

**Communication-Efficient Optimization Algorithms
for Decentralized Machine Learning**

*Submitted in partial fulfillment of the requirements for
the degree of
Doctor of Philosophy
in
Electrical and Computer Engineering*

Boyue Li

B.S., Electronic Engineering, Tsinghua University
M.S., Language Technologies, Carnegie Mellon University

Carnegie Mellon University
Pittsburgh, PA

May 2023

© Boyue Li, 2023
All rights reserved.

Acknowledgements

I would like to thank the generous support by ONR under N00014-18-1-2142 and N00014-19-1-2404, ARO under W911NF-18-1-0303, AFOSR/AFRL under FA8750-20-2-0504, NSF under ECCS-1818571, CCF-1806154, CCF-1901199 and CCF-2007911, as well as the Wei Shen and Xuehong Zhang Presidential Fellowship at Carnegie Mellon University.

I am deeply grateful for my advisor and the chair of the committee, Dr. Yuejie Chi, for her mentorship and support through my PhD journey. Your passion for research and dedication to students have been inspiring me since the first day we work together. I would like to extend my sincere thanks to my other committee members: Dr. Guannan Gu, Dr. Mingyi Hong and Dr. Soumya Kar, for giving invaluable feedback to help shaping my thesis.

I would like to acknowledge everyone in Yuejie's group: Dr. Yuanxin Li, Dr. Maxime Ferreira Da Costa, Dr. Harlin Lee, Dr. Vincent J. Monardo, Dr. Tian Tong, Shicong Cen, Laixi Shi, Diogo Cardoso, Pedro Valdeira, Jiin Woo, Harry Dong, Dr. Zhize Li, Lingjing Kong, Zixin Wen and He Wang, my life in Porter Hall basement was brighter because of all of you.

Finally, special thanks to my lovely fiancée, it would be impossible for me to get through the COVID-19 pandemic without you. And thanks to my parents for your unconditional support over the years.

Boyue Li

Abstract

Emerging applications in multi-agent environments such as internet-of-things, networked sensing, large-scale machine learning and federated learning, have attracting increasing attention for decentralized optimization algorithms that are resource efficient in both computation and communication while being able to protect data privacy. This thesis considers the prototypical setting where the agents work collaboratively to minimize the sum of local loss functions by only communicating with their neighbors over a predetermined network topology. We propose four decentralized optimization algorithms, with the intertwined goals of achieving communication efficiency, computation efficiency, as well as data privacy through carefully designed update procedures. For all algorithms, we provide theoretical convergence guarantees and perform extensive numerical experiments to support the analyses.

First, we propose a Newton-type algorithm called *Network-DANE* for decentralized problems with strongly convex objectives, which utilizes gradient tracking and extra mixing (i.e., multiple mixing rounds per iteration) to extend the celebrated server/client algorithm DANE to the decentralized setting. Our analysis shows that, similar to DANE, *Network-DANE* achieves linear convergence guarantees for general smooth strongly convex and quadratic objective functions, and can provably harness data similarity across agents to accelerate convergence, which highlights its communication efficiency. We further extend *Network-DANE* by allowing a nonsmooth penalty term for composite optimization problems, and by using stochastic variance-reduced local updates for computation efficiency.

Next, for more general decentralized nonconvex empirical risk minimization (ERM) problems, we propose *DESTRESS*, which matches the optimal incremental first-order oracle (IFO) complexity of centralized algorithms for finding first-order stationary points, while maintaining communication efficiency. In addition to gradient tracking and extra mixing, *DESTRESS* also incorporates randomly activated stochastic recursive mini-batch gradient updates to avoid unnecessary computations, which allows the improvement upon prior decentralized algorithms over a wide range of parameter regimes.

Then, we consider communication compression to further improve communication efficiency for decentralized nonconvex optimization problems, which leads to the development of *BEER*. This algorithm

also leverages stochastic gradient tracking, and in addition incorporates communication compression together with error feedback to improve communication quality, which allows it to maintain communication efficiency even when data distribution over agents is highly heterogeneous.

Finally, based on BEER, we propose PORTER for decentralized nonconvex optimization, which can provably converge to global first-order stationary points while preserving each agent's privacy under the notion of differential privacy. PORTER utilizes stochastic gradient tracking, communication compression together with error feedback as BEER does, and further leverages Gaussian perturbation with gradient clipping to preserve privacy for arbitrary objective functions.

In summary, our work emphasizes 1) the effectiveness of gradient tracking in estimating global gradients, 2) by using extra mixing, communication compression and error feedback, the overall efficiency can be substantially improved, and 3) privacy can be preserved thorough gradient clipping and Gaussian perturbation. The key algorithm design ideas can also be applied, in a systematic manner, to design new resource-efficient decentralized optimization algorithms.

Keywords: decentralized optimization, communication and computation efficiency, gradient tracking, variance reduction, communication compression, error feedback, differential privacy

Contents

Acknowledgements	iii
Abstract	iv
Contents	vi
List of Tables	ix
List of Figures	x
1 Introduction	1
1.1 Problem formulation	2
1.2 Contributions	4
1.3 Related works	13
1.4 Notation	15
1.5 Thesis organization	16
2 Decentralized Newton-style algorithm	17
2.1 Preliminaries	18
2.2 The Network-DANE algorithm	20
2.3 Convergence guarantees	22
2.4 Extension to nonsmooth composite optimization	27
2.5 Extension with variance reduction	27
2.6 Numerical experiments	29
3 Decentralized stochastic recursive gradient algorithm	35
3.1 The DESTRESS algorithm	35
3.2 Convergence guarantees	37

3.3 Numerical experiments	39
4 Decentralized stochastic algorithm with communication compression	44
4.1 Preliminaries	44
4.2 The BEER algorithm	46
4.3 Convergence guarantees	48
4.4 Numerical experiments	50
5 Decentralized private stochastic algorithm with communication compression	54
5.1 Preliminaries	54
5.2 The PORTER algorithm	57
5.3 Local differential privacy guarantee	58
5.4 Convergence with bounded gradient assumption	58
5.5 Convergence without bounded gradient assumption	60
5.6 Numerical experiments	61
6 Conclusions	64
A Appendix for Chapter 2	66
A.1 Derivation of Equation (2.6)	66
A.2 Proof of Theorem 1 and Theorem 2	66
A.3 Proofs of Theorem 3 and Theorem 4	71
A.4 Proof of Theorem 5	74
A.5 Proof of Lemma 1	75
A.6 Proof of Lemma 2	81
A.7 Proof of Lemma 3	84
B Appendix for Chapter 3	90
B.1 Experiment details	90
B.2 Proof of Theorem 6	90
B.3 Proof of Corollary 3	93
B.4 Proof of Lemma 5	94
B.5 Proof of Lemma 6	101
C Appendix for Chapter 4	107
C.1 Technical lemmas	107

C.2 Recursive relations of main errors	109
C.3 Proof of Theorem 7	113
C.4 Proof of Theorem 8	115
D Appendix for Chapter 5	117
D.1 Proof of Theorem 9	117
D.2 Proof of Theorem 10	118
D.3 Proof of Theorem 11	126
Bibliography	132

List of Tables

1.1	Proposed algorithms, corresponding chapters and contribution highlights	5
1.2	Complexities of Network-DANE	6
1.3	Complexities of stochastic variance-reduced algorithms	8
1.4	Complexities of BEER	10
1.5	Utilities and iterations complexities of private optimization algorithms	12
3.1	Communication complexities of DESTRESS	39
3.2	Settings for logistic regression with nonconvex regularization experiments for DESTRESS	40
3.3	Settings for neural network training experiments for DESTRESS	41

List of Figures

1.1	Illustration of two distributed settings	2
2.1	Convergence of Network-DANE and Network-SVRG for linear regression experiments	30
2.2	Convergence of Network-DANE under different rounds of mixing for linear regression	31
2.3	Number of communication rounds till converge with respect to different numbers of local iterations for Network-SVRG	32
2.4	Convergence of Network-DANE under different network topologies for linear regression	32
2.5	Convergence of Network-DANE under different rounds of mixing for linear regression with ℓ_1 -norm regularization	33
2.6	Convergence of Network-DANE and Network-SVRG under different rounds of mixing for logistic regression	34
2.7	Convergence of Network-DANE and Network-SVRG for neural networks	34
3.1	Convergence of DESTRESS under different network topologies for logistic regression with non-convex regularization	41
3.2	Convergence of DESTRESS under different network topologies for one-hidden-layer neural network training	42
3.3	The convergence precision $1/\epsilon^2$ with respect to gradient evaluations	42
4.1	Convergence of BEER for logistic regression with nonconvex regularization	51
4.2	Convergence of BEER for 1-hidden-layer neural network training	52
4.3	Convergence of BEER under different network topologies for logistic regression with nonconvex regularization	53
4.4	Convergence of BEER under different compression schemes for logistic regression with nonconvex regularization	53
5.1	Illustration of input norm and clipped norm for clipping operators	57

5.2	Convergence of PORTER for logistic regression with nonconvex regularization when guaranteeing $(10^{-2}, 10^{-3})$ -LDP	62
5.3	Convergence of PORTER for logistic regression with nonconvex regularization when guaranteeing $(10^{-1}, 10^{-3})$ -LDP	62
5.4	Convergence of PORTER for neural network training when guaranteeing $(10^{-2}, 10^{-3})$ -LDP . . .	63
5.5	Convergence of PORTER for neural network training when guaranteeing $(10^{-1}, 10^{-3})$ -LDP . . .	63

Chapter 1

Introduction

Distributed optimization has been a classic topic [BT89] yet is attracting significant attention recently in machine learning due to its numerous applications such as distributed training [BPC⁺11], multi-agent learning [NOP10], and federated learning [KMR15, KMY⁺16, MMR⁺17]. At least two facts contribute towards this resurgence of interest: (1) the scale of modern datasets has oftentimes far exceeded the **computation** capacity of a single machine and requires coordination across multiple machines; (2) **communication** and **privacy** constraints disfavor information sharing in a centralized manner and necessitates distributed infrastructures.

Broadly speaking, in terms of communication patterns, there are two distributed settings that have received the most interest as illustrated in Figure 1.1: 1) the *server/client* setting, which assumes the existence of a central parameter server that can aggregate and share information with all agents; and 2) the *decentralized* setting, where each agent is only permitted to communicate with its neighbors over a locally connected network specified by a communication graph. For both settings, each agent only has access to a disjoint subset of the data samples and aims to work collaboratively to optimize the global objective function $f(x)$, by only exchanging information with the parameter server or its neighbors (in other words, no centralized coordination is present). It is in general more challenging to developing algorithms for the latter setting, which is the focus of this thesis.

For a typical decentralized optimization algorithm, agents generally alternate between (1) communication, which propagates local information and enforces consensus and (2) computation, which updates individual optimization variables and improves convergence using information received from neighbors. Compared to the server/client setting, [LZZ⁺17] suggests decentralized algorithms can effectively avoid communication jams on busy agents, e.g., the parameter server, and be more efficient in wall-clock time. We can roughly break down an algorithm's total running cost C_{total} to the total number of iterations T

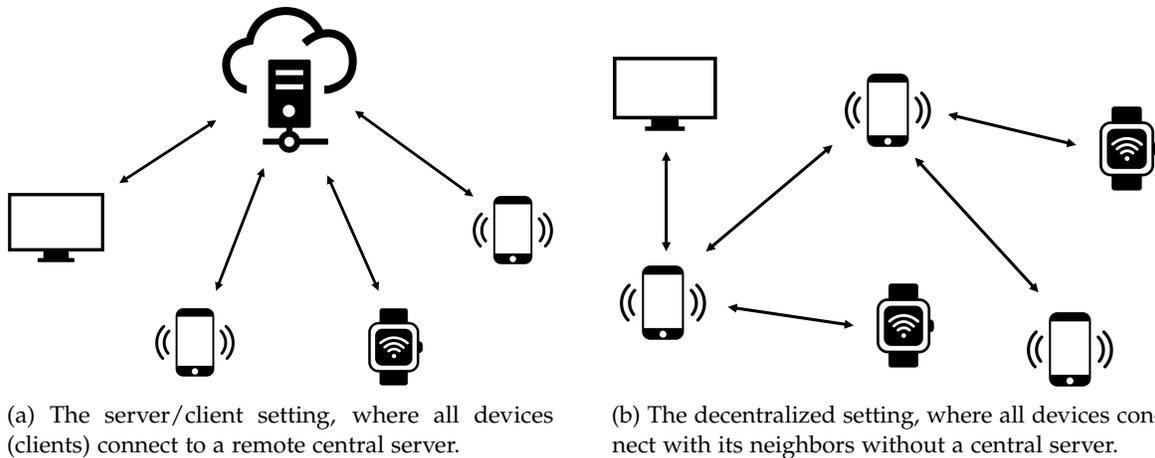


Figure 1.1: Illustration of two distributed settings: (a) the server/client setting, and (b) the decentralized setting. This thesis focuses on the decentralized setting.

multiplied by the sum of per-iteration communication cost C_{comm} and per-iteration computation cost C_{comp} , i.e.

$$C_{\text{total}} = T \times (C_{\text{comm}} + C_{\text{comp}}). \quad (1.1)$$

For example, when reducing total time spent to reach certain accuracy is the top priority, the running cost can be defined as wall-clock time. Horizontally scaling the system by adding more agents to run in parallel may reduce T , but it will increase the communication cost C_{comm} in the meantime, which may lead to a worse total cost C_{total} in some cases. Therefore, achieving a desired level of resource efficiency for a decentralized algorithm often requires careful and delicate trade-offs between computation and communication costs, as these objectives are often conflicting in nature.

Privacy is another crucial constraint in large-scale decentralized optimization, because these applications often involve transferring computation outcomes on sensitive data that may contain personal or confidential information to external agents. Without carefully designed privacy-preserving protocols, systems may risk to directly reveal sensitive data or be vulnerable to other forms of attacks, e.g., linkage attacks [DR14]. A notable example of linkage attacks is the identification of the medical records of the governor of Massachusetts from anonymized medical data with voter registration data [Lam01], which emphasizes the importance of privacy-preserving.

1.1 Problem formulation

In this thesis, we investigate decentralized minimization problems, with the aim of achieving communication and computation efficiency simultaneously, as well as satisfying privacy-preserving constraints. In this section, we formally define the problem and related concepts.

Objective functions Consider the following minimization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \quad f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad (1.2)$$

where $\mathbf{x} \in \mathbb{R}^d$ denotes the optimization variable, $f_i(\cdot)$ denotes the local objective function at the i -th agent ($1 \leq i \leq n$) of n agents. In many machine learning applications, for example, empirical risk minimization (ERM) problems, local objective functions often have a finite-sum structure that is defined as

$$f_i(\mathbf{x}) := \frac{1}{m} \sum_{\mathbf{z} \in \mathcal{M}_i} \ell(\mathbf{x}; \mathbf{z}), \quad (1.3)$$

where $\ell(\mathbf{x}; \mathbf{z})$ denotes the sample loss of the sample \mathbf{z} at \mathbf{x} , \mathcal{M}_i denotes the dataset at agent i , $m = |\mathcal{M}_i|$ denotes the number of data samples at each agent.¹ In addition, we define the full dataset $\mathcal{M} = \cup_{i=1}^n \mathcal{M}_i$ and total sample size $N = |\mathcal{M}| = mn$. The communication pattern of the network is specified by an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where \mathcal{V} denotes the set of all agents, and two agents can exchange information if and only if there is an edge in \mathcal{E} connecting them.

General assumptions We assume Assumptions 1 and 2 for all analysis in this thesis. Assumption 1 guarantees the existence of nontrivial solution(s) of the global optimization problem. Assumption 2 limits the change of sample gradients, which also implies the local objective functions $f_i(\cdot)$ and global objective function $f(\cdot)$ have Lipschitz gradients.

Assumption 1 (Bounded global objective function). *The global objective function $f(\mathbf{x})$ is bounded below, i.e.,*

$$f^* = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) > -\infty.$$

Assumption 2 (L -smooth sample loss function). *A sample loss function $\ell(\mathbf{x}; \mathbf{z})$ is L -smooth if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\forall \mathbf{z} \in \mathcal{Z}$, the following inequality holds:*

$$\|\nabla \ell(\mathbf{x}; \mathbf{z}) - \nabla \ell(\mathbf{y}; \mathbf{z})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2.$$

Information mixing The information mixing between agents is conducted by updating the local information via a weighted sum of information from neighbors, which is characterized by a mixing (gossiping) matrix. Concerning this matrix is an important quantity called the mixing rate, defined in Definition 1.

Definition 1 (Mixing matrix and mixing rate). *The mixing matrix is a matrix $\mathbf{W} = [w_{ij}] \in \mathbb{R}^{n \times n}$, such that $w_{ij} = 0$ if agent i and j are not connected according to the communication graph \mathcal{G} . Furthermore, $\mathbf{W}\mathbf{1}_n = \mathbf{1}_n$ and $\mathbf{W}^\top \mathbf{1}_n = \mathbf{1}_n$. The mixing rate of a mixing matrix \mathbf{W} is defined as*

$$\alpha := \|\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top\|_{\text{op}} \in [0, 1). \quad (1.4)$$

¹We assuming the data is distributed equally among all agents, but it can be easily generalized to the unequal splitting case.

The mixing rate indicates the speed of information shared across the network. For example, for a fully connected network, choosing $W = \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top$ leads to $\alpha = 0$. For general networks and mixing matrices, [NOR18, Proposition 5] provides comprehensive bounds on $1 - \alpha$ —also known as the spectral gap—for various graphs. For instance, one has $\alpha \asymp 1$ with high probability in an Erdős-Rényi random graph if the graph is connected. In practice, FDLA matrices [XB04] are more favorable because it can achieve a much smaller mixing rate, but they usually contain negative elements and are not symmetric. Our analysis can handle arbitrary mixing matrices as long as their row/column sums equal to one and its mixing rate is smaller than 1.

Convergence metrics For strongly convex optimization algorithms, a global optimum is guaranteed to exist. Thus, we consider the distance of output to the global optimum defined in Definition 2, where a smaller ν means better performance.

Definition 2 (ν -accurate solution). *The output of a deterministic optimization algorithm $\mathbf{x} \in \mathbb{R}^d$ is a ν -solution of a differentiable strongly convex function $f(\mathbf{x})$ if*

$$\|\mathbf{x} - \mathbf{x}^*\|_2^2 \leq \nu^2,$$

where $\mathbf{x}^* = \arg \min f(\mathbf{x})$.

For nonconvex optimization algorithms, we consider the norm of the output's gradient as the quality metric. Definition 3 defines a ν -first-order stationary point, where a smaller ν indicates better convergence.

Definition 3 (ν -first-order stationary point). *The output of an optimization algorithm $\mathbf{x} \in \mathbb{R}^d$ is a ν -first-order stationary point of a differentiable function $f(\mathbf{x})$ if*

- $\|\nabla f(\mathbf{x})\|_2^2 \leq \nu^2$ for deterministic optimization algorithms.
- $\mathbb{E}\|\nabla f(\mathbf{x})\|_2^2 \leq \nu^2$ for stochastic optimization algorithms.

1.2 Contributions

Our main contribution is the design and analysis of four resource-efficient primal-only decentralized optimization algorithms: Network-DANE (cf. Algorithm 1), DESTRESS (cf. Algorithm 4), BEER (cf. Algorithm 5) and PORTER (cf. Algorithm 7) as summarized in Table 1.1. While all algorithms share a similar framework based on gradient tracking, each one focuses on a unique perspective of decentralized ERM problems, resulting in distinct algorithm designs that can achieve overall resource efficiency and address problem-specific constraints. Proposed algorithms encompass a variety of methods, such as approximate

Algorithm	Chapter	Communication efficiency	Computation efficiency	Communication compression	Differential privacy
Network-DANE	2	✓			
DESTRESS	3	✓	✓		
BEER	4	✓	✓	✓	
PORTER	5	✓	✓	✓	✓

Table 1.1: Proposed algorithms, corresponding chapters and contribution highlights.

Newton-type methods, stochastic variance-reduced methods, communication compression, error feedback and privacy perturbation.

In addition, our work indicates that by performing a judiciously chosen amount of local communication and computation per iteration, the overall efficiency can be remarkably boosted while simultaneously meeting problem-specific constraints. Extensive numerical experiments are provided to corroborate our theoretical findings, and to demonstrate the practical efficacy of the proposed algorithms over competitive baselines. All code can be accessed at

<https://github.com/liboyue/Network-Distributed-Algorithm>.

1.2.1 Decentralized Newton-style algorithm

This subsection highlights our contributions of designing the decentralized Newton-type method `Network-DANE` (cf. Chapter 2) [LCCC20], which converges linearly for smooth strongly convex objectives and quadratic objectives, and can be extended to nonsmooth composite objectives using a proximal operator.

Algorithm development We start by studying an approximate Newton-type method called DANE [SSZ14], which is among the most popular communication-efficient server/client optimization algorithms for solving ERM problems. Then, we develop `Network-DANE`, which generalizes DANE to the decentralized (network) setting. The main challenge in developing such an algorithm is to track and update a faithful estimate of the global gradient at each agent, despite the lack of centralized information aggregation. Towards this end, we leverage *dynamic average consensus* (originally proposed in the control literature [ZM10] and later adopted in decentralized optimization [QL18, NOS17, DLS16]) to track and correct locally aggregated gradients at each agent — a scheme commonly referred to as *gradient tracking*. We then employ the corrected gradient in local computations, according to the subroutine adapted from DANE. This simple idea allows one to adapt approximate Newton-type methods to network-distributed optimization, without communicating Hessian matrices.

Algorithm	Communication rounds	Extra averaging	Loss functions	β
EXTRA [SLWY15a]	$\kappa^2 \log(1/\epsilon)$	✗	Strongly convex	Arbitrary
DGD [QL18]	$\frac{\kappa^2 \log(1/\epsilon)}{(1-\alpha)^2}$	✗	Strongly convex	
Network-DANE (Algorithm 1)	$\frac{\kappa(\beta/\sigma+1) \log(1/\epsilon)}{(1-\alpha)^2}$	✗	Quadratic	Arbitrary
	$\log \kappa \cdot \frac{(\beta^2/\sigma^2+1) \log(1/\epsilon)}{(1-\alpha)^{1/2}}$	✓		
	$\frac{\kappa^2 \log(1/\epsilon)}{(1-\alpha)^2}$	✗	Strongly convex	
	$\log \kappa \cdot \frac{\kappa(\beta/\sigma+1) \log(1/\epsilon)}{(1-\alpha)^{1/2}}$	✓		
Network-SVRG (Algorithm 3)	$\log \kappa \cdot \frac{\log(1/\epsilon)}{(1-\alpha)^{1/2}}$	✓	Strongly convex	$\beta \leq \sigma/200$

Table 1.2: Communication complexity of the proposed algorithms for quadratic and strongly convex losses to reach ϵ -accuracy. Here, σ , L and $\kappa = L/\sigma$ are the strong convexity, smoothness, and condition number of the local loss functions. In addition, $\beta \leq L$ is the homogeneity parameter gauging the similarities of the local loss functions defined in (4), and α is the mixing rate defined in (1.4). Here, we assume the extra averaging step is implemented via Chebyshev acceleration scheme [AS14]. EXTRA [SLWY15a] and DGD [QL18] are listed as baselines. The big- O notation (defined in Section 1.4) is omitted for simplicity.

Our ideas for designing Network-DANE can be extended, in a systematic manner, to obtain decentralized versions of other algorithms developed for the server/client setting, by modifying the local computation step properly. As a notable example, we develop Network-SVRG, which performs variance-reduced stochastic optimization locally to enable further computational savings [JZ13]. We also demonstrate that Network-DANE can be extended to the proximal setting for nonsmooth composite optimization in a straightforward manner.

Efficiency analysis Network-DANE achieves an intriguing trade-off between communication and computation efficiency. During every iteration, each agent only communicates the parameter and gradient estimate to its neighbors, and is therefore communication-efficient globally; moreover, the local subproblem at each agent can be solved efficiently with accelerated gradient methods, and is thus computation-efficient locally. When the network exhibits a high degree of locality, we show that by allowing multiple rounds of local mixing within each iteration, an improved overall communication complexity can be achieved as it accelerates the rate of convergence. Theoretically, we establish the linear convergence of Network-DANE for strongly convex and quadratic losses, and show that incorporating extra mixing leads to great improvements in communication complexity. For Network-SVRG, we establish its linear convergence for the case of smooth strongly convex losses with extra rounds of averaging. Our analysis is highly nontrivial, as it needs to deal with the tight couplings of optimization and network consensus errors through a carefully designed linear

system of Lyapunov functions, especially in the context of approximate Newton-type methods which are known to be harder to handle than simple gradient-type methods.

Table 1.2 summarizes the convergence rates of the proposed algorithm and baseline algorithms. Let σ , L , $\kappa = L/\sigma$ and β denote the strong convexity, smoothness, condition number of the local loss functions and homogeneity parameter. For general strongly convex losses, Network-DANE matches the results of both EXTRA [SLWY15a] and DGD [QL18] without extra mixing, but achieves a significantly improved communication complexity when incorporating extra mixing by a factor of $O((1 - \alpha)^{3/2}(\beta/\sigma + 1)\kappa^{-1} \log \kappa)$. The resulting convergence rate attains the optimal network dependency $O((1 - \alpha)^{-1/2})$, and improves the condition number dependency from $O(\kappa^2)$ to $O((\beta/\sigma + 1)\kappa \log \kappa)$, which can be a significant improvement in the case that data similarity parameter β is small. For example, in the big data scenario, each agent can access a large number of i.i.d. data samples, which results in a small homogeneity parameter. Moreover, Network-DANE reaches a condition number free (up to a log factor) communication complexity for quadratic losses, which emphasizes the potential of the Newton-type method and extra mixing. Network-SVRG reaches a condition number free (up to a log factor) communication complexity as well for strongly convex objectives, providing β is small enough, at a much lower computation complexity than Network-DANE due to its stochastic nature. Our results shed light on the impacts of data homogeneity and network connectivity upon the rate of convergence.

1.2.2 Decentralized stochastic recursive gradient algorithm

This subsection highlights our contributions of the development of DEcentralized STochastic REcurSive gradient methodS (DESTRESS, cf. Chapter 3) [LLC22], which is a resource-efficient algorithm for decentralized nonconvex ERM problems. DESTRESS provably finds first-order stationary points of the global objective function $f(x)$ with the optimal incremental first-order (IFO) oracle complexity, i.e., the complexity of evaluating sample gradients, but at a much lower communication complexity compared to existing decentralized algorithms over a wide range of parameter regimes.

Algorithm development DESTRESS tries to improve computation and communication simultaneously to achieve overall resource efficiency. To reduce local computation, DESTRESS harnesses the finite-sum structure of the empirical risk function by performing stochastic variance-reduced recursive gradient updates [NvP⁺22, FLLZ18, WJZ⁺19, Li19, LR21b, LBZR21, ZXG18]—an approach that is shown to be optimal in terms of IFO complexity in the centralized setting—in a randomly activated manner to further improve computational efficiency when the local sample size is limited. To reduce communication, DESTRESS employs gradient tracking [ZM10] with a few mixing rounds per iteration, which helps accelerate the

Algorithms	Setting	Per-agent IFO Complexity	Communication Rounds
SVRG [AZH16, RHS+16]	centralized	$N + \frac{N^{2/3}L}{\epsilon^2}$	-
SCSG/SVRG+ [LJCJ17, LL18]	centralized	$N + \frac{N^{2/3}L}{\epsilon^2}$	-
SNVRG [ZXG18]	centralized	$N + \frac{N^{1/2}L}{\epsilon^2}$	-
SARAH/SPIDER/SpiderBoost [NvP+22, FLLZ18, WJZ+19]	centralized	$N + \frac{N^{1/2}L}{\epsilon^2}$	-
SSRGD/ZeroSARAH/PAGE [Li19, LR21b, LBZR21]	centralized	$N + \frac{N^{1/2}L}{\epsilon^2}$	-
D-GET [SLH20]	decentralized	$m + \frac{1}{(1-\alpha)^2} \cdot \frac{m^{1/2}L}{\epsilon^2}$	Same as IFO
GT-SARAH [XKK22a]	decentralized	$m + \max\left(\frac{1}{(1-\alpha)^2}, \left(\frac{m}{n}\right)^{1/2}, \frac{(m/n+1)^{1/3}}{1-\alpha}\right) \cdot \frac{L}{\epsilon^2}$	Same as IFO
DESTRESS (Algorithm 4)	decentralized	$m + \frac{(m/n)^{1/2}L}{\epsilon^2}$	$\frac{1}{(1-\alpha)^{1/2}} \cdot \left((mn)^{1/2} + \frac{L}{\epsilon^2}\right)$

Table 1.3: The per-agent IFO complexities and communication complexities to find ϵ -first-order stationary points by stochastic variance-reduced algorithms for nonconvex ERM problems. The first five algorithms are designed for the centralized setting, and the remaining D-GET, GT-SARAH and DESTRESS are for the decentralized setting. m, n, L are defined in Section 1.1 and α is the mixing rate defined in (1.4). The big- O notation (defined in Section 1.4) and logarithmic terms are omitted for simplicity.

convergence through better information sharing [LCCC20]; the extra mixing scheme can be implemented using Chebyshev acceleration [AS14] (detailed in Section 2.1.2) to further improve the communication efficiency.

Efficiency analysis In a nutshell, to find an ϵ -first-order stationary point (cf. Definition 3), where x^{out} is the output of DESTRESS, and the expectation is taken with respect to the randomness of the algorithm, DESTRESS requires:

- $O\left(m + (m/n)^{1/2}L/\epsilon^2\right)$ per-agent IFO calls, which is *network-independent*; and
- $O\left(\frac{\log\left((n/m)^{1/2}+2\right)}{(1-\alpha)^{1/2}} \cdot \left((mn)^{1/2} + \frac{L}{\epsilon^2}\right)\right)$ rounds of communication,

where L is the smoothness parameter of the sample loss, $\alpha \in [0, 1)$ is the mixing rate of the mixing matrix. DESTRESS is resource-efficient for it reaches optimal computation complexity with state-of-the-art communication complexity.

Table 1.3 summarizes convergence guarantees of representative stochastic variance-reduced algorithms for finding first-order stationary points across centralized and decentralized communication settings.

- In terms of the computation complexity, the overall IFO complexity of DESTRESS—when summed

over all agents—becomes

$$n \cdot O\left(m + (m/n)^{1/2}L/\epsilon^2\right) = O\left(mn + (mn)^{1/2}L/\epsilon^2\right) = O\left(N + N^{1/2}L/\epsilon^2\right),$$

matching the optimal IFO complexity of centralized algorithms (e.g., SPIDER [FLLZ18], PAGE [LBZR21]) and distributed server/client algorithms (e.g., D-ZeroSARAH [LR21b]). However, the state-of-the-art decentralized algorithm GT-SARAH [XKK22a] is not able to achieve this optimal IFO complexity for all situations (see Table 1.3). To the best of our knowledge, DESTRESS is the first algorithm to achieve the optimal IFO complexity for the decentralized setting regardless of network topology and sample size.

- When it comes to the communication complexity, it is observed that the communication rounds of DESTRESS can be decomposed into the sum of an ϵ -independent term and an ϵ -dependent term (up to a logarithmic factor), i.e.,

$$\underbrace{\frac{1}{(1-\alpha)^{1/2}} \cdot (mn)^{1/2}}_{\epsilon\text{-independent}} + \underbrace{\frac{1}{(1-\alpha)^{1/2}} \cdot \frac{L}{\epsilon^2}}_{\epsilon\text{-dependent}};$$

similar decompositions also apply to competing decentralized algorithms. DESTRESS significantly improves the ϵ -dependent term of D-GET and GT-SARAH by at least a factor of $\frac{1}{(1-\alpha)^{3/2}}$, and therefore, saves more communications over poorly connected networks. Further, the ϵ -independent term of DESTRESS is also smaller than that of D-GET/GT-SARAH as long as the local sample size is sufficiently large, i.e., $m = \Omega\left(\frac{n}{1-\alpha}\right)$, which also holds for a wide variety of application scenarios.

In sum, DESTRESS harnesses the ideas of random client activation, variance reduction, gradient tracking and extra mixing in a sophisticated manner to achieve a scalable decentralized algorithm for nonconvex empirical risk minimization that is competitive in both computation and communication over existing approaches.

1.2.3 Decentralized stochastic algorithm with communication compression

This subsection highlights our contributions of the development of Better comprEssion for dEcentRalized optimization (BEER, cf. Chapter 4) [ZLL⁺22], which is a communication-efficient algorithm for decentralized nonconvex optimization problems. Communication compression is a well-established method to improve communication efficiency for server/client distributed optimization algorithms. However, for the decentralized setting, most existing algorithms, e.g. [KSJ19, TLQ⁺19, SDGD21], require strong assumptions as bounded gradient assumption or bounded dissimilarity assumption to guarantee convergence. BEER

Algorithm	Communication rounds	Per-agent IFO complexity	Extra assumption
SQuARM-SGD [SDGD21]	$\frac{n\tau^2}{\epsilon^2} + \frac{\sigma^2}{bne^4}$	$\frac{n\tau^2}{\epsilon^2} + \frac{\sigma^2}{ne^4}$	Bounded gradient ⁽¹⁾
DeepSqueeze [TLQ ⁺ 19]	$\frac{\beta}{\epsilon^3} + \frac{\sigma^2}{bne^4}$	$\frac{\beta}{\epsilon^3} + \frac{\sigma^2}{ne^4}$	Bounded dissimilarity ⁽²⁾
CHOCO-SGD [KSJ19]	$\frac{\tau}{\epsilon^3} + \frac{\sigma^2}{bne^4}$	$\frac{\tau}{\epsilon^3} + \frac{\sigma^2}{ne^4}$	Bounded gradient ⁽¹⁾
BEER (Algorithm 5)	$\frac{1}{\epsilon^2}$	$\frac{1}{\epsilon^2} + \frac{\sigma^2}{\epsilon^4}$	-

Table 1.4: Comparison of iteration (communication) complexity, per-agent IFO complexity, and extra assumptions in addition to BEER’s assumptions, for existing decentralized stochastic nonconvex optimization algorithms to reach ϵ -first-order stationary points. σ^2 and b denote the local gradient variance and batch size, respectively.

⁽¹⁾ The bounded gradient assumption is defined as $\forall \mathbf{x} \in \mathbb{R}^d, \mathbb{E}_{z_i \sim \mathcal{Z}_i} \|\nabla \ell(\mathbf{x}, z_i)\|_2^2 \leq \tau^2$.

⁽²⁾ The bounded dissimilarity assumption is defined as $\forall \mathbf{x} \in \mathbb{R}^d, \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|_2^2 \leq \beta^2$.

removes these strong assumptions, which enables it to converge under arbitrary data heterogeneity where other algorithms fail. BEER enjoys a faster convergence rate to find first-order stationary points of the global objective function $f(\mathbf{x})$ than existing algorithms that use communication compression. We establish convergence analysis for BEER using stochastic gradient for nonconvex optimization problem, and further extend this convergence to problems satisfying the Polyak-Łojasiewicz (PL) condition.

Algorithm development We develop BEER, in hope to match the convergence rate of centralized stochastic gradient descent (SGD) without aforementioned strong assumptions and under arbitrary data distributions, by leveraging gradient tracking and error feedback [RSF21]. For each variable that needs to be communicated, BEER tracks and maintains a control sequence that serves as compressed surrogate by communicating and accumulating a compressed error from the variable and the control sequence, which leads to an improved convergence rate.

Efficiency analysis To find an ϵ -first-order stationary point, using both full gradient and stochastic gradient (with proper batch size), BEER takes

- $O(\epsilon^{-2})$ iterations for optimizing nonconