, 2018.
- [51] Jianyu Wang and Gauri Joshi. Cooperative sgd: A unified framework for the design and analysis of communication-efficient sgd algorithms. *arXiv preprint arXiv:1808.07576*, 2018.
- [52] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [53] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Distributed dual averaging in networks. In *NIPS*, pages 550–558. Citeseer, 2010.

- [54] Alekh Agarwal and Leon Bottou. A lower bound for the optimization of finite sums. *arXiv preprint arXiv:1410.0723*, 2014.
- [55] Yossi Arjevani and Ohad Shamir. Communication complexity of distributed convex learning and optimization. In *Advances in neural information processing systems*, pages 1756–1764, 2015.
- [56] Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method. *Mathematical programming*, 171(1-2):167–215, 2018.
- [57] Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, pages 689–699, 2018.
- [58] Yossi Arjevani and Ohad Shamir. Oracle complexity of second-order methods for finite-sum problems. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 205–213. JMLR. org, 2017.
- [59] Zeyuan Allen-Zhu. How to make the gradients small stochastically: Even faster convex and nonconvex sgd. In *Advances in Neural Information Processing Systems*, pages 1157–1167, 2018.
- [60] Dylan Foster, Ayush Sekhari, Ohad Shamir, Nathan Srebro, Karthik Sridharan, and Blake Woodworth. The complexity of making the gradient small in stochastic convex optimization. *arXiv preprint arXiv:1902.04686*, 2019.
- [61] Jelena Diakonikolas and Cristóbal Guzmán. Lower bounds for parallel and randomized convex optimization. *arXiv preprint arXiv:1811.01903*, 2018.
- [62] Eric Balkanski and Yaron Singer. Parallelization does not accelerate convex optimization: Adaptivity lower bounds for non-smooth convex minimization. *arXiv preprint arXiv:1808.03880*, 2018.
- [63] Quoc Tran-Dinh, Ahmet Alacaoglu, Olivier Fercoq, and Volkan Cevher. An adaptive primal-dual framework for nonsmooth convex minimization. *Mathematical Programming Computation*, pages 1–41, 2019.
- [64] Igor Colin, Ludovic Dos Santos, and Kevin Scaman. Theoretical limits of pipeline parallel optimization and application to distributed deep learning. In *Advances in Neural Information Processing Systems*, pages 12350–12359, 2019.
- [65] Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points ii: First-order methods. *arXiv preprint arXiv:1711.00841*, 2017.
- [66] Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points i. *Mathematical Programming*, pages 1–50, 2019.
- [67] Dongruo Zhou and Quanquan Gu. Lower bounds for smooth nonconvex finite-sum optimization. *arXiv preprint arXiv:1901.11224*, 2019.
- [68] Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Woodworth. Lower bounds for non-convex stochastic optimization. *arXiv preprint arXiv:1912.02365*, 2019.
- [69] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [70] Kevin Scaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *international conference on machine learning*, pages 3027–3036. PMLR, 2017.
- [71] Ji Liu and Ce Zhang. Distributed learning systems with first-order methods. *arXiv preprint arXiv:2104.05245*, 2021.

- [72] Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Mike Rabbat. Stochastic gradient push for distributed deep learning. In *International Conference on Machine Learning*, pages 344–353. PMLR, 2019.
- [73] Saeed Ghadimi and Guanhui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [74] Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. In *Doklady an ussr*, volume 269, pages 543–547, 1983.
- [75] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [76] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- [77] Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. *arXiv preprint arXiv:1806.01811*, 2018.
- [78] Matthew D Zeiler. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- [79] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of markov chain monte carlo*. CRC press, 2011.
- [80] Balázs Gerencsér. Markov chain mixing time on cycles. *Stochastic processes and their applications*, 121(11):2553–2570, 2011.
- [81] David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [82] Leonidas Georgopoulos. Definitive consensus for distributed data inference. Technical report, EPFL, 2011.
- [83] Chih-Kai Ko. *On matrix factorization and scheduling for finite-time average-consensus*. PhD thesis, California Institute of Technology, 2010.
- [84] Julien M Hendrickx, Raphaël M Jungers, Alexander Olshevsky, and Guillaume Vankeerberghen. Graph diameter, eigenvalues, and minimum-time consensus. *Automatica*, 50(2):635–640, 2014.
- [85] Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi. Don’t use large mini-batches, use local sgd. *arXiv preprint arXiv:1808.07217*, 2018.
- [86] Raphaël Berthier, Francis Bach, and Pierre Gaillard. Accelerated gossip in networks of given dimension using jacobi polynomial iterations. *SIAM Journal on Mathematics of Data Science*, 2(1):24–47, 2020.
- [87] Ji Liu and A Stephen Morse. Accelerated linear iterations for distributed averaging. *Annual Reviews in Control*, 35(2):160–165, 2011.
- [88] Haishan Ye, Luo Luo, Ziang Zhou, and Tong Zhang. Multi-consensus decentralized accelerated gradient descent. *arXiv preprint arXiv:2005.00797*, 2020.

Supplementary Material

A Experimental Details

A.1 Hyperparameter Tuning

In the experiment of training LeNet on CIFAR10, we tune the step size using grid search inside the following range: $\{5\text{e-}3, 1\text{e-}3, 5\text{e-}4, 2.5\text{e-}4, 1\text{e-}4, 5\text{e-}5\}$. Note that this range is in general smaller than the one chosen in [42], since here we are working with unshuffled data, and we found original range in baselines causes algorithms to diverge easily. Following [24], we let each run warm up for 10 epochs with step size $1\text{e-}5$. For DeTAG, we further tune the accelerated gossip parameter η within $\{0, 0.1, 0.2, 0.4\}$ and phase length R within $\{1, 2, 3\}$. We fix the momentum term to be 0.9 and weight decay to be $1\text{e-}4$.

In the experiment of training Resnet20 on CIFAR100, we tune the step size using grid search inside the following range: $\{0.5, 0.1, 0.05, 0.01, 0.005\}$. For DeTAG, we further tune the accelerated gossip parameter η within $\{0, 0.1, 0.2, 0.4\}$ and phase length R within $\{1, 2, 3\}$. We fix the momentum term to be 0.9 and weight decay to be $5\text{e-}4$.

The hyperparameters adopted for each runs are shown in Table 3 and Table 4.

A.2 Techniques of Running DeTAG

We can see in the main loop of DeTAG, several gradient queries are made at the same point. This essentially is equivalent to a large mini-batch size. In practice, however, we can modify this to use local-steps and get better empirical results [85]. Another technique is to use warm-up epochs when data is decentralized. We observe it ensures a smooth convergence in practice. Last but not least, since at first the noise in the algorithms is generally large, we can use a dynamic phase length to obtain better results. That is, we start from phase length 1 for the first few epochs, and let DeTAG follow the special case of DSGT. Then we can gradually increase the phase length following given policies. The intuition is that as algorithm converges, we would need less noise from communication, and thus a longer phase length can benefit.

Table 3: (Initial) Step size α used for each experiments.

| Experiment | Setting | Algorithm | | | |
|-------------------|-----------------|-----------|----------------|--------|-------|
| | | D-PSGD | D ² | DSGT | DeTAG |
| LeNet/CIFAR10 | 100% Shuffled | 5e-3 | 5e-3 | 5e-3 | 5e-3 |
| | 50% Shuffled | 5e-5 | 2.5e-4 | 2.5e-4 | 5e-4 |
| | 25% Shuffled | 5e-5 | 1e-4 | 2.5e-4 | 5e-4 |
| | 0% Shuffled | 5e-5 | 1e-4 | 2.5e-4 | 5e-4 |
| Resnet20/CIFAR100 | $\kappa = 1$ | 0.5 | 0.5 | 0.5 | 0.5 |
| | $\kappa = 0.1$ | 0.5 | 0.5 | 0.5 | 0.5 |
| | $\kappa = 0.05$ | 0.5 | 0.5 | 0.5 | 0.5 |
| | $\kappa = 0.01$ | 0.5 | 0.5 | 0.5 | 0.5 |

Table 4: DeTAG-specific hyperparameters used for each experiments.

| Experiment | Setting | Accelerate Factor η | Phase Length R |
|-------------------|-----------------|--------------------------|------------------|
| LeNet/CIFAR10 | 100% Shuffled | 0 | 1 |
| | 50% Shuffled | 0.2 | 2 |
| | 25% Shuffled | 0.2 | 2 |
| | 0% Shuffled | 0.2 | 2 |
| Resnet20/CIFAR100 | $\kappa = 1$ | 0 | 1 |
| | $\kappa = 0.1$ | 0.2 | 2 |
| | $\kappa = 0.05$ | 0.2 | 2 |
| | $\kappa = 0.01$ | 0.4 | 2 |

B Technical Proof

B.1 Proof to Theorem 1

Proof To prove this theorem, it suffices for us to provide two examples, each has a (set of) loss function $f \in \mathcal{F}_{\Delta, L}$, a set of underlying oracles $O \in \mathcal{O}_{\sigma^2}$, a graph $G \in \mathcal{G}_{n, D}$, such that $\inf_{A \in \mathcal{A}_B} T_\epsilon(A, f, O, G)$ is lower bounded by $\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$ and $\Omega\left(\frac{\Delta L D}{\epsilon^2}\right)$ iterations on these two examples, respectively. Then we will obtain the final bound as $\max\left\{\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right), \Omega\left(\frac{\Delta L D}{\epsilon^2}\right)\right\}$, i.e., $\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4} + \frac{\Delta L D}{\epsilon^2}\right)$ as desired. For simplicity, we denote $\mathbf{z}^{(i)}$ as the i -th coordinate of vector $\mathbf{z} \in \mathbb{R}^d$.

For each setting, our constructions contain three main steps.

