Proximal Stochastic Recursive Momentum Methods for Nonconvex Composite Decentralized Optimization

Gabriel Mancino-Ball¹, Shengnan Miao¹, Yangyang Xu¹, Jie Chen²

¹Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, NY 12180

²MIT IBM-Watson AI Lab, IBM Research, Cambridge, MA 02142

mancig@rpi.edu, snmiao236@gmail.com, xuy21@rpi.edu, chenjie@us.ibm.com

Abstract

Consider a network of N decentralized computing agents collaboratively solving a nonconvex stochastic composite problem. In this work, we propose a single-loop algorithm, called DEEPSTORM, that achieves optimal sample complexity for this setting. Unlike double-loop algorithms that require a large batch size to compute the (stochastic) gradient once in a while, DEEPSTORM uses a small batch size, creating advantages in occasions such as streaming data and online learning. This is the first method achieving optimal sample complexity for decentralized nonconvex stochastic composite problems, requiring $\mathcal{O}(1)$ batch size. We conduct convergence analysis for DEEPSTORM with both constant and diminishing step sizes. Additionally, under proper initialization and a small enough desired solution error, we show that DEEPSTORM with a constant step size achieves a network-independent sample complexity, with an additional linear speed-up with respect to N over centralized methods. All codes are made available at https://github.com/gmancino/DEEPSTORM.

1 Introduction

Recent years have seen an increase in designing efficient algorithms for solving large-scale machine learning problems, over a network of N computing agents connected by a communication graph $\mathcal{G}=(\mathcal{V},\mathcal{E})$. Agents collaboratively solve the following composite problem:

$$\min_{\mathbf{x}} \frac{1}{N} \sum_{i=1}^{N} \left\{ \phi_i(\mathbf{x}) \triangleq f_i(\mathbf{x}) + r(\mathbf{x}) \right\}, \tag{1}$$

where the decision variable $\mathbf{x} \in \mathbb{R}^{1 \times p}$ is treated as a row vector; f_i is a smooth, possibly nonconvex function known only to agent i; and r is a convex, possibly non-smooth regularizer common to all agents. Agents i and j can communicate only if $(i,j) \in \mathcal{E}$. Many real-world applications in machine learning (Vogels et al. 2021; Ying et al. 2021; Yuan et al. 2021; Chamideh, Tärneberg, and Kihl 2021) and reinforcement learning (Zhang et al. 2018; Qu et al. 2019) fit the form of (1). Such scenarios differ from the centralized setting (McMahan et al. 2017; T. Dinh, Tran, and Nguyen 2020), where the agents are assumed to be able to communicate with one another globally via either a parameter server or a collective communication protocol. This setting arises naturally when data is distributed over a large geographic

region or when a centralized communication structure is too costly (Xin, Khan, and Kar 2021a).

Utilizing the communication topology induced by \mathcal{G} , we reformulate (1) into the following equivalent *decentralized* consensus optimization problem:

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_N} \frac{1}{N} \sum_{i=1}^N \phi_i(\mathbf{x}_i), \text{ s.t. } \mathbf{x}_i = \mathbf{x}_j, \ \forall (i, j) \in \mathcal{E}.$$
 (2)

Problem (2) allows for agents to maintain and update a local copy of the decision variable by locally computing gradients and performing neighbor communications.

The existence of a non-smooth regularizer r renders many decentralized optimization methods for a smooth objective inappropriate. We assume that r admits an easily computable (e.g. closed form) proximal mapping. Moreover, we are interested in the case where each local function f_i takes the following expectation form:

$$f_i(\mathbf{x}) \triangleq \mathbb{E}_{\varepsilon \sim \mathcal{D}_i} \left[f_i(\mathbf{x}; \xi) \right],$$
 (3)

with a slight abuse of notation for ease of exposition. In such a case, agents locally compute stochastic gradients of f_i . We adapt ideas from recent advances of stochastic optimization to the decentralized setting, by combining variance reduction techniques (Johnson and Zhang 2013; Nguyen et al. 2017; Allen-Zhu 2018; Wang et al. 2019; Cutkosky and Orabona 2019; Tran-Dinh et al. 2022) with gradient tracking (Lorenzo and Scutari 2016; Nedic, Olshevsky, and Shi 2017; Lu et al. 2019; Zhang and You 2020; Koloskova, Lin, and Stich 2021), to produce an algorithmic framework that achieves the optimal sample complexity bounds established in (Arjevani et al. 2022) for nonconvex stochastic methods.

Our framework, coined DEEPSTORM, is a single-loop algorithm with an attractive property that, besides the initial iteration, each agent only needs $m=\mathcal{O}(1)$ stochastic samples to compute a gradient estimate. Further, when a diminishing step size is used, even the first iteration does not need a large batch, at the expense of an additional logarithmic factor in the sample complexity result. Intuitively, DEEPSTORM utilizes a momentum based variance reduction technique (Cutkosky and Orabona 2019; Xu and Xu 2023; Levy, Kavis, and Cevher 2021; Tran-Dinh et al. 2022) to guarantee convergence under a small batch size. The use of momentum simultaneously accelerates the computation and communication complexities over non-momentum

based methods in the small batch setting; see Table 1 for a comparison. The recent ProxGT-SR-O/E (Xin et al. 2021) method can also achieve optimal sample complexity for solving (2), but at the expense of performing a double-loop which requires a large (stochastic) gradient computation every time the inner loop is completed. In scenarios where the batch size is uncontrollable, such as streaming or online learning, DEEPSTORM is advantageous.

When discussing sample complexity, it is paramount to specify the impact of the communication graph \mathcal{G} . With a constant step size, we show that under a sufficient amount of initial, or *transient*, iterations and proper initialization, DEEPSTORM behaves similarly to its centralized counterparts (Cutkosky and Orabona 2019; Levy, Kavis, and Cevher 2021; Tran-Dinh et al. 2022), while enjoying a linear speedup with respect to N.

We summarize the contributions of this work below:

- We propose a novel decentralized framework, DEEP-STORM, for nonconvex stochastic composite optimization problems. We show that DEEPSTORM achieves the optimal sample complexity with respect to solution accuracy, where each agent needs only O(1) samples to compute a local stochastic gradient. To the best of our knowledge, this is the first decentralized method that achieves optimal sample complexity for solving stochastic composite problems by using only small batches.
- Additionally, we establish convergence guarantees of DEEPSTORM with both constant and diminishing step sizes. When a constant step size is used, we show that under sufficiently many transient iterations and proper initialization, DEEPSTORM achieves a linear speed-up with respect to N, signifying an advantage over analogous centralized variance reduction methods (Cutkosky and Orabona 2019; Levy, Kavis, and Cevher 2021; Tran-Dinh et al. 2022).

2 Related works

A rich body of literature exists for solving the problem (2) in the decentralized setting. We discuss related works below.

Nonconvex decentralized methods. Of particular relevance to this work are methods for nonconvex f_i 's. When f_i takes the finite-sum form, deterministic methods (with full gradient computation) such as DGD (Zeng and Yin 2018), Near-DGD (Iakovidou and Wei 2021), Prox-PDA (Hong, Hajinezhad, and Zhao 2017), xFILTER (Sun and Hong 2019), and SONATA (Scutari and Sun 2019) converge to an ε -stationary point in $\mathcal{O}\left(\varepsilon^{-1}\right)$ iterations. They all work for the case $r\equiv 0$ only, except SONATA. For stochastic methods, we summarize a few representative ones in Table 1, including the information of whether they handle $r\not\equiv 0$. Note that D-PSGD (Lian et al. 2017) extends the convergence results of DGD; D^2 (Tang et al. 2018b) further improves over D-PSGD by relaxing a dissimilarity assumption.

Gradient tracking (Lorenzo and Scutari 2016; Nedic, Olshevsky, and Shi 2017) has been introduced as a tool to track the gradient of the global objective and has been studied extensively in the nonconvex and stochastic setting, under different names (Zhang and You 2020; Lu et al.

2019; Koloskova, Lin, and Stich 2021; Xin, Khan, and Kar 2021b). Many works now utilize this technique to improve the performance of their methods; those that mimic the SARAH (Nguyen et al. 2017) and Spider (Wang, Yin, and Zeng 2019) updates have become popular for their improved theoretical convergence rates. D-SPIDER-SFO (Pan, Liu, and Wang 2020) and D-GET (Sun, Lu, and Hong 2020) are two such methods. When f_i takes the finite-sum form, GT-SARAH (Xin, Khan, and Kar 2022) and DESTRESS (Li, Li, and Chi 2022) improve the analysis of D-GET by obtaining an optimal sample complexity and an optimal communication complexity, respectively. All these methods require computing a stochastic gradient with a large batch size every few iterations.

GT-HSGD (Xin, Khan, and Kar 2021a) can be considered a special case of our method. It uses a stochastic gradient estimator of the form proposed in (Cutkosky and Orabona 2019; Levy, Kavis, and Cevher 2021), requiring a large initial batch size, followed by $\mathcal{O}(1)$ batch size subsequently. The convergence analysis of GT-HSGD requires $r\equiv 0$; hence part of our work is to extend it to the case of $r\not\equiv 0$. Similar extensions have been proposed for other methods; for example, ProxGT-SR-O/E (Xin et al. 2021) extends D-GET, GT-SARAH, and DESTRESS. Additionally, the primal-dual method SPPDM (Wang et al. 2021) is shown to converge in $\mathcal{O}\left(\varepsilon^{-1}\right)$ communications, but it requires a large batch size proportional to ε^{-1} . Using such a batch size can negatively impact the performance on machine learning problems (Keskar et al. 2017).

Other decentralized methods. Several other decentralized methods exist for scenarios differing from that considered here. They include methods that work for convex problems only, such as DGD (Yuan, Ling, and Yin 2016), EXTRA (Shi et al. 2015), ADMM (Shi et al. 2014), DIGing (Nedic, Olshevsky, and Shi 2017), Acc-DNGD (Qu and Li 2019), MSDA (Scaman et al. 2017), DPAG (Ye et al. 2020), Flex-PD (Mansoori and Wei 2021), IDEAL (Arjevani et al. 2020), PUDA (Alghunaim et al. 2021), PMGT-VR (Ye, Xiong, and Zhang 2020), and DPSVRG (Li et al. 2021); asynchronous methods, such as AD-PSGD (Lian et al. 2018), the Asynchronous Primal-Dual method (Wu et al. 2017), APPG (Zhang and You 2021), asynchronous ADMM (Wei and Ozdaglar 2013; Hong 2018), and AD-OGP (Jiang et al. 2021); methods that operate under a timevarying network topology, such as Acc-GT (Li and Lin 2021) and ADOM (Kovalev et al. 2021); and methods that focus on providing convergence guarantees when communication compression is used, such as DCD-PSGD (Tang et al. 2018a), SQuARM-SGD (Singh et al. 2021), and the Primal-Dual method developed in (Chen et al. 2021).

3 DEEPSTORM framework

We first state the assumed conditions of each ϕ_i and the communication graph \mathcal{G} . They are standard in variance reduction (Cutkosky and Orabona 2019; Xu and Xu 2023; Tran-Dinh et al. 2022) and decentralized methods (Lian et al. 2017; Sun, Lu, and Hong 2020; Xin, Khan, and Kar 2021b).

Assumption 1 The following conditions hold.

Method	$r \not\equiv 0$	Batch size	Sample complexity (per agent)	
D-PSGD (Lian et al. 2017)	Х	$\mathcal{O}\left(1\right)$	$\mathcal{O}\left(\max\left\{rac{1}{Narepsilon^2},rac{N^2}{(1- ho)^2arepsilon} ight\} ight)$	
DSGT (Xin, Khan, and Kar 2021b)	X	$\mathcal{O}\left(1\right)$	$\mathcal{O}\left(\max\left\{rac{1}{Narepsilon^2},rac{N^2}{(1- ho)^2arepsilon} ight\} ight) \ \mathcal{O}\left(\max\left\{rac{1}{Narepsilon^2},rac{ ho N}{(1- ho)^3arepsilon} ight\} ight)$	
D-GET (Sun, Lu, and Hong 2020)	X	$\mathcal{O}\left(\frac{1}{arepsilon}\right)$ or $\mathcal{O}\left(\frac{1}{arepsilon^{0.5}}\right)$	$\mathcal{O}\left(\frac{1}{(1- ho)^aarepsilon^{1.5}} ight)$	
GT-HSGD (Xin, Khan, and Kar 2021a)	X	$\mathcal{O}\left(\frac{1}{arepsilon^{0.5}}\right)$ then $\mathcal{O}\left(1\right)$	$\mathcal{O}\left(\max\left\{\frac{1}{N\varepsilon^{1.5}}, \frac{\rho^4}{N(1-\rho)^3\varepsilon}, \frac{\rho^{1.5}N^{0.5}}{(1-\rho)^{2.25}\varepsilon^{0.75}}\right\}\right)$	
SPPDM (Wang et al. 2021)	✓	$\Omega(rac{N}{arepsilon})$	$\mathcal{O}\left(rac{1}{(1- ho)^barepsilon^2} ight)$	
ProxGT-SR-O/E (Xin et al. 2021)	✓	$\mathcal{O}\left(\frac{1}{arepsilon}\right)$ or $\mathcal{O}\left(\frac{1}{arepsilon^{0.5}}\right)$	$\mathcal{O}\left(rac{1}{Narepsilon^{1.5}} ight)^{\dagger}$	
Theorem 1	1	$\mathcal{O}\left(\frac{1}{\varepsilon^{0.5}}\right)$ then $\mathcal{O}\left(1\right)$	$\mathcal{O}\left(\max\left\{\frac{1}{N\varepsilon^{1.5}}, \frac{1}{(1-\rho)^2\varepsilon}, \frac{N^{0.5}}{\varepsilon^{0.75}}\right\}\right)^{\frac{1}{4}}$	
Theorem 2	✓	$\mathcal{O}\left(1\right)$	$ ilde{\mathcal{O}}\left(rac{1}{arepsilon^{1.5}} ight)$	

Table 1: Comparison between DEEPSTORM (bottom two rows) and representative decentralized stochastic nonconvex methods. The sample complexity takes into account both the stationarity and consensus violation. Since D-GET and SPPDM do not show the dependence on ρ , we use unspecified powers a and b, following the practice of (Xin, Khan, and Kar 2021a). †The sample complexity of ProxGT-SR-O/E is independent of ρ by *requiring* multiple communications per update; this is similar to our result in Theorem 2. ‡With multiple communications and $\varepsilon \leq N^{-2}$, Theorem 1 guarantees our algorithm attains the optimal $\mathcal{O}\left(N^{-1}\varepsilon^{-1.5}\right)$ sample complexity, but with a smaller batch size than ProxGT-SR-O/E.

- (i) The regularizer function r is convex and admits an easily computable proximal mapping.
- (ii) Each component function f_i is mean-squared L-smooth; i.e. there exists a constant $0 < L < \infty$ such that $\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^{1 \times p}$ and $\forall i = 1, \dots, N$,

$$\mathbb{E}_{\xi} \|\nabla f_i(\mathbf{a}; \xi) - \nabla f_i(\mathbf{b}; \xi)\|_2^2 \le L^2 \|\mathbf{a} - \mathbf{b}\|_2^2.$$
 (4)

(iii) There exists $\sigma > 0$ such that $\forall \mathbf{a} \in \mathbb{R}^{1 \times p}$,

$$\mathbb{E}_{\xi}[\nabla f_i(\mathbf{a}; \xi)] = \nabla f_i(\mathbf{a}),$$

$$\mathbb{E} \|\nabla f_i(\mathbf{a}; \xi) - \nabla f_i(\mathbf{a})\|_2^2 \le \sigma^2.$$
(5)

(iv) The global function $\phi = \frac{1}{N} \sum_{i=1}^{N} \phi_i$ is lower bounded; i.e. there exists a constant ϕ^* such that

$$-\infty < \phi^* < \phi(\mathbf{a}), \ \forall \mathbf{a} \in \mathbb{R}^{1 \times p}.$$
 (6)

Assumption 2 The graph \mathcal{G} is connected and undirected. It can be represented by a mixing matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$ such that:

- (i) (Decentralized property) $w_{ij} > 0$ if $(i, j) \in \mathcal{E}$ and $w_{ij} = 0$ otherwise;
- (ii) (Symmetric property) $W = W^{\top}$;
- (iii) (Null-space property) $\operatorname{null}(\mathbf{I} \mathbf{W}) = \operatorname{span}\{\mathbf{e}\},$ where $\mathbf{e} \in \mathbb{R}^N$ is the vector of all ones; and
- (iv) (Spectral property) the eigenvalues of \mathbf{W} lie in the range (-1,1] with

$$\rho \triangleq \left\| \mathbf{W} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right\|_{2} < 1. \tag{7}$$

Note that the entry values of ${\bf W}$ can be flexibly designed as long as Assumption 2 holds. One example is ${\bf W}={\bf I}-{\bf L}/\tau$, where ${\bf L}$ is the combinatorial Laplacian of ${\cal G}$ and τ is a value greater than half of its largest eigenvalue. It is not hard

to see that the consensus constraint $\mathbf{x}_i = \mathbf{x}_j$ for all $(i, j) \in \mathcal{E}$ in (2) is equivalent to $\mathbf{W}\mathbf{X} = \mathbf{X}$, where the *i*-th row of \mathbf{X} is \mathbf{x}_i . The value ρ in (7) indicates the connectedness of the graph. The quantity $1 - \rho$ is sometimes referred to as the *spectral gap*; a higher value suggests that the graph is more connected and consensus of the \mathbf{x}_i 's is easier to achieve.

Under Assumptions 1 and 2, we now present the DEEP-STORM framework. We start with the basic algorithm and later generalize the simple communication (using W) with a more general communication operator, denoted by \mathcal{W}_T .

Basic algorithm. Let $\mathbf{x}_i^{(k)}$ be the k-th iterate for agent i, and let the matrix $\mathbf{X}^{(k)}$ contain all the k-th iterates among agents, stacked as a matrix. We will similarly use such vector and matrix notations for other variables. Our **DEcE**ntralized **Proximal STO**chastic **Recursive Momentum** framework, DEEPSTORM, uses a variance reduction variable $\mathbf{d}_i^{(k)}$ and a gradient tracking variable $\mathbf{y}_i^{(k)}$ to improve the convergence of $\mathbf{x}_i^{(k)}$. DEEPSTORM contains the following steps in each iteration k:

1. Communicate the local variables:

$$\mathbf{Z}^{(k)} = \mathbf{W}\mathbf{X}^{(k)}.\tag{8}$$

2. Update each local variable (by using, e.g., proximal mappings):

$$\mathbf{x}_{i}^{(k+1)} = \underset{\mathbf{x}_{i}}{\operatorname{argmin}} \left\{ \alpha_{k} r(\mathbf{x}_{i}) + \frac{1}{2} \left\| \mathbf{x}_{i} - \left(\mathbf{z}_{i}^{(k)} - \alpha_{k} \mathbf{y}_{i}^{(k)} \right) \right\|^{2} \right\}.$$

3. Update the variance reduction variable:

$$\mathbf{d}_{i}^{(k+1)} = (1 - \beta_{k}) \left(\mathbf{d}_{i}^{(k)} + \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right) + \beta_{k} \tilde{\mathbf{v}}_{i}^{(k+1)},$$
(10)

where

$$\mathbf{v}_{i}^{(k+1)} = \frac{1}{m} \sum_{\xi \in B_{i}^{(k+1)}} \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}; \xi),$$

$$\mathbf{u}_{i}^{(k+1)} = \frac{1}{m} \sum_{\xi \in B_{i}^{(k+1)}} \nabla f_{i}(\mathbf{x}_{i}^{(k)}; \xi).$$
(11)

Here, $\boldsymbol{B}_{i}^{(k+1)}$ is a batch of m samples at the current iteration. Note that while $\mathbf{v}_i^{(k+1)}$ is evaluated at the current iterate, $\mathbf{u}_{i}^{(k+1)}$ is evaluated at the previous iterate. We make the assumption that for all k and all agents i and j, $B_i^{(k+1)}$ and $B_j^{(k+1)}$ contain independent and mutually independent random variables. The part $\tilde{\mathbf{v}}_i^{(k+1)}$ can be any unbiased estimate of $\nabla f_i(\mathbf{x}_i^{(k+1)})$ with bounded variance; its details will be elaborated soon.

4. Update the gradient tracking variable via communication:

$$\mathbf{Y}^{(k+1)} = \mathbf{W} \left(\mathbf{Y}^{(k)} + \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right).$$
 (12)

The step that updates the variance reduction variable, (10), is motivated by Hybrid-SGD (Tran-Dinh et al. 2022), which allows for a single-loop update. Intuitively, this variable is a convex combination of the SARAH (Nguyen et al. 2017) update and $\tilde{\mathbf{v}}_i^{(k+1)}$, allowing for strong variance reduction and meanwhile flexibility in design. By doing so, a constant batch size m suffices for convergence. This is a useful property in scenarios of online learning and real-time decision making, where it is unrealistic to obtain and store mega batches for training (Xu and Xu 2023; Xin, Khan, and Kar 2021a).

Examples of $\tilde{\mathbf{v}}_i^{(k+1)}$. The vector $\tilde{\mathbf{v}}_i^{(k+1)}$ in (10) can be any unbiased local gradient estimate. In this work, we consider two cases: either $\tilde{\mathbf{v}}_i^{(k+1)}$ is evaluated on another set of samples $\tilde{B}_i^{(k+1)}$, defined analogously to $B_i^{(k+1)}$ that is used to compute $\mathbf{v}_i^{(k+1)}$ in (11), such that

$$\begin{split} \tilde{B}_i^{(k+1)} &\text{ is independent of } B_i^{(k+1)} \text{ with} \\ \mathbb{E} \left\| \tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right\|_2^2 \leq \hat{\sigma}^2; \end{split} \tag{v1}$$

$$\tilde{\mathbf{v}}_i^{(k+1)} = \mathbf{v}_i^{(k+1)} \text{ with } \mathbb{E} \left\| \mathbf{v}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right\|_2^2 \le \hat{\sigma}^2,$$

for some $\hat{\sigma} > 0$. Two possible unbiased estimators that sat-

$$\tilde{\mathbf{v}}_i^{(k+1)} = \frac{1}{m} \sum_{\tilde{\epsilon} \in \tilde{R}^{(k+1)}} \nabla f_i(\mathbf{x}_i^{(k+1)}; \tilde{\xi}), \tag{v1-SG}$$

$$\begin{split} \tilde{\mathbf{v}}_{i}^{(k+1)} = & \frac{1}{m} \sum_{\tilde{\xi} \in \tilde{B}_{i}^{(k+1)}} \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}; \tilde{\xi}) \\ & + \frac{1}{m} \sum_{\tilde{\xi} \in \tilde{B}_{i}^{(\tau_{k+1})}} \nabla f_{i}(\mathbf{x}_{i}^{(\tau_{k+1})}; \tilde{\xi}) - \frac{1}{m} \sum_{\tilde{\xi} \in \tilde{B}_{i}^{(k+1)}} \nabla f_{i}(\mathbf{x}_{i}^{(\tau_{k+1})}; \tilde{\xi}), \end{split}$$

$$(v1-SVRG)$$

for some $\tau_{k+1} < k+1$. The first estimator is a standard one, evaluated by using a batch $\tilde{B}_{i}^{(k+1)}$ independent of $B_{\cdot}^{(k+1)}$. The second estimator, which introduces further variance reduction, uses an additional past-time iterate $\mathbf{x}_i^{(au_{k+1})}$ and a batch $\tilde{\tilde{B}}_i^{(au_{k+1})}$, whose size is generally greater than m. Such an update is inspired by the SVRG method (Johnson and Zhang 2013). Here, we have $\hat{\sigma}^2 = m^{-1}\sigma^2$ for the estimators (v1-SG) and (v2); while $\hat{\sigma}^2 =$ $\left(3m^{-1}+6\left|\tilde{\tilde{B}}_{i}^{(\tau_{k+1})}\right|^{-1}\right)\sigma^{2}$ for (v1-SVRG), where we recall that σ^2 comes from (5). Note that beyond the two exam-

ples, our proof techniques hold for any unbiased estimator satisfying (v1), leaving more open designs.

Generalized communication. Steps (8) and (12) use the mixing matrix to perform weighted averaging of neighbor information. The closer \mathbf{W} is to $\frac{1}{N}\mathbf{e}\mathbf{e}^{\mathsf{T}}$, the more uniform the rows of $\mathbf{X}^{(k+1)}$ are, implying agents are closer to consensus. Hence, to improve convergence, we can apply multiple mixing rounds in each iteration. To this end, we generalize the network communication by using an operator \mathcal{W}_T , which is a degree-T polynomial in W that must satisfy Assumption 2 parts (ii)-(iv). We adopt Chebyshev acceleration (Auzinger and Melenk 2011; Scaman et al. 2017; Xin et al. 2021; Li, Li, and Chi 2022), which defines for any input matrix \mathbf{B}_0 , $\mathbf{B}_T = \mathcal{W}_T(\mathbf{B}_0)$, where $\mathbf{B}_1 = \mathbf{W}\mathbf{B}_0$, $\mu_0 = 1$, $\mu_1 = \frac{1}{\rho}$ for ρ defined in (7), and recursively,

$$\mu_{t+1} = \frac{2}{\rho} \mu_t - \mu_{t-1} \text{ and}$$

$$\mathbf{B}_{t+1} = \frac{2\mu_t}{\rho \mu_{t+1}} \mathbf{W} \mathbf{B}_t - \frac{\mu_{t-1}}{\mu_{t+1}} \mathbf{B}_{t-1}, \text{ for } t \le T - 1.$$
(13)

It is not hard to see that e is an eigenvector of \mathcal{W}_T , associated to eigenvalue 1, whose algebraic multiplicity is 1. Therefore,

$$\tilde{\rho} \triangleq \left\| \mathbf{W}_T - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right\|_2 < 1.$$
 (14)

Moreover, $\tilde{\rho}$ converges to zero exponentially with T, bringing W_T rather close to an averaging operator (for details, see Appendix B). Notice with $T=1, \mathcal{W}_T$ reduces to W.

We summarize the overall algorithm in Algorithm 1, by replacing W in (8) and (12) with \mathcal{W}_T . Additionally, see the discussions after Theorems 1 and 2 regarding the probability distribution for choosing the output of Algorithm 1.

Convergence results

For the convergence of DEEPSTORM, we start with the following standard definitions (Xu and Xu 2023; Xin et al. 2021).

Definition 1 Given $\mathbf{x} \in \text{dom}(r)$, \mathbf{y} , and $\eta > 0$, define the proximal gradient mapping of y at x to be

$$P(\mathbf{x}, \mathbf{y}, \eta) \triangleq \frac{1}{\eta} (\mathbf{x} - \text{prox}_{\eta r} (\mathbf{x} - \eta \mathbf{y})),$$
 (15)

where prox denotes the proximal operator $prox_a(\mathbf{v}) =$ $\operatorname{argmin}_{\mathbf{u}} \{ g(\mathbf{u}) + \frac{1}{2} ||\mathbf{u} - \mathbf{v}||_{2}^{2} \}.$

Algorithm 1: DEEPSTORM

Input: Initial $\mathbf{X}^{(0)}$, mixing rounds T_0, T , iteration K, and $\{\alpha_k\}, \{\beta_k\}$

- 1: Compute $\mathbf{d}_{i}^{(0)} = \frac{1}{m_{0}} \sum_{\xi \in B_{i}^{(0)}} \nabla f_{i}(\mathbf{x}_{i}^{(0)}; \xi) \ \forall i$
- 2: Communicate to obtain $\mathbf{Y}^{(0)} = \mathcal{W}_{T_0}(\mathbf{D}^{(0)})$
- 3: **for** $k = 0, \dots, K 1$ **do**
- 4: Communicate to obtain $\mathbf{Z}^{(k)} = \mathcal{W}_T(\mathbf{X}^{(k)})$
- 5: Update local decision variables by (9)
- 6: Obtain local gradient estimator by (10)
- 7: Communicate to update gradient tracking variable $\mathbf{Y}^{(k+1)} = \mathbf{W}_T(\mathbf{Y}^{(k)} + \mathbf{D}^{(k+1)} \mathbf{D}^{(k)})$
- 8: end for

Output: $\mathbf{Z}^{(\tau)}$ with τ chosen randomly from $\{0,\dots,K-1\}$

Definition 2 A stochastic matrix $\mathbf{X} \in \mathbb{R}^{N \times p}$ is called a stochastic ε -stationary point of (2) if

$$\mathbb{E}\left[\frac{1}{N}\sum_{i=1}^{N}\|P\left(\mathbf{x}_{i},\nabla f(\mathbf{x}_{i}),\eta\right)\|_{2}^{2} + \frac{L^{2}}{N}\|\mathbf{X}_{\perp}\|_{F}^{2}\right] \leq \varepsilon,$$
(16)

where $\eta > 0$, $\nabla f \triangleq \frac{1}{N} \sum_{j=1}^{N} \nabla f_j$, \mathbf{x}_i is the *i*-th row of \mathbf{X} , and $\mathbf{X}_{\perp} \triangleq \mathbf{X} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{X}$ is the difference between all \mathbf{x}_i and their average $\frac{1}{N} \sum_{j=1}^{N} \mathbf{x}_j$.

Our analyses rely on the construction of two novel Lyapunov functions as indicated by Theorems 1 and 2 below. These Lyapunov functions guarantee convergence through the careful design of function coefficients which result from solving non-linear systems of inequalities in either the constant or diminishing step size case. We first consider the use of a constant step size. The convergence rate result is given in the following theorem. Its proof is given in Appendix C.2.

Theorem 1 Under Assumptions 1 and 2, let $\left\{ \left(\mathbf{X}^{(k)}, \mathbf{D}^{(k)}, \mathbf{Y}^{(k)}, \mathbf{Z}^{(k)} \right) \right\}$ be obtained by Algorithm 1 via (9), (12), and (10) such that $\tilde{\mathbf{v}}_i^{(k+1)}$ is any unbiased gradient estimator that satisfies either (v1) or (v2). Further, let α_k and β_k be chosen as

$$\alpha_{k} = \frac{\alpha}{K^{\frac{1}{3}}}, \ \beta_{k} = \frac{144L^{2}\alpha^{2}}{NK^{\frac{2}{3}}}, \ \text{with}$$

$$\alpha \le \min\left\{\frac{K^{\frac{1}{3}}}{32L}, \frac{(1-\tilde{\rho})^{2}K^{\frac{1}{3}}}{64L}\right\},$$
(17)

for all k = 0, ..., K - 1. Then, it holds that $\beta_k \in (0, 1)$ for all $k \ge 0$ and that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left(\frac{1}{N} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k} \right) \right\|_{2}^{2} + \frac{L^{2}}{N} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2} \right) \\
\leq \frac{512}{\alpha K^{\frac{2}{3}}} \left(\Phi^{(0)} - \phi^{*} \right) + \left(\frac{2048}{L(1 - \tilde{\rho})^{2} K} \right) \frac{144^{2} L^{4} \alpha^{3} \hat{\sigma}^{2}}{N^{2}} \\
+ \left(\frac{128}{3L^{2} \alpha K^{\frac{2}{3}}} + \frac{8192\alpha}{K^{\frac{4}{3}}} + \frac{2048\alpha}{NK^{\frac{4}{3}}} \right) \frac{144^{2} L^{4} \alpha^{3} \hat{\sigma}^{2}}{N^{2}}, \tag{18}$$

for some $\Phi^{(0)} > \phi^*$ that depends on the initialization. Note that $\Phi^{(k)}$ is defined in (C.43) in Appendix C for any $k \geq 0$.

Network-independent sample complexity, linear speed-up, and communication complexity. Theorem 1 establishes convergence based on the sequence $\{\mathbf{Z}^{(k)}\}$ defined in (8). As a consequence, if we let each agent start with the same initial variable $\mathbf{x}^{(0)}$, set $\alpha = \frac{N^{\frac{2}{3}}}{64L}$ and the initial batch size $m_0 = \sqrt[3]{NK}$, and choose initial communication rounds $T_0 = \mathcal{O}\left((1-\rho)^{-0.5}\right)$ for $\mathbf{Y}^{(0)}$, then for all $K \geq \frac{N^2}{(1-\hat{\rho})^6}$, DEEPSTORM achieves stochastic ε -stationarity for some iterate $\mathbf{Z}^{(\tau)}$, where τ is selected uniformly from $\{0,\ldots,K-1\}$, by using

$$\mathcal{O}\left(\max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^2}{(1-\tilde{\rho})^2\varepsilon}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right) \quad (19)$$

local stochastic gradient computations. For the formal statement, see Corollary 1 in Appendix C.2. Here, $\Delta = \Phi^{(0)} - \phi^*$ denotes an initial function gap, which is independent of $\tilde{\rho},$ N, and K. Moreover, when $\varepsilon \leq N^{-2}(1-\tilde{\rho})^4$, we see that $\mathcal{O}\left(N^{-1}\varepsilon^{-1.5}\right)$ dominates in (C.60); hence, this result manifests a linear speed-up with respect to N over the centralized counterparts (Cutkosky and Orabona 2019; Tran-Dinh et al. 2022) of DEEPSTORM. Furthermore, if the number of Chebyshev mixing rounds is $T = \lceil \frac{2}{\sqrt{1-\rho}} \rceil$, we have $(1-\tilde{\rho}) \geq \frac{1}{\sqrt{2}}$, which suggests that ε does not need to be small for the linear speed-up to hold. For details, see Lemma B.1 and Remark C.2 in the Appendix. The communication cost is $\mathcal{O}\left(T_0 + TK\right)$.

In parallel, we state a result for the case of diminishing step size. Its proof is given in Appendix C.3.

Theorem 2 *Under the same assumptions as Theorem 1, let* α_k *and* β_k *be chosen as*

$$\alpha_{k} = \frac{\alpha}{(k+k_{0})^{\frac{1}{3}}}, \quad \beta_{k} = 1 - \frac{\alpha_{k+1}}{\alpha_{k}} + 48L^{2}\alpha_{k+1}^{2}, \text{ with}$$

$$\alpha \leq \min\left\{\frac{k_{0}^{\frac{1}{3}}}{32L}, \frac{(1-\tilde{\rho})^{2}k_{0}^{\frac{1}{3}}}{64L}\right\},$$
(20)

for all k = 0, ..., K-1, where $k_0 \ge \lceil \frac{2}{1-\tilde{\rho}^3} \rceil$. Then, it holds that $\beta_k \in (0,1)$ for all $k \ge 0$ and that

$$\sum_{k=0}^{K-1} c\alpha_k \mathbb{E}\left(\frac{1}{N} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_i^{(k)}, \nabla f(\mathbf{z}_i^{(k)}), \alpha_k\right) \right\|_2^2 + \frac{L^2}{N} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_F^2 \right)$$

$$\leq 12 \left(\hat{\Phi}^{(0)} - \phi^*\right) + \sum_{k=0}^{K-1} \left(\frac{1}{L^2 \alpha_{k+1}} + \frac{48}{L(1-\tilde{\rho})^2}\right) \beta_k^2 \hat{\sigma}^2,$$
(21)

for some $\hat{\Phi}^{(0)} > \phi^*$ that depends on initialization and $c \triangleq \frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} > \frac{1}{4}$. Note that $\hat{\Phi}^{(k)}$ is defined in (C.69) in Appendix C for any $k \geq 0$.

Sample complexity. Theorem 2 establishes the convergence rate of DEEPSTORM with diminishing step sizes. If we choose $k_0 = \lceil \frac{2}{(1-\bar{\rho})^6} \rceil$ in (20), then DEEPSTORM achieves stochastic ε -stationarity for some iterate $\mathbf{Z}^{(\tau)}$, where τ is chosen according to (C.87), by us-

ing $\tilde{\mathcal{O}}\left((1-\tilde{\rho})^{-3}\varepsilon^{-1.5}\right)$ local stochastic gradient computations; this sample complexity is network-dependent. However, by using an initialization technique similar to the case of constant step sizes above and letting the initial batch size be $\mathcal{O}(1)$, we can set the Chebyshev mixing rounds to be $T=\lceil\frac{2}{\sqrt{1-\tilde{\rho}}}\rceil$, so that $(1-\tilde{\rho})^{-1}\leq\sqrt{2}$. This leads to the network-independent sample complexity reported in Table 1. For a full statement of the complexity results, see Corollary 2 in Appendix C.3 and Remark C.4.

5 Experiments

In this section, we empirically validate the convergence theory of DEEPSTORM and demonstrate its effectiveness in comparison with representative decentralized methods. We compare all versions of DEEPSTORM with DSGT (Lu et al. 2019; Zhang and You 2020; Koloskova, Lin, and Stich 2021; Xin, Khan, and Kar 2021b), SPPDM (Wang et al. 2021), and ProxGT-SR-O/E (Xin et al. 2021). DSGT uses gradient tracking but it is not designed for non-smooth objectives; nevertheless, it outperforms strong competitors (e.g., D-PSGD (Lian et al. 2017) and D² (Tang et al. 2018b)) in practice (Zhang and You 2020; Xin, Khan, and Kar 2021b). SP-PDM is a primal-dual method, but it does not utilize gradient tracking and its convergence theory requires a large batch size. ProxGT-SR-O/E is a double-loop algorithm, which requires using a mega-batch to compute the (stochastic) gradient at each outer iteration. All experiments are conducted using the AiMOS ¹ supercomputer with eight NVIDIA Tesla V100 GPUs in total, with code implemented in PyTorch (v1.6.0) and OpenMPI (v3.1.4).

Problems. We conduct tests on three classification problems. Each local agent i has the objective $\phi_i(\mathbf{x}_i) = \frac{1}{M} \sum_{j=1}^M \ell\left(g\left(\mathbf{x}_i, \mathbf{a}_j\right), \mathbf{b}_j\right) + \lambda \left\|\mathbf{x}_i\right\|_1$, where $g(\mathbf{x}, \mathbf{a})$ is the output of a neural network with parameters \mathbf{x} on data \mathbf{a} , and ℓ is the cross-entropy loss function between the output and the true label \mathbf{b} . The data is uniformly randomly split among the agents, each obtaining M training examples. The L_1 regularization promotes sparsity of the trained network. The regularization strength λ is set to 0.0001 following general practice.

Data sets and neural networks. The three data sets we experiment with are summarized in Table 2 in Appendix A. Two of them are tabular data and we use the standard multilayer perceptron for g (one hidden layer with 64 units). The other data set contains images; thus, we use a convolutional neural network. Both neural networks use the tanh activation to satisfy the smoothness condition of the objective function.

Communication graphs. Each data set is paired with a different communication graph, indicated by, and visualized in, Table 2 in Appendix A. For the ladder and random graphs, the mixing matrix is set as $\mathbf{W} = \mathbf{I} - \gamma \mathbf{L}$, where γ is reciprocal of the maximum eigenvalue of the combinatorial Laplacian \mathbf{L} . For the ring graph, self-weighting and neighbor weights are set to be $\frac{1}{2}$.

Performance metrics. We evaluate on four metrics: training loss, stationarity violation, solution sparsity, and test

accuracy. Further, we compare the methods with respect to data passes and algorithm iterations, which reflect the sample complexity and communication complexity, respectively. Note that for each iteration, all methods except SP-PDM communicate two variables. For the training loss, stationarity violation, and test accuracy, we evaluate on the average solution $\bar{\mathbf{x}}$. The stationarity violation is defined as $\|\bar{\mathbf{x}} - \mathrm{prox}_r (\bar{\mathbf{x}} - \nabla f(\bar{\mathbf{x}}))\|_2^2 + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|_2^2$, which measures both optimality and consensus. For sparsity, we use the average percentage of non-zeros in each \mathbf{x}_i prior to local communication.

Protocols. For hyperparameter selection, see Appendix A. We perform ten runs with different starting points for each dataset. In several runs for the MNIST dataset, DSGT and SPPDM converge to solutions with $\ll 1\%$ nonzero entries, but the training loss and test accuracy are not competitive at all. We remove these runs and keep only the five best runs for reporting the (averaged) performance.

Results. Figure 1 summarizes the results for all performance metrics, by using the same number of data passes for all methods when convergence has been observed. For a9a and MiniBooNE, the results are averaged over passes 80 to 100; while for MNIST, over passes 180 to 200. Figure 2 compares different methods by using the same number of algorithm iterations.

Overall, we see that DEEPSTORM (all variants) generally yields a lower training loss and significantly fewer nonzeros in the solution than the other decentralized algorithms. This observation suggests that DEEPSTORM indeed solves the optimization problem (2) much more efficiently in terms of both data passes and iterations. Moreover, the test accuracy is also highly competitive, concluding the practical usefulness of DEEPSTORM.

6 Conclusion

We have presented a novel decentralized algorithm for solving the nonconvex stochastic composite problem (2) by leveraging variance reduction and gradient tracking. It is the first such work that achieves optimal sample complexity for this class of problems by using $\mathcal{O}(1)$ batch sizes. Our algorithm is a framework with an open term (see (10)), for which we analyze two examples that allow the framework to achieve network-independent complexity bounds, suggesting no sacrifice over centralized variance reduction methods. Our proof technique can be used to analyze more designs of the open term. While our work is one of the few studies on the nonconvex stochastic composite problem (2), our analysis is for the synchronous setting with a static communication graph. Analysis (or adaptation of the algorithm) for asynchronous or time-varying settings is an avenue of future investigation.

7 Acknowledgments

This work was supported by the Rensselaer-IBM AI Research Collaboration, part of the IBM AI Horizons Network, NSF grants DMS-2053493 and DMS-2208394, and the ONR award N00014-22-1-2573.

¹See: https://cci.rpi.edu/aimos

Method	Train loss	Stationarity	% Non-zeros	Test accuracy					
a9a									
DSGT	0.3308±1.272e-4	0.0003±1.819e-4	74.18±160.09e-4	84.89±271.02e-4					
SPPDM	$0.5457 \pm 20.014e$ -4	$0.001\pm2.99e-4$	$46.19\pm51.04e-4$	$76.38 \pm 0.0e-4$					
ProxGT-SR-E	$0.545\pm85.017e-4$	$0.0491\pm64.099e-4$	98.04±15.035e-4	$76.38 \pm 0.0e-4$					
DEEPSTORM v1-SG	$0.3306 \pm 9.46e-4$	$0.0002\pm1.292e-4$	$2.99\pm60.066e-4$	$84.96 \pm 1235.0e-4$					
DEEPSTORM v1-SVRG	$0.3308 \pm 7.689e - 4$	0.0001 ±0.21278e-4	$2.86 \pm 45.018e-4$	$84.94 \pm 929.04e-4$					
DEEPSTORM v2	0.3277 ±7.461e-4	0.0001 ±0.8179e-4	1.92 ±53.073e-4	85.11 ±478.03e-4					
MiniBooNE									
DSGT	0.3735±3.844e-4	0.0003±2.076e-4	81.83±227.0e-4	84.24±202.07e-4					
SPPDM	$0.5699 \pm 61.016e-4$	$0.0025\pm5.565e-4$	$35.32 \pm 77.02e-4$	$72.02 \pm 0.0e-4$					
ProxGT-SR-E	$0.5663\pm32.027e-4$	$0.0115\pm7.57e-4$	97.88±17.017e-4	$72.02\pm0.0e-4$					
DEEPSTORM v1-SG	0.3637 ±19.015e-4	$0.0002 \pm 0.6464e-4$	$4.34 \pm 60.07e-4$	$84.24 \pm 1902.0e-4$					
DEEPSTORM v1-SVRG	$0.3653\pm23.054e-4$	$\overline{0.0002} \pm 0.9716e-4$	$4.42\pm65.068e-4$	$\overline{84.15} \pm 1974.0e-4$					
DEEPSTORM v2	0.3637 ±18.046e-4	$0.0001 \pm 0.4136e-4$	4.2 ±61.073e-4	84.25 ±1752.0e-4					
MNIST									
DSGT	0.1055±24.03e-4	0.0024±3.554e-4	51.05±896.0e-4	97.61±1346.0e-4					
SPPDM	$0.1851\pm55.065e-4$	$\overline{0.0051} \pm 2.058e-4$	66.81±616.03e-4	95.55±1488.0e-4					
ProxGT-SR-E	$1.699 \pm 903.07e-4$	$0.21299 \pm 268.0e-4$	91.4±70.087e-4	52.25±41480.0e-4					
DEEPSTORM v1-SG	$0.081\pm33.014e-4$	$0.0027 \pm 5.376 e$ -4	$10.31 \pm 70.031e-4$	97.97±1261.0e-4					
DEEPSTORM v1-SVRG	$0.078\pm34.022e-4$	$0.0031 \pm 7.366e-4$	$\overline{10.99} \pm 82.095e-4$	98.08±1485.0e-4					
DEEPSTORM v2	0.0768 ±29.095e-4	0.0016 ±1.83e-4	7.36 ±50.07e-4	98.15 ±659.04e-4					

Figure 1: Comparisons of different methods by running them with the same number of data passes. Bold values indicate the best results and underlined values indicate the second best.

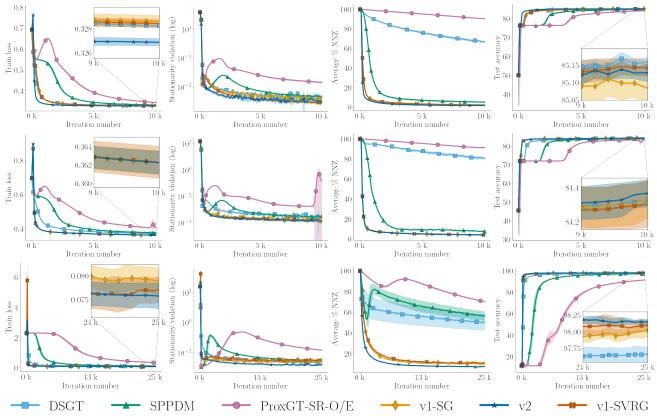


Figure 2: Comparison of different methods by running them with the same number of iterations. From top to bottom: a9a, MiniBooNE, MNIST. From left to right: training loss, stationarity violation, average percentage non-zeros, testing accuracy. The shaded regions indicate standard deviations (with some being small and unnoticeable).

References

- Alghunaim, S. A.; Ryu, E. K.; Yuan, K.; and Sayed, A. H. 2021. Decentralized Proximal Gradient Algorithms With Linear Convergence Rates. *IEEE Transactions on Automatic Control*, 66(6): 2787–2794.
- Allen-Zhu, Z. 2018. Katyusha: The First Direct Acceleration of Stochastic Gradient Methods. *Journal of Machine Learning Research*, 18(221): 1–51.
- Arjevani, Y.; Bruna, J.; Can, B.; Gurbuzbalaban, M.; Jegelka, S.; and Lin, H. 2020. IDEAL: Inexact DEcentralized Accelerated Augmented Lagrangian Method. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M. F.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 20648–20659. Curran Associates, Inc.
- Arjevani, Y.; Carmon, Y.; Duchi, J. C.; Foster, D. J.; Srebro, N.; and Woodworth, B. 2022. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*.
- Auzinger, W.; and Melenk, J. M. 2011. Iterative Solution of Large Linear Systems. *TU Wien, Lecture Notes*.
- Chamideh, S.; Tärneberg, W.; and Kihl, M. 2021. Evaluation of Decentralized Algorithms for Coordination of Autonomous Vehicles at Intersections. In 2021 IEEE International Intelligent Transportation Systems Conference (ITSC), 1954–1961.
- Chen, C.; Zhang, J.; Shen, L.; Zhao, P.; and Luo, Z. 2021. Communication Efficient Primal-Dual Algorithm for Nonconvex Nonsmooth Distributed Optimization. In Banerjee, A.; and Fukumizu, K., eds., *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, 1594–1602. PMLR.
- Cutkosky, A.; and Orabona, F. 2019. Momentum-Based Variance Reduction in Non-Convex SGD. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d' Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Ghadimi, S.; Lan, G.; and Zhang, H. 2016. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1-2): 267–305.
- Hong, M. 2018. A Distributed, Asynchronous, and Incremental Algorithm for Nonconvex Optimization: An ADMM Approach. *IEEE Transactions on Control of Network Systems*, 5(3): 935–945.
- Hong, M.; Hajinezhad, D.; and Zhao, M.-M. 2017. Prox-PDA: The Proximal Primal-Dual Algorithm for Fast Distributed Nonconvex Optimization and Learning Over Networks. In Precup, D.; and Teh, Y. W., eds., *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 1529–1538. International Convention Centre, Sydney, Australia: PMLR.
- Iakovidou, C.; and Wei, E. 2021. On the Convergence of NEAR-DGD for Nonconvex Optimization with Second Order Guarantees. In 2021 60th IEEE Conference on Decision and Control (CDC), 259–264.
- Jiang, J.; Zhang, W.; GU, J.; and Zhu, W. 2021. Asynchronous Decentralized Online Learning. In Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*.
- Johnson, R.; and Zhang, T. 2013. Accelerating Stochastic Gradient Descent using Predictive Variance Reduction. In Burges, C. J. C.; Bottou, L.; Welling, M.; Ghahramani, Z.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.

- Keskar, N. S.; Mudigere, D.; Nocedal, J.; Smelyanskiy, M.; and Tang, P. T. P. 2017. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. Open-Review.net.
- Koloskova, A.; Lin, T.; and Stich, S. U. 2021. An Improved Analysis of Gradient Tracking for Decentralized Machine Learning. In Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*.
- Kovalev, D.; Shulgin, E.; Richtarik, P.; Rogozin, A. V.; and Gasnikov, A. 2021. ADOM: Accelerated Decentralized Optimization Method for Time-Varying Networks. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, 5784–5793. PMLR.
- Levy, K. Y.; Kavis, A.; and Cevher, V. 2021. STORM+: Fully Adaptive SGD with Recursive Momentum for Nonconvex Optimization. In Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*.
- Li, B.; Li, Z.; and Chi, Y. 2022. DESTRESS: Computation-Optimal and Communication-Efficient Decentralized Nonconvex Finite-Sum Optimization. *SIAM Journal on Mathematics of Data Science*, 4(3): 1031–1051.
- Li, H.; and Lin, Z. 2021. Accelerated Gradient Tracking over Timevarying Graphs for Decentralized Optimization. *arXiv preprint arXiv:2104.02596*.
- Li, X.; Xu, Y.; Wang, J. H.; Wang, X.; and Lui, J. C. S. 2021. Decentralized Stochastic Proximal Gradient Descent with Variance Reduction over Time-varying Networks. *arXiv preprint arXiv:2112.10389*.
- Lian, X.; Zhang, C.; Zhang, H.; Hsieh, C.-J.; Zhang, W.; and Liu, J. 2017. Can Decentralized Algorithms Outperform Centralized Algorithms? A Case Study for Decentralized Parallel Stochastic Gradient Descent. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30, 5330–5340. Curran Associates, Inc.
- Lian, X.; Zhang, W.; Zhang, C.; and Liu, J. 2018. Asynchronous Decentralized Parallel Stochastic Gradient Descent. *Proceedings of the 35th International Conference on Machine Learning*, 80: 3043–3052.
- Lorenzo, P. D.; and Scutari, G. 2016. NEXT: In-Network Nonconvex Optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2): 120–136.
- Lu, S.; Zhang, X.; Sun, H.; and Hong, M. 2019. GNSD: a Gradient-Tracking Based Nonconvex Stochastic Algorithm for Decentralized Optimization. In 2019 IEEE Data Science Workshop (DSW), 315–321.
- Mancino-Ball, G.; Xu, Y.; and Chen, J. 2021. A Decentralized Primal-Dual Framework for Non-convex Smooth Consensus Optimization. *arXiv preprint arXiv:2107.11321*.
- Mansoori, F.; and Wei, E. 2021. FlexPD: A Flexible Framework of First-Order Primal-Dual Algorithms for Distributed Optimization. *IEEE Transactions on Signal Processing*, 69: 3500–3512.
- McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and Arcas, B. A. y. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Singh, A.; and Zhu, J., eds., *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, 1273–1282. PMLR.

- Nedic, A.; Olshevsky, A.; and Shi, W. 2017. Achieving Geometric Convergence for Distributed Optimization Over Time-Varying Graphs. *SIAM Journal on Optimization*, 27: 2597 2633.
- Nguyen, L. M.; Liu, J.; Scheinberg, K.; and Takáč, M. 2017. SARAH: A Novel Method for Machine Learning Problems Using Stochastic Recursive Gradient. In Precup, D.; and Teh, Y. W., eds., *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 2613–2621. International Convention Centre, Sydney, Australia: PMLR.
- Pan, T.; Liu, J.; and Wang, J. 2020. D-SPIDER-SFO: A Decentralized Optimization Algorithm with Faster Convergence Rate for Nonconvex Problems. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 1619–1626. AAAI Press.
- Qu, C.; Mannor, S.; Xu, H.; Qi, Y.; Song, L.; and Xiong, J. 2019. Value Propagation for Decentralized Networked Deep Multi-agent Reinforcement Learning. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d' Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Qu, G.; and Li, N. 2019. Accelerated Distributed Nesterov Gradient Descent. *IEEE Transactions on Automatic Control*.
- Scaman, K.; Bach, F.; Bubeck, S.; Lee, Y. T.; and Massoulié, L. 2017. Optimal Algorithms for Smooth and Strongly Convex Distributed Optimization in Networks. In Precup, D.; and Teh, Y. W., eds., *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 3027–3036. International Convention Centre, Sydney, Australia: PMLR.
- Scutari, G.; and Sun, Y. 2019. Distributed nonconvex constrained optimization over time-varying digraphs. *Mathematical Programming*, 176(1): 497–544.
- Shi, W.; Ling, Q.; Wu, G.; and Yin, W. 2015. EXTRA: An Exact First-Order Algorithm for Decentralized Consensus Optimization. *SIAM Journal on Optimization*, 25: 944 966.
- Shi, W.; Ling, Q.; Yuan, K.; Wu, G.; and Yin, W. 2014. On the Linear Convergence of the ADMM in Decentralized Consensus Optimization. *IEEE Transactions on Signal Processing*, 62(7): 1750–1761.
- Singh, N.; Data, D.; George, J.; and Diggavi, S. 2021. SQuARM-SGD: Communication-Efficient Momentum SGD for Decentralized Optimization. In 2021 IEEE International Symposium on Information Theory (ISIT), 1212–1217.
- Sun, H.; and Hong, M. 2019. 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.
- Sun, H.; Lu, S.; and Hong, M. 2020. Improving the Sample and Communication Complexity for Decentralized Non-Convex Optimization: Joint Gradient Estimation and Tracking. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 9217–9228. Virtual: PMLR.
- T. Dinh, C.; Tran, N.; and Nguyen, J. 2020. Personalized Federated Learning with Moreau Envelopes. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M. F.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 21394–21405. Curran Associates, Inc.

- Tang, H.; Gan, S.; Zhang, C.; Zhang, T.; and Liu, J. 2018a. Communication Compression for Decentralized Training. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- Tang, H.; Lian, X.; Yan, M.; Zhang, C.; and Liu, J. 2018b. D^2 : Decentralized Training over Decentralized Data. In Dy, J.; and Krause, A., eds., *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 4848–4856. Stockholmsmässan, Stockholm Sweden: PMLR.
- Tran-Dinh, Q.; Pham, N. H.; Phan, D. T.; and Nguyen, L. M. 2022. A hybrid stochastic optimization framework for composite non-convex optimization. *Mathematical Programming*, 191(2): 1005–1071.
- Vogels, T.; He, L.; Koloskova, A.; Karimireddy, S. P.; Lin, T.; Stich, S. U.; and Jaggi, M. 2021. RelaySum for Decentralized Deep Learning on Heterogeneous Data. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 28004–28015. Curran Associates, Inc.
- Wang, Y.; Yin, W.; and Zeng, J. 2019. Global Convergence of ADMM in Nonconvex Nonsmooth Optimization. *Journal of Scientific Computing*, 78(1): 29–63.
- Wang, Z.; Ji, K.; Zhou, Y.; Liang, Y.; and Tarokh, V. 2019. SpiderBoost and Momentum: Faster Variance Reduction Algorithms. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d' Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Wang, Z.; Zhang, J.; Chang, T.-H.; Li, J.; and Luo, Z.-Q. 2021. Distributed Stochastic Consensus Optimization With Momentum for Nonconvex Nonsmooth Problems. *IEEE Transactions on Signal Processing*, 69: 4486–4501.
- Wei, E.; and Ozdaglar, A. 2013. On the O(1/k) convergence of asynchronous distributed alternating Direction Method of Multipliers. In 2013 IEEE Global Conference on Signal and Information Processing, 551–554.
- Wu, T.; Yuan, K.; Ling, Q.; Yin, W.; and Sayed, A. 2017. Decentralized Consensus Optimization With Asynchrony and Delays. *IEEE Transactions on Signal and Information Processing over Networks*, 4: 293 307.
- Xin, R.; Das, S.; Khan, U. A.; and Kar, S. 2021. A Stochastic Proximal Gradient Framework for Decentralized Non-Convex Composite Optimization: Topology-Independent Sample Complexity and Communication Efficiency. arXiv preprint arXiv:2110.01594.
- Xin, R.; Khan, U.; and Kar, S. 2021a. A Hybrid Variance-Reduced Method for Decentralized Stochastic Non-Convex Optimization. In Meila, M.; and Zhang, T., eds., *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, 11459–11469. PMLR.
- Xin, R.; Khan, U. A.; and Kar, S. 2021b. An Improved Convergence Analysis for Decentralized Online Stochastic Non-Convex Optimization. *IEEE Transactions on Signal Processing*, 69: 1842–1858.
- Xin, R.; Khan, U. A.; and Kar, S. 2022. Fast Decentralized Nonconvex Finite-Sum Optimization with Recursive Variance Reduction. *SIAM Journal on Optimization*, 32(1): 1–28.
- Xu, Y.; and Xu, Y. 2023. Momentum-Based Variance-Reduced Proximal Stochastic Gradient Method for Composite Nonconvex Stochastic Optimization. *Journal of Optimization Theory and Applications*, 196(1): 266–297.

- Ye, H.; Xiong, W.; and Zhang, T. 2020. PMGT-VR: A decentralized proximal-gradient algorithmic framework with variance reduction. *arXiv* preprint arXiv:2012.15010.
- Ye, H.; Zhou, Z.; Luo, L.; and Zhang, T. 2020. Decentralized Accelerated Proximal Gradient Descent. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M. F.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 18308–18317. Curran Associates, Inc.
- Ying, B.; Yuan, K.; Chen, Y.; Hu, H.; PAN, P.; and Yin, W. 2021. Exponential Graph is Provably Efficient for Decentralized Deep Training. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 13975–13987. Curran Associates, Inc.
- Yuan, K.; Chen, Y.; Huang, X.; Zhang, Y.; Pan, P.; Xu, Y.; and Yin, W. 2021. DecentLaM: Decentralized Momentum SGD for Large-Batch Deep Training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 3029–3039.
- Yuan, K.; Ling, Q.; and Yin, W. 2016. On the Convergence of Decentralized Gradient Descent. *SIAM Journal on Optimization*, 26: 1835 1854.
- Zeng, J.; and Yin, W. 2018. On Nonconvex Decentralized Gradient Descent. *IEEE Transactions on Signal Processing*, 66: 2834 2848
- Zhang, J.; and You, K. 2020. Decentralized Stochastic Gradient Tracking for Non-convex Empirical Risk Minimization. *arXiv* preprint arXiv:1909.02712.
- Zhang, J.; and You, K. 2021. Fully Asynchronous Distributed Optimization with Linear Convergence in Directed Networks. *arXiv* preprint arXiv:1901.08215.
- Zhang, K.; Yang, Z.; Liu, H.; Zhang, T.; and Basar, T. 2018. Fully Decentralized Multi-Agent Reinforcement Learning with Networked Agents. In Dy, J.; and Krause, A., eds., *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, 5872–5881. PMLR.

Reproducibility

Data sets and communication graphs. The data sets and communication graphs are summarized/visualized in Table 2.

Dataset	Train	Test	Features	Model	Graph
a9a	32,561	16,281	123	MLP	Ladder
MiniBooNE	100,000	30,064	50	MLP	Ring
MNIST	60,000	10,000	784	LENET	Random

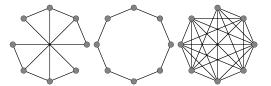


Table 2: Summary of data sets; all are downloaded from https://www.openml.org. Graphs from left to right: ladder, ring, and random.

Code. Code for all numerical experiments is available at https://github.com/gmancino/DEEPSTORM.

Hyperparameter selection. We choose the batch size according to theoretical guidance, while observing reasonably good performance. For all methods but SPPDM, we set the batch size to be 64, 128, and 64 for a9a, MiniBooNE, and MNIST, respectively. For SPPDM, the respective sizes are 512, 1024, and 128.

For all variants of DEEPSTORM, we set the number of communication rounds to be T=1 such that $\mathcal{W}_1=\mathbf{W}$. We use a diminishing step size as in (20) with $k_0 = \lceil \frac{2}{1-\rho^3} \rceil$ and $\beta_k = 1 - \frac{\alpha_{k+1}}{\alpha_k} + \beta \alpha_{k+1}^2$ with $\beta < \frac{1}{\alpha_0 \alpha_1}$, such that $\beta_0 < 1$. Such a choice ensures that $\beta_k \in (0,1)$ for all k. For the (v1-SVRG) variant, we compute the snapshot gradient every four passes, by using all local data for a9a and MiniBooNE; whereas for MNIST, we compute the snapshot gradient at the end of every pass, by using 20% of the local data.

For DSGT, we set the step size to be $\alpha_k = \frac{\alpha}{\sqrt{k+1}}$ for all k, according to (Lu et al. 2019; Xin, Khan, and Kar 2021b). For SPPDM, we follow the choices of many hyperparameters used in the original paper and only tune $c \in \{0.1, 1\}$ and α . For ProxGT-SR-O/E, we tune the step size α and the frequency of communicating the full local gradient, q. We find that q=32yields the most stable results for a9a and MiniBooNE and q=64 performs the best for MNIST. For all these methods, α is tuned from $\{10.0, 5.0, 1.0, 0.1, 0.01, 0.005, 0.001\}$. We choose the Pareto optimal α that balances a small stationarity violation and a high test accuracy.

Chebyshev acceleration

The Chebyshev mixing protocol (Auzinger and Melenk 2011) can be summarized in the following pseudo-code.

Algorithm B.1: Chebyshev mixing protocol $\mathcal{W}_T(\mathbf{B})$

Input: Mixing matrix W, input B, rounds T

- 1: Let $\mathbf{B}_0 = \mathbf{B}$ and $\mathbf{B}_1 = \mathbf{W}\mathbf{B}_0$
- 2: Compute step sizes $\mu_0 = 1, \mu_1 = \frac{1}{a}$

- 3: **for** t = 0, ..., T 1 **do**4: $\mu_{t+1} \leftarrow \frac{2}{\rho} \mu_t \mu_{t-1}$ 5: $\mathbf{B}_{t+1} \leftarrow \frac{2\mu_t}{\rho\mu_{t+1}} \mathbf{W} \mathbf{B}_t \frac{\mu_{t-1}}{\mu_{t+1}} \mathbf{B}_{t-1}$
- 6: end for

Output: $\mathbf{B}_T = \mathcal{W}_T(\mathbf{B}_0)$

This method is accompanied with the following convergence result, relating the spectrum of \mathcal{W}_T to the spectrum of W. For a proof, see (Mancino-Ball, Xu, and Chen 2021).

Lemma B.1 The output of Algorithm B.1 can be denoted as $B_T = \mathcal{W}_T B$, where \mathcal{W}_T is a degree-T polynomial of W and it satisfies Assumptions 2(ii)-(iv). Additionally, we use the bar notation to mean replacing each row of a matrix by the average of its rows; that is, $\bar{\mathbf{B}} = \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{B}$. Then, $\bar{\mathbf{B}}_T = \bar{\mathbf{B}}$ for all T and

$$\|\mathbf{B}_T - \bar{\mathbf{B}}\|_F \le 2\left(1 - \sqrt{1 - \rho}\right)^T \|\mathbf{B} - \bar{\mathbf{B}}\|_F.$$
 (B.1)

The analysis in Appendix C uses a sufficiently large degree T such that $\tilde{\rho}$ as defined in (14) is bounded by a constant, independent of the communication graph. For this, by Corollary 6.1 in (Auzinger and Melenk 2011), it holds that

$$\tilde{\rho} \le 2\left(1 - \sqrt{1 - \rho}\right)^T. \tag{B.2}$$

Hence, by the proof of Theorem 4 in (Mancino-Ball, Xu, and Chen 2021), we see that when $T = \lceil \frac{2}{\sqrt{1-\rho}} \rceil$, we obtain

$$(1 - \tilde{\rho})^2 \ge \frac{1}{2}.\tag{B.3}$$

C Convergence results

We denote the global objective function and the corresponding smooth part to be

$$\phi \triangleq \frac{1}{N} \sum_{i=1}^{N} \phi_i$$
 and $f \triangleq \frac{1}{N} \sum_{i=1}^{N} f_i$ (C.1)

respectively. Crucially, our analysis relies on bounding the difference between the local first-order estimators given in (10) and the true local gradient; namely we define

$$\mathbf{r}_{i}^{(k)} \triangleq \mathbf{d}_{i}^{(k)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}). \tag{C.2}$$

Additionally, we define the following matrix terms to be used throughout the analysis,

$$\bar{\mathbf{A}} \triangleq \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A}, \ \forall \mathbf{A} \in \mathbb{R}^{N \times p},$$
 (C.3)

$$\mathbf{X}_{\perp} \triangleq \mathbf{X} - \bar{\mathbf{X}},\tag{C.4}$$

$$\mathbf{Y}_{\perp} \triangleq \mathbf{Y} - \bar{\mathbf{Y}},$$
 (C.5)

$$\mathbf{R} \triangleq \mathbf{D} - \nabla F(\mathbf{X}),\tag{C.6}$$

where ∇F is the gradient of the smooth part of the objective function written in the following matrix form

$$\nabla F(\mathbf{X}) \triangleq \begin{bmatrix} \nabla f_1(\mathbf{x}_1) \\ \vdots \\ \nabla f_N(\mathbf{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times p}.$$
 (C.7)

Before beginning with the analysis, we present two preparatory Lemmas. The first is standard in the literature (Ghadimi, Lan, and Zhang 2016).

Lemma C.1 Let $r: \mathbb{R}^{1 \times p} \to \mathbb{R}$ be a closed, convex function, then for any $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{1 \times p}$, it holds that

$$\|\operatorname{prox}_{r}(\mathbf{a}) - \operatorname{prox}_{r}(\mathbf{b})\|_{2} \le \|\mathbf{a} - \mathbf{b}\|_{2}.$$
 (C.8)

Lemma C.2 For all $k \ge 0$,

$$\bar{\mathbf{y}}^{(k)} = \bar{\mathbf{d}}^{(k)}.\tag{C.9}$$

Proof We proceed by induction. Notice \mathcal{W}_T is a degree-T polynomial of \mathbf{W} , so $\mathbf{e}^{\top} \mathcal{W}_T = \mathbf{e}^{\top}$ and thus $\bar{\mathbf{y}}^{(0)} = \frac{1}{N} \mathbf{e}^{\top} \mathbf{Y}^{(0)} = \frac{1}{N} \mathbf{e}^{\top} \mathbf{W}_T \mathbf{D}^{(0)} = \frac{1}{N} \mathbf{e}^{\top} \mathbf{D}^{(0)} = \bar{\mathbf{d}}^{(0)}$. For $k \geq 0$, we have

$$\bar{\mathbf{y}}^{(k)} = \frac{1}{N} \mathbf{e}^{\top} \mathbf{Y}^{(k)} \stackrel{\text{(12)}}{=} \frac{1}{N} \mathbf{e}^{\top} \mathcal{W}_{T} \left(\mathbf{Y}^{(k-1)} + \mathbf{D}^{(k)} - \mathbf{D}^{(k-1)} \right) \\
= \frac{1}{N} \mathbf{e}^{\top} \left(\mathbf{Y}^{(k-1)} + \mathbf{D}^{(k)} - \mathbf{D}^{(k-1)} \right) \\
= \bar{\mathbf{y}}^{(k-1)} + \bar{\mathbf{d}}^{(k)} - \bar{\mathbf{d}}^{(k-1)} \\
= \bar{\mathbf{d}}^{(k)}$$

where in the last step we used the inductive hypothesis, $\bar{\mathbf{y}}^{(k-1)} = \bar{\mathbf{d}}^{(k-1)}$.

C.1 Building blocks for constant and diminishing step size convergence.

Our analysis begins by building a non-increasing Lyapunov function by relating changes in X and Y to various quantities.

Lemma C.3 For all $k \ge 0$ and for all i = 1, ..., N,

$$r(\mathbf{x}_{i}^{(k+1)}) - r(\bar{\mathbf{x}}^{(k)}) + \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$\leq -\frac{1}{2\alpha_{k}} \left(\left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} - \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right)$$
(C.10)

Proof By (9), we have

$$\mathbf{0} \in \alpha_k \partial r(\mathbf{x}_i^{(k+1)}) + \mathbf{x}_i^{(k+1)} - \left(\mathbf{z}_i^{(k)} - \alpha_k \mathbf{y}_i^{(k)}\right).$$

Thus, for some $\tilde{\nabla}r(\mathbf{x}_i^{(k+1)}) \in \partial r(\mathbf{x}_i^{(k+1)})$, and for any $\mathbf{x}_i \in \mathbb{R}^{1 \times p}$

$$\left\langle \mathbf{x}_{i}^{(k+1)} - \mathbf{x}_{i}, \tilde{\nabla}r(\mathbf{x}_{i}^{(k+1)}) + \frac{1}{\alpha_{k}} \left(\mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right) + \mathbf{y}_{i}^{(k)} \right\rangle = 0.$$
 (C.11)

By the convexity of r, it holds for any $\mathbf{x}_i \in \mathbb{R}^{1 \times p}$,

$$r(\mathbf{x}_{i}^{(k+1)}) - r(\mathbf{x}_{i}) + \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$\leq \left\langle \mathbf{x}_{i}^{(k+1)} - \mathbf{x}_{i}, \tilde{\nabla}r(\mathbf{x}_{i}^{(k+1)}) + \mathbf{y}_{i}^{(k)} \right\rangle$$

$$\stackrel{\text{(C.11)}}{=} -\frac{1}{\alpha_{k}} \left\langle \mathbf{x}_{i}^{(k+1)} - \mathbf{x}_{i}, \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\rangle$$

$$\stackrel{(a)}{=} -\frac{1}{2\alpha_{k}} \left(\left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{x}_{i} \right\|_{2}^{2} + \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} - \left\| \mathbf{x}_{i} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right),$$

where (a) follows from $\langle \mathbf{a}, \mathbf{b} \rangle = \frac{1}{2} \left(\|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2 - \|\mathbf{a} - \mathbf{b}\|_2^2 \right)$. Letting $\mathbf{x}_i = \bar{\mathbf{x}}^{(k)}$ completes the proof.

Lemma C.4 For all $k \ge 0$.

$$\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)}) \\
\leq \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle \\
- \frac{1}{2N\alpha_{k}} \sum_{i=1}^{N} \left(\left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} - \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right), \tag{C.12}$$

where $\nabla f(\bar{\mathbf{x}}^{(k)}) \triangleq \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(\bar{\mathbf{x}}^{(k)})$ comes from (C.1).

Proof From the L-smoothness of each f_i and the convexity of r, we have

$$\begin{split} &\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)}) \\ &= \frac{1}{N} \sum_{i=1}^{N} \left(f_{i}(\bar{\mathbf{x}}^{(k+1)}) + r(\bar{\mathbf{x}}^{(k+1)}) \right) - \frac{1}{N} \sum_{i=1}^{N} \left(f_{i}(\bar{\mathbf{x}}^{(k)}) + r(\bar{\mathbf{x}}^{(k)}) \right) \\ &= \frac{1}{N} \sum_{i=1}^{N} \left(f_{i}(\bar{\mathbf{x}}^{(k+1)}) - f_{i}(\bar{\mathbf{x}}^{(k)}) \right) + r(\bar{\mathbf{x}}^{(k+1)}) - r(\bar{\mathbf{x}}^{(k)}) \\ &\leq \frac{1}{N} \sum_{i=1}^{N} \left(f_{i}(\bar{\mathbf{x}}^{(k+1)}) - f_{i}(\bar{\mathbf{x}}^{(k)}) \right) + \frac{1}{N} \sum_{i=1}^{N} r(\mathbf{x}_{i}^{(k+1)}) - r(\bar{\mathbf{x}}^{(k)}) \\ &\leq \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\langle \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\bar{\mathbf{x}}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle + \frac{1}{N} \sum_{i=1}^{N} r(\mathbf{x}_{i}^{(k+1)}) - r(\bar{\mathbf{x}}^{(k)}) \\ &\leq \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\langle \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\bar{\mathbf{x}}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle \\ &- \frac{1}{2N\alpha_{k}} \sum_{i=1}^{N} \left(\left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} - \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right). \end{split}$$

Utilizing (C.1) to have $\nabla f(\bar{\mathbf{x}}^{(k)}) = \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(\bar{\mathbf{x}}^{(k)})$ completes the proof.

Lemma C.5 For all $k \ge 0$, the following equality holds,

$$\left\langle \nabla f(\bar{\mathbf{x}}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$= \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle$$

$$+ \frac{1}{N} \sum_{i=1}^{N} \left\langle \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)}, \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle,$$
(C.13)

where $\bar{\mathbf{r}}^{(k)} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{r}_{i}^{(k)}$ for all k.

Proof We have,

$$\left\langle \nabla f(\bar{\mathbf{x}}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$= \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}), \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$\stackrel{(a)}{=} \frac{1}{N} \sum_{i=1}^{N} \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}), \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle - \frac{1}{N} \sum_{i=1}^{N} \left\langle \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}, \mathbf{y}_{i}^{(k)} \right\rangle$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \bar{\mathbf{y}}^{(k)} + \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)}, \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle$$

$$\stackrel{(b)}{=} \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \bar{\mathbf{d}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle + \frac{1}{N} \sum_{i=1}^{N} \left\langle \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)}, \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle$$
(C.14)

where (a) utilizes the linearity of the inner product and (b) comes from Lemma C.2 in conjunction with the linearity of the inner product. Now,

$$\left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \bar{\mathbf{d}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle = \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(\mathbf{x}_i^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle + \left\langle \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(\mathbf{x}_i^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \mathbf{d}_i^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle. \tag{C.15}$$

Plugging (C.15) into (C.14) and utilizing (C.2) completes the proof.

Lemma C.6 For all k > 0, the following inequality holds,

$$\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)}) \leq -\frac{1}{2N} \left(\frac{1}{\alpha_k} - 3L \right) \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_F^2 - \left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle
+ \frac{1}{2N\alpha_k} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{Z}^{(k)} \right\|_F^2 - \frac{1}{2N\alpha_k} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2
+ \frac{L}{2N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \frac{1}{2NL} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_F^2.$$
(C.16)

Proof From (C.12), we use (C.13) to have

$$\begin{split} & \phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)}) \\ \leq & \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} - \left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \\ & + \left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle + \frac{1}{N} \sum_{i=1}^{N} \left\langle \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)}, \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \\ & - \frac{1}{2N\alpha_{k}} \sum_{i=1}^{N} \left(\left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} - \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right). \end{split}$$

We bound terms individually. By Jensen's inequality, we have

$$\left\|\bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)}\right\|_{2}^{2} \le \frac{1}{N} \sum_{i=1}^{N} \left\|\mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)}\right\|_{2}^{2}.$$
(C.17)

By the Peter-Paul inequality, we have

$$\frac{1}{N} \sum_{i=1}^{N} \left\langle \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)}, \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \leq \frac{1}{N} \sum_{i=1}^{N} \left(\frac{L}{2} \left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \frac{1}{2L} \left\| \bar{\mathbf{y}}^{(k)} - \mathbf{y}_{i}^{(k)} \right\|_{2}^{2} \right)$$

and

$$\left\langle \nabla f(\bar{\mathbf{x}}^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \\
\leq \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \frac{1}{2L} \left\| \nabla f(\bar{\mathbf{x}}^{(k)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2} \\
\leq \frac{L}{2} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \frac{1}{2NL} \sum_{i=1}^{N} \left\| \nabla f_{i}(\bar{\mathbf{x}}^{(k)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2} \\
\leq \frac{L}{2N} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \frac{L}{2N} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2},$$

where the second inequality also uses Jensen's inequality. Combining like terms results in

$$\begin{split} & \phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)}) \\ \leq & -\frac{1}{2N} \left(\frac{1}{\alpha_k} - 3L \right) \sum_{i=1}^N \left\| \mathbf{x}_i^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_2^2 - \left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \\ & + \frac{1}{2N\alpha_k} \sum_{i=1}^N \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_i^{(k)} \right\|_2^2 - \frac{1}{2N\alpha_k} \sum_{i=1}^N \left\| \mathbf{x}_i^{(k+1)} - \mathbf{z}_i^{(k)} \right\|_2^2 \\ & + \frac{L}{2N} \sum_{i=1}^N \left\| \mathbf{x}_i^{(k)} - \bar{\mathbf{x}}^{(k)} \right\|_2^2 + \frac{1}{2NL} \sum_{i=1}^N \left\| \bar{\mathbf{y}}^{(k)} - \mathbf{y}_i^{(k)} \right\|_2^2. \end{split}$$

We complete the proof by writing the summations of the 2-norms into the equivalent Frobenius norm expressions.

Lemma C.7 For all $k \ge 0$, the followings hold,

$$\left\| \mathbf{X}_{\perp}^{(k+1)} \right\|_{F}^{2} \leq \tilde{\rho} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{\alpha_{k}^{2}}{1 - \tilde{\rho}} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2}, \tag{C.18}$$

and

$$\mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k+1)} \right\|_{F}^{2} \leq \tilde{\rho} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{1}{1 - \tilde{\rho}} \left(8L^{2} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{F}^{2} + 4\beta_{k}^{2} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} + 4N\beta_{k}^{2} \hat{\sigma}^{2} \right), \tag{C.19}$$

where $\tilde{\rho}$ is defined in (14) and $\hat{\sigma}^2 > 0$ is defined in (v1).

Proof We first prove (C.18). First, we use the following identity

$$\begin{aligned} \operatorname{prox}_{\alpha_k R} \left(\bar{\mathbf{X}}^{(k)} - \alpha_k \bar{\mathbf{Y}}^{(k)} \right) &\triangleq \begin{bmatrix} \operatorname{prox}_{\alpha_k r} \left(\bar{\mathbf{x}}^{(k)} - \alpha_k \bar{\mathbf{y}}^{(k)} \right) \\ \vdots \\ \operatorname{prox}_{\alpha_k r} \left(\bar{\mathbf{x}}^{(k)} - \alpha_k \bar{\mathbf{y}}^{(k)} \right) \end{bmatrix} \in \mathbb{R}^{N \times p} \\ &= \frac{1}{N} \mathbf{e} \mathbf{e}^\top \operatorname{prox}_{\alpha_k R} \left(\bar{\mathbf{X}}^{(k)} - \alpha_k \bar{\mathbf{Y}}^{(k)} \right), \end{aligned}$$

since each row is identical. Then by (9) we have

$$\begin{split} \left\| \mathbf{X}_{\perp}^{(k+1)} \right\|_{F}^{2} &= \left\| \operatorname{prox}_{\alpha_{k}R} \left(\mathbf{W}_{T}(\mathbf{X}^{(k)}) - \alpha_{k} \mathbf{Y}^{(k)} \right) - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \operatorname{prox}_{\alpha_{k}R} \left(\mathbf{W}_{T}(\mathbf{X}^{(k)}) - \alpha_{k} \mathbf{Y}^{(k)} \right) \right\|_{F}^{2} \\ &\stackrel{(a)}{=} \left\| \operatorname{prox}_{\alpha_{k}R} \left(\mathbf{W}_{T}(\mathbf{X}^{(k)}) - \alpha_{k} \mathbf{Y}^{(k)} \right) - \operatorname{prox}_{\alpha_{k}R} \left(\bar{\mathbf{X}}^{(k)} - \alpha_{k} \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &- \left\| \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \left(\operatorname{prox}_{\alpha_{k}R} \left(\bar{\mathbf{X}}^{(k)} - \alpha_{k} \bar{\mathbf{Y}}^{(k)} \right) - \operatorname{prox}_{\alpha_{k}R} \left(\mathbf{W}_{T}(\mathbf{X}^{(k)}) - \alpha_{k} \mathbf{Y}^{(k)} \right) \right) \right\|_{F}^{2} \\ &\stackrel{(\mathbf{C}.8)}{\leq} \left\| \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} - \alpha_{k} \left(\mathbf{Y}^{(k)} - \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &= \left\| \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \left\| \alpha_{k} \left(\mathbf{Y}^{(k)} - \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &- 2 \left\langle \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \left\| \alpha_{k} \left(\mathbf{Y}^{(k)} - \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &\stackrel{(b)}{\leq} \left\| \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \left\| \alpha_{k} \left(\mathbf{Y}^{(k)} - \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &+ \delta \left\| \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \frac{1}{\delta} \left\| \alpha_{k} \left(\mathbf{Y}^{(k)} - \bar{\mathbf{Y}}^{(k)} \right) \right\|_{F}^{2} \\ &\stackrel{(c)}{=} (1 + \delta) \left\| \left(\mathbf{W}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + (1 + \frac{1}{\delta}) \alpha_{k}^{2} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2}, \end{split}$$

where (a) uses that $\frac{1}{N} \mathbf{e} \mathbf{e}^{\top}$ is a projection operator to have, for any matrix $\mathbf{A} \in \mathbb{R}^{N \times p}$,

$$\begin{split} \left\| \mathbf{A} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\|_{F}^{2} &= \left\| \mathbf{A} \right\|_{F}^{2} - 2 \left\langle \mathbf{A}, \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\rangle + \left\| \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\|_{F}^{2} \\ &= \left\| \mathbf{A} \right\|_{F}^{2} - 2 \left\| \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\|_{F}^{2} + \left\| \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\|_{F}^{2} \\ &= \left\| \mathbf{A} \right\|_{F}^{2} - \left\| \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathbf{A} \right\|_{F}^{2}, \end{split}$$

(b) uses the Peter-Paul inequality with $\delta>0$, and (c) uses $\left(\mathcal{W}_T-\frac{1}{N}\mathbf{e}\mathbf{e}^\top\right)=\left(\mathcal{W}_T-\frac{1}{N}\mathbf{e}\mathbf{e}^\top\right)\left(\mathbf{I}-\frac{1}{N}\mathbf{e}\mathbf{e}^\top\right)$. Choosing $\delta=\frac{1-\tilde{\rho}}{\tilde{\rho}}$ with $\tilde{\rho}$ defined in (14) and using the compatibility of the Frobenius norm and the 2-norm to have

$$\left\| \left(\boldsymbol{\mathcal{W}}_T - \frac{1}{N} \mathbf{e} \mathbf{e}^\top \right) \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \leq \left\| \boldsymbol{\mathcal{W}}_T - \frac{1}{N} \mathbf{e} \mathbf{e}^\top \right\|_2^2 \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \stackrel{\text{(14)}}{=} \tilde{\rho}^2 \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2$$

yields (C.18).

To prove (C.19), we use Assumption 2 parts (ii) and (iii) to have

$$\begin{aligned} & \left\| \mathbf{Y}_{\perp}^{(k+1)} \right\|_{F}^{2} \\ & \stackrel{(12)}{=} \left\| \mathbf{W}_{T} \left(\mathbf{Y}^{(k)} + \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right) - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \left(\mathbf{Y}^{(k)} + \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right) \right\|_{F}^{2} \\ & \leq (1 + c_{1}) \left\| \left(\mathbf{W}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{Y}^{(k)} \right\|_{F}^{2} + (1 + \frac{1}{c_{1}}) \left\| \left(\mathbf{W}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \left(\mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right) \right\|_{F}^{2} \\ & \leq (1 + c_{1}) \tilde{\rho}^{2} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + (1 + \frac{1}{c_{1}}) \tilde{\rho}^{2} \left\| \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right\|_{F}^{2} \\ & \leq (1 + c_{1}) \tilde{\rho}^{2} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + (1 + \frac{1}{c_{1}}) \left\| \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right\|_{F}^{2}, \end{aligned}$$

where (a) utilizes $\left(\mathbf{W} - \frac{1}{N}\mathbf{e}\mathbf{e}^{\top}\right) = \left(\mathbf{W} - \frac{1}{N}\mathbf{e}\mathbf{e}^{\top}\right)\left(\mathbf{I} - \frac{1}{N}\mathbf{e}\mathbf{e}^{\top}\right)$ coupled with part (iv) of Assumption 2 and (b) uses $\tilde{\rho}^2 < 1$

and $\left\|\mathbf{I} - \frac{1}{N}\mathbf{e}\mathbf{e}^{\top}\right\|_{2} \leq 1$. Next, by Young's inequality we have

$$\begin{split} & \left\| \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right\|_{F}^{2} \\ &= \left\| (1 - \beta_{k}) \left(\mathbf{V}^{(k+1)} - \mathbf{U}^{(k+1)} \right) + \beta_{k} \left(\tilde{\mathbf{V}}^{(k+1)} - \mathbf{D}^{(k)} \right) \right\|_{F}^{2} \\ &\leq 4 \left((1 - \beta_{k})^{2} \left\| \mathbf{V}^{(k+1)} - \mathbf{U}^{(k+1)} \right\|_{F}^{2} + \beta_{k}^{2} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \right) \\ &+ 4 \left(\beta_{k}^{2} \left\| \tilde{\mathbf{V}}^{(k+1)} - \nabla F(\mathbf{X}^{(k+1)}) \right\|_{F}^{2} + \beta_{k}^{2} \left\| \nabla F(\mathbf{X}^{(k+1)}) - \nabla F(\mathbf{X}^{(k)}) \right\|_{F}^{2} \right) \\ &\leq 4 \left(\left\| \mathbf{V}^{(k+1)} - \mathbf{U}^{(k+1)} \right\|_{F}^{2} + \beta_{k}^{2} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \right) \\ &+ 4 \left(\beta_{k}^{2} \left\| \tilde{\mathbf{V}}^{(k+1)} - \nabla F(\mathbf{X}^{(k+1)}) \right\|_{F}^{2} + \left\| \nabla F(\mathbf{X}^{(k+1)}) - \nabla F(\mathbf{X}^{(k)}) \right\|_{F}^{2} \right), \end{split}$$
(C.20)

where the last inequality comes from the assumption that $\beta_k \in (0,1)$. Hence we have

$$\begin{split} & \left\| \mathbf{Y}_{\perp}^{(k+1)} \right\|_{F}^{2} \\ & \stackrel{\text{(C.20)}}{\leq} (1+c_{1}) \tilde{\rho}^{2} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 4(1+\frac{1}{c_{1}}) \left(\left\| \mathbf{V}^{(k+1)} - \mathbf{U}^{(k+1)} \right\|_{F}^{2} + \beta_{k}^{2} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \right) \\ & + 4(1+\frac{1}{c_{1}}) \left(\beta_{k}^{2} \left\| \tilde{\mathbf{V}}^{(k+1)} - \nabla F(\mathbf{X}^{(k+1)}) \right\|_{F}^{2} + \left\| \nabla F(\mathbf{X}^{(k+1)}) - \nabla F(\mathbf{X}^{(k)}) \right\|_{F}^{2} \right). \end{split}$$

Letting $c_1 = \frac{1}{\bar{\rho}} - 1 > 0$ and then first taking the expectation with respect to the samples and utilizing (4) and (v1) on the above two inequalities and then taking the full expectation, completes the proof.

Our analysis relies on bounding the gradient error term defined in (C.6). Hence, we present the following two Lemmas which define a recursive error bound given either (v1) or (v2) holds for the unbiased estimator $\tilde{\mathbf{v}}_i^{(k+1)}$ in (10).

Lemma C.8 Suppose $\{\mathbf{d}_i^{(t)}\}_{t=0}^{(k)}$ is updated by (10) such that $\tilde{\mathbf{v}}_i$ satisfies (v1) for all iterates $t=0,\ldots,k$, for each agent $i=1,\ldots,N$. Then at iteration k+1, the following bound holds

$$\mathbb{E} \left\| \mathbf{R}^{(k+1)} \right\|_{F}^{2} \leq N \beta_{k}^{2} \hat{\sigma}^{2} + (1 - \beta_{k})^{2} L^{2} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{F}^{2} + (1 - \beta_{k})^{2} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2}. \tag{C.21}$$

Proof The proof follows the same logic as the proof of Lemmas 3 and 4 in (Tran-Dinh et al. 2022), but is included here for the sake of completeness. For sake of brevity, define $B \triangleq B_i^{(k+1)}$ and $\tilde{B} \triangleq \tilde{B}_i^{(k+1)}$. Then for each agent i, by (v1) and the definition of $\mathbf{r}_i^{(k+1)}$ in (C.2), it holds that

$$\mathbb{E}_{(B,\tilde{B})} \left\| \mathbf{r}_{i}^{(k+1)} \right\|_{2}^{2} \\
\stackrel{(10)}{=} \mathbb{E}_{(B,\tilde{B})} \left\| (1-\beta_{k}) \mathbf{d}_{i}^{(k)} + (1-\beta_{k}) \left(\mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right) + \beta_{k} \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right\|_{2}^{2} \\
= \mathbb{E}_{(B,\tilde{B})} \left((1-\beta_{k})^{2} \left\| \mathbf{d}_{i}^{(k)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2} + \beta_{k}^{2} \left\| \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right\|_{2}^{2} \right) \\
+ \mathbb{E}_{(B,\tilde{B})} \left((1-\beta_{k})^{2} \left\| \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right\|_{2}^{2} \right) \\
+ 2(1-\beta_{k})^{2} \mathbb{E}_{(B,\tilde{B})} \left\langle \mathbf{d}_{i}^{(k)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right\rangle \\
+ 2\beta_{k}(1-\beta_{k}) \mathbb{E}_{(B,\tilde{B})} \left\langle \tilde{\mathbf{d}}_{i}^{(k)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right\rangle \\
+ 2\beta_{k}(1-\beta_{k}) \mathbb{E}_{(B,\tilde{B})} \left\langle \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}), \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right\rangle, \tag{C.22}$$

where the second equality comes from adding and subtracting $\beta_k \nabla f_i(\mathbf{x}_i^{(k+1)})$ and $(1-\beta_k) \nabla f_i(\mathbf{x}_i^{(k)})$ and expanding the norm squared. The first two inner products evaluate to zero by the unbiasedness in Assumption 1 (iv). Next, since all $\xi \in B$ are independent from all $\tilde{\xi} \in \tilde{B}$, it holds by $\mathbb{E}_{(B,\tilde{B})}[\cdot] = \mathbb{E}_{\tilde{B}}\left[\mathbb{E}_B\left[\cdot\middle|\tilde{B}\right]\right]$ that the final inner product is zero. Using the unbiasedness assumption of $\mathbf{v}_i^{(k+1)}$ and $\mathbf{u}_i^{(k+1)}$, we have

$$\mathbb{E}_{(B,\tilde{B})} \left((1 - \beta_{k})^{2} \left\| \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right\|_{2}^{2} \right) \\
= \mathbb{E}_{(B,\tilde{B})} (1 - \beta_{k})^{2} \left(\left\| \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right\|_{2}^{2} + \left\| \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2} \right) \\
- \mathbb{E}_{(B,\tilde{B})} \left[2(1 - \beta_{k})^{2} \left\langle \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)}, \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\rangle \right] \\
= (1 - \beta_{k})^{2} \left(\mathbb{E}_{B} \left\| \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right\|_{2}^{2} - \left\| \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2} \right) \\
\leq (1 - \beta_{k})^{2} \mathbb{E}_{B} \left\| \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right\|_{2}^{2}. \tag{C.23}$$

Summing over the agents i = 1, ..., N, utilizing (4), (5), (v1), and taking the full expectation completes the proof.

Lemma C.9 Suppose $\{\mathbf{d}_i^{(t)}\}_{t=0}^{(k)}$ is updated by (10) with (v2) for each agent $i=1,\ldots,N$. Then at iteration k+1, the following bound holds

$$\mathbb{E} \left\| \mathbf{R}^{(k+1)} \right\|_F^2 \le 2N\beta_k^2 \hat{\sigma}^2 + 2(1-\beta_k)^2 L^2 \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_F^2 + (1-\beta_k)^2 \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_F^2. \tag{C.24}$$

Proof The proof follows from Lemma 2 of (Xu and Xu 2023), but is included here for sake of completeness. Using (10) with (v2) and defining $B \triangleq B_i^{(k+1)}$, for each agent i we have

$$\mathbb{E}_{B} \left[\left\| \mathbf{r}_{i}^{(k+1)} \right\|^{2} \right]$$

$$= \mathbb{E}_{B} \left[\left\| \mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) + (1 - \beta_{k}) \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) + (1 - \beta_{k}) \mathbf{r}_{i}^{(k)} \right\|_{2}^{2} \right]$$

$$= \mathbb{E}_{B} \left[\left\| \mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) + (1 - \beta_{k}) \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) \right\|_{2}^{2} + (1 - \beta_{k})^{2} \left\| \mathbf{r}_{i}^{(k)} \right\|_{2}^{2},$$

where we have used

$$\mathbb{E}_{B}\left[\left\langle \mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}), \mathbf{r}_{i}^{(k)} \right\rangle\right] = 0 \text{ and } \mathbb{E}_{B}\left[\left\langle \mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \mathbf{r}_{i}^{(k)} \right\rangle\right] = 0$$

by the definition of the \mathbf{v}_i and \mathbf{u}_i and the unbiasedness in Assumption 1 (iv). Adding and subtracting $\beta_k \left(\mathbf{v}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right)$ inside of the norm of the first term and using Young's inequality results in

$$\mathbb{E}_{B} \left[\left\| \mathbf{r}_{i}^{(k+1)} \right\|^{2} \right] \leq 2\beta_{k}^{2} \mathbb{E}_{B} \left\| \mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right\|_{2}^{2} + (1 - \beta_{k})^{2} \left\| \mathbf{r}_{i}^{(k)} \right\|^{2} \\
+ 2(1 - \beta_{k})^{2} \mathbb{E}_{B} \left\| \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) + \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) \right\|_{2}^{2}.$$

Applying (C.23) to the last term above, summing over all agents $i=1,\ldots,N,$ utilizing both (4) and (v2), and taking the full expectation completes the proof.

Lemma C.10 Suppose $\{\mathbf{d}_i^{(t)}\}_{t=0}^{(k)}$ is updated by (10) such that $\tilde{\mathbf{v}}_i$ satisfies (v1) for each agent $i=1,\ldots,N$. Then at iteration k+1, the following bound holds

$$\mathbb{E}\left\|\bar{\mathbf{r}}^{(k+1)}\right\|_{2}^{2} \le (1-\beta_{k})^{2} \mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_{2}^{2} + \frac{(1-\beta_{k})^{2} L^{2}}{N^{2}} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + \frac{\beta_{k}^{2} \hat{\sigma}^{2}}{N}.$$
(C.25)

Proof The proof follows from Lemma 3 of (Xin, Khan, and Kar 2021a), but is included here for sake of completeness.

Following similar notation to the proof of Lemma C.8, by (10), it holds that

$$\begin{split} \mathbf{r}_{i}^{(k+1)(\mathrm{C.2})} & \stackrel{\cdot}{=} (1-\beta_{k}) \left(\mathbf{d}_{i}^{(k)} + \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right) + \beta_{k} \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \\ & = \beta_{k} \left(\tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) + (1-\beta_{k}) \left(\mathbf{d}_{i}^{(k)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) + \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right) \\ & + (1-\beta_{k}) \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right). \end{split}$$

Taking the average results in

$$\bar{\mathbf{r}}^{(k+1)} = (1 - \beta_k)\bar{\mathbf{r}}^{(k)} + \frac{\beta_k}{N} \sum_{i=1}^{N} \left(\tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right) + \frac{(1 - \beta_k)}{N} \sum_{i=1}^{N} \left(\mathbf{v}_i^{(k+1)} - \mathbf{u}_i^{(k+1)} + \nabla f_i(\mathbf{x}_i^{(k)}) - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right).$$
(C.26)

Defining $B \triangleq \left\{ \xi : \xi \in \bigcup_{i=1}^N B_i^{(k+1)} \right\}$ and $\tilde{B} \triangleq \left\{ \xi : \xi \in \bigcup_{i=1}^N \tilde{B}_i^{(k+1)} \right\}$ we take the norm squared and compute $\mathbb{E}_{(B,\tilde{B})}$, resulting in

$$\mathbb{E}_{(B,\tilde{B})} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} \\
= \mathbb{E}_{(B,\tilde{B})} \left((1 - \beta_{k})^{2} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{\beta_{k}^{2}}{N^{2}} \left\| \sum_{i=1}^{N} \left(\tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) \right\|_{2}^{2} \right) \\
+ \mathbb{E}_{(B,\tilde{B})} \left(\frac{(1 - \beta_{k})^{2}}{N^{2}} \left\| \sum_{i=1}^{N} \left(\left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right) \right\|_{2}^{2} \right) \\
+ \frac{2(1 - \beta_{k})^{2}}{N} \mathbb{E}_{(B,\tilde{B})} \sum_{i=1}^{N} \left\langle \bar{\mathbf{r}}^{(k)}, \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right\rangle \\
+ \frac{2\beta_{k}(1 - \beta_{k})}{N} \mathbb{E}_{(B,\tilde{B})} \sum_{i=1}^{N} \left\langle \bar{\mathbf{r}}^{(k)}, \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right\rangle \\
+ \frac{2\beta_{k}(1 - \beta_{k})}{N^{2}} \mathbb{E}_{(B,\tilde{B})} \sum_{i=1}^{N} \sum_{j=1}^{N} \left\langle \tilde{\mathbf{v}}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}), \left(\mathbf{v}_{j}^{(k+1)} - \nabla f_{j}(\mathbf{x}_{j}^{(k+1)}) \right) - \left(\mathbf{u}_{j}^{(k+1)} - \nabla f_{j}(\mathbf{x}_{j}^{(k)}) \right) \right\rangle. \quad (C.27)$$

Similar to the proof of Lemma C.8, we have the first two inner products evaluate to zero by the unbiasedness in Assumption 1 (iv). Next, since all $\xi \in B$ are independent from all $\tilde{\xi} \in \tilde{B}$, it holds by $\mathbb{E}_{(B,\tilde{B})}[\cdot] = \mathbb{E}_{\tilde{B}}\left[\mathbb{E}_{B}\left[\cdot\middle|\tilde{B}\right]\right]$ that the final inner product is zero. Define

$$\hat{\nabla}_i^{(k+1)} \triangleq \left(\mathbf{v}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right) - \left(\mathbf{u}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k)}) \right)$$

for all $i=1,\ldots,N.$ Hence, by using $\xi\in B$ are independent from $\tilde{\xi}\in \tilde{B},$

$$\mathbb{E}_{(B,\tilde{B})} \left(\frac{(1-\beta_{k})^{2}}{N^{2}} \left\| \sum_{i=1}^{N} \left(\left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right) \right\|_{2}^{2} \right) \\
= \mathbb{E}_{(B,\tilde{B})} \left(\frac{(1-\beta_{k})^{2}}{N^{2}} \sum_{i=1}^{N} \left\| \left(\left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right) \right\|_{2}^{2} \right) \\
+ \frac{(1-\beta_{k})^{2}}{N^{2}} \sum_{i\neq j} \mathbb{E}_{(B,\tilde{B})} \left\langle \hat{\nabla}_{i}^{(k+1)}, \hat{\nabla}_{j}^{(k+1)} \right\rangle \\
= \mathbb{E}_{(B,\tilde{B})} \left(\frac{(1-\beta_{k})^{2}}{N^{2}} \sum_{i=1}^{N} \left\| \left(\left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) - \left(\mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right) \right) \right\|_{2}^{2} \right) \\
\leq \frac{(C.23)}{N^{2}} \sum_{i=1}^{N} \mathbb{E}_{B} \left\| \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right\|_{2}^{2}. \tag{C.28}$$

By similar logic,

$$\mathbb{E}_{\tilde{B}} \frac{\beta_k^2}{N^2} \left\| \sum_{i=1}^N \left(\tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right) \right\|_2^2 \\
= \frac{\beta_k^2}{N^2} \sum_{i=1}^N \mathbb{E}_{\tilde{B}} \left\| \left(\tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right) \right\|_2^2 \\
+ \frac{\beta_k^2}{N^2} \sum_{i \neq j} \mathbb{E}_{\tilde{B}} \left\langle \tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}), \tilde{\mathbf{v}}_j^{(k+1)} - \nabla f_j(\mathbf{x}_j^{(k+1)}) \right\rangle \\
= \frac{\beta_k^2}{N^2} \sum_{i=1}^N \mathbb{E}_{\tilde{B}} \left\| \left(\tilde{\mathbf{v}}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)}) \right) \right\|_2^2. \tag{C.29}$$

Plugging (C.28) to (C.29) into (C.27) and taking the full expectation yields

$$\mathbb{E} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} \\
\leq (1 - \beta_{k})^{2} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{(1 - \beta_{k})^{2}}{N^{2}} \sum_{i=1}^{N} \mathbb{E} \left\| \mathbf{v}_{i}^{(k+1)} - \mathbf{u}_{i}^{(k+1)} \right\|_{2}^{2} + \frac{\beta_{k}^{2} \hat{\sigma}^{2}}{N} \\
\stackrel{(4)}{\leq} (1 - \beta_{k})^{2} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{(1 - \beta_{k})^{2} L^{2}}{N^{2}} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{E}^{2} + \frac{\beta_{k}^{2} \hat{\sigma}^{2}}{N},$$

where we have used (5) and the equivalence of the Frobenius norm to the sum of the squared 2-norms. This completes the proof. \Box

Lemma C.11 Suppose $\{\mathbf{d}_i^{(t)}\}_{t=0}^{(k)}$ is updated by (10) such that $\tilde{\mathbf{v}}_i$ satisfies (v2) for each agent $i=1,\ldots,N$. Then at iteration k+1, the following bound holds

$$\mathbb{E} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} \le (1 - \beta_{k})^{2} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{2(1 - \beta_{k})^{2} L^{2}}{N^{2}} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{F}^{2} + \frac{2\beta_{k}^{2} \hat{\sigma}^{2}}{N}. \tag{C.30}$$

Proof The proof follows from Lemma 3 of (Xin, Khan, and Kar 2021a), but is included here for sake of completeness.

Following similar notation as the proof of Lemma C.10, by (v2), we define $B \triangleq \left\{ \xi : \xi \in \bigcup_{i=1}^N B_i^{(k+1)} \right\}$ to have

$$\begin{split} & \mathbb{E}_{B} \left[\left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} \right] \\ & \stackrel{(10)}{=} \mathbb{E}_{B} \left[\left\| \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) + (1 - \beta_{k}) \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) \right) + (1 - \beta_{k}) \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} \right] \\ & = \mathbb{E}_{B} \left[\left\| \frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) + (1 - \beta_{k}) \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) \right) \right\|_{2}^{2} \right] + (1 - \beta_{k})^{2} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} \end{split}$$

where we have used, for all i = 1, ..., N,

$$\mathbb{E}_{B}\left[\left\langle \mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}), \bar{\mathbf{r}}^{(k)} \right\rangle\right] = 0 \text{ and } \mathbb{E}_{B}\left[\left\langle \mathbf{u}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k)}), \bar{\mathbf{r}}^{(k)} \right\rangle\right] = 0$$

by the definition of the \mathbf{v}_i and \mathbf{u}_i and the unbiasedness in Assumption 1 (iv). Adding and subtracting $\frac{\beta_k}{N}\sum_{i=1}^N \left(\mathbf{v}_i^{(k+1)} - \nabla f_i(\mathbf{x}_i^{(k+1)})\right)$ inside of the norm of the first term and using Young's inequality results in

$$\mathbb{E}_{B} \left[\left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} \right] \leq \frac{2\beta_{k}^{2}}{N^{2}} \mathbb{E}_{B} \left\| \sum_{i=1}^{N} \left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) \right\|_{2}^{2} + (1 - \beta_{k})^{2} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{2(1 - \beta_{k})^{2}}{N^{2}} \mathbb{E}_{B} \left\| \sum_{i=1}^{N} \left(\left(\mathbf{v}_{i}^{(k+1)} - \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) \right) + \left(\nabla f_{i}(\mathbf{x}_{i}^{(k)}) - \mathbf{u}_{i}^{(k+1)} \right) \right) \right\|_{2}^{2}.$$

Applying (C.28) and (C.29), utilizing both (4) and (v2), and taking the full expectation completes the proof.

Before presenting our convergence results, we give the following Lemma which relates relevant terms to a stochastic ε -stationary point as defined in Definition 2.

Lemma C.12 For all $k \ge 0$, the following bound holds,

$$\frac{1}{N} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{L^{2}}{N} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2}
\leq \frac{6}{N} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 6 \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{32L^{2}}{N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{2}{N\alpha_{k}^{2}} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2}.$$
(C.31)

Proof Notice that by (15), for all i = 1, ..., N, we have

$$\left\| P\left(\mathbf{z}_{i}^{(k)}, \mathbf{y}_{i}^{(k)}, \alpha_{k}\right) \right\|_{2}^{2} = \left\| \frac{1}{\alpha_{k}} \left(\mathbf{z}_{i}^{(k)} - \operatorname{prox}_{\alpha_{k}r} \left(\mathbf{z}_{i}^{(k)} - \alpha_{k} \mathbf{y}_{i}^{(k)}\right)\right) \right\|_{2}^{2}$$

$$\stackrel{(9)}{=} \frac{2}{\alpha_{k}^{2}} \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2}$$
(C.32)

and by (C.8), we further have

$$\begin{aligned}
& \left\| P\left(\mathbf{z}_{i}^{(k)}, \mathbf{y}_{i}^{(k)}, \alpha_{k}\right) - P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2} \\
&= \frac{1}{\alpha_{k}^{2}} \left\| \operatorname{prox}_{\alpha_{k}r} \left(\mathbf{z}_{i}^{(k)} - \alpha_{k} \mathbf{y}_{i}^{(k)}\right) - \operatorname{prox}_{\alpha_{k}r} \left(\mathbf{z}_{i}^{(k)} - \alpha_{k} \nabla f(\mathbf{z}_{i}^{(k)})\right) \right\|_{2}^{2} \\
&\leq \left\| \mathbf{y}_{i}^{(k)} - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2}.
\end{aligned} (C.33)$$

By Young's inequality and (C.32), we have

$$\frac{1}{2} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2} \leq \frac{1}{\alpha_{k}^{2}} \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} + \left\| P\left(\mathbf{z}_{i}^{(k)}, \mathbf{y}_{i}^{(k)}, \alpha_{k}\right) - P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2}.$$
(C.34)

Plugging (C.33) into (C.34) and summing over i = 1, ..., N yields

$$\frac{1}{2} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2}$$

$$\leq \frac{1}{\alpha_{k}^{2}} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} + \sum_{i=1}^{N} \left\| \mathbf{y}_{i}^{(k)} - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2}$$

$$= \frac{1}{\alpha_{k}^{2}} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} + \sum_{i=1}^{N} \left\| \mathbf{y}_{i}^{(k)} - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2}, \tag{C.35}$$

where we have utilized the definition of the Frobenius norm. Next, we bound

$$\sum_{i=1}^{N} \left\| \mathbf{y}_{i}^{(k)} - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2} \\
\stackrel{\text{(C.9)}}{=} \sum_{i=1}^{N} \left\| \mathbf{y}_{i}^{(k)} - \bar{\mathbf{y}}^{(k)} + \bar{\mathbf{d}}^{(k)} - \frac{1}{N} \sum_{j=1}^{N} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) + \frac{1}{N} \sum_{j=1}^{N} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2} \\
\leq 3 \sum_{i=1}^{N} \left(\left\| \mathbf{y}_{i}^{(k)} - \bar{\mathbf{y}}^{(k)} \right\|_{2}^{2} + \left\| \bar{\mathbf{d}}^{(k)} - \frac{1}{N} \sum_{j=1}^{N} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) \right\|_{2}^{2} + \left\| \frac{1}{N} \sum_{j=1}^{N} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2} \right). \quad (C.36)$$

Looking at terms individually, we have

$$\left\| \bar{\mathbf{d}}^{(k)} - \frac{1}{N} \sum_{j=1}^{N} \nabla f_j(\mathbf{x}_j^{(k)}) \right\|_2^2 = \left\| \frac{1}{N} \sum_{j=1}^{N} \left(\mathbf{d}_j^{(k)} - \nabla f_j(\mathbf{x}_j^{(k)}) \right) \right\|_2^2 = \left\| \bar{\mathbf{r}}^{(k)} \right\|_2^2.$$
 (C.37)

By Jensen's inequality, we have

$$\left\| \frac{1}{N} \sum_{j=1}^{N} \nabla f_{j}(\mathbf{x}_{j}^{(k)}) - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2}$$

$$= \left\| \frac{1}{N} \sum_{j=1}^{N} \left(\nabla f_{j}(\mathbf{x}_{j}^{(k)}) - \nabla f_{j}(\bar{\mathbf{x}}^{(k)}) + \nabla f_{j}(\bar{\mathbf{x}}^{(k)}) - \nabla f_{j}(\mathbf{z}_{i}^{(k)}) \right) \right\|_{2}^{2}$$

$$\leq \frac{1}{N} \sum_{j=1}^{N} \left\| \nabla f_{j}(\mathbf{x}_{j}^{(k)}) - \nabla f_{j}(\bar{\mathbf{x}}^{(k)}) + \nabla f_{j}(\bar{\mathbf{x}}^{(k)}) - \nabla f_{j}(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2}$$

$$\stackrel{(4)}{\leq} \frac{2L^{2}}{N} \sum_{j=1}^{N} \left(\left\| \mathbf{x}_{j}^{(k)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} + \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2} \right)$$

$$= \frac{2L^{2}}{N} \left\| \mathbf{X}^{(k)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \frac{2L^{2}}{N} \sum_{j=1}^{N} \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2}$$

$$= \frac{2L^{2}}{N} \left\| \mathbf{X}^{(k)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + 2L^{2} \left\| \bar{\mathbf{x}}^{(k)} - \mathbf{z}_{i}^{(k)} \right\|_{2}^{2}.$$
(C.38)

Plugging (C.37) and (C.38) into (C.36) yields

$$\sum_{i=1}^{N} \left\| \mathbf{y}_{i}^{(k)} - \nabla f(\mathbf{z}_{i}^{(k)}) \right\|_{2}^{2} \\
\leq 3 \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 3N \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + 6L^{2} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + 6L^{2} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{Z}^{(k)} \right\|_{F}^{2}.$$
(C.39)

Adding $\frac{L^2}{2} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2$ to both sides of (C.35), applying (C.39), and dividing by N results in

$$\frac{1}{2N} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2} + \frac{L^{2}}{2N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \\
\leq \frac{3}{N} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 3 \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{13L^{2}}{2N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{1}{N\alpha_{L}^{2}} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} + \frac{6L^{2}}{N} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{Z}^{(k)} \right\|_{F}^{2}.$$
(C.40)

Notice that by line 3 of Algorithm 1 and Assumption 2 (iv), $\bar{\mathbf{Z}}^{(k)} = \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \mathcal{W}_T(\mathbf{X}^{(k)}) = \bar{\mathbf{X}}^{(k)}$. Hence

$$\left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2} \stackrel{(8)}{=} \left\| \mathbf{W}_{T}(\mathbf{X}^{(k)}) - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} = \left\| \left(\mathbf{W}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \stackrel{(14)}{\leq} \tilde{\rho}^{2} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} < \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2}, \tag{C.41}$$

so adding $\frac{L^2}{2N} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{E}^2$ to both sides of (C.40) and using (C.41) results in

$$\begin{split} & \frac{1}{2N} \sum_{i=1}^{N} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k} \right) \right\|_{2}^{2} + \frac{L^{2}}{2N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{L^{2}}{2N} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2} \\ & \leq & \frac{3}{N} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 3 \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{16L^{2}}{N} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{1}{N\alpha_{L}^{2}} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2}. \end{split}$$

Finally, we multiply both sides by 2, which completes the proof.

We are now in position to define a lower bounded Lyapunov function and use this to show the convergence of DEEPSTORM v1 and v2. Notice that until this point, the analyses of v1 and v2 of our method only differ slightly, i.e. in terms of the constants involved in Lemmas C.8, C.9, C.10, and C.11. Since the bound established in (C.24) is larger than (C.21), we upper bound (C.21) by (C.24). Additionally, we notice that Lemma C.11 provides an upper bound on the results from Lemma C.10 and hence use (C.30) in the following Lemma.

C.2 Constant step size

Lemma C.13 For all $k \ge 0$, the following inequality holds

$$\underbrace{\left(\frac{1}{4N\alpha_{k}} - \frac{3L}{2N}\right)}_{(A)} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{1}{2N\alpha_{k}} - \frac{72L^{2}\alpha_{k}}{N} - \frac{8L^{2}\alpha_{k}(1-\beta_{k})^{2}}{N^{2}} - \frac{16L^{2}\gamma_{1}^{(k)}}{1-\tilde{\rho}} - (\gamma_{3}^{(k)} + \frac{\gamma_{4}^{(k)}}{N^{2}})4L^{2}(1-\beta_{k})^{2}\right)}_{(B)} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\gamma_{2}^{(k-1)} - \tilde{\rho}\gamma_{2}^{(k)} - \frac{L}{2N} - \frac{\tilde{\rho}^{2}}{2N\alpha_{k}}\right)}_{(C)} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(-(9 + \frac{(1-\beta_{k})^{2}}{N})\frac{32L^{2}\alpha_{k}}{N} - \frac{64L^{2}\gamma_{1}^{(k)}}{1-\tilde{\rho}} - (\gamma_{3}^{(k)} + \frac{\gamma_{4}^{(k)}}{N^{2}})16L^{2}(1-\beta_{k})^{2}\right)}_{(C)} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\gamma_{1}^{(k-1)} - \tilde{\rho}\gamma_{1}^{(k)} - \frac{1}{2NL} - \frac{\gamma_{2}^{(k)}\alpha_{k}^{2}}{1-\tilde{\rho}}\right)}_{(D)}_{(E)} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\gamma_{3}^{(k-1)} - \frac{16\alpha_{k}\beta_{k}^{2}}{N} - \frac{4\gamma_{1}^{(k)}\beta_{k}^{2}}{1-\tilde{\rho}} - \gamma_{3}^{(k)}(1-\beta_{k})^{2}\right)}_{(E)} \mathbb{E} \left\| \mathbf{x}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\gamma_{4}^{(k-1)} - \gamma_{4}^{(k)}(1-\beta_{k})^{2} - 2\alpha_{k}(1-\beta_{k})^{2}\right)}_{(E)}}_{(E)} \mathbb{E} \left\| \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2} \\
\leq \mathbb{E} \left[\Phi^{(k)} - \Phi^{(k+1)} \right] + \underbrace{\left(\frac{4\gamma_{1}^{(k)}N}{1-\tilde{\rho}} + 2N\gamma_{3}^{(k)} + \frac{2\gamma_{4}^{(k)}}{N} + 16\alpha_{k} + \frac{4\alpha_{k}}{N}\right)}_{A} \hat{\sigma}^{2}\beta_{k}^{2},$$

where $\gamma_1^{(k)}, \gamma_2^{(k)}, \gamma_3^{(k)}, \gamma_4^{(k)}$ are strictly positive values and

$$\Phi^{(k)} \triangleq \phi(\bar{\mathbf{x}}^{(k)}) + \gamma_1^{(k-1)} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_F^2 + \gamma_2^{(k-1)} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \gamma_3^{(k-1)} \left\| \mathbf{R}^{(k)} \right\|_F^2 + \gamma_4^{(k-1)} \left\| \bar{\mathbf{r}}^{(k)} \right\|_2^2 \tag{C.43}$$

is a lower bounded Lyapunov function.

Proof We start by using part (iv) of Assumption 2 and (8) to note that

$$\frac{1}{2N\alpha_{k}} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} = \frac{1}{2N\alpha_{k}} \left\| \left(\boldsymbol{\mathcal{W}}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{X}^{(k)} \right\|_{F}^{2} \\
\stackrel{(14)}{\leq} \frac{\tilde{\rho}^{2}}{2N\alpha_{k}} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2}.$$
(C.44)

Next, we utilize the Peter-Paul inequality and Jensen's inequality to have

$$-\left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle$$

$$= \left\langle -\bar{\mathbf{r}}^{(k+1)} + \bar{\mathbf{r}}^{(k+1)} - \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle$$

$$\leq \alpha_{k} \left\| -\bar{\mathbf{r}}^{(k+1)} + \bar{\mathbf{r}}^{(k+1)} - \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{1}{4\alpha_{k}} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2}$$

$$\leq 2\alpha_{k} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} + 2\alpha_{k} \left\| \bar{\mathbf{r}}^{(k+1)} - \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{1}{4N\alpha_{k}} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2}$$

$$= 2\alpha_{k} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_{2}^{2} + 2\alpha_{k} \left\| \bar{\mathbf{r}}^{(k+1)} - \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{1}{4N\alpha_{k}} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2}. \tag{C.45}$$

Further, by Young's inequality it holds that

$$2\alpha_{k} \left\| \bar{\mathbf{r}}^{(k+1)} - \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2}$$

$$= 2\alpha_{k} \left\| \bar{\mathbf{d}}^{(k+1)} - \bar{\mathbf{d}}^{(k)} + \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) - \frac{1}{N} \sum_{i=1}^{N} \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2}$$

$$\leq \frac{4\alpha_{k}}{N} \sum_{i=1}^{N} \left\| \mathbf{d}_{i}^{(k+1)} - \mathbf{d}_{i}^{(k)} \right\|_{2}^{2} + \frac{4\alpha_{k}}{N} \sum_{i=1}^{N} \left\| \nabla f_{i}(\mathbf{x}_{i}^{(k+1)}) - \nabla f_{i}(\mathbf{x}_{i}^{(k)}) \right\|_{2}^{2}$$

$$\stackrel{(4)}{\leq} \frac{4\alpha_{k}}{N} \left\| \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right\|_{F}^{2} + \frac{4L^{2}\alpha_{k}}{N} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k+1)} - \mathbf{x}_{i}^{(k)} \right\|_{2}^{2}$$

$$= \frac{4\alpha_{k}}{N} \left\| \mathbf{D}^{(k+1)} - \mathbf{D}^{(k)} \right\|_{F}^{2} + \frac{4L^{2}\alpha_{k}}{N} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{F}^{2}. \tag{C.46}$$

Taking the expectation conditioned on the local samples and then taking the full expectation yields

$$\frac{4\alpha_{k}}{N}\mathbb{E}\left\|\mathbf{D}^{(k+1)} - \mathbf{D}^{(k)}\right\|_{F}^{2} \leq \frac{4\alpha_{k}}{N}\left(8L^{2}\mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + 4\beta_{k}^{2}\mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} + 4N\beta_{k}^{2}\hat{\sigma}^{2}\right),\tag{C.47}$$

where we have also used (4) and (v1). Plugging (C.47) into (C.46) and using (C.45) yields

$$-\mathbb{E}\left\langle\bar{\mathbf{r}}^{(k)},\bar{\mathbf{x}}^{(k+1)}-\bar{\mathbf{x}}^{(k)}\right\rangle$$

$$\leq 2\alpha_{k}\mathbb{E}\left\|\bar{\mathbf{r}}^{(k+1)}\right\|_{2}^{2}+\frac{1}{4N\alpha_{k}}\mathbb{E}\left\|\mathbf{X}^{(k+1)}-\bar{\mathbf{X}}^{(k)}\right\|_{F}^{2}$$

$$+\frac{36L^{2}\alpha_{k}}{N}\mathbb{E}\left\|\mathbf{X}^{(k+1)}-\mathbf{X}^{(k)}\right\|_{F}^{2}+\frac{16\alpha_{k}\beta_{k}^{2}}{N}\mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2}+16\alpha_{k}\beta_{k}^{2}\hat{\sigma}^{2}.$$
(C.48)

Using (C.44) and (C.48) in (C.16) and taking the full expectation results in

$$\mathbb{E}\left[\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)})\right] \leq -\frac{1}{2N} \left(\frac{1}{2\alpha_k} - 3L\right) \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_F^2 \\
+ \left(\frac{L}{2N} + \frac{\tilde{\rho}^2}{2N\alpha_k}\right) \mathbb{E}\left\|\mathbf{X}^{(k)}_{\perp}\right\|_F^2 - \frac{1}{2N\alpha_k} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_F^2 \\
+ \frac{1}{2NL} \mathbb{E}\left\|\mathbf{Y}^{(k)}_{\perp}\right\|_F^2 + 2\alpha_k \mathbb{E}\left\|\bar{\mathbf{r}}^{(k+1)}\right\|_2^2 \\
+ \frac{36L^2\alpha_k}{N} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_F^2 + \frac{16\alpha_k\beta_k^2}{N} \mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_F^2 + 16\alpha_k\beta_k^2\hat{\sigma}^2. \tag{C.49}$$

Noticing that the right-hand side of (C.24) is larger than the right-hand side of (C.21), we add $\gamma_1^{(k)} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k+1)} \right\|_F^2$, $\gamma_2^{(k)} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k+1)} \right\|_F^2$, $\gamma_3^{(k)} \mathbb{E} \left\| \mathbf{R}^{(k+1)} \right\|_F^2$, $\gamma_4^{(k)} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k+1)} \right\|_2^2$ to both sides of the above inequality and use the results from Lemmas C.7, C.9, and C.11 with (C.43), and subtract $\gamma_1^{(k-1)} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_F^2$, $\gamma_2^{(k-1)} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2$, $\gamma_3^{(k-1)} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_F^2$, $\gamma_4^{(k-1)} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_2^2$

from both sides of the above inequality yields

$$\mathbb{E}\left[\Phi^{(k+1)} - \Phi^{(k)}\right] \\
\leq -\frac{1}{2N}\left(\frac{1}{2\alpha_{k}} - 3L\right) \mathbb{E}\left\|\mathbf{X}_{\perp}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_{F}^{2} \\
+ \left(\frac{L}{2N} + \frac{\tilde{\rho}^{2}}{2N\alpha_{k}}\right) \mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} - \frac{1}{2N\alpha_{k}} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_{F}^{2} \\
+ \frac{1}{2NL} \mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2} + 2\alpha_{k} \mathbb{E}\left\|\bar{\mathbf{r}}^{(k+1)}\right\|_{2}^{2} \\
+ \frac{36L^{2}\alpha_{k}}{N} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + \frac{16\alpha_{k}\beta_{k}^{2}}{N} \mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} + 16\alpha_{k}\beta_{k}^{2}\hat{\sigma}^{2} \\
+ \gamma_{1}^{(k)}\left(\tilde{\rho}\mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2} + \frac{1}{1-\tilde{\rho}}\left(8L^{2}\mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + 4\beta_{k}^{2}\mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} + 4N\beta_{k}^{2}\hat{\sigma}^{2}\right)\right) \\
+ \gamma_{2}^{(k)}\left(\tilde{\rho}\mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} + \frac{\alpha_{k}^{2}}{1-\tilde{\rho}}\mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2}\right) \\
+ \gamma_{3}^{(k)}\left(2N\beta_{k}^{2}\hat{\sigma}^{2} + 2(1-\beta_{k})^{2}L^{2}\mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + (1-\beta_{k})^{2}\mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2}\right) \\
+ \gamma_{4}^{(k)}\left((1-\beta_{k})^{2}\mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_{2}^{2} + \frac{2(1-\beta_{k})^{2}L^{2}}{N^{2}}\mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} + \frac{2\beta_{k}^{2}\hat{\sigma}^{2}}{N}\right) \\
- \gamma_{1}^{(k-1)}\mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2} - \gamma_{2}^{(k-1)}\mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} - \gamma_{3}^{(k-1)}\mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} - \gamma_{4}^{(k-1)}\mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_{2}^{2}. \tag{C.50}$$

Grouping like terms in (C.50) results in

$$\begin{split} & \mathbb{E}\left[\Phi^{(k+1)} - \Phi^{(k)}\right] \\ & \leq -\frac{1}{2N}\left(\frac{1}{2\alpha_k} - 3L\right) \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_F^2 - \frac{1}{2N\alpha_k} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_2^{(k-1)} + \tilde{\rho}\gamma_2^{(k)} + \frac{L}{2N} + \frac{\tilde{\rho}^2}{2N\alpha_k}\right) \mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_1^{(k-1)} + \tilde{\rho}\gamma_1^{(k)} + \frac{1}{2NL} + \frac{\gamma_2^{(k)}\alpha_k^2}{1 - \tilde{\rho}}\right) \mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_3^{(k-1)} + \gamma_3^{(k)}(1 - \beta_k)^2 + \frac{16\alpha_k\beta_k^2}{N} + \frac{4\gamma_1^{(k)}\beta_k^2}{1 - \tilde{\rho}}\right) \mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_4^{(k-1)} + \gamma_4^{(k)}(1 - \beta_k)^2\right) \mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_2^2 \\ & + \left(\frac{36L^2\alpha_k}{N} + \frac{8L^2\gamma_1^{(k)}}{1 - \tilde{\rho}} + 2L^2\gamma_3^{(k)}(1 - \beta_k)^2 + \frac{2L^2\gamma_4^{(k)}(1 - \beta_k)^2}{N^2}\right) \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_F^2 \\ & + \left(16\alpha_k + \frac{4N\gamma_1^{(k)}}{1 - \tilde{\rho}} + 2N\gamma_3^{(k)} + \frac{2\gamma_4^{(k)}}{N}\right) \beta_k^2\hat{\sigma}^2 \\ & + 2\alpha_k \mathbb{E}\left\|\bar{\mathbf{r}}^{(k+1)}\right\|_2^2. \end{split}$$

Next, we use that the right-hand side of (C.30) is larger than the right-hand side of (C.25) to have,

$$\begin{split} & \mathbb{E}\left[\Phi^{(k+1)} - \Phi^{(k)}\right] \\ & \stackrel{\text{(C.30)}}{\leq} - \frac{1}{2N}\left(\frac{1}{2\alpha_k} - 3L\right) \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_F^2 - \frac{1}{2N\alpha_k} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_2^{(k-1)} + \tilde{\rho}\gamma_2^{(k)} + \frac{L}{2N} + \frac{\tilde{\rho}^2}{2N\alpha_k}\right) \mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_1^{(k-1)} + \tilde{\rho}\gamma_1^{(k)} + \frac{1}{2NL} + \frac{\gamma_2^{(k)}\alpha_k^2}{1 - \tilde{\rho}}\right) \mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_3^{(k-1)} + \gamma_3^{(k)}(1 - \beta_k)^2 + \frac{16\alpha_k\beta_k^2}{N} + \frac{4\gamma_1^{(k)}\beta_k^2}{1 - \tilde{\rho}}\right) \mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_F^2 \\ & + \left(-\gamma_4^{(k-1)} + \gamma_4^{(k)}(1 - \beta_k)^2\right) \mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_2^2 \\ & + \left(\frac{36L^2\alpha_k}{N} + \frac{8L^2\gamma_1^{(k)}}{1 - \tilde{\rho}} + 2L^2\gamma_3^{(k)}(1 - \beta_k)^2 + \frac{2L^2\gamma_4^{(k)}(1 - \beta_k)^2}{N^2}\right) \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_F^2 \\ & + \left(16\alpha_k + \frac{4N\gamma_1^{(k)}}{1 - \tilde{\rho}} + 2N\gamma_3^{(k)} + \frac{2\gamma_4^{(k)}}{N}\right)\beta_k^2\hat{\sigma}^2 \\ & + 2\alpha_k\left((1 - \beta_k)^2\mathbb{E}\left\|\bar{\mathbf{r}}^{(k)}\right\|_2^2 + \frac{2(1 - \beta_k)^2L^2}{N^2}\mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_F^2 + \frac{2\beta_k^2\hat{\sigma}^2}{N}\right). \end{split}$$

Further using

$$\begin{aligned} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_F^2 &= \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} + \mathbf{Z}^{(k)} - \mathbf{X}^{(k)} \right\|_F^2 \\ &\leq 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 2 \left\| \mathbf{W}_T(\mathbf{X}^{(k)}) - \mathbf{W}_T \bar{\mathbf{X}}^{(k)} + \mathbf{W}_T \bar{\mathbf{X}}^{(k)} - \mathbf{X}^{(k)} \right\|_F^2 \\ &= 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 2 \left\| (\mathbf{I} - \mathbf{W}_T) \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \\ &\leq 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 8 \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \end{aligned}$$

and combining like terms completes the proof.

Proof of Theorem 1 Proof We approach this proof in phases; first, we note that $\beta_k \in (0,1)$ by $0 < \alpha \le \frac{K^{\frac{1}{3}}}{32L}$. Second, let

$$\gamma_{1}^{(k)} \triangleq \frac{1}{NL(1-\tilde{\rho})}, \gamma_{2}^{(k)} \triangleq \frac{16K^{\frac{1}{3}}}{N(1-\tilde{\rho})\alpha}, \gamma_{3}^{(k)} \triangleq \frac{K^{\frac{1}{3}}}{48NL^{2}\alpha}, \text{ and } \gamma_{4}^{(k)} \triangleq \frac{NK^{\frac{1}{3}}}{48L^{2}\alpha}$$
(C.51)

in (C.42) to have

$$\underbrace{\left(\frac{K^{\frac{1}{3}}}{4N\alpha} - \frac{3L}{2N}\right)}_{(A')} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{K^{\frac{1}{3}}}{2N\alpha} - \frac{72L^{2}\alpha}{NK^{\frac{1}{3}}} - \frac{8L^{2}\alpha(1-\beta_{k})^{2}}{N^{2}K^{\frac{1}{3}}} - \frac{16L}{N(1-\tilde{\rho})^{2}} - \frac{(1-\beta_{k})^{2}K^{\frac{1}{3}}}{6N\alpha} \right)}_{(B')} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{16K^{\frac{1}{3}}}{N\alpha} - \frac{L}{2N} - \frac{\tilde{\rho}^{2}K^{\frac{1}{3}}}{2N\alpha} - \frac{288L^{2}\alpha}{NK^{\frac{1}{3}}} - \frac{32L^{2}\alpha(1-\beta_{k})^{2}}{N^{2}K^{\frac{1}{3}}} - \frac{64L}{N(1-\tilde{\rho})^{2}} - \frac{2(1-\beta_{k})^{2}K^{\frac{1}{3}}}{3N\alpha} \right)}_{(C')} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{1}{NL} - \frac{1}{2NL} - \frac{16\alpha}{N(1-\tilde{\rho})^{2}K^{\frac{1}{3}}}\right)}_{(D')} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{K^{\frac{1}{3}}}{48NL^{2}\alpha} - \frac{(1-\beta_{k})^{2}K^{\frac{1}{3}}}{48NL^{2}\alpha} - \frac{16\alpha\beta_{k}^{2}}{NK^{\frac{1}{3}}} - \frac{4\beta_{k}^{2}}{NL(1-\tilde{\rho})^{2}}\right)}_{(E')} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{NK^{\frac{1}{3}}}{48L^{2}\alpha} - \frac{(1-\beta_{k})^{2}NK^{\frac{1}{3}}}{48L^{2}\alpha} - \frac{2\alpha(1-\beta_{k})^{2}}{K^{\frac{1}{3}}}\right)}_{(E')} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} \\
\leq \mathbb{E} \left[\Phi^{(k)} - \Phi^{(k+1)} \right] + \underbrace{\left(\frac{4}{L(1-\tilde{\rho})^{2}} + \frac{K^{\frac{1}{3}}}{12L^{2}\alpha} + \frac{16\alpha}{K^{\frac{1}{3}}} + \frac{4\alpha}{NK^{\frac{1}{3}}}\right)}_{C} \hat{\sigma}^{2} \beta_{k}^{2}.$$
(C.52)

Next, we lower bound (A') - (E'). For (A'), we have

$$(A') = \frac{K^{\frac{1}{3}}}{4N\alpha} - \frac{3L}{2N} \ge \frac{K^{\frac{1}{3}}}{4N\alpha} - \frac{3K^{\frac{1}{3}}}{64N\alpha}$$
$$> \frac{K^{\frac{1}{3}}}{5N\alpha},$$

where the first inequality uses $\alpha \leq \frac{K^{\frac{1}{3}}}{32L}$. For (B'), we use $(1-\beta_k)^2 < 1$ and $\alpha \leq \min\{\frac{K^{\frac{1}{3}}}{32L}, \frac{(1-\tilde{\rho})^2K^{\frac{1}{3}}}{64L}\}$ to have

$$\begin{split} (B') = & \frac{K^{\frac{1}{3}}}{2N\alpha} - \frac{72L^{2}\alpha}{NK^{\frac{1}{3}}} - \frac{8L^{2}\alpha(1-\beta_{k})^{2}}{N^{2}K^{\frac{1}{3}}} - \frac{16L}{N(1-\tilde{\rho})^{2}} - \frac{(1-\beta_{k})^{2}K^{\frac{1}{3}}}{6N\alpha} \\ = & \frac{K^{\frac{1}{3}}}{N\alpha} \left(\frac{1}{2} - \frac{72L^{2}\alpha^{2}}{K^{\frac{2}{3}}} - \frac{8L^{2}\alpha^{2}(1-\beta_{k})^{2}}{NK^{\frac{2}{3}}} - \frac{16L\alpha}{(1-\tilde{\rho})^{2}K^{\frac{1}{3}}} - \frac{(1-\beta_{k})^{2}}{6} \right) \\ > & \frac{K^{\frac{1}{3}}}{N\alpha} \left(\frac{1}{2} - \frac{72}{1024} - \frac{8}{1024N} - \frac{16}{64} - \frac{1}{6} \right) \\ > & \frac{K^{\frac{1}{3}}}{256N\alpha}. \end{split}$$

For (C'), we again use $(1 - \beta_k) < 1$ and $\tilde{\rho}^2 < 1$ to have

$$\begin{split} (C') = & \frac{16K^{\frac{1}{3}}}{N\alpha} - \frac{L}{2N} - \frac{\tilde{\rho}^2K^{\frac{1}{3}}}{2N\alpha} - \frac{288L^2\alpha}{NK^{\frac{1}{3}}} - \frac{32L^2\alpha(1-\beta_k)^2}{N^2K^{\frac{1}{3}}} - \frac{64L}{N(1-\tilde{\rho})^2} - \frac{2(1-\beta_k)^2K^{\frac{1}{3}}}{3N\alpha} \\ > & \frac{K^{\frac{1}{3}}}{N\alpha} \left(16 - \frac{L\alpha}{2K^{\frac{1}{3}}} - \frac{1}{2} - \frac{288L^2\alpha^2}{K^{\frac{2}{3}}} - \frac{32L^2\alpha^2}{NK^{\frac{2}{3}}} - \frac{64L\alpha}{(1-\tilde{\rho})^2K^{\frac{1}{3}}} - \frac{2}{3} \right) \\ \ge & \frac{K^{\frac{1}{3}}}{N\alpha} \left(16 - \frac{1}{64} - \frac{1}{2} - \frac{288}{1024} - \frac{32}{1024N} - 1 - \frac{2}{3} \right) \\ > & \frac{12K^{\frac{1}{3}}}{N\alpha}, \end{split}$$

where the second to last inequality uses $\alpha \leq \min\{\frac{K^{\frac{1}{3}}}{32L}, \frac{(1-\tilde{\rho})^2K^{\frac{1}{3}}}{64L}\}$. For (D'), it holds that

$$(D') = \frac{1}{NL} - \frac{1}{2NL} - \frac{16\alpha}{N(1-\tilde{\rho})^2 K^{\frac{1}{3}}} = \frac{1}{2NL} - \frac{16\alpha}{N(1-\tilde{\rho})^2 K^{\frac{1}{3}}}$$
$$\geq \frac{1}{2NL} - \frac{1}{4NL}$$
$$= \frac{1}{4NL},$$

where we have used $\alpha \leq \frac{(1-\tilde{\rho})^2K^{\frac{1}{3}}}{64L}$. For (E'), we expand $1-(1-\beta_k)^2=2\beta_k-\beta_k^2>\beta_k$ and use (17) and $\beta_k^2<\beta_k$ to have

$$\begin{split} (E') = & \frac{K^{\frac{1}{3}}}{48NL^{2}\alpha} - \frac{(1-\beta_{k})^{2}K^{\frac{1}{3}}}{48NL^{2}\alpha} - \frac{16\alpha\beta_{k}^{2}}{NK^{\frac{1}{3}}} - \frac{4\beta_{k}^{2}}{NL(1-\tilde{\rho})^{2}} \\ > & \frac{3\alpha}{N^{2}K^{\frac{1}{3}}} - \frac{16\cdot144^{2}L^{4}\alpha^{5}}{N^{3}K^{\frac{5}{3}}} - \frac{4\cdot144^{2}L^{3}\alpha^{4}}{N^{3}(1-\tilde{\rho})^{2}K^{\frac{4}{3}}} \\ \geq & \frac{\alpha}{N^{2}K^{\frac{1}{3}}} \left(3 - \frac{1}{2} - \frac{4}{3}\right) \\ > & \frac{\alpha}{N^{2}K^{\frac{1}{3}}}, \end{split}$$

where the second inequality uses $\alpha^4 \leq \frac{K^{\frac{4}{3}}}{32^4L^4}$ and $\alpha^3 \leq \frac{(1-\tilde{\rho})^2K}{64\cdot 32^2L^3}$. For (F'), we also expand $1-(1-\beta_k)^2=2\beta_k-\beta_k^2>\beta_k$ and use (17) to have

$$(F') = \frac{NK^{\frac{1}{3}}}{48L^{2}\alpha} - \frac{(1-\beta_{k})^{2}NK^{\frac{1}{3}}}{48L^{2}\alpha} - \frac{2\alpha(1-\beta_{k})^{2}}{K^{\frac{1}{3}}}$$
$$> \frac{3\alpha}{K^{\frac{1}{3}}} - \frac{2\alpha}{K^{\frac{1}{3}}}$$
$$= \frac{\alpha}{K^{\frac{1}{3}}}.$$

Next, we sum (C.52) over k = 0 to K - 1 and divide by K; using the established lower bounds to have

$$\frac{1}{K} \sum_{k=0}^{K-1} \left(\frac{K^{\frac{1}{3}}}{5N\alpha} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \frac{K^{\frac{1}{3}}}{256N\alpha} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} + \frac{\alpha}{K^{\frac{1}{3}}} \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} \right) \\
+ \frac{1}{K} \sum_{k=0}^{K-1} \left(\frac{12K^{\frac{1}{3}}}{N\alpha} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{1}{4NL} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{\alpha}{N^{2}K^{\frac{1}{3}}} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \right) \\
\leq \frac{1}{K} \left(\Phi^{(0)} - \phi^{*} \right) + \left(\frac{4}{L(1-\tilde{\rho})^{2}K^{\frac{4}{3}}} + \frac{1}{12L^{2}\alpha K} + \frac{16\alpha}{K^{\frac{5}{3}}} + \frac{4\alpha}{NK^{\frac{5}{3}}} \right) \frac{144^{2}L^{4}\alpha^{4}\hat{\sigma}^{2}}{N^{2}}, \tag{C.53}$$

where we have used $\phi^* \leq \Phi^{(k)}$ for any $k \geq 0$.

The final phase of the proof uses Lemma C.12 to provide a concise convergence statement. To do so, multiply both sides

of (C.31) by $\frac{\alpha}{K^{\frac{1}{3}}}$, sum from $k=0,\ldots,K-1$ and divide by K, and take the expectation to have

$$\frac{\alpha}{K^{\frac{1}{3}} \cdot K} \sum_{k=0}^{K-1} \left(\frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k} \right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2} \right) \\
\leq \frac{\alpha}{K^{\frac{1}{3}} \cdot K} \sum_{k=0}^{K-1} \left(\frac{6}{N} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + 6 \mathbb{E} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{32L^{2}}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{2}{N\alpha_{k}^{2}} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} \right). \tag{C.54}$$

We relate each of the terms on the right-hand side of (C.54) to 512 times of the left-hand side of (C.53). Since $\alpha \leq \frac{K^{\frac{1}{3}}}{32L}$ it holds that

$$\frac{\alpha}{K^{\frac{1}{3}}} \cdot \frac{6}{NK} \sum_{k=0}^{K-1} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} \le \frac{128}{NL} \cdot \frac{1}{K} \sum_{k=0}^{K-1} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2}. \tag{C.55}$$

Using $\alpha \leq \frac{K^{\frac{1}{3}}}{32L}$ we have,

$$\frac{\alpha}{K^{\frac{1}{3}}} \cdot \frac{32L^2}{NK} \sum_{k=0}^{K-1} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \le \frac{6144K^{\frac{1}{3}}}{\alpha} \cdot \frac{1}{NK} \sum_{k=0}^{K-1} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2. \tag{C.56}$$

Combining (C.55) and (C.56) in conjunction with 512 times of (C.53) and (C.54) yields:

$$\begin{split} &\frac{\alpha}{K^{\frac{1}{3}} \cdot K} \sum_{k=0}^{K-1} \frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(k)}, \nabla f(\mathbf{z}_{i}^{(k)}), \alpha_{k}\right) \right\|_{2}^{2} \\ &+ \frac{\alpha}{K^{\frac{1}{3}} \cdot K} \sum_{k=0}^{K-1} \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \frac{\alpha}{K^{\frac{1}{3}} \cdot K} \sum_{k=0}^{K-1} \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_{F}^{2} \\ &\leq \frac{512}{K} \left(\Phi^{(0)} - \phi^{*} \right) + \left(\frac{2048}{L(1 - \tilde{\rho})^{2} K^{\frac{4}{3}}} + \frac{128}{3L^{2} \alpha K} + \frac{8192 \alpha}{K^{\frac{5}{3}}} + \frac{2048 \alpha}{N K^{\frac{5}{3}}} \right) \frac{144^{2} L^{4} \alpha^{4} \hat{\sigma}^{2}}{N^{2}}. \end{split}$$

Multiplying both sides by $\frac{K^{\frac{1}{3}}}{\alpha}$ and using $\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} \geq 0$ and $\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} \geq 0$ for all $k \geq 0$ completes the proof.

Complexity analysis Before presenting the complexity analysis for DEEPSTORM with a constant step size, we first present a preparatory Lemma.

Lemma C.14 For any real numbers $x \in (0,1)$, it holds that

$$\frac{1}{x} \ge \frac{-1}{\ln(1-x)}.\tag{C.57}$$

Proof We prove

$$g(x) \triangleq x + \ln(1 - x) \le 0,\tag{C.58}$$

which is equivalent to (C.57). Computing the first derivative results in

$$\frac{d}{dx}g(x) = 1 - \frac{1}{1-x} < 0, \ \forall x \in (0,1)$$

since (1-x) < 1 for all $x \in (0,1)$. Hence g(x) is decreasing on (0,1). Computing g(0+) = 0, we have (C.58) and hence (C.57).

We make the following remark in order to aid in the discussion of presenting final complexity results for Algorithm 1 with a constant step size.

Remark C.1 Notice that the convergence of Algorithm 1 depends upon $\Phi^{(0)}$ (see (C.43)) which in turn depends upon $\gamma_2^{(0)}$, $\gamma_3^{(0)}$, and $\gamma_4^{(0)}$, all of which are $\mathcal{O}\left(K^{\frac{1}{3}}\right)$. In order to obtain the best possible convergence rate which is also independent of the communication network, we need $\|\mathbf{R}^{(0)}\|_F^2 = \mathcal{O}\left(K^{-\frac{1}{3}}\right)$, which occurs when the initial batch size, denoted as m_0 , is large enough; by Jensen's inequality, this will in turn make $\|\bar{\mathbf{r}}^{(0)}\|_2^2 = \mathcal{O}\left(K^{-\frac{1}{3}}\right)$. Additionally, we make a standard assumption (Lian

et al. 2017; Tang et al. 2018b; Xin, Khan, and Kar 2021a) that $\mathbf{x}_i^{(0)} = \mathbf{x}_j^{(0)}$ for all i and j; this eliminates the $\left\|\mathbf{X}_{\perp}^{(0)}\right\|_F^2$ error. For the gradient tracking term $\left\|\mathbf{Y}_{\perp}^{(0)}\right\|_F^2$, we need to perform sufficiently many initial communications so that $\gamma_1^{(0)}\left\|\mathbf{Y}_{\perp}^{(0)}\right\|_F^2 = \mathcal{O}\left(1\right)$ independent of $\tilde{\rho}$. Notice that $\frac{1}{NL(1-\tilde{\rho})} = \mathcal{O}\left(1\right)$ if NL is sufficiently large and $\tilde{\rho}$ is not too close to 1; in the following Corollary, we assume the worst case so that $\frac{1}{NL(1-\tilde{\rho})} = \mathcal{O}\left(\frac{1}{(1-\tilde{\rho})}\right)$.

Corollary 1 Let $\varepsilon > 0$ be given and assume that $L \ge 1$. Under the same conditions as in Theorem 1, let $\mathbf{x}_i^{(0)} = \mathbf{x}_j^{(0)}$ for all $i, j = 1, \ldots, N$, let the initial batch size $m_0 = \sqrt[3]{NK}$ for all $i = 1, \ldots, N$, and perform $T_0 = \left\lceil \frac{-2\ln(1-\tilde{\rho})}{\sqrt{1-\rho}} \right\rceil$ communications by Algorithm B.1 for the initial gradient tracking update in line 1 of Algorithm 1. Let $\tilde{\mathbf{v}}_i^{(k+1)}$ be any unbiased gradient estimator such that either (v1) or (v2) holds, let the local batch size $m = \mathcal{O}(1)$ for all remaining iterations, and choose α such that

$$\alpha = \frac{N^{\frac{2}{3}}}{64L}.\tag{C.59}$$

Then, provided $K \geq \frac{N^2}{(1-\tilde{\rho})^6}$, Algorithm 1 produces a stochastic ε -stationary point as defined in Definition 2 in

$$K = \mathcal{O}\left(\max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^2}{(1 - \tilde{\rho})^2\varepsilon}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right)$$
(C.60)

local stochastic gradient computations and $T_0 + T(K-1)$ neighbor communications for any $T \ge 1$, where $\Delta \triangleq \Phi^{(0)} - \phi^*$, with $\Phi^{(0)}$ defined in (C.43).

Proof First, notice that if $K \geq \frac{N^2}{(1-\tilde{\rho})^6}$, then $K^{\frac{1}{3}} \geq \frac{N^{\frac{2}{3}}}{(1-\tilde{\rho})^2}$, so the choice of α in (C.59) satisfies $\alpha \leq \frac{(1-\rho)^2K^{\frac{1}{3}}}{64L}$ since

$$\alpha = \frac{N^{\frac{2}{3}}}{64L} \le \frac{(1-\rho)^2 K^{\frac{1}{3}}}{64L}.$$

Next, by (C.43) and (C.51), we have

$$\Phi^{(0)} = \phi(\bar{\mathbf{x}}^{(0)}) + \frac{1}{NL(1-\tilde{\rho})} \left\| \mathbf{Y}_{\perp}^{(0)} \right\|_{F}^{2} + \frac{1024LK^{\frac{1}{3}}}{N^{\frac{5}{3}}(1-\tilde{\rho})} \left\| \mathbf{X}_{\perp}^{(0)} \right\|_{F}^{2} + \frac{4K^{\frac{1}{3}}}{3N^{\frac{5}{3}}L} \left\| \mathbf{R}^{(0)} \right\|_{F}^{2} + \frac{4N^{\frac{1}{3}}K^{\frac{1}{3}}}{3L} \left\| \bar{\mathbf{r}}^{(0)} \right\|_{2}^{2}.$$
 (C.61)

Notice that

$$\frac{1}{N} \mathbb{E} \left\| \mathbf{R}^{(0)} \right\|_{F}^{2} = \frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| \mathbf{d}_{i}^{(0)} - \nabla f_{i}(\mathbf{x}_{i}^{(0)}) \right\|_{2}^{2} \le \frac{\sigma^{2}}{m_{0}}$$

by (5) since $\mathbf{d}_{i}^{(0)} = \frac{1}{m_0} \sum_{\xi \in B_{i}^{(0)}} \nabla f_i(\mathbf{x}_{i}^{(0)}, \xi)$ for all i = 1, ..., N. Hence with $m_0 = \sqrt[3]{NK}$, it holds that

$$\frac{4K^{\frac{1}{3}}}{3N^{\frac{5}{3}}L} \left\| \mathbf{R}^{(0)} \right\|_{F}^{2} \leq \frac{4\sigma^{2}}{3NL} \leq \frac{4\sigma^{2}}{3L}$$

which is independent of K. By Jensen's inequality, we further have

$$\frac{4N^{\frac{1}{3}}K^{\frac{1}{3}}}{3L} \left\| \bar{\mathbf{r}}^{(0)} \right\|_{2}^{2} \le \frac{4K^{\frac{1}{3}}}{N^{\frac{2}{3}}3L} \left\| \mathbf{R}^{(0)} \right\|_{F}^{2} \le \frac{4\sigma^{2}}{3L},$$

which is also independent of K. Next, notice that $\mathbf{Y}^{(0)} = \mathcal{W}_{T_0}(\mathbf{D}^{(0)})$ by line 1 in Algorithm 1; hence it holds

$$\frac{1}{NL(1-\tilde{\rho})} \left\| \mathbf{Y}_{\perp}^{(0)} \right\|_{F}^{2} \stackrel{\text{(C.9)}}{=} \frac{1}{NL(1-\tilde{\rho})} \left\| \boldsymbol{\mathcal{W}}_{T_{0}}(\mathbf{D}^{(0)}) - \bar{\mathbf{D}}^{(0)} \right\|_{F}^{2} \stackrel{\text{(B.1)}}{\leq} \frac{4\left(1-\sqrt{1-\rho}\right)^{2T_{0}}}{\left(1-\tilde{\rho}\right)} \left\| \mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)} \right\|_{F}^{2}$$
(C.62)

where we have also used $N, L \ge 1$. By the choice of $T_0 = \left\lceil \frac{-2\ln(1-\tilde{\rho})}{\sqrt{1-\rho}} \right\rceil$, we have

$$T_{0} = \left\lceil \frac{-2\ln(1-\tilde{\rho})}{\sqrt{1-\rho}} \right\rceil \ge \frac{-2\ln(1-\tilde{\rho})^{\text{C.57}}}{\sqrt{1-\rho}} \ge \frac{2\ln(1-\tilde{\rho})}{\ln(1-\sqrt{1-\rho})}$$
(C.63)

where we have used Lemma (C.14) with $x = \sqrt{1-\rho}$ and $\ln(1-\tilde{\rho}) \le 0$. By (C.63), it holds that

$$\frac{4\left(1-\sqrt{1-\rho}\right)^{2T_0}}{(1-\tilde{\rho})} \le \frac{4\left(1-\sqrt{1-\rho}\right)^{\frac{4\ln(1-\tilde{\rho})}{\ln(1-\sqrt{1-\rho})}}}{(1-\tilde{\rho})} \le 4 \tag{C.64}$$

since $4\ln(1-\tilde{\rho}) \leq \ln(1-\tilde{\rho})$ as $\tilde{\rho} \in [0,1)$. Thus, by $\mathbf{x}_i^{(0)} = \mathbf{x}_j^{(0)}$ for all $i,j=1,\ldots,N$, we have that Δ is independent of $\tilde{\rho},N$, and K. Hence, for τ chosen uniformly at random from $\{0,\ldots,K-1\}$, we have

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(\tau)}, \nabla f(\mathbf{z}_{i}^{(\tau)}), \alpha_{\tau}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(\tau)} \right\|_{F}^{2} \leq \varepsilon,$$

provided

$$K = \mathcal{O}\left(\max\left\{\frac{\Delta^{\frac{3}{2}}}{(\alpha\varepsilon)^{\frac{3}{2}}}, \frac{L^{3}\hat{\sigma}^{3}\alpha^{3}}{N^{3}\varepsilon^{\frac{3}{2}}}, \frac{L^{3}\alpha^{3}\hat{\sigma}^{2}}{N^{2}(1-\tilde{\rho})^{2}\varepsilon}, \frac{L^{3}\alpha^{3}\hat{\sigma}^{\frac{3}{2}}}{N^{\frac{3}{2}}\varepsilon^{\frac{3}{4}}}\right\}\right).$$
(C.65)

Plugging $\alpha = \frac{N^{\frac{2}{3}}}{64L}$ into (C.65) results in

$$K = \mathcal{O}\left(\max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^2}{(1-\tilde{\rho})^2\varepsilon}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right).$$

Hence, the number of gradient evaluations is

$$\mathcal{K} \triangleq m(K-1) = \mathcal{O}\left(\max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^2}{(1-\tilde{\rho})^2\varepsilon}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right),$$

which yields a total number of gradient evaluations $\lceil \mathcal{K} + \sqrt[3]{NK} \rceil$, provided $K \geq \frac{N^2}{(1-\hat{\rho})^6}$. Since $\mathcal{K} = \Omega\left(\sqrt[3]{NK}\right)$, we drop the $\sqrt[3]{NK}$ and obtain (C.60).

Remark C.2 Similar to other works (Lian et al. 2017; Xin, Khan, and Kar 2021b,a), we have a minimum requirement on the number of iterations, called transient iterations, in order to achieve the complexity results in (C.60). Further, we notice that if the connectivity of the original network is poor, i.e. $\rho \approx 1$ for ρ defined in (7), and we only perform one neighbor communication during lines 3 and 6 in Algorithm 1 so that $\tilde{\rho} = \rho$, then it could be that $\frac{\hat{\sigma}^2}{(1-\rho)^2\varepsilon}$ dominates in (C.60), meaning DEEPSTORM is network-dependent. Similar to (Xin, Khan, and Kar 2021a), we can place a requirement that $\varepsilon \leq N^{-2}(1-\rho)^4$, in which case DEEPSTORM achieves the optimal complexity result and is independent of the communication network. In order to relax this requirement to $\varepsilon \leq N^{-2}$ (which can be significantly greater than $N^{-2}(1-\rho)^4$), we perform Algorithm B.1 during the neighbor communications (lines 3 and 6 in Algorithm 1) such that for $T = \lceil \frac{2}{\sqrt{1-\rho}} \rceil$, $\tilde{\rho} \leq \frac{1}{2}$ by (B.3), so that the number of local gradient computations becomes

$$\mathcal{O}\left(\max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right) \tag{C.66}$$

which is independent of ρ . Additionally, the number of local neighbor communications becomes

$$T_0 + TK = \mathcal{O}\left(\frac{1}{\sqrt{1-\rho}} \max\left\{\frac{(L\Delta)^{\frac{3}{2}} + \hat{\sigma}^3}{N\varepsilon^{\frac{3}{2}}}, \frac{\sqrt{N}\hat{\sigma}^{\frac{3}{2}}}{\varepsilon^{\frac{3}{4}}}\right\}\right)$$
(C.67)

which is optimal in terms of the dependence upon ρ (Scaman et al. 2017).

C.3 Diminishing step size

The Lyapunov function and relation defined in (C.13) are specially designed for the constant step size proof. Here, we make analogous designs for the diminishing step size proof.

Lemma C.15 For all $k \ge 0$, the following inequality holds

$$\mathbb{E}\left[\hat{\Phi}^{(k+1)} - \hat{\Phi}^{(k)}\right] \leq \underbrace{\left(4\gamma_{3}^{(k)}L^{2}(1-\beta_{k})^{2} - \frac{1}{2N}\left(\frac{1}{2\alpha_{k}} - 3L\right)\right)}_{(A)} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_{F}^{2} + \underbrace{\left(\frac{16\gamma_{1}^{(k)}L^{2}}{1-\tilde{\rho}} - \frac{1}{2N\alpha_{k}}\right)}_{(B)} \mathbb{E}\left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_{F}^{2} + \underbrace{\left(\frac{L}{2N} + \frac{\tilde{\rho}^{2}}{2N\alpha_{k}} + \frac{64\gamma_{1}^{(k)}L^{2}}{1-\tilde{\rho}} + \gamma_{2}^{(k)}\tilde{\rho} + 4\gamma_{3}^{(k)}L^{2}(1-\beta_{k})^{2} - \gamma_{2}^{(k-1)}\right)}_{(C)} \mathbb{E}\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} + \underbrace{\left(\frac{1}{2NL} + \gamma_{1}^{(k)}\tilde{\rho} + \gamma_{2}^{(k)}\frac{\alpha_{k}^{2}}{1-\tilde{\rho}} - \gamma_{1}^{(k-1)}\right)}_{(D)} \mathbb{E}\left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2} + \underbrace{\left(\frac{\alpha_{k}}{N} + \frac{4\gamma_{1}^{(k)}\beta_{k}^{2}}{1-\tilde{\rho}} + \gamma_{3}^{(k)}(1-\beta_{k})^{2} - \gamma_{3}^{(k-1)}\right)}_{(E)} \mathbb{E}\left\|\mathbf{R}^{(k)}\right\|_{F}^{2} + \underbrace{\left(\gamma_{3}^{(k)} + \frac{2\gamma_{1}^{(k)}}{1-\tilde{\rho}}\right)2\hat{\sigma}^{2}N\beta_{k}^{2}}_{2}}_{2}\right\}$$

where $\gamma_1^{(k)}, \gamma_2^{(k)}, \gamma_3^{(k)}$ are strictly positive values and

$$\hat{\Phi}^{(k)} \triangleq \phi(\bar{\mathbf{x}}^{(k)}) + \gamma_1^{(k-1)} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{E}^{2} + \gamma_2^{(k-1)} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{E}^{2} + \gamma_3^{(k-1)} \left\| \mathbf{R}^{(k)} \right\|_{E}^{2}$$
(C.69)

is a lower bounded Lyapunov function.

Proof We start by using part (iv) of Assumption 2 and (8) to note that

$$\frac{1}{2N\alpha_{k}} \left\| \bar{\mathbf{X}}^{(k)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} = \frac{1}{2N\alpha_{k}} \left\| \left(\boldsymbol{\mathcal{W}}_{T} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \left(\mathbf{I} - \frac{1}{N} \mathbf{e} \mathbf{e}^{\top} \right) \mathbf{X}^{(k)} \right\|_{F}^{2} \\
\stackrel{(14)}{\leq} \frac{\tilde{\rho}^{2}}{2N\alpha_{k}} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2}.$$
(C.70)

Next, we utilize the Peter-Paul inequality and Jensen's inequality to have

$$-\left\langle \bar{\mathbf{r}}^{(k)}, \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\rangle \leq \alpha_{k} \left\| \bar{\mathbf{r}}^{(k)} \right\|_{2}^{2} + \frac{1}{4\alpha_{k}} \left\| \bar{\mathbf{x}}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2}$$

$$\leq \frac{\alpha_{k}}{N} \sum_{i=1}^{N} \left\| \mathbf{r}_{i}^{(k)} \right\|_{2}^{2} + \frac{1}{4N\alpha_{k}} \sum_{i=1}^{N} \left\| \mathbf{x}_{i}^{(k+1)} - \bar{\mathbf{x}}^{(k)} \right\|_{2}^{2}$$

$$= \frac{\alpha_{k}}{N} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} + \frac{1}{4N\alpha_{k}} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2}. \tag{C.71}$$

Applying (C.70) and (C.71) to (C.16) results in

$$\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)})
\leq -\frac{1}{2N} \left(\frac{1}{\alpha_k} - 3L - \frac{1}{2\alpha_k} \right) \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_F^2 - \frac{1}{2N\alpha_k} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2
+ \left(\frac{L}{2N} + \frac{\tilde{\rho}^2}{2N\alpha_k} \right) \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \frac{1}{2NL} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_F^2 + \frac{\alpha_k}{N} \left\| \mathbf{R}^{(k)} \right\|_F^2.$$
(C.72)

Next, from Lemma (C.9), we use Young's inequality to have

$$\mathbb{E} \left\| \mathbf{R}^{(k+1)} \right\|_{F}^{2} \\
\stackrel{\text{(C.24)}}{\leq} 2N\beta_{k}^{2} \hat{\sigma}^{2} + 2(1-\beta_{k})^{2} L^{2} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_{F}^{2} + (1-\beta_{k})^{2} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \\
\leq 2N\beta_{k}^{2} \hat{\sigma}^{2} + 4(1-\beta_{k})^{2} L^{2} \left(\mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \right) + (1-\beta_{k})^{2} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2}. \tag{C.73}$$

Adding $\gamma_1^{(k)} \left\| \mathbf{Y}_{\perp}^{(k+1)} \right\|_F^2$, $\gamma_2^{(k)} \left\| \mathbf{X}_{\perp}^{(k+1)} \right\|_F^2$, $\gamma_3^{(k)} \left\| \mathbf{R}^{(k+1)} \right\|_F^2$ to both sides of (C.72) and taking the full expectation in conjunction with the results from Lemmas C.7 and (C.73) yields

$$\begin{split} & \mathbb{E}\left[\phi(\bar{\mathbf{x}}^{(k+1)}) - \phi(\bar{\mathbf{x}}^{(k)})\right] + \gamma_{1}^{(k)} \mathbb{E} \left\|\mathbf{Y}_{\perp}^{(k+1)}\right\|_{F}^{2} + \gamma_{2}^{(k)} \mathbb{E} \left\|\mathbf{X}_{\perp}^{(k+1)}\right\|_{F}^{2} + \gamma_{3}^{(k)} \mathbb{E} \left\|\mathbf{R}^{(k+1)}\right\|_{F}^{2} \\ & \leq \left(4\gamma_{3}^{(k)} L^{2} (1 - \beta_{k})^{2} - \frac{1}{2N} \left(\frac{1}{2\alpha_{k}} - 3L\right)\right) \mathbb{E} \left\|\mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)}\right\|_{F}^{2} \\ & + \frac{8\gamma_{1}^{(k)} L^{2}}{1 - \tilde{\rho}} \mathbb{E} \left\|\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}\right\|_{F}^{2} - \frac{1}{2N\alpha_{k}} \mathbb{E} \left\|\mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)}\right\|_{F}^{2} \\ & + \left(\frac{L}{2N} + \frac{\tilde{\rho}^{2}}{2N\alpha_{k}} + \gamma_{2}^{(k)} \tilde{\rho} + 4\gamma_{3}^{(k)} L^{2} (1 - \beta_{k})^{2}\right) \mathbb{E} \left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} \\ & + \left(\frac{1}{2NL} + \gamma_{1}^{(k)} \tilde{\rho} + \gamma_{2}^{(k)} \frac{\alpha_{k}^{2}}{1 - \tilde{\rho}}\right) \mathbb{E} \left\|\mathbf{Y}_{\perp}^{(k)}\right\|_{F}^{2} \\ & + \left(\frac{\alpha_{k}}{N} + \frac{4\gamma_{1}^{(k)} \beta_{k}^{2}}{1 - \tilde{\rho}} + \gamma_{3}^{(k)} (1 - \beta_{k})^{2}\right) \mathbb{E} \left\|\mathbf{R}^{(k)}\right\|_{F}^{2} \\ & + \left(\gamma_{3}^{(k)} + \frac{2\gamma_{1}^{(k)}}{1 - \tilde{\rho}}\right) 2\hat{\sigma}^{2} N \beta_{k}^{2}. \end{split}$$

Next we apply the following bound to the above relation,

$$\begin{aligned} \left\| \mathbf{X}^{(k+1)} - \mathbf{X}^{(k)} \right\|_F^2 &= \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} + \mathbf{Z}^{(k)} - \mathbf{X}^{(k)} \right\|_F^2 \\ &\leq 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 2 \left\| \mathbf{W}_T(\mathbf{X}^{(k)}) - \mathbf{W}_T \bar{\mathbf{X}}^{(k)} + \mathbf{W}_T \bar{\mathbf{X}}^{(k)} - \mathbf{X}^{(k)} \right\|_F^2 \\ &= 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 2 \left\| (\mathbf{I} - \mathbf{W}_T) \mathbf{X}_{\perp}^{(k)} \right\|_F^2 \\ &\leq 2 \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 + 8 \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2, \end{aligned}$$

where we have used Assumption 2 part (iv) to bound $\|\mathbf{I} - \mathbf{W}_T\|_2 \leq 2$. Finally, we subtract $\gamma_1^{(k-1)} \mathbb{E} \|\mathbf{Y}_{\perp}^{(k)}\|_F^2$, $\gamma_2^{(k-1)} \mathbb{E} \|\mathbf{X}_{\perp}^{(k)}\|_F^2$, and $\gamma_3^{(k-1)} \mathbb{E} \|\mathbf{R}^{(k)}\|_F^2$ from both sides to complete the proof. The lower boundedness of (C.69) follows from the non-negativity of the Frobenius norm and Assumption 1 (iv).

Lemma C.16 Let $x \in [0,1)$. Then for any $y \ge \lceil \frac{2}{1-x^3} \rceil$, it holds that

$$\left(\frac{y-1}{y}\right)^{\frac{1}{3}} - x > \frac{1-x}{2}.$$
 (C.74)

Proof The proof begins by analyzing $h(x) \triangleq x^3 - x^2 - x + 1$ for $x \in [0,1)$. Notice that h(x) is decreasing on [0,1) by $h'(x) = 3x^2 - 2x - 1 < 0$. Since h(0) = 1 and h(1-) = 0, it holds that h(x) > 0 for $x \in [0,1)$. Hence

$$3 + 3x^3 > 3x + 3x^2.$$

Adding $1 + x^3$ to both sides results in

$$4 + 4x^3 > (1+x)^3$$
.

Dividing by 4, adding 1 to both sides, and rearranging results in

$$2 - \frac{1}{4}(1+x)^3 > 1 - x^3.$$

Since x < 1, we divide both sides by $(1 - x^3)$ and use $\lceil a \rceil \ge a$ for any $a \in \mathbb{R}$ to have

$$y\left(1-\left(\frac{1+x}{2}\right)^3\right) \ge \left\lceil\frac{2}{1-x^3}\right\rceil \left(1-\left(\frac{1+x}{2}\right)^3\right) \ge \left(\frac{2}{1-x^3}\right) \left(1-\left(\frac{1+x}{2}\right)^3\right) > 1$$

where we have used $y \ge \lceil \frac{2}{1-x^3} \rceil$. Rearranging

$$y\left(1-\left(\frac{1+x}{2}\right)^3\right) > 1$$

results in

$$\frac{y-1}{y} > \left(\frac{1+x}{2}\right)^3.$$

Taking the cube-root and subtracting x from both sides completes the proof.

Proof of Theorem 2 Proof The proof follows similar steps as the proof for Theorem 1. We frequently use (20) to have $\alpha_k \leq \min\{\frac{1}{32L}, \frac{(1-\tilde{\rho})^2}{64L}\}$. First, we show $\beta_k \in (0,1)$ for all $k \geq 0$. By $\alpha \leq \frac{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}}{64L} < \frac{k_0^{\frac{1}{3}}}{\sqrt{480}L}$ it holds

$$L^{2}\alpha_{k+1}^{2} < \frac{1}{480} = \frac{1}{10 \cdot 48} < \frac{1}{48} \left(\frac{1}{4} + \left(\frac{2}{3} \right)^{\frac{1}{3}} - 1 \right)$$

$$\leq \frac{1}{48} \left(\frac{1}{4} + \left(\frac{k_{0}}{k_{0} + 1} \right)^{\frac{1}{3}} - 1 \right)$$

$$< \frac{1}{48} \left(\frac{(k_{0} + 1)^{\frac{1}{3}}}{2k_{0}^{\frac{1}{3}} + (k_{0} + 1)^{\frac{1}{3}}} + \frac{\alpha_{k+1}}{\alpha_{k}} - 1 \right)$$
(C.75)

where the first inequality uses $k_0 \ge \lceil \frac{2}{1-\tilde{\rho}^3} \rceil \ge 2$ and the last uses $2k_0^{\frac{1}{3}} < 3(k_0+1)^{\frac{1}{3}}$ for all $k_0 \ge 2$. Rearranging (C.75) results in

$$\beta_k < \frac{(k_0 + 1)^{\frac{1}{3}}}{2k_3^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} < 1 \tag{C.76}$$

where the right most inequality uses $2k_0^{\frac{1}{3}} > 0$ by $k_0 \ge 2$. Since $1 - \frac{\alpha_{k+1}}{\alpha_k} > 0$ by $\alpha_k > \alpha_{k+1}$, we also have $0 < \beta_k$. Second, let

$$\gamma_1^{(k)} \triangleq \frac{1}{NL(1-\tilde{\rho})}, \gamma_2^{(k)} \triangleq \frac{16}{N(1-\tilde{\rho})\alpha_k}, \text{ and } \gamma_3^{(k)} \triangleq \frac{1}{24NL^2\alpha_{k+1}} \tag{C.77}$$

in (C.68) to have

$$\underbrace{\left(\frac{1}{4N\alpha_{k}} - \frac{(1-\beta_{k})^{2}}{6N\alpha_{k+1}} - \frac{3L}{2N}\right)}_{(A')} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{1}{2N\alpha_{k}} - \frac{16L}{N(1-\tilde{\rho})^{2}}\right)}_{(B')} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{16}{N(1-\tilde{\rho})\alpha_{k-1}} - \frac{16\tilde{\rho}}{N(1-\tilde{\rho})\alpha_{k}} - \frac{L}{2N} - \frac{\tilde{\rho}^{2}}{2N\alpha_{k}} - \frac{64L}{N(1-\tilde{\rho})^{2}} - \frac{(1-\beta_{k})^{2}}{6N\alpha_{k+1}}\right)}_{(C')} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{1}{NL} - \frac{1}{2NL} - \frac{16\alpha_{k}}{N(1-\tilde{\rho})^{2}}\right)}_{(D')} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} \\
+ \underbrace{\left(\frac{1}{24NL^{2}\alpha_{k}} - \frac{(1-\beta_{k})^{2}}{24NL^{2}\alpha_{k+1}} - \frac{4\beta_{k}^{2}}{NL(1-\tilde{\rho})^{2}} - \frac{\alpha_{k}}{N}\right)}_{(E')} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \\
\leq \mathbb{E} \left[\hat{\Phi}^{(k)} - \hat{\Phi}^{(k+1)}\right] + \underbrace{\left(\frac{1}{24NL^{2}\alpha_{k+1}} + \frac{2}{NL(1-\tilde{\rho})^{2}}\right)}_{(C.78)} 2N\beta_{k}^{2}\hat{\sigma}^{2}. \tag{C.78}$$

Next we lower bound (A') - (E'). For (A'), since $(1 - \beta_k)^2 < 1$ we have

$$(A') = \frac{1}{4N\alpha_k} - \frac{(1-\beta_k)^2}{6N\alpha_{k+1}} - \frac{3L}{2N} > \frac{1}{N\alpha_k} \left(\frac{1}{4} - \frac{1}{6} \left(\frac{k+k_0+1}{k+k_0}\right)^{\frac{1}{3}} - \frac{3L\alpha}{2k_0^{\frac{1}{3}}}\right)$$
$$\geq \frac{1}{N\alpha_k} \left(\frac{1}{4} - \frac{1}{6} \left(\frac{3}{2}\right)^{\frac{1}{3}} - \frac{3}{64}\right)$$
$$> \frac{1}{96N\alpha_k}$$

where we have used $k_0 \ge 2$ and $\alpha \le \frac{k_0^{\frac{1}{3}}}{32L}$. For (B'), we use $\alpha \le \frac{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}}{64L}$ so $\alpha_k \le \frac{(1-\tilde{\rho})^2}{64L}$ to have

$$(B') = \frac{1}{2N\alpha_k} - \frac{16L}{N(1-\tilde{\rho})^2} \ge \frac{1}{N\alpha_k} \left(\frac{1}{2} - \frac{1}{4}\right) = \frac{1}{4N\alpha_k}.$$

For (C'), we have

$$\begin{split} (C') = & \frac{16}{N(1-\tilde{\rho})\alpha_{k-1}} - \frac{16\tilde{\rho}}{N(1-\tilde{\rho})\alpha_k} - \frac{L}{2N} - \frac{\tilde{\rho}^2}{2N\alpha_k} - \frac{64L}{N(1-\tilde{\rho})^2} - \frac{(1-\beta_k)^2}{6N\alpha_{k+1}} \\ = & \frac{1}{N\alpha_k} \left(\frac{16}{(1-\tilde{\rho})} \left(\frac{\alpha_k}{\alpha_{k-1}} - \tilde{\rho} \right) - \frac{L\alpha_k}{2} - \frac{\tilde{\rho}^2}{2} - \frac{64L\alpha_k}{(1-\tilde{\rho})^2} - \frac{(1-\beta_k)^2\alpha_k}{6\alpha_{k+1}} \right) \\ \geq & \frac{1}{N\alpha_k} \left(\frac{16}{(1-\tilde{\rho})} \left(\left(\frac{k_0-1}{k_0} \right)^{\frac{1}{3}} - \tilde{\rho} \right) - \frac{L\alpha_k}{2} - \frac{\tilde{\rho}^2}{2} - \frac{64L\alpha_k}{(1-\tilde{\rho})^2} - \frac{(1-\beta_k)^2\alpha_k}{6\alpha_{k+1}} \right). \end{split}$$

Next, using Lemma C.16, since $\tilde{\rho} \in [0,1)$ and $k_0 \geq \lceil \frac{2}{1-\tilde{\rho}^3} \rceil$, it holds that

$$\left(\frac{k_0-1}{k_0}\right)^{\frac{1}{3}}-\tilde{\rho}\geq\frac{1-\tilde{\rho}}{2}.$$

Thus, using $(1 - \beta_k)^2 < 1$ and $\alpha_k \le \min\{\frac{1}{32L}, \frac{(1 - \tilde{\rho})^2}{64L}\}$, it holds

$$\begin{split} (C') > & \frac{1}{N\alpha_k} \left(8 - \frac{L\alpha_k}{2} - \frac{\tilde{\rho}^2}{2} - \frac{64L\alpha_k}{(1 - \tilde{\rho})^2} - \frac{(1 - \beta_k)^2 \alpha_k}{6\alpha_{k+1}} \right) \\ > & \frac{1}{N\alpha_k} \left(8 - \frac{1}{64} - \frac{1}{2} - 1 - \frac{1}{6} \left(\frac{k + k_0 + 1}{k + k_0} \right)^{\frac{1}{3}} \right) \\ > & \frac{1}{N\alpha_k} \left(8 - \frac{1}{64} - \frac{1}{2} - 1 - \frac{1}{6} \left(\frac{3}{2} \right)^{\frac{1}{3}} \right) \\ > & \frac{4}{N\alpha_k}. \end{split}$$

For (D'), we use $\alpha \leq \frac{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}}{64L}$ so $\alpha_k \leq \frac{(1-\tilde{\rho})^2}{64L}$ to have

$$(D') = \frac{1}{NL} - \frac{1}{2NL} - \frac{16\alpha_k}{N(1-\tilde{\rho})^2} \ge \frac{1}{2NL} - \frac{16\alpha}{N(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}} \ge \frac{1}{4NL}.$$

For (E'), we factor out $\frac{1}{24NL^2\alpha_{k+1}}$ and expand $(1-\beta_k)^2$ to have

$$\begin{split} (E') = & \frac{1}{24NL^2\alpha_{k+1}} \left(\frac{\alpha_{k+1}}{\alpha_k} - 1 + 2\beta_k - \beta_k^2 - 24L^2\alpha_{k+1}\alpha_k - \frac{96L\alpha_{k+1}\beta_k^2}{(1-\tilde{\rho})^2} \right) \\ \geq & \frac{1}{24NL^2\alpha_{k+1}} \left(\frac{\alpha_{k+1}}{\alpha_k} - 1 + 2\beta_k - \beta_k^2 - 48L^2\alpha_{k+1}^2 - \frac{96L\alpha_{k+1}\beta_k^2}{(1-\tilde{\rho})^2} \right), \end{split}$$

where we have used $\alpha_k \leq 2\alpha_{k+1}$. Plugging in the definition of β_k from (20) and using $\beta_k \geq 48L^2\alpha_{k+1}^2$ gives

$$\begin{split} (E') \geq & \frac{\beta_k}{24NL^2\alpha_{k+1}} \left(1 - \beta_k \left(1 + \frac{96L\alpha_{k+1}}{(1 - \tilde{\rho})^2} \right) \right) \\ \geq & \frac{\beta_k}{24NL^2\alpha_{k+1}} \left(1 - \beta_k \left(1 + \frac{96}{64} \left(\frac{k_0}{k_0 + 1} \right)^{\frac{1}{3}} \right) \right) \\ \stackrel{\text{(C.76)}}{\geq 24NL^2\alpha_{k+1}} \left(1 - \frac{(k_0 + 1)^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} \left(1 + \frac{3}{2} \left(\frac{k_0}{k_0 + 1} \right)^{\frac{1}{3}} \right) \right) \\ = & \frac{\beta_k}{24NL^2\alpha_{k+1}} \left(\frac{1}{2} \cdot \frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} \right) \\ \geq & \frac{2\alpha_{k+1}}{N} \left(\frac{1}{2} \cdot \frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} \right) \\ \geq & \frac{\alpha_k}{2N} \left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}} \right) \end{split}$$

where the second inequality uses $\alpha_{k+1} \leq \frac{\alpha}{(k_0+1)^{\frac{1}{3}}}$ and $\alpha \leq \frac{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}}{64L}$ and the last inequality uses $\alpha_k \leq 2\alpha_{k+1}$. Next, we sum (C.78) over k=0 to K-1; using the established lower bounds to have

$$\sum_{k=0}^{K-1} \frac{1}{96N\alpha_{k}} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \bar{\mathbf{X}}^{(k)} \right\|_{F}^{2} + \sum_{k=0}^{K-1} \frac{1}{4N\alpha_{k}} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_{F}^{2} + \sum_{k=0}^{K-1} \frac{4}{N\alpha_{k}} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_{F}^{2} + \sum_{k=0}^{K-1} \frac{1}{4NL} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_{F}^{2} + \sum_{k=0}^{K-1} \left(\frac{k_{0}^{\frac{1}{3}}}{2k_{0}^{\frac{1}{3}} + (k_{0} + 1)^{\frac{1}{3}}} \right) \frac{\alpha_{k}}{2N} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_{F}^{2} \\
\leq \left(\hat{\Phi}^{(0)} - \phi^{*} \right) + \sum_{k=0}^{K-1} \left(\frac{1}{12L^{2}\alpha_{k+1}} + \frac{4}{L(1 - \tilde{\rho})^{2}} \right) \beta_{k}^{2} \hat{\sigma}^{2}. \tag{C.79}$$

where we have used $\phi^* \leq \hat{\Phi}^{(k)}$ for any $k \geq 0$. The final phase of the proof uses Lemma C.12 to provide a concise convergence statement. To do so, we use Jensen's inequality to have

$$\left\|\bar{\mathbf{r}}^{(k)}\right\|_2^2 \leq \frac{1}{N} \sum_{i=1}^N \left\|\mathbf{r}_i^{(k)}\right\|_2^2 = \frac{1}{N} \left\|\mathbf{R}^{(k)}\right\|_F^2.$$

Applying this to (C.31), multiply both sides of (C.31) by $\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}}+(k_0+1)^{\frac{1}{3}}}\right)\alpha_k$, sum from $k=0,\ldots,K-1$, and take the expectation to have

$$\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E} \left\| P\left(\mathbf{z}_i^{(k)}, \nabla f(\mathbf{z}_i^{(k)}), \alpha_k\right) \right\|_2^2 \right)
+ \left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{L^2}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \frac{L^2}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_F^2 \right)
\leq \left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{6}{N} \mathbb{E} \left\| \mathbf{Y}_{\perp}^{(k)} \right\|_F^2 + \frac{6}{N} \mathbb{E} \left\| \mathbf{R}^{(k)} \right\|_F^2 \right)
+ \left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{32L^2}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \frac{2}{N\alpha_k^2} \mathbb{E} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 \right).$$
(C.80)

We relate each of the terms on the right-hand side of (C.80) to 12 times of the left-hand side of (C.79). Since $\alpha_k \leq \frac{1}{32L}$ and $\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \leq 1$, it holds for all $k \geq 0$,

$$\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \frac{6\alpha_k}{N} \left\|\mathbf{Y}_{\perp}^{(k)}\right\|_F^2 \le \frac{3}{NL} \left\|\mathbf{Y}_{\perp}^{(k)}\right\|_F^2.$$
(C.81)

Next, by $\alpha_k \leq \frac{1}{32L}$ and $\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}}+(k_0+1)^{\frac{1}{3}}}\right) \leq 1$, we have

$$\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \frac{32L^2\alpha_k^2}{N} \left\|\mathbf{X}_{\perp}^{(k)}\right\|_F^2 \le \frac{48}{N} \left\|\mathbf{X}_{\perp}^{(k)}\right\|_F^2.$$
(C.82)

Additionally, since $\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}}+(k_0+1)^{\frac{1}{3}}}\right)\leq 1$, it holds that

$$\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \frac{2}{N\alpha_k} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2 \le \frac{3}{N\alpha_k} \left\| \mathbf{X}^{(k+1)} - \mathbf{Z}^{(k)} \right\|_F^2.$$
(C.83)

Combining (C.81) - (C.83) in conjunction with 12 times of (C.79) and (C.80) results in

$$\left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{1}{N} \sum_{i=1}^N \mathbb{E} \left\| P\left(\mathbf{z}_i^{(k)}, \nabla f(\mathbf{z}_i^{(k)}), \alpha_k\right) \right\|_2^2\right) \\
+ \left(\frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0 + 1)^{\frac{1}{3}}}\right) \sum_{k=0}^{K-1} \alpha_k \left(\frac{L^2}{N} \mathbb{E} \left\| \mathbf{X}_{\perp}^{(k)} \right\|_F^2 + \frac{L^2}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(k)} \right\|_F^2\right) \\
\leq 12 \left(\hat{\Phi}^{(0)} - \phi^*\right) + \sum_{k=0}^{K-1} \left(\frac{1}{L^2 \alpha_{k+1}} + \frac{48}{L(1 - \tilde{\rho})^2}\right) \beta_k^2 \hat{\sigma}^2.$$

Using $\left\|\mathbf{X}_{\perp}^{(k)}\right\|_{F}^{2} \geq 0$ completes the proof.

Complexity analysis

Remark C.3 Similar to Remark C.1, the convergence of Algorithm 1 depends upon $\hat{\Phi}^{(0)}$ (see (C.69)), but in this setting we do not need a big initial batch, in terms of dependence upon ε . The initial variables must be equal for all agents and the number of initial communications to have $\mathbf{Y}^{(0)}$ must be sufficiently large, as in Corollary 1.

Corollary 2 Let $\varepsilon > 0$ be given and assume that $L \ge 1$. Under the same conditions as in Theorem 2, let $\mathbf{x}_i^{(0)} = \mathbf{x}_j^{(0)}$ for all $i, j = 1, \ldots, N$, let the local batch size $m = \mathcal{O}(1)$ for all iterations, choose $k_0 = \lceil \frac{2}{(1-\tilde{\rho})^6} \rceil$, and perform $T_0 = \lceil \frac{-2\ln(1-\tilde{\rho})}{\sqrt{1-\rho}} \rceil$ communications by Algorithm B.1 for the initial gradient tracking update in line 1 of Algorithm 1. Then choose α such that

$$\alpha = \frac{1}{64L}.\tag{C.84}$$

Then for all

$$K \ge 2^{\frac{3}{2}} k_0,\tag{C.85}$$

Algorithm 1 produces a stochastic ε -stationary point as defined in Definition 2 in

$$K = \mathcal{O}\left(\max\left\{\frac{(L\delta)^{\frac{3}{2}} + \hat{\sigma}^3 + k_0^{\frac{1}{2}}\sigma^3}{\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^3 \left(|\ln\varepsilon| + |\ln\hat{\sigma}|\right)^{\frac{3}{2}}}{\varepsilon^{\frac{3}{2}}}\right\}\right)$$
(C.86)

local stochastic gradient computations and $T_0 + T(K-1)$ neighbor communications for any $T \ge 1$.

Proof First, notice that $k_0 = \lceil \frac{2}{(1-\tilde{\rho})^6} \rceil \ge \lceil \frac{2}{(1-\tilde{\rho}^3)} \rceil$ by $(1-\tilde{\rho})^6 \le (1-\tilde{\rho})^3 \le 1-\tilde{\rho}^3$ for all $\tilde{\rho} \in [0,1)$. Hence k_0 satisfies the requirements of Theorem 2. Next, we have $c \triangleq \frac{k_0^{\frac{1}{3}}}{2k_0^{\frac{1}{3}} + (k_0+1)^{\frac{1}{3}}} > \frac{1}{4}$ for all $k_0 \ge 2$. Hence it holds that

$$c \sum_{k=0}^{K-1} \alpha_k = c\alpha \sum_{k=0}^{K-1} \frac{1}{(k+k_0)^{\frac{1}{3}}}$$

$$\geq c\alpha \int_0^K \frac{1}{(x+k_0)^{\frac{1}{3}}} dx$$

$$= \frac{3c\alpha}{2} \left((K+k_0)^{\frac{2}{3}} - k_0^{\frac{2}{3}} \right)$$

$$> \frac{3\alpha}{8} \left((K+k_0)^{\frac{2}{3}} - k_0^{\frac{2}{3}} \right)$$

Notice $(K + k_0)^{\frac{2}{3}} - k_0^{\frac{2}{3}} > \frac{K^{\frac{2}{3}}}{2}$ as long as K satisfies (C.85). Hence for some iterate $\tau \in \{0, 1, \dots, K-1\}$ chosen with probability

$$\operatorname{Prob}(\tau = k) = \frac{\frac{c\alpha}{(k+k_0)^{\frac{1}{3}}}}{\sum_{j=0}^{K-1} \frac{c\alpha}{(j+k_0)^{\frac{1}{3}}}}, \quad k = 0, \dots, K-1,$$
(C.87)

it holds that

$$\begin{split} & \frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(\tau)}, \nabla f(\mathbf{z}_{i}^{(\tau)}), \alpha_{\tau}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(\tau)} \right\|_{F}^{2} \\ \leq & \frac{64 \hat{\Delta}}{\alpha K^{\frac{2}{3}}} + \frac{16}{3\alpha K^{\frac{2}{3}}} \sum_{k=0}^{K-1} \left(\frac{1}{L^{2} \alpha_{k+1}} + \frac{48}{L(1-\tilde{\rho})^{2}} \right) \beta_{k}^{2} \hat{\sigma}^{2}. \end{split}$$

Expanding the summation on the right yields

$$\sum_{k=0}^{K-1} \left(\frac{1}{L^2 \alpha_{k+1}} + \frac{48}{L(1-\tilde{\rho})^2} \right) \beta_k^2 \hat{\sigma}^2$$

$$= \hat{\sigma}^2 \sum_{k=0}^{K-1} \left(\frac{1}{L^2 \alpha_{k+1}} + \frac{48}{L(1-\tilde{\rho})^2} \right) \left(1 - \frac{\alpha_{k+1}}{\alpha_k} + 48L^2 \alpha_{k+1}^2 \right)^2$$

$$\leq 2\hat{\sigma}^2 \sum_{k=0}^{K-1} \left(\frac{1}{L^2 \alpha_{k+1}} + \frac{48}{L(1-\tilde{\rho})^2} \right) \left(\left(1 - \frac{\alpha_{k+1}}{\alpha_k} \right)^2 + \left(48L^2 \alpha_{k+1}^2 \right)^2 \right) \tag{C.88}$$

where the last inequality uses $(a+b)^2 \le 2a^2 + 2b^2$ for any $a,b \in \mathbb{R}$. Utilizing $a^3 - b^3 = (a-b)(a^2 + ab + b^2)$ for any $a,b \in \mathbb{R}$, it holds that

$$1 - \frac{\alpha_{k+1}}{\alpha_k} = 1 - \frac{(k+k_0)^{\frac{1}{3}}}{(k+k_0+1)^{\frac{1}{3}}} = \frac{(k+k_0+1)^{-\frac{1}{3}}}{(k+k_0+1)^{\frac{2}{3}} + (k+k_0+1)^{\frac{1}{3}}(k+k_0)^{\frac{1}{3}} + (k+k_0)^{\frac{2}{3}}}.$$
 (C.89)

Notice that for all $k \ge 1$ and $k_0 \ge 2$, it holds that

$$2(k+k_0+1)^{\frac{2}{3}} \le (k+k_0+1)^{\frac{2}{3}} + (k+k_0+1)^{\frac{1}{3}}(k+k_0)^{\frac{1}{3}} + (k+k_0)^{\frac{2}{3}}.$$

Squaring both sides of the above inequality and rearranging results in

$$\frac{1}{\left((k+k_0+1)^{\frac{2}{3}}+(k+k_0+1)^{\frac{1}{3}}(k+k_0)^{\frac{1}{3}}+(k+k_0)^{\frac{2}{3}}\right)^2} \le \frac{1}{4(k+k_0+1)^{\frac{4}{3}}}$$
(C.90)

Utilizing (C.89) and multiplying both sides of (C.90) by $(k + k_0 + 1)^{-\frac{1}{3}}$, we further bound

$$2\hat{\sigma}^{2} \sum_{k=0}^{K-1} \frac{1}{L^{2}\alpha_{k+1}} \left(1 - \frac{\alpha_{k+1}}{\alpha_{k}}\right)^{2}$$

$$= \frac{2\hat{\sigma}^{2}}{L^{2}\alpha} \sum_{k=0}^{K-1} \frac{(k+k_{0}+1)^{-\frac{1}{3}}}{\left((k+k_{0}+1)^{\frac{2}{3}} + (k+k_{0}+1)^{\frac{1}{3}}(k+k_{0})^{\frac{1}{3}} + (k+k_{0})^{\frac{2}{3}}\right)^{2}}$$

$$\leq \frac{\hat{\sigma}^{2}}{2L^{2}\alpha} \sum_{k=0}^{K-1} \frac{1}{(k+k_{0}+1)^{\frac{5}{3}}}$$

$$\leq \frac{\hat{\sigma}^{2}}{2L^{2}\alpha} \int_{-1}^{K-1} \frac{1}{(x+k_{0}+1)^{\frac{5}{3}}} dx$$

$$\leq \frac{\hat{\sigma}^{2}}{2L^{2}\alpha} \cdot \frac{3}{2k_{0}^{\frac{2}{3}}}$$

$$= \frac{3\hat{\sigma}^{2}}{4L^{2}\alpha k_{0}^{\frac{2}{3}}}$$
(C.91)

Again, utilizing (C.89) and multiplying both sides of (C.90) by $(k + k_0 + 1)^{-\frac{2}{3}}$, we have

$$2\hat{\sigma}^2 \sum_{k=0}^{K-1} \frac{48}{L(1-\tilde{\rho})^2} \left(1 - \frac{\alpha_{k+1}}{\alpha_k}\right)^2 \le \frac{24\hat{\sigma}^2}{L(1-\tilde{\rho})^2} \sum_{k=0}^{K-1} \frac{1}{(k+k_0+1)^2}$$

$$\le \frac{24\hat{\sigma}^2}{L(1-\tilde{\rho})^2 k_0}, \tag{C.92}$$

where we have also upper bounded the summation by the corresponding integral, since both $(x+k_0+1)^{-2}$ and $(x+k_0+1)^{-\frac{5}{3}}$ are decreasing for all x>0. Next, we bound

$$2\hat{\sigma}^{2} \sum_{k=0}^{K-1} \frac{1}{L^{2}\alpha_{k+1}} \left(48L^{2}\alpha_{k+1}^{2}\right)^{2} = 4608\hat{\sigma}^{2}L^{2}\alpha^{3} \sum_{k=0}^{K-1} \frac{1}{k+k_{0}+1}$$

$$\leq 4608\hat{\sigma}^{2}L^{2}\alpha^{3} \left(\ln(K+k_{0}) - \ln(k_{0})\right)$$

$$\leq 4608\hat{\sigma}^{2}L^{2}\alpha^{3} \ln(K+k_{0}), \tag{C.93}$$

and for $b = 2 \cdot 48^3$ we have.

$$2\hat{\sigma}^{2} \sum_{k=0}^{K-1} \frac{48}{L(1-\tilde{\rho})^{2}} \left(48L^{2}\alpha_{k+1}^{2}\right)^{2} = \frac{b\hat{\sigma}^{2}L^{3}\alpha^{4}}{(1-\tilde{\rho})^{2}} \sum_{k=0}^{K-1} \frac{1}{(k+k_{0}+1)^{\frac{4}{3}}}$$

$$\leq \frac{3b\hat{\sigma}^{2}L^{3}\alpha^{4}}{(1-\tilde{\rho})^{2}k_{0}^{\frac{1}{3}}}.$$
(C.94)

Plugging (C.91) - (C.94) into (C.88) results in an inequality of the form

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(\tau)}, \nabla f(\mathbf{z}_{i}^{(\tau)}), \alpha_{\tau}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(\tau)} \right\|_{F}^{2} \\
\leq \frac{64\hat{\Delta}}{\alpha K^{\frac{2}{3}}} + \frac{16}{3\alpha K^{\frac{2}{3}}} \left(\frac{3\hat{\sigma}^{2}}{4L^{2}\alpha k_{0}^{\frac{2}{3}}} + \frac{24\hat{\sigma}^{2}}{L(1-\tilde{\rho})^{2}k_{0}} \right) \\
+ \frac{16}{3\alpha K^{\frac{2}{3}}} \left(4608\hat{\sigma}^{2}L^{2}\alpha^{3} \ln(K+k_{0}) + \frac{3b\hat{\sigma}^{2}L^{3}\alpha^{4}}{(1-\tilde{\rho})^{2}k_{0}^{\frac{1}{3}}} \right). \tag{C.95}$$

Next, by (C.69) and (C.77), we have

$$\hat{\Phi}^{(0)} \triangleq \phi(\bar{\mathbf{x}}^{(0)}) + \frac{1}{NL(1-\tilde{\rho})} \left\| \mathbf{Y}_{\perp}^{(0)} \right\|_{F}^{2} + \frac{16k_{0}^{\frac{1}{3}}}{N(1-\tilde{\rho})\alpha} \left\| \mathbf{X}_{\perp}^{(0)} \right\|_{F}^{2} + \frac{k_{0}^{\frac{1}{3}}}{24NL^{2}\alpha} \left\| \mathbf{R}^{(0)} \right\|_{F}^{2}, \tag{C.96}$$

where we have defined $\alpha_{-1} \triangleq \alpha_0$ for the $\gamma_2^{(-1)}$ term in (C.69). Similar to the proof of Corollary 1, we bound each of the terms on the right-hand side of (C.96) by the initialization from Algorithm 1. We have

$$\frac{k_0^{\frac{1}{3}}}{24NL^2\alpha} \left\| \mathbf{R}^{(0)} \right\|_F^2 \le \frac{k_0^{\frac{1}{3}}\sigma^2}{24L^2\alpha m_0} \tag{C.97}$$

by (5) since $\mathbf{d}_{i}^{(0)} = \frac{1}{m_0} \sum_{\xi \in B_{i}^{(0)}} \nabla f_{i}(\mathbf{x}_{i}^{(0)}, \xi)$ for all $i = 1, \dots, N$. Notice that $\mathbf{Y}^{(0)} = \mathcal{W}_{T_0}(\mathbf{D}^{(0)})$. Hence (C.62) still holds, so by $T_0 = \begin{bmatrix} \frac{-2 \ln(1-\bar{\rho})}{\sqrt{1-\bar{\rho}}} \end{bmatrix}$, we have (C.64). Thus, by $\mathbf{x}_{i}^{(0)} = \mathbf{x}_{j}^{(0)}$ for all $i, j = 1, \dots, N$, we have that $\hat{\Delta} \leq \phi(\bar{\mathbf{x}}^{(0)}) + 4 \|\mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)}\|_F^2 + \frac{k_0^{1/2}\sigma^2}{24L^2\alpha m_0} - \phi^*$. By $k_0 = \lceil \frac{2}{(1-\bar{\rho})^6} \rceil$, the α in (C.84) satisfies $\alpha \leq \frac{(1-\bar{\rho})^2 k_0^{1/2}}{64L}$. Hence,

$$\frac{\hat{\Delta}}{\alpha} \le 256L \left\| \mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)} \right\|_F^2 + \frac{172k_0^{\frac{1}{3}}\sigma^2}{m_0} + 64L \left(\phi(\bar{\mathbf{x}}^{(0)}) - \phi^* \right). \tag{C.98}$$

Additionally we further bound the two terms on the right-hand side of (C.95) that contain $(1 - \tilde{\rho})^{-2}$. By the choice of k_0 , we have $k_0 \ge \frac{1}{(1-\tilde{\rho})^6}$, thus it holds that

$$\frac{1}{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}} \leq 1 \text{ and } \frac{1}{(1-\tilde{\rho})^2 k_0} \leq (1-\tilde{\rho})^4 \leq 1.$$

Hence, by recalling $\alpha = \frac{1}{64L}$ and $b = 2 \cdot 48^3$, we have

$$\frac{24\hat{\sigma}^2}{\alpha L(1-\tilde{\rho})^2 k_0} \le 1536\hat{\sigma}^2,\tag{C.99}$$

$$\frac{3b\hat{\sigma}^2 L^3 \alpha^3}{(1-\tilde{\rho})^2 k_0^{\frac{1}{3}}} \le 3\hat{\sigma}^2. \tag{C.100}$$

Further, it holds that $\frac{3\hat{\sigma}^2}{4L^2\alpha^2k_0^{\frac{2}{3}}} \leq 3072\hat{\sigma}^2$; using this and plugging (C.98), (C.99), and (C.100) into (C.95) yields

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(\tau)}, \nabla f(\mathbf{z}_{i}^{(\tau)}), \alpha_{\tau}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(\tau)} \right\|_{F}^{2} \\
\leq \frac{64}{K^{\frac{2}{3}}} \left(256L \left\| \mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)} \right\|_{F}^{2} + \frac{172k_{0}^{\frac{1}{3}}\sigma^{2}}{m_{0}} + 64L \left(\phi(\bar{\mathbf{x}}^{(0)}) - \phi^{*} \right) \right) \\
+ \frac{16}{3K^{\frac{2}{3}}} \left(3072\hat{\sigma}^{2} + 3\hat{\sigma}^{2} + 1536\hat{\sigma}^{2} \right) + \frac{6\hat{\sigma}^{2} \ln\left(K + k_{0}\right)}{K^{\frac{2}{3}}} \\
\leq \frac{4096L \left(4 \left\| \mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)} \right\|_{F}^{2} + \phi(\bar{\mathbf{x}}^{(0)}) - \phi^{*} \right)}{K^{\frac{2}{3}}} \\
+ \frac{\hat{\sigma}^{2}}{K^{\frac{2}{3}}} \left(24592 + 32\ln\left(K + k_{0}\right) + \frac{11008k_{0}^{\frac{1}{3}}\sigma^{2}}{m_{0}\hat{\sigma}^{2}} \right). \tag{C.101}$$

Finally, for τ chosen according to (C.87), we have

$$\frac{1}{N} \sum_{i=1}^{N} \mathbb{E} \left\| P\left(\mathbf{z}_{i}^{(\tau)}, \nabla f(\mathbf{z}_{i}^{(\tau)}), \alpha_{\tau}\right) \right\|_{2}^{2} + \frac{L^{2}}{N} \mathbb{E} \left\| \mathbf{Z}_{\perp}^{(\tau)} \right\|_{F}^{2} \leq \varepsilon,$$

provided

$$K = \mathcal{O}\left(\max\left\{\frac{(L\delta)^{\frac{3}{2}} + \hat{\sigma}^3}{\varepsilon^{\frac{3}{2}}}, \frac{k_0^{\frac{1}{2}}\sigma^3}{m_0^{\frac{3}{2}}\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^3\left(\left|\ln\varepsilon\right| + \left|\ln\hat{\sigma}\right|\right)^{\frac{3}{2}}}{\varepsilon^{\frac{3}{2}}}\right\}\right)$$
(C.102)

where $\delta \triangleq \|\mathbf{D}^{(0)} - \bar{\mathbf{D}}^{(0)}\|_F^2 + \phi(\bar{\mathbf{x}}^{(0)}) - \phi^*$. Choosing the initial batch size $m_0 = m = \mathcal{O}(1)$ yields the total number of gradient computations mK in (C.86), provided K satisfies (C.85).

Remark C.4 Similar to the discussion provided in Remark C.2, we note that Chebyshev acceleration can be utilized to perform the neighbor communications (lines 3 and 6 in Algorithm 1). Since $k_0^{\frac{1}{2}} = \mathcal{O}\left((1-\tilde{\rho})^{-3}\right)$, this number can dominate in (C.86), indicating that the sample complexity result is network-dependent. In order to have the complexity result as indicated in Table 1, we perform $T = \lceil \frac{2}{\sqrt{1-\rho}} \rceil$ Chebyshev communications rounds by Algorithm B.1 so that $\tilde{\rho} \leq \frac{1}{2}$ by (B.3) and hence the sample complexity cost is independent of ρ and $\tilde{\rho}$. In this regime, the number of local neighbor communications is

$$T_0 + TK = \mathcal{O}\left(\frac{1}{\sqrt{1-\rho}}\max\left\{\frac{(L\delta)^{\frac{3}{2}} + \hat{\sigma}^3 + \sigma^3}{\varepsilon^{\frac{3}{2}}}, \frac{\hat{\sigma}^3\left(|\ln\varepsilon| + |\ln\hat{\sigma}|\right)^{\frac{3}{2}}}{\varepsilon^{\frac{3}{2}}}\right\}\right)$$

which is optimal in terms of dependence upon ρ (Scaman et al. 2017). Alternatively, if we let the initial batch size be $m_0 = \mathcal{O}\left(\frac{k_0^{\frac{1}{3}}}{\sigma^2}\right)$ independent of ε , then the middle term in (C.102) can be dominated by the first term, in which case the sample complexity is network-independent, after the initial iteration. For cases where the original communication network is not too sparse, e.g. ρ is not too close to 1, this may be preferred over performing Chebyshev acceleration.