(1) The first step is to follow the construction of a zero chain function model [65, 66]. Following [68] and define

$$\text{prog}(\mathbf{z}) = \max\{i \geq 0 | \mathbf{z}^{(i)} \neq 0\}, \forall \mathbf{z} \in \mathbb{R}^d. \quad (18)$$

A zero chain function f has the following property:

$$\text{prog}(\nabla f(\mathbf{x})) \leq \text{prog}(\mathbf{x}) + 1, \quad (19)$$

that means, for a model start from $\mathbf{x} = \mathbf{0}$, a single gradient evaluation can only make at most one more coordinate to be non-zero. The name of "chain" comes from the fact that the adjacent coordinates are linked like a chain and only if the previous coordinate becomes non-zero that the current coordinate can become non-zero via a gradient update. Consider a model with d dimension, if we show that $\|\nabla f(\mathbf{x})\| \geq \epsilon$ for any $\mathbf{x} \in \mathbb{R}^d$ with $\mathbf{x}^{(d)} = 0$, we will obtain d as a lower bound on the gradient calls to obtain the ϵ -stationary point. We refer such sequential lower bound as T_0 .

(2) Step two is to construct a graph $G \in \mathcal{G}_{n, D}$ and a set of oracle $O \in \mathcal{O}_{\sigma^2}$. To do this, our basic idea is to follow [68] and introduce randomness on the $\text{prog}(\mathbf{x})$, and thus the whole chain only make progress with probability p . As will be shown later, this requires $\Omega(T_0/p)$ iterations in total.

(3) The third and last step is to rescale the function and distribution so as to make it belong to the function and oracle classes we consider. In other words, this step is to guarantee the result is shown in terms of Δ , L , σ , n and D .

We start from a smooth and (potentially) non-convex zero chain function \hat{f} [66] as defined below:

$$\hat{f}(\mathbf{x}) = -\Psi(1)\Phi(\mathbf{x}^{(1)}) + \sum_{i=1}^{T-1} [\Psi(-\mathbf{x}^{(i)})\Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)})\Phi(\mathbf{x}^{(i+1)})], \quad (20)$$

where for $\forall \mathbf{z} \in \mathbb{R}$

$$\Psi(z) = \begin{cases} 0 & z \leq 1/2 \\ \exp\left(1 - \frac{1}{(2z-1)^2}\right) & z > 1/2 \end{cases}, \quad \Phi(z) = \sqrt{e} \int_{-\infty}^z e^{\frac{1}{2}t^2} dt. \quad (21)$$

This function, as shown in previous works [66, 68], is a zero-chain function and thus is generally "hard" to optimize: it costs at least T gradient evaluations to find a stationary point. We summarize some properties of Equation (20) as the following (Proof can be found in Lemma 2 in [68]):

1. $\hat{f}(\mathbf{x}) - \inf_{\mathbf{x}} \hat{f}(\mathbf{x}) \leq \Delta_0 T$, $\forall \mathbf{x} \in \mathbb{R}^d$, where $\Delta_0 = 12$.
2. \hat{f} is l_1 -smooth, where $l_1 = 152$.
3. $\forall \mathbf{x} \in \mathbb{R}^T$, $\|\nabla \hat{f}(\mathbf{x})\|_{\infty} \leq G_{\infty}$, where $G_{\infty} = 23$.
4. $\forall \mathbf{x} \in \mathbb{R}^T$, if $\text{prog}(\mathbf{x}) < T$, then $\|\hat{f}(\mathbf{x})\|_{\infty} \geq 1$.

(Setting 1) Next we discuss the first setting with lower bound $\Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right)$. (Setting 1, Step 1) The loss functions are defined as

$$\hat{f}_i(\mathbf{x}) = \hat{f}(\mathbf{x}), \quad (22)$$

note that $1/n \sum_{i=1}^n \hat{f}_i = \hat{f}$. It can be seen from Property 2 that all the \hat{f}_i are l_1 -smooth. (Setting 1, Step 2) For this setting we consider complete graph. We construct the oracle on worker i as the following:

$$[\hat{g}_i(\mathbf{x})]_j = \nabla_j \hat{f}_i(\mathbf{x}) \cdot \left(1 + \mathbb{1}\{j > \text{prog}(\mathbf{x})\} \left(\frac{z}{p} - 1\right)\right), \quad (23)$$

where $z \sim \text{Bernoulli}(p)$. It can be seen that

$$\mathbb{E}[\hat{g}_i(\mathbf{x})] = \nabla \hat{f}_i(\mathbf{x}), \quad (24)$$

and from Property 3 we know

$$\mathbb{E}\|\hat{g}_i(\mathbf{x}) - \nabla \hat{f}_i(\mathbf{x})\|^2 = |\nabla_{\text{prog}(\mathbf{x})+1} \hat{f}(\mathbf{x})|^2 \mathbb{E}\left(\frac{z}{p} - 1\right)^2 \leq \frac{\|\nabla \hat{f}_i(\mathbf{x})\|_{\infty}^2 (1-p)}{p} \leq \frac{\|\nabla \hat{f}(\mathbf{x})\|_{\infty}^2 (1-p)}{p} \leq \frac{G_{\infty}^2 (1-p)}{p}.$$

(Setting 1, Step 3) Finally we rescale each function as $f_i = L\lambda^2/l_1 \hat{f}_i(\mathbf{x}/\lambda)$ where λ is a parameter subject to change. For L : note that all f_i are $\frac{L}{l_1} \cdot l_1 = L$ -smooth. For the Δ ,

$$f - f^* = \frac{L\lambda^2}{l_1} (\hat{f} - \hat{f}^*) = \frac{L\lambda^2 \Delta_0 T}{l_1} \leq \Delta. \quad (25)$$

For the oracle, to be consistent with f_i , we rescale it as $g_i(\mathbf{x}) = L\lambda/l_1 \hat{g}_i(\mathbf{x}/\lambda)$, and we have

$$\mathbb{E}\|g_i(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2 \leq \frac{L^2 \lambda^2}{l_1^2} \mathbb{E}\left\|g_i\left(\frac{\mathbf{x}}{\lambda}\right) - \nabla \hat{f}_i\left(\frac{\mathbf{x}}{\lambda}\right)\right\|^2 \leq \frac{L^2 \lambda^2 G_{\infty}^2 (1-p)}{l_1^2 p} \leq \sigma^2. \quad (26)$$

We assign $\lambda = 2l_1\epsilon/L$, then Equation (25) and (26) are fulfilled with

$$T = \left\lfloor \frac{\Delta}{\Delta_0 l_1 (2\epsilon)^2} \right\rfloor, \\ p = \min\{(2G_{\infty}\epsilon)^2/\sigma^2, 1\}.$$

Take $\delta = 1/2$ in Lemma 2, we have for probability at least $1/2$, $\|\nabla f(\hat{\mathbf{x}}^{(t)})\| \geq \epsilon$ for all $t \leq \frac{T + \log(\delta)}{\min\{nBp, 1\}(e-1)}$. Use Property 4, for any $\mathbf{x} \in \mathbb{R}^T$ such that $\text{prog}(\mathbf{x}) < T$ it holds that $\|\nabla f(\mathbf{x})\| \geq 2\epsilon$, therefore,

$$\mathbb{E}\|\nabla f(\hat{\mathbf{x}}_T)\| > \epsilon. \quad (27)$$

Then with small ϵ it follows that

$$T_{\epsilon}(A, f, O, G) \geq \frac{T-1}{nBp(e-1)} \geq \Omega\left(\frac{\Delta L \sigma^2}{n B \epsilon^4}\right), \quad (28)$$

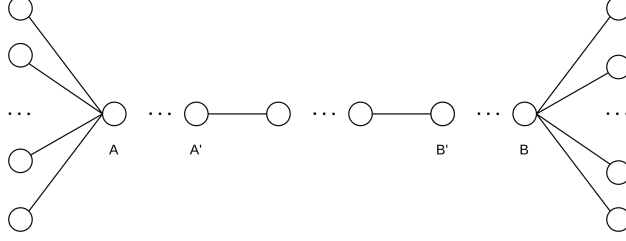


Figure 4: Illustration graph for setting 2 to in the proof of Theorem 1.

and that completes the proof for setting 1.

(Setting 2) We proceed to prove second bound $\Omega\left(\frac{\Delta L D}{\epsilon^2}\right)$.

(Setting 2 Step 1 & Step 2) We assign all the workers with index from 1 to n , we first define two indices set

$$\begin{aligned} I_0 &= \{1, \dots, |I_0|\}, \\ I_1 &= \{n, n-1, \dots, n - |I_1| + 1\}. \end{aligned} \quad (29)$$

where $|\cdot|$ denotes a cardinality of a set. Consider the construction of G in Figure 4:

If $D \geq n - 2\lceil n/3 \rceil + 2$, then it implies the number of nodes between A and B is larger than $\lceil n/3 \rceil$. In this case, denote A' and B' as a sub linear graph where its number of nodes is exactly $\lceil n/3 \rceil$. Let all the nodes on the left of A' be in I_0 and all the nodes on the right of B' be I_1 . We define all the local functions on such graph as following:

$$\hat{f}_i(\mathbf{x}) = \begin{cases} -\frac{2n}{n-\lceil n/3 \rceil} \Psi(1) \Phi(\mathbf{x}^{(1)}) + \sum_{i=2k, k \in \{1, 2, \dots\}, i < T} \frac{2n}{n-\lceil n/3 \rceil} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_0, \\ \sum_{i=2k-1, k \in \{1, 2, \dots\}, i < T} \frac{2n}{n-\lceil n/3 \rceil} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_1, \\ 0 & i \notin I_0, I_1. \end{cases} \quad (30)$$

If $D < n - 2\lceil n/3 \rceil + 2$, the distance between node A and node B is $D - 2$ and the sub linear graph whose end points are A and B contains $D - 1$ nodes. We let the number of nodes on the left of A be $\lceil \frac{n-D+1}{2} \rceil$, we denote the set of indices of all such nodes as I_0 ; and then we let the number of nodes on the right of B be $\lfloor \frac{n-D+1}{2} \rfloor$, we denote the set of indices of all such nodes as I_1 . Since $D < n - 2\lceil n/3 \rceil + 2$, this implies $|I_0|, |I_1| > n/3$. We define all the local functions on such graph as following:

$$\hat{f}_i(\mathbf{x}) = \begin{cases} -\frac{n}{|I_0|} \Psi(1) \Phi(\mathbf{x}^{(1)}) + \sum_{i=2k, k \in \{1, 2, \dots\}, i < T} \frac{n}{|I_0|} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_0, \\ \sum_{i=2k-1, k \in \{1, 2, \dots\}, i < T} \frac{n}{|I_1|} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_1, \\ 0 & i \notin I_0, I_1. \end{cases} \quad (31)$$

In both cases discussed based on D , we can see that $\hat{f}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \hat{f}_i(\mathbf{x})$, and we are splitting hard zero-chain function into two main different part: the even components of the chain and the odd components of the chain. It is easy to see that for the zero chain function to make progress, it takes at least $\lceil n/3 \rceil$, i.e., $\Omega(D)$ number of iterations in the first case (since here $D = \tilde{\gamma}n$ for some $\tilde{\gamma} > 1/3$) and D number of iterations in the second case. Then the total number of iterations is lower bounded by $\Omega(TD)$.

For the oracle, we let oracle on worker i as

$$[\hat{g}_i(\mathbf{x})]_j = \nabla_j \hat{f}_i(\mathbf{x}). \quad (32)$$

(Setting 2, Step 3) The last step is to rescale the parameters. Compared to setting 1, we know here all the \hat{f}_i are $3l_1$ -smooth, as before we let

$$f_i(\mathbf{x}) = \frac{L\lambda^2}{3l_1} \hat{f}_i\left(\frac{\mathbf{x}}{\lambda}\right), \quad \lambda = \frac{6l_1\epsilon}{L}. \quad (33)$$

For the Δ bound we have

$$L\lambda^2\Delta_0T/3l_1 \leq \Delta \quad (34)$$

to fulfill this it suffices to set

$$T = \left\lfloor \frac{\Delta L}{\Delta_0 l_1 (12\epsilon)^2} \right\rfloor. \quad (35)$$

It also can be seen that f is L -smooth. So in this setting,

$$T_\epsilon(A, f, O, G) \geq \Omega(TD) = \Omega\left(\frac{\Delta LD}{\epsilon^2}\right). \quad (36)$$

Combining Setting 1 and 2 we complete the proof.

Lemma 2 In setting 1 in the proof of Theorem 1, with probability at least $1 - \delta$, $\|\nabla f(\mathbf{x}_t)\| \geq \epsilon$ for all $t \leq \frac{T + \log(\delta)}{\min\{nBp, 1\}(e-1)}$.

Proof Define a filtration at iteration t as the sigma field of all the previous events happened before iteration t . Let $i_j^{(t)} = \text{prog}(\mathbf{x}_{t,j}), \forall j \in [n]$ and $i^{(t)} = \max_j i_j^{(t)}$. And we denote $\mathcal{E}^{(t,m,j)}$ as the event of the $i_m^{(t)} + 1$ -th coordinate of output of j -th query on worker m at iteration t is non-zero. Based on the independent sampling, these events are independent. Thus we know:

$$\mathbb{P}[i^{(t+1)} - i^{(t)} = 1 | \mathcal{U}^{(t)}] = \mathbb{P}\left[\bigcup_{\substack{i \in [n] \\ j \leq B}} \mathcal{E}^{(t,i,j)} | \mathcal{U}^{(t)}\right] \leq \sum_{i \in [n], j \leq B} \mathbb{P}[\mathcal{E}^{(t,i,j)} | \mathcal{U}^{(t)}] \leq \min\{nBp, 1\}. \quad (37)$$

Let $q^{(t)} = i^{(t+1)} - i^{(t)}$, with Chernoff bound, we obtain

$$\mathbb{P}[i^{(t)} \geq T] = \mathbb{P}[e^{\sum_{j=0}^{t-1} q^{(j)}} \geq e^T] \leq e^{-T} \mathbb{E}[e^{\sum_{j=0}^{t-1} q^{(j)}}]. \quad (38)$$

For the expectation term we know that

$$\mathbb{E}[e^{\sum_{j=0}^{t-1} q^{(j)}}] = \mathbb{E}\left[\prod_{j=0}^{t-1} \mathbb{E}[e^{q^{(j)}} | \mathcal{U}^{(j)}]\right] \leq (1 - \min\{nBp, 1\} + \min\{nBp, 1\}e)^t \leq e^{\min\{nBp, 1\}t(e-1)}. \quad (39)$$

Thus we know

$$\mathbb{P}[i^{(t)} \geq T] \leq e^{(e-1)\min\{nBp, 1\}t-T} \leq \delta, \quad (40)$$

for every $t \leq \frac{T + \log(\delta)}{\min\{nBp, 1\}(e-1)}$.

B.2 Proof to Corollary

Proof The proof of the Corollary consists of two parts: In the first part, we first prove given any n , how to construct the graph and the gossip matrix with $\lambda = 0, \cos(\pi/n)$. Then we proceed to discuss how the other $\lambda \in (0, \cos(\pi/n))$ can be achieved.

We start from the first part of proving the lower bound. We propose two settings of construction. The first setting is the same as the one shown in the proof of Theorem 1 as the complete graph. For setting 2, consider the two special cases of dumbbell graph in Figure 4: if the graph is fully connected, then we can use the average consensus matrix \mathbf{W} that fulfills $\lambda = 0$, and this can be seen as the special case where $D = 1$; on the other hand, if the graph is a linear graph, we can first prove the lower bound using the diameter as follows:

(Linear graph, Step 1) We first let $|I_0| = |I_1| = \lceil n/3 \rceil$ in the proof of Theorem 1, meaning I_0 denotes the first $\lceil n/3 \rceil$ workers and I_1 denotes the last $\lceil n/3 \rceil$ workers. We define all the local functions as following:

$$\hat{f}_i(\mathbf{x}) = \begin{cases} -\frac{n}{\lceil n/3 \rceil} \Psi(1) \Phi(\mathbf{x}^{(1)}) + \sum_{i=2k, k \in \{1, 2, \dots\}, i < T} \frac{n}{\lceil n/3 \rceil} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_0, \\ \sum_{i=2k-1, k \in \{1, 2, \dots\}, i < T} \frac{n}{\lceil n/3 \rceil} [\Psi(-\mathbf{x}^{(i)}) \Phi(-\mathbf{x}^{(i+1)}) - \Psi(\mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i+1)})] & i \in I_1, \\ 0 & i \notin I_0, I_1. \end{cases} \quad (41)$$

we can see that $\hat{f}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \hat{f}_i(\mathbf{x})$. (Linear graph, Step 2) We consider linear graph in this setting and from one end to the other, the worker's index is 1 to n , without the loss of generality. It is easy to see that for the zero chain function to make progress, it takes at least $n - 2\lceil n/3 \rceil + 1$ number of iterations. Note that in linear graph $n - 1 = D$, the total number of iterations is at least

$$\Omega(TD). \quad (42)$$

For the oracle, we let oracle on worker i as

$$[\hat{g}_i(\mathbf{x})]_j = \nabla_j \hat{f}_i(\mathbf{x}) \quad (43)$$

(Linear graph, Step 3) The last step is to rescale the parameters. Compared to setting 1, we know here all the \hat{f}_i are $3l_1$ -smooth, as before we let

$$f_i(\mathbf{x}) = \frac{L\lambda^2}{3l_1} \hat{f}_i\left(\frac{\mathbf{x}}{\lambda}\right), \quad \lambda = \frac{6l_1\epsilon}{L}. \quad (44)$$

For the Δ bound we have

$$L\lambda^2\Delta_0T/3l_1 \leq \Delta, \quad (45)$$

to fulfill this it suffices to set

$$T = \left\lfloor \frac{\Delta L}{\Delta_0 l_1 (12\epsilon)^2} \right\rfloor. \quad (46)$$

It also can be seen that f is L -smooth. So in this setting,

$$T_\epsilon(A, f, O, G) \geq \Omega(TD) \geq \Omega\left(\frac{\Delta LD}{\epsilon^2}\right). \quad (47)$$

Given the bound, we use two additional results on linear graph as [86]: the random walk matrix \mathbf{W}_{rw} on linear graph with λ fulfilling

$$\frac{1}{\sqrt{1-\lambda}} = O(D). \quad (48)$$

Then we can rewrite the lower bound in the form of λ as shown in Corollary 1.

So far we've proved the two boundary cases, now the dumbbell graph can be seen as the intermediates between the two boundary cases. Take any dumbbell graph as shown in Figure 4 and performs a random walk, we can easily see the second largest eigenvalue for that random walk is some $0 < \lambda' < \lambda = \cos(\pi/n)$, where $\cos(\pi/n)$ is the second largest eigenvalue for the linear graph for $n \geq 2$. Denote all these λ' associated with different dumbbell graph as $\{\tilde{\lambda}_1, \dots, \tilde{\lambda}_K\}$, then for any $0 \leq \tilde{\lambda}_i < \lambda < \tilde{\lambda}_{i+1} \leq \cos(\pi/n)$, a corresponding matrix \mathbf{W} with second largest eigenvalue λ can always be achieved by performing linear combination of the matrices with second largest eigenvalue $\tilde{\lambda}_i$ and $\tilde{\lambda}_{i+1}$.

Finally, using the conclusion of $\lambda = \cos(\pi/n)$ for $n \in \{2, 3, \dots\}$ on linear graph we complete the proof.

B.3 Proof to Theorem 2

Proof As (partially) discussed in the paper, DeFacto is statistically equivalent to centralized SGD. Specifically, it conduct $K = T/2R$ gradient steps where each step contains a mini-batch of R at the point of $\mathbf{x}_{k,i}, \forall i \in [n]$. Take the well-known convergence rate for centralized SGD:

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\hat{\mathbf{x}})\|^2 \leq O\left(\frac{\Delta L\sigma}{\sqrt{nBT}} + \frac{\Delta L}{T}\right). \quad (49)$$

The convergence rate of DeFacto can be expressed as:

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(\hat{\mathbf{x}})\|^2 \leq O\left(\frac{\Delta L\sigma/\sqrt{R}}{\sqrt{nBK}} + \frac{\Delta L}{K}\right) = O\left(\frac{\Delta L\sigma}{\sqrt{nBT}} + \frac{\Delta LR}{T}\right) = O\left(\frac{\Delta L\sigma}{\sqrt{nBT}} + \frac{\Delta LD}{T}\right), \quad (50)$$

then we obtain for DeFacto, when $T = O(\Delta L \sigma^2 (nB\epsilon^4)^{-1} + \Delta LD\epsilon^{-2})$,

$$\min_{t=0,1,\dots,T-1} \mathbb{E} \|\nabla f(\hat{\mathbf{x}})\| \leq \sqrt{\min_{t=0,1,\dots,T-1} \mathbb{E} \|\nabla f(\hat{\mathbf{x}})\|^2} \leq \sqrt{O\left(\frac{\Delta L \sigma}{\sqrt{nBT}} + \frac{\Delta LD}{T}\right)} \leq \epsilon, \quad (51)$$

that completes the proof.

B.4 Proof to Theorem 3

Proof In this proof, we adopt an updated version of notation: we denote at the beginning of phase k , the three quantities of interests are \mathbf{X}_k , \mathbf{Y}_k and $\tilde{\mathbf{G}}_k$, and the update rule becomes:

$$\mathbf{Y}_{k+1} = \mathcal{M}(\mathbf{Y}_k + \tilde{\mathbf{G}}_k - \tilde{\mathbf{G}}_{k-1}), \quad (52)$$

$$\mathbf{X}_{k+1} = \mathcal{M}(\mathbf{X}_k - \alpha \mathbf{Y}_k), \quad (53)$$

with

$$\tilde{\mathbf{G}}_{k+1} = [\nabla \tilde{f}_1(\mathbf{x}_{k,1}), \dots, \nabla \tilde{f}_n(\mathbf{x}_{k,n})] \in \mathbb{R}^{d \times n}, \quad (54)$$

$$\mathbf{G}_{k+1} = [\nabla f_1(\mathbf{x}_{k,1}), \dots, \nabla f_n(\mathbf{x}_{k,n})] \in \mathbb{R}^{d \times n}, \quad (55)$$

$$\mathbf{X}_k = [\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,n}] \in \mathbb{R}^{d \times n}, \quad (56)$$

$$\mathbf{Y}_k = [\mathbf{y}_{k,1}, \dots, \mathbf{y}_{k,n}] \in \mathbb{R}^{d \times n}, \quad (57)$$

where $\nabla \tilde{f}_i$ denotes the stochastic gradient oracle on worker i , and ∇f_i denotes the full gradient oracle on worker i . We use $\bar{\mathbf{X}}$ denote $\mathbf{X} \frac{1}{n}$ for any matrix \mathbf{X} with appropriate shape. We use $\lambda_i(\mathbf{W})$ to denote the i -th general largest eigenvalue of matrix \mathbf{W} . Under such notation, λ in the main paper is equivalent to $\lambda_2(\mathbf{W})$. We use $\mathcal{M}(\cdot)$ to denote the R -step accelerated gossip which has the following property [87]:

$$\|\mathcal{M}(\mathbf{X}) - \bar{\mathbf{X}}\| \leq \rho \|\mathbf{X} - \bar{\mathbf{X}}\|; \quad \mathcal{M}(\mathbf{X}) \frac{1}{n} = \mathbf{X} \frac{1}{n}, \quad (58)$$

where $\rho = \left(1 - \sqrt{1 - \lambda_2(\mathbf{W})}\right)^R$. The proof to the statement of Equation (58) can be found in [88].

For the stochastic oracle, based on the oracle class assumption, we have

$$\mathbb{E} \|\nabla \tilde{f}_i(\mathbf{x}) - \nabla f_i(\mathbf{x})\|^2 \leq \sigma^2, \quad (59)$$

and we denote $\tilde{\sigma}^2 = \frac{\sigma^2}{BR}$ as the variance of mini-batch of R .

First, from the update rule of DeTAG,

$$\bar{\mathbf{Y}}_k = \mathcal{M}(\mathbf{Y}_{k-1} + \tilde{\mathbf{G}}_{k-1} - \tilde{\mathbf{G}}_{k-2}) \frac{1}{n} = \bar{\mathbf{Y}}_{k-1} + \bar{\tilde{\mathbf{G}}}_{k-1} - \bar{\tilde{\mathbf{G}}}_{k-2} = \bar{\mathbf{Y}}_{-1} + \sum_{j=-1}^{k-1} (\bar{\tilde{\mathbf{G}}}_j - \bar{\tilde{\mathbf{G}}}_{j-1}) = \bar{\tilde{\mathbf{G}}}_{k-1} \quad (60)$$

and

$$\bar{\mathbf{X}}_{k+1} = \mathcal{M}(\mathbf{X}_k - \alpha \mathbf{Y}_k) \frac{1}{n} = \bar{\mathbf{X}}_k - \alpha \bar{\mathbf{Y}}_k. \quad (61)$$

By Taylor Theorem, we obtain

$$\mathbb{E} f(\bar{\mathbf{X}}_{k+1}) = \mathbb{E} f(\bar{\mathbf{X}}_k - \alpha \bar{\mathbf{Y}}_k) \quad (62)$$

$$\leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{Y}}_k \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{Y}}_k\|^2 \quad (63)$$

$$\stackrel{(60)}{=} \mathbb{E} f(\bar{\mathbf{X}}_k) - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{G}}_{k-1} \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2. \quad (64)$$

For the last term, we have

$$\mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 = \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \mathbb{E} \|\bar{\mathbf{G}}_{k-1} - \tilde{\mathbf{G}}_{k-1}\|^2 + 2\mathbb{E} \langle \bar{\mathbf{G}}_{k-1}, \bar{\mathbf{G}}_{k-1} - \tilde{\mathbf{G}}_{k-1} \rangle \quad (65)$$

$$= \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \mathbb{E} \|\bar{\mathbf{G}}_{k-1} - \tilde{\bar{\mathbf{G}}}_{k-1}\|^2 \quad (66)$$

$$= \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \|\mathbf{G}_{k-1} \mathbf{e}_i - \tilde{\mathbf{G}}_{k-1} \mathbf{e}_i\|^2 \quad (67)$$

$$\leq \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \frac{\tilde{\sigma}^2}{n}, \quad (68)$$

where in the second step, we use the fact that the sampling noise is independent of the gradient itself. Putting it back we obtain

$$\mathbb{E} f(\bar{\mathbf{X}}_{k+1}) \leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{G}}_{k-1} \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \frac{\alpha^2 \tilde{\sigma}^2 L}{2n} \quad (69)$$

$$= \mathbb{E} f(\bar{\mathbf{X}}_k) - \frac{\alpha}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 - \frac{\alpha - \alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-1}\|^2 + \frac{\alpha^2 \tilde{\sigma}^2 L}{2n} + \frac{\alpha}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-1} - \nabla f(\bar{\mathbf{X}}_k)\|^2, \quad (70)$$

where the last step we use $2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2$. Expand the last term, we obtain

$$\mathbb{E} \|\bar{\mathbf{G}}_{k-1} - \nabla f(\bar{\mathbf{X}}_k)\|^2 \quad (71)$$

$$\leq 2\mathbb{E} \|\bar{\mathbf{G}}_{k-1} - \bar{\mathbf{G}}_{k+1}\|^2 + 2\mathbb{E} \|\bar{\mathbf{G}}_{k+1} - \nabla f(\bar{\mathbf{X}}_k)\|^2 \quad (72)$$

$$= 2\mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_{k,i}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_{k-2,i}) \right\|^2 + 2\mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_{k,i}) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\bar{\mathbf{X}}_k) \right\|^2 \quad (73)$$

$$\leq \frac{2}{n} \sum_{i=1}^n \mathbb{E} \|\nabla f_i(\mathbf{x}_{k,i}) - \nabla f_i(\mathbf{x}_{k-2,i})\|^2 + \frac{2}{n} \sum_{i=1}^n \mathbb{E} \|\nabla f_i(\mathbf{x}_{k,i}) - \nabla f_i(\bar{\mathbf{X}}_k)\|^2 \quad (74)$$

$$\leq \frac{2L^2}{n} \mathbb{E} \|\mathbf{X}_k - \mathbf{X}_{k-2}\|_F^2 + \frac{2L^2}{n} \mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2. \quad (75)$$

Denote $f(\mathbf{0}) - f^* \leq \Delta$, we obtain

$$\sum_{k=0}^{K-1} \alpha(1 - \alpha L) \|\bar{\mathbf{G}}_k\|^2 + \sum_{k=0}^{K-1} \alpha \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \quad (76)$$

$$\leq 2\Delta + \frac{\alpha^2 \tilde{\sigma}^2 LK}{n} + \frac{2\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + \frac{2\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \mathbf{X}_{k-2}\|_F^2 \quad (77)$$

$$\leq 2\Delta + \frac{\alpha^2 \tilde{\sigma}^2 LK}{n} + \frac{16\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + \frac{6\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{X}}_k \mathbf{1}_n^\top - \bar{\mathbf{X}}_{k-2} \mathbf{1}_n^\top\|_F^2, \quad (78)$$

where in the last step we use

$$\frac{2\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \mathbf{X}_{k-2}\|_F^2 \quad (79)$$

$$\leq \frac{6\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + \frac{6\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_{k-2} - \bar{\mathbf{X}}_{k-2} \mathbf{1}_n^\top\|_F^2 + \frac{6\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{X}}_k \mathbf{1}_n^\top - \bar{\mathbf{X}}_{k-2} \mathbf{1}_n^\top\|_F^2. \quad (80)$$

In addition, for the last term we have

$$\frac{6\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{X}}_k \mathbf{1}_n^\top - \bar{\mathbf{X}}_{k-2} \mathbf{1}_n^\top\|_F^2 = \frac{6\alpha L^2 n}{n} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{X}}_k - \bar{\mathbf{X}}_{k-2}\|^2 \quad (81)$$

$$\stackrel{(61)}{=} \frac{24\alpha^3 L^2 n}{n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \bar{\mathbf{G}}_k \right\|^2 \quad (82)$$

$$\stackrel{(65)}{\leq} 24\alpha^3 L^2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \bar{\mathbf{G}}_k \right\|^2 + \frac{24\alpha^3 \tilde{\sigma}^2 L^2}{n}. \quad (83)$$

Push it back we have

$$\sum_{k=0}^{K-1} \alpha(1 - \alpha L - 24\alpha^2 L^2) \left\| \bar{\mathbf{G}}_k \right\|^2 + \sum_{k=0}^{K-1} \alpha \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2 \quad (84)$$

$$\leq 2\Delta + \frac{\alpha^2 \tilde{\sigma}^2 L K}{n} + \frac{16\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{24\alpha^3 \tilde{\sigma}^2 L^2}{n}. \quad (85)$$

The rest of the proof is to bound $\frac{16\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2$.

We start from

$$\left\| \mathbf{X}_{k+1} - \bar{\mathbf{X}}_{k+1} \mathbf{1}_n^\top \right\|_F^2 \quad (86)$$

$$\stackrel{(61)}{=} \left\| \mathcal{M}(\mathbf{X}_k - \alpha \mathbf{Y}_k) - (\bar{\mathbf{X}}_k - \alpha \bar{\mathbf{Y}}_k) \mathbf{1}_n^\top \right\|_F^2 \quad (87)$$

$$= \left\| \mathcal{M}(\mathbf{X}_k) - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2 - 2\alpha \langle \mathcal{M}(\mathbf{X}_k) - \bar{\mathbf{X}}_k \mathbf{1}_n^\top, \mathcal{M}(\mathbf{Y}_k) - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \rangle + \alpha^2 \left\| \mathcal{M}(\mathbf{Y}_k) - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 \quad (88)$$

$$\stackrel{(58)}{\leq} \rho^2 \left\| \mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{\rho^2(1 - \rho^2)}{1 + \rho^2} \left\| \mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{\rho^2(1 + \rho^2)\alpha^2}{1 - \rho^2} \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 \quad (89)$$

$$+ \alpha^2 \rho^2 \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 \quad (90)$$

$$= \frac{2\rho^2}{(1 + \rho^2)} \left\| \mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{2\rho^2\alpha^2}{1 - \rho^2} \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2, \quad (91)$$

where in the third step we use

$$-2\langle \mathbf{a}, \mathbf{b} \rangle \leq \frac{1 - \rho^2}{1 + \rho^2} \|\mathbf{a}\|^2 + \frac{1 + \rho^2}{1 - \rho^2} \|\mathbf{b}\|^2. \quad (92)$$

Similarly, for \mathbf{Y}_{k+1} , we obtain

$$\mathbb{E} \left\| \mathbf{Y}_{k+1} - \bar{\mathbf{Y}}_{k+1} \mathbf{1}_n^\top \right\|_F^2 \quad (93)$$

$$\stackrel{(60)}{=} \mathbb{E} \left\| \mathcal{M}(\mathbf{Y}_k + \tilde{\mathbf{G}}_k - \tilde{\mathbf{G}}_{k-1}) - (\bar{\mathbf{Y}}_k + \bar{\tilde{\mathbf{G}}}_k - \bar{\tilde{\mathbf{G}}}_{k-1}) \mathbf{1}_n^\top \right\|_F^2 \quad (94)$$

$$= \mathbb{E} \left\| \mathcal{M}(\mathbf{Y}_k) - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 + \mathbb{E} \left\| \mathcal{M}(\tilde{\mathbf{G}}_k - \tilde{\mathbf{G}}_{k-1}) - (\bar{\tilde{\mathbf{G}}}_k - \bar{\tilde{\mathbf{G}}}_{k-1}) \mathbf{1}_n^\top \right\|_F^2 \quad (95)$$

$$+ 2\mathbb{E} \left\langle \mathcal{M}(\mathbf{Y}_k) - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top, \mathcal{M}(\tilde{\mathbf{G}}_k - \tilde{\mathbf{G}}_{k-1}) - (\bar{\tilde{\mathbf{G}}}_k - \bar{\tilde{\mathbf{G}}}_{k-1}) \mathbf{1}_n^\top \right\rangle \quad (96)$$

$$\stackrel{(92)(58)}{\leq} \rho^2 \mathbb{E} \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 + \rho^2 \mathbb{E} \left\| \mathbf{G}_k - \mathbf{G}_{k-1} - (\bar{\mathbf{G}}_k - \bar{\mathbf{G}}_{k-1}) \mathbf{1}_n^\top \right\|_F^2 \quad (97)$$

$$+ \frac{(1 - \rho^2)\rho^2}{1 + \rho^2} \mathbb{E} \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{(1 + \rho^2)\rho^2}{1 - \rho^2} \mathbb{E} \left\| \mathbf{G}_k - \mathbf{G}_{k-1} - (\bar{\mathbf{G}}_k - \bar{\mathbf{G}}_{k-1}) \mathbf{1}_n^\top \right\|_F^2 \quad (98)$$

$$+ 2\rho^2 \mathbb{E} \left\| \mathbf{G}_k - \tilde{\mathbf{G}}_k \right\|_F^2 + 2\rho^2 \mathbb{E} \left\| \mathbf{G}_{k-1} - \tilde{\mathbf{G}}_{k-1} \right\|_F^2 + 2\rho^2 \mathbb{E} \left\| \bar{\mathbf{G}}_k \mathbf{1}_n^\top - \bar{\tilde{\mathbf{G}}}_k \mathbf{1}_n^\top \right\|_F^2 + 2\rho^2 \mathbb{E} \left\| \bar{\mathbf{G}}_{k-1} \mathbf{1}_n^\top - \bar{\tilde{\mathbf{G}}}_{k-1} \mathbf{1}_n^\top \right\|_F^2 \quad (99)$$

$$\leq \frac{2\rho^2}{1 + \rho^2} \mathbb{E} \left\| \mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top \right\|_F^2 + \frac{4\rho^2}{1 - \rho^2} \mathbb{E} \left\| \mathbf{G}_{k+2} - \mathbf{G}_{k+1} \right\|_F^2 \quad (100)$$

$$+ \frac{4\rho^2}{1 - \rho^2} \mathbb{E} \left\| \mathbf{G}_{k+2} - \mathbf{G}_{k+1} - \mathbf{G}_k + \mathbf{G}_{k-1} \right\|_F^2 + 8n\rho^2 \tilde{\sigma}^2, \quad (101)$$

where in the last step we use $\|I - \frac{\mathbf{1}\mathbf{1}^\top}{n}\| \leq 1$ and $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \|\mathbf{B}\|$.

For the second term, we have

$$\mathbb{E} \|\mathbf{G}_{k+2} - \mathbf{G}_{k+1}\|_F^2 \quad (102)$$

$$= \sum_{i=1}^n \mathbb{E} \|\nabla f(\mathbf{x}_{k+1,i}) - \nabla f(\mathbf{x}_{k,i})\|^2 \quad (103)$$

$$\leq L^2 \sum_{i=1}^n \mathbb{E} \|\mathbf{x}_{k+1,i} - \mathbf{x}_{k,i}\|^2 \quad (104)$$

$$= L^2 \mathbb{E} \|\mathbf{X}_{k+1} - \mathbf{X}_k\|_F^2 \quad (105)$$

$$\stackrel{(61)}{=} L^2 \mathbb{E} \|\mathcal{M}(\mathbf{X}_k) - \mathbf{X}_k - \alpha \mathcal{M}(\mathbf{Y}_k)\|_F^2 \quad (106)$$

$$= L^2 \mathbb{E} \|\mathcal{M}(\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top) - (\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top) - \alpha \mathcal{M}(\mathbf{Y}_k)\|_F^2 \quad (107)$$

$$\leq 4L^2 \mathbb{E} \|\mathcal{M}(\mathbf{X}_k) - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + 4L^2 \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + 4\alpha^2 L^2 \mathbb{E} \|\mathcal{M}(\mathbf{Y}_k) - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|_F^2 \quad (108)$$

$$+ 4\alpha^2 n L^2 \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2 \quad (109)$$

$$\leq 4(1 + \rho^2) L^2 \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + 4\alpha^2 \rho^2 L^2 \mathbb{E} \|\mathbf{Y}_k - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|_F^2 + 4\alpha^2 n L^2 \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2. \quad (110)$$

Putting it back we obtain

$$\mathbb{E} \|\mathbf{Y}_{k+1} - \overline{\mathbf{Y}}_{k+1} \mathbf{1}_n^\top\|_F^2 \quad (111)$$

$$\leq \left(\frac{2\rho^2}{1 + \rho^2} + \frac{16\alpha^2 \rho^4 L^2}{1 - \rho^2} \right) \mathbb{E} \|\mathbf{Y}_k - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|_F^2 + \frac{16\rho^2(1 + \rho^2)L^2}{1 - \rho^2} \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + \frac{16\alpha^2 \rho^2 n L^2}{1 - \rho^2} \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2 \quad (112)$$

$$\frac{4\rho^2}{1 - \rho^2} \mathbb{E} \|\mathbf{G}_{k+2} - \mathbf{G}_{k+1} - \mathbf{G}_k + \mathbf{G}_{k-1}\|_F^2 + 8n\rho^2 \tilde{\sigma}^2. \quad (113)$$

Combining Equation (91) and Equation (100), we have

$$\begin{bmatrix} \mathbb{E} \|\mathbf{X}_{k+1} - \overline{\mathbf{X}}_{k+1} \mathbf{1}_n^\top\|_F^2 \\ \mathbb{E} \|\mathbf{Y}_{k+1} - \overline{\mathbf{Y}}_{k+1} \mathbf{1}_n^\top\|_F^2 \end{bmatrix} \preceq \begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{bmatrix} \begin{bmatrix} \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 \\ \mathbb{E} \|\mathbf{Y}_k - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|_F^2 \end{bmatrix} \quad (114)$$

$$+ \begin{bmatrix} 0 \\ \frac{4\rho^2}{1 - \rho^2} \mathbb{E} \|\mathbf{U}_k\|_F^2 + \frac{16\alpha^2 \rho^2 n L^2}{1 - \rho^2} \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2 + 8n\rho^2 \tilde{\sigma}^2 \end{bmatrix}, \quad (115)$$

where

$$\mathbf{P}_{11} = \frac{2\rho^2}{(1 + \rho^2)} \quad (116)$$

$$\mathbf{P}_{12} = \frac{2\rho^2 \alpha^2}{1 - \rho^2} \quad (117)$$

$$\mathbf{P}_{21} = \frac{16\rho^2(1 + \rho^2)L^2}{1 - \rho^2} \quad (118)$$

$$\mathbf{P}_{22} = \frac{2\rho^2}{1 + \rho^2} + \frac{16\alpha^2 \rho^4 L^2}{1 - \rho^2} \quad (119)$$

$$\mathbf{U}_k = \mathbf{G}_{k+2} - \mathbf{G}_{k+1} - \mathbf{G}_k + \mathbf{G}_{k-1}. \quad (120)$$

For simplicity, define

$$\mathbf{z}_k = \begin{bmatrix} \mathbb{E} \|\mathbf{X}_{k+1} - \overline{\mathbf{X}}_{k+1} \mathbf{1}_n^\top\|_F^2 \\ \mathbb{E} \|\mathbf{Y}_{k+1} - \overline{\mathbf{Y}}_{k+1} \mathbf{1}_n^\top\|_F^2 \end{bmatrix} \quad (121)$$

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{11} & \mathbf{P}_{12} \\ \mathbf{P}_{21} & \mathbf{P}_{22} \end{bmatrix} \quad (122)$$

$$\mathbf{u}_k = \begin{bmatrix} 0 \\ \frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_k\|_F^2 + \frac{16\alpha^2\rho^2 n L^2}{1-\rho^2} \mathbb{E} \|\bar{\mathbf{Y}}_k\|^2 + 8n\rho^2 \tilde{\sigma}^2 \end{bmatrix} \quad (123)$$

then we can write this linear system as

$$\mathbf{z}_k \preceq \mathbf{P} \mathbf{z}_{k-1} + \mathbf{u}_{k-1} \preceq \mathbf{P}^k \mathbf{z}_0 + \sum_{t=0}^{k-1} \mathbf{P}^{k-t} \mathbf{u}_t, \quad (124)$$

for simplicity.

Let $\lambda_1(\mathbf{P}), \lambda_2(\mathbf{P})$ denote the two eigenvalues of \mathbf{P} (without the loss of generality, we denote $\lambda_1(\mathbf{P}) < \lambda_2(\mathbf{P})$), define

$$\Psi = \sqrt{(\mathbf{P}_{11} - \mathbf{P}_{22})^2 + 4\mathbf{P}_{12}\mathbf{P}_{21}}, \quad (125)$$

then with eigendecomposition, we obtain

$$\lambda_1(\mathbf{P}) = \frac{\mathbf{P}_{11} + \mathbf{P}_{22} - \Psi}{2} \quad (126)$$

$$\lambda_2(\mathbf{P}) = \frac{\mathbf{P}_{11} + \mathbf{P}_{22} + \Psi}{2} = \frac{2\rho^2}{1+\rho^2} + \frac{8\alpha^2\rho^4 L^2}{1-\rho^2} + \frac{16\alpha\rho^2 L \sqrt{\alpha^2\rho^4 L^2 + (1+\rho^2)}}{1-\rho^2} \quad (127)$$

$$\mathbf{P}^k \preceq \begin{bmatrix} \frac{\lambda_1^k(\mathbf{P}) + \lambda_2^k(\mathbf{P})}{2} + \frac{(\mathbf{P}_{11} - \mathbf{P}_{22})(\lambda_2^k(\mathbf{P}) - \lambda_1^k(\mathbf{P}))}{2\Psi} & \frac{\mathbf{P}_{12}}{\Psi}(\lambda_2^k(\mathbf{P}) - \lambda_1^k(\mathbf{P})) \\ \frac{\mathbf{P}_{21}}{\Psi}(\lambda_2^k(\mathbf{P}) - \lambda_1^k(\mathbf{P})) & \frac{\lambda_1^k(\mathbf{P}) + \lambda_2^k(\mathbf{P})}{2} + \frac{(\mathbf{P}_{11} - \mathbf{P}_{22})(\lambda_1^k(\mathbf{P}) - \lambda_2^k(\mathbf{P}))}{2\Psi} \end{bmatrix}, \quad (128)$$

when the step size is small enough such that

$$\alpha L < \frac{(1-\rho)^2}{32}, \quad (129)$$

it can be verified that $\lambda_2(\mathbf{P}) \leq \frac{\sqrt{\rho} + \rho}{1+\rho}$, and then we can compute the $\mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|^2$ and $\mathbb{E} \|\mathbf{Y}_k - \bar{\mathbf{Y}}_k \mathbf{1}_n^\top\|^2$. We use $\mathbf{X}[1 :]$ to denote the first row of matrix \mathbf{X} . First for \mathbf{X}_k , we obtain:

$$\mathbf{P}^k \mathbf{z}_0[1 :] \leq \mathbf{P}_{12} k \lambda_2^{k-1}(\mathbf{P}) \mathbb{E} \|\mathbf{Y}_0 - \bar{\mathbf{Y}}_0 \mathbf{1}_n^\top\|_F^2 = \frac{2\rho^2 \alpha^2 k}{1-\rho^2} \lambda_2^{k-1}(\mathbf{P}) \mathbb{E} \|\mathbf{Y}_0 - \bar{\mathbf{Y}}_0 \mathbf{1}_n^\top\|_F^2 \quad (130)$$

where we use the property that $\lambda_2^k(\mathbf{P}) - \lambda_1^k(\mathbf{P}) = (\lambda_2(\mathbf{P}) - \lambda_1(\mathbf{P})) \sum_{l=0}^{k-1} \lambda_2(\mathbf{P})^l \lambda_1(\mathbf{P})^{k-1-l} = \Psi k \lambda_2^{k-1}(\mathbf{P})$ and, similarly

$$\mathbf{P}^{k-t} \mathbf{u}_t[1 :] \quad (131)$$

$$\leq \frac{2\rho^2 \alpha^2 (k-t)}{1-\rho^2} \lambda_2^{k-t-1}(\mathbf{P}) \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2 \rho^2 n L^2}{1-\rho^2} \mathbb{E} \|\bar{\mathbf{Y}}_t\|^2 + 8n\rho^2 \tilde{\sigma}^2 \right) \quad (132)$$

$$= \frac{2\rho^2 \alpha^2 (k-t)}{1-\rho^2} \lambda_2^{k-t-1}(\mathbf{P}) \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2 \rho^2 n L^2}{1-\rho^2} \mathbb{E} \|\bar{\bar{\mathbf{G}}}_{t-1}\|^2 + 8n\rho^2 \tilde{\sigma}^2 \right) \quad (133)$$

$$\stackrel{(65)}{\leq} \frac{2\rho^2 \alpha^2 (k-t)}{1-\rho^2} \lambda_2^{k-t-1}(\mathbf{P}) \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2 \rho^2 n L^2}{1-\rho^2} \mathbb{E} \|\bar{\mathbf{G}}_{t-1}\|^2 + \frac{16\alpha^2 \rho^2 \tilde{\sigma}^2 L^2}{1-\rho^2} + 8n\rho^2 \tilde{\sigma}^2 \right), \quad (134)$$

then we obtain

$$\mathbb{E} \|\mathbf{X}_k - \bar{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 \quad (135)$$

$$\leq \frac{2\rho^2 \alpha^2 k}{1-\rho^2} \lambda_2^{k-1}(\mathbf{P}) \mathbb{E} \|\mathbf{Y}_0 - \bar{\mathbf{Y}}_0 \mathbf{1}_n^\top\|_F^2 \quad (136)$$

$$+ \sum_{t=0}^{k-1} \frac{2\rho^2\alpha^2(k-t)}{1-\rho^2} \lambda_2^{k-t-1}(\mathbf{P}) \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{G}}_{t-1}\|^2 + \frac{16\alpha^2\rho^2\tilde{\sigma}^2L^2}{1-\rho^2} + 8n\rho^2\tilde{\sigma}^2 \right). \quad (137)$$

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

$$\sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 \quad (138)$$

$$\leq \frac{2\rho^2\alpha^2}{(1-\rho^2)(1-\lambda_2(\mathbf{P}))^2} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{Y}_0 - \overline{\mathbf{Y}}_0 \mathbf{1}_n^\top\|_F^2 \quad (139)$$

$$+ \frac{2\rho^2\alpha^2}{(1-\rho^2)(1-\lambda_2(\mathbf{P}))^2} \sum_{k=0}^{K-1} \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_k\|_F^2 + \frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{16\alpha^2\rho^2\tilde{\sigma}^2L^2}{1-\rho^2} + 8n\rho^2\tilde{\sigma}^2 \right) \quad (140)$$

$$\leq \frac{2\rho^2\alpha^2(1+\rho)nK\zeta_0^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{32\rho^4\alpha^4nL^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{8\rho^4\alpha^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{U}_k\|_F^2 \quad (141)$$

$$+ \frac{32\rho^4\alpha^4\tilde{\sigma}^2L^2K}{(1-\rho)^2(1-\sqrt{\rho})^2} + \frac{8\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)K}{(1-\rho)(1-\sqrt{\rho})^2} \quad (142)$$

$$\leq \frac{2\rho^2\alpha^2(1+\rho)nK\zeta_0^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{32\rho^4\alpha^4nL^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{8\rho^4\alpha^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{U}_k\|_F^2 \quad (143)$$

$$+ \frac{16\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)K}{(1-\rho)(1-\sqrt{\rho})^2}, \quad (144)$$

where in the second step we used $\frac{1}{1-\lambda_2(\mathbf{P})} < \frac{1+\rho}{1-\sqrt{\rho}}$ since $\lambda_2(\mathbf{P}) \leq \frac{\sqrt{\rho}+\rho}{1+\rho}$. And the third step holds due to Equation (129). We proceed to analyze the case in \mathbf{Y}_k : we first have

$$[\mathbf{P}^k]_{22} = \frac{\lambda_1^k(\mathbf{P}) + \lambda_2^k(\mathbf{P})}{2} + \frac{(\mathbf{P}_{11} - \mathbf{P}_{22})(\lambda_1^k(\mathbf{P}) - \lambda_2^k(\mathbf{P}))}{2\Psi} \quad (145)$$

$$\leq \lambda_2^k(\mathbf{P}) + \frac{8\alpha^2\rho^4L^2k\lambda_2^{k-1}(\mathbf{P})}{1-\rho^2}, \quad (146)$$

then we can have

$$\mathbf{P}^k \mathbf{z}_0[2:] \leq \left(\lambda_2^k(\mathbf{P}) + \frac{8\alpha^2\rho^4L^2k\lambda_2^{k-1}(\mathbf{P})}{1-\rho^2} \right) \mathbb{E} \|\mathbf{Y}_0 - \overline{\mathbf{Y}}_0 \mathbf{1}_n^\top\|_F^2, \quad (147)$$

and

$$\mathbf{P}^{k-t} \mathbf{u}_t[2:] \quad (148)$$

$$\leq \left(\lambda_2^{k-t}(\mathbf{P}) + \frac{8\alpha^2\rho^4L^2(k-t)\lambda_2^{k-t-1}(\mathbf{P})}{1-\rho^2} \right) \cdot \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{Y}}_t\|^2 + 8n\rho^2\tilde{\sigma}^2 \right) \quad (149)$$

$$\leq \left(\lambda_2^{k-t}(\mathbf{P}) + \frac{8\alpha^2\rho^4L^2(k-t)\lambda_2^{k-t-1}(\mathbf{P})}{1-\rho^2} \right) \cdot \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_t\|_F^2 + \frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{G}}_{t-1}\|^2 + \frac{16\alpha^2\rho^2\tilde{\sigma}^2L^2}{1-\rho^2} + 8n\rho^2\tilde{\sigma}^2 \right). \quad (150)$$

Summing over $k = 0$ to $K - 1$, we obtain

$$\sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{Y}_k - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|^2 \quad (151)$$

$$\begin{aligned}
&\leq \frac{(1+\rho)nK\zeta_0^2}{1-\sqrt{\rho}} + \frac{8\alpha^2\rho^4(1+\rho)L^2nK\zeta_0^2}{(1-\rho)(1-\sqrt{\rho})^2} \\
&\quad + \left(\frac{1+\rho}{1-\sqrt{\rho}} + \frac{8\alpha^2\rho^4(1+\rho)L^2}{(1-\rho)(1-\sqrt{\rho})^2} \right) \sum_{k=0}^{K-1} \left(\frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_k\|_F^2 + \frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{16\alpha^2\rho^2\tilde{\sigma}^2L^2}{1-\rho^2} + 8n\rho^2\tilde{\sigma}^2 \right).
\end{aligned} \tag{152}$$

$$\tag{153}$$

We next solve $\|\mathbf{U}_k\|_F$, from the definition of \mathbf{U}_k we obtain that

$$\sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{U}_k\|_F^2 = \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{G}_{k+2} - \mathbf{G}_{k+1} - \mathbf{G}_k + \mathbf{G}_{k-1}\|_F^2 \tag{154}$$

$$\leq 2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{G}_{k+2} - \mathbf{G}_{k+1}\|_F^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{G}_k - \mathbf{G}_{k-1}\|_F^2 \tag{155}$$

$$\leq 4 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{G}_{k+2} - \mathbf{G}_{k+1}\|_F^2 \tag{156}$$

$$= 4 \sum_{k=0}^{K-1} \sum_{i=1}^n \mathbb{E} \|\nabla f(\mathbf{x}_{k+1,i}) - \nabla f(\mathbf{x}_{k,i})\|^2 \tag{157}$$

$$\leq 4L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n \mathbb{E} \|\mathbf{x}_{k+1,i} - \mathbf{x}_{k,i}\|^2 \tag{158}$$

$$= 4L^2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_{k+1} - \mathbf{X}_k\|_F^2. \tag{159}$$

Fit in the derivation from Equation (110) we obtain

$$\sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{U}_k\|_F^2 \tag{160}$$

$$\leq 4L^2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_{k+1} - \mathbf{X}_k\|_F^2 \tag{161}$$

$$\leq 16(1+\rho^2)L^2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + 16\alpha^2\rho^2L^2 \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{Y}_k - \overline{\mathbf{Y}}_k \mathbf{1}_n^\top\|_F^2 + 16\alpha^2nL^2 \sum_{k=0}^{K-1} \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2 \tag{162}$$

$$\leq \frac{32\rho^2\alpha^2(1+\rho)^2nK\zeta_0^2L^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{512\rho^4(1+\rho)\alpha^4nL^4}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{256\rho^4(1+\rho)\alpha^2L^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{U}_k\|_F^2 \tag{163}$$

$$+ \frac{256\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)^2KL^2}{(1-\rho)(1-\sqrt{\rho})^2} \tag{164}$$

$$+ \frac{16\alpha^2\rho^2(1+\rho)nK\zeta_0^2L^2}{1-\sqrt{\rho}} + \frac{128\alpha^4\rho^6(1+\rho)L^4nK\zeta_0^2}{(1-\rho)(1-\sqrt{\rho})^2} \tag{165}$$

$$+ \left(\frac{16\alpha^2\rho^2(1+\rho)L^2}{1-\sqrt{\rho}} + \frac{128\alpha^4\rho^6(1+\rho)L^4}{(1-\rho)(1-\sqrt{\rho})^2} \right) \sum_{k=0}^{K-1} \frac{4\rho^2}{1-\rho^2} \mathbb{E} \|\mathbf{U}_k\|_F^2 \tag{166}$$

$$+ \left(\frac{16\alpha^2\rho^2(1+\rho)L^2}{1-\sqrt{\rho}} + \frac{128\alpha^4\rho^6(1+\rho)L^4}{(1-\rho)(1-\sqrt{\rho})^2} \right) \sum_{k=0}^{K-1} \left(\frac{16\alpha^2\rho^2nL^2}{1-\rho^2} \mathbb{E} \|\overline{\mathbf{G}}_k\|^2 + \frac{16\alpha^2\rho^2\tilde{\sigma}^2L^2}{1-\rho^2} + 8n\rho^2\tilde{\sigma}^2 \right) \tag{167}$$

$$+ 16\alpha^2nL^2 \sum_{k=0}^{K-1} \mathbb{E} \|\overline{\mathbf{Y}}_k\|^2 \tag{168}$$

$$\leq \frac{64\rho^2\alpha^2(1+\rho)^2nK\varsigma_0^2L^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{512\rho^4(1+\rho)\alpha^2L^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E}\|U_k\|_F^2 + \frac{512\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)^2KL^2}{(1-\rho)(1-\sqrt{\rho})^2} \quad (169)$$

$$+ 32\alpha^2nL^2 \sum_{k=0}^{K-1} \mathbb{E}\|\overline{\mathbf{G}}_k\|^2, \quad (170)$$

where in the third step we use the derivation from Equation (138) and (151), in the fourth step we repeatedly use Equation (129) and Equation (65), solve it we obtain

$$\sum_{k=0}^{K-1} \mathbb{E}\|U_k\|_F^2 \leq \frac{128\rho^2\alpha^2(1+\rho)^2nK\varsigma_0^2L^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{1024\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)^2KL^2}{(1-\rho)(1-\sqrt{\rho})^2} + 64\alpha^2nL^2 \sum_{k=0}^{K-1} \mathbb{E}\|\overline{\mathbf{G}}_k\|^2, \quad (171)$$

where again we use Equation (129), combine it with Equation (138) we obtain

$$\sum_{k=0}^{K-1} \mathbb{E}\|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 \quad (172)$$

$$\leq \frac{4\rho^2\alpha^2(1+\rho)nK\varsigma_0^2}{(1-\rho)(1-\sqrt{\rho})^2} + \frac{544\rho^4\alpha^4nL^2}{(1-\rho)^2(1-\sqrt{\rho})^2} \sum_{k=0}^{K-1} \mathbb{E}\|\overline{\mathbf{G}}_k\|^2 + \frac{32\rho^4\alpha^2n\tilde{\sigma}^2(1+\rho)K}{(1-\rho)(1-\sqrt{\rho})^2}, \quad (173)$$

where we use Equation (129).

Recall from Equation (84) that

$$\sum_{k=0}^{K-1} \alpha(1-\alpha L - 24\alpha^2L^2) \|\overline{\mathbf{G}}_k\|^2 + \sum_{k=0}^{K-1} \alpha \mathbb{E}\|\nabla f(\overline{\mathbf{X}}_k)\|^2 \quad (174)$$

$$\leq 2\Delta + \frac{\alpha^2\tilde{\sigma}^2LK}{n} + \frac{16\alpha L^2}{n} \sum_{k=0}^{K-1} \mathbb{E}\|\mathbf{X}_k - \overline{\mathbf{X}}_k \mathbf{1}_n^\top\|_F^2 + \frac{24\alpha^3\tilde{\sigma}^2L^2}{n}. \quad (175)$$

Combine Equation (129) and (173), we obtain

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}\|\nabla f(\overline{\mathbf{X}}_k)\|^2 \leq O\left(\frac{\Delta}{\alpha K} + \frac{\alpha\tilde{\sigma}^2L}{n} + \frac{\rho^2\alpha^2L^2\varsigma_0^2}{(1-\rho)^3} + \frac{\rho^4\alpha^2\tilde{\sigma}^2L^2}{(1-\rho)^3} + \frac{\alpha^2\tilde{\sigma}^2L^2}{nK}\right), \quad (176)$$

where we omit the numerical constants. Set

$$\alpha = \frac{1}{\tilde{\sigma}\sqrt{KL/n\Delta} + \frac{\rho^{\frac{2}{3}}L^{\frac{2}{3}}\varsigma_0^{\frac{2}{3}}K^{\frac{1}{3}}}{\Delta^{\frac{1}{3}}(1-\rho)} + \frac{32L}{(1-\rho)^2}}, \quad (177)$$

we obtain

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}\|\nabla f(\overline{\mathbf{X}}_k)\|^2 \leq O\left(\frac{\sqrt{\Delta L}\tilde{\sigma}}{\sqrt{nK}} + \frac{(\rho\Delta L\varsigma_0)^{\frac{2}{3}}}{(1-\rho)K^{\frac{2}{3}}} + \frac{\rho^2n\Delta L}{(1-\rho)^3K} + \frac{\Delta L}{(1-\rho)^2K}\right). \quad (178)$$

Fit in $T = KR$ and $\tilde{\sigma}^2 = \sigma^2/BR$, we obtain

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}\|\nabla f(\overline{\mathbf{X}}_k)\|^2 \leq O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{(\rho\Delta L\varsigma_0R)^{\frac{2}{3}}}{(1-\rho)T^{\frac{2}{3}}} + \frac{\rho^2nR\Delta L}{(1-\rho)^3T} + \frac{R\Delta L}{(1-\rho)^2T}\right), \quad (179)$$

set

$$R = \frac{1}{\sqrt{1-\lambda_2(\mathbf{W})}} \max\left(\frac{1}{2}\log(n), \frac{1}{2}\log\left(\frac{\varsigma_0^2T}{\Delta L}\right)\right),$$

we first have $\rho \leq 1/\sqrt{2}$ since

$$R \geq \frac{\log(n)}{2\sqrt{1-\lambda_2(\mathbf{W})}} \geq \frac{-\log(n)}{2\log(1-\sqrt{1-\lambda_2(\mathbf{W})})} \Rightarrow \left(1 - \sqrt{1-\lambda_2(\mathbf{W})}\right)^R \leq \frac{1}{\sqrt{n}} \leq \frac{1}{\sqrt{2}}, \quad (180)$$

this implies

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \leq O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{(\rho\Delta L\varsigma_0 R)^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{\rho^2 n R \Delta L}{T} + \frac{R\Delta L}{T}\right) \quad (181)$$

$$\leq O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{(\rho\varsigma_0\sqrt{T}R/\sqrt{\Delta L})^{\frac{2}{3}}\Delta L}{T} + \frac{\rho^2 n R \Delta L}{T} + \frac{R\Delta L}{T}\right), \quad (182)$$

with the assignment of R , $\rho^2 n < 1$ and $\rho\varsigma_0\sqrt{T}/\sqrt{\Delta L} < 1$, so since it also holds that $R \geq 1$ (and so $R^{2/3} \leq R$),

$$\min_t \|\nabla f(\bar{\mathbf{X}}_t)\|^2 \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \quad (183)$$

$$\leq O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{(R)^{\frac{2}{3}}\Delta L}{T} + \frac{R\Delta L}{T} + \frac{R\Delta L}{T}\right) \quad (184)$$

$$\leq O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{R\Delta L}{T}\right) \quad (185)$$

$$= O\left(\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} + \frac{\Delta L}{T\sqrt{1-\lambda_2(\mathbf{W})}} \cdot \max\left(\log(n), \log\left(\frac{\varsigma_0^2 T}{\Delta L}\right)\right)\right), \quad (186)$$

when

$$T \leq O\left(\frac{\Delta L\sigma^2}{nB\epsilon^4}\right), \quad (187)$$

we have

$$\frac{\sqrt{\Delta L}\sigma}{\sqrt{nBT}} \leq O(\epsilon^2). \quad (188)$$

On the other hand, when

$$T \leq O\left(\max\left(\frac{\log(n)\Delta L}{\epsilon^2\sqrt{1-\lambda_2(\mathbf{W})}}, \frac{\Delta L}{\epsilon^2\sqrt{1-\lambda_2(\mathbf{W})}} \log\left(\frac{\varsigma_0^2}{\epsilon^2\Delta L}\right)\right)\right), \quad (189)$$

we have

$$\frac{\Delta L}{T\sqrt{1-\lambda_2(\mathbf{W})}} \cdot \max\left(\log(n), \log\left(\frac{\varsigma_0^2 T}{\Delta L}\right)\right) \leq O(\epsilon^2), \quad (190)$$

to see this, note that

$$\frac{\Delta L}{T\sqrt{1-\lambda_2(\mathbf{W})}} \log\left(\frac{\varsigma_0^2 T}{\Delta L}\right) = \epsilon^2 \frac{\log\left(\frac{\varsigma_0^2}{\epsilon^2\Delta L} \log\left(\frac{\varsigma_0^2}{\epsilon^2\Delta L}\right)\right)}{\log\left(\frac{\varsigma_0^2}{\epsilon^2\Delta L}\right)} \leq O(\epsilon^2). \quad (191)$$

Finally, we can obtain the upper bound

$$T \leq O\left(\frac{\Delta L\sigma^2}{nB\epsilon^4} + \max\left(\frac{\log(n)\Delta L}{\epsilon^2\sqrt{1-\lambda_2(\mathbf{W})}}, \frac{\Delta L}{\epsilon^2\sqrt{1-\lambda_2(\mathbf{W})}} \log\left(\frac{\varsigma_0^2}{\epsilon^2\Delta L}\right)\right)\right) \quad (192)$$

$$= O\left(\frac{\Delta L\sigma^2}{nB\epsilon^4} + \frac{\Delta L}{\epsilon^2\sqrt{1-\lambda_2(\mathbf{W})}} \log\left(n + \frac{\varsigma_0 n}{\epsilon\sqrt{\Delta L}}\right)\right), \quad (193)$$

as desired.

C Details to footnotes

C.1 Asynchronous Algorithm (Footnote 2)

In the full paper, we focus on the synchronous algorithms, i.e., we assume the existence of a synchronization process among workers between two adjacent iterations. We now extend our formulation to asynchronous algorithms. Since workers now update and communicate asynchronously, we define any gradient update that took place on a randomly chosen worker as one iteration. This randomness depends on system implementation, stochastic events, etc. This is a commonly adopted definition in the analysis of (decentralized) asynchronous algorithms [13]. To obtain a lower bound in such case, consider the two settings as shown in the proof of Theorem 1. In setting 1, it can be easily verified that the lower bound for sample complexity is

$$\Omega\left(\frac{\Delta L \sigma^2}{B \epsilon^4}\right). \quad (194)$$

This holds because in the extreme case, only one worker is making contributions to the optimization. And since we have not made any assumption on how workers are sampled to conduct the next iteration, this is a valid bound for arbitrary distribution. On the other hand, considering setting 2, the lower bound is still $\Omega(T_0 D)$ where $T_0 = \Omega(\Delta L \epsilon^{-2})$ is the lower bound in the sequential case, since the systems need at least $\Omega(D)$ iterations for the workers in I_0 and I_2 to contact. The lower bound for communication complexity is then

$$\Omega\left(\frac{\Delta L D}{\epsilon^2}\right). \quad (195)$$

Combining them together, we can get the final lower bound as:

$$\Omega\left(\frac{\Delta L \sigma^2}{B \epsilon^4} + \frac{\Delta L D}{\epsilon^2}\right). \quad (196)$$

Note that this bound holds with probability 1. It is possible to propose finer-grained assumption on how workers are chosen (e.g. uniformly random) and use concentration inequalities (e.g. Hoeffding's inequality) to get tighter bounds, we leave this as future work.

C.2 Relax zero-respecting assumption (Footnote 3)

To relax the zero-respecting assumption, we can use the technique proposed by [66] (See their proofs to Proposition 1 and 2). The basic idea is that to adversarially construct the loss function and rotate the non-zero coordinates in t -th iterations, such that when the algorithm operates on the rotated function, the first t iterations match with that of the old function. However, the new rotated function is still zero-respecting to the algorithm after t -th iteration so is generally hard to optimize. The details can be found in [66].

C.3 Specific algorithm for Average Consensus (Footnote 8)

Many algorithms have been proposed on solving the Average Consensus problem, readers can find details in many previous works on graph theory such as [82, 83, 84]. A straightforward algorithm is the Minimum Spanning Tree, that is, we first generate a spanning tree of the graph, and then the workers send and receive message using propagation on the tree. Specifically, starting from the leaves, all the children nodes of the tree send its accumulated value to the parents and the root compute the averaged value after gathering the information from the graph. And then reversely, the parent nodes send the value back to the child nodes and eventually all the nodes will get the averaged value. This algorithm is also known as the **GATHER-PROPAGATE** algorithm as discussed in [83], section 3. We include the detailed pseudo-code¹¹ in Algorithm 4.

¹¹This code is proposed by Ko [83], we do not intend to take credit for this.

Algorithm 4 GATHER-PROPAGATE (Spanning Tree) for a single coordinate

Input: communication graph G , a single coordinate on workers (all the coordinates follow the same instructions) to be communicated $\mathbf{X} \in \mathbb{R}^n$.

- 1: $\mathbf{d} \leftarrow$ vector of 1's indexed by $V(G)$ (vertices set of graph G).
- 2: $\mathcal{T} \leftarrow$ a spanning tree of G with root r arbitrarily picked.
- 3: **for** $v \in V(\mathcal{T})$ **do**
- 4: $l_v \leftarrow \bar{D}(r, v)$ (the distance between r and v)
- 5: **end for**
- 6: **for** $\alpha = \max_v l_v, \dots, 1$ **do**
- 7: **for** v with $l_v = \alpha$ **do**
- 8: v gives all its value onto its parents u :

$$\begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_v \end{bmatrix} \leftarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_v \end{bmatrix}$$

- 9: $\mathbf{d}_u \leftarrow \mathbf{d}_u + \mathbf{d}_v$
- 10: **end for**
- 11: **end for**
- 12: **for** $\alpha = 0, \dots, \max_v l_v - 1$ **do**
- 13: **for** u with $l_u = \alpha$ **do**
- 14: $\{v_1, \dots, v_\beta\} \leftarrow$ set of children of u
- 15: re-distribute the results:

$$\begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_{v_1} \\ \vdots \\ \mathbf{X}_{v_\beta} \end{bmatrix} \leftarrow \frac{1}{\mathbf{d}_u} \begin{bmatrix} \mathbf{d}_u - \mathbf{d}_u - \dots - \mathbf{d}_{v_\beta} \\ \mathbf{d}_{v_1} \\ \vdots \\ \mathbf{d}_{v_\beta} \end{bmatrix} \mathbf{X}_u$$

- 16: **end for**
 - 17: **end for**
 - 18: **return** $\mathbf{X} \frac{\mathbf{1}\mathbf{1}^\top}{n}$
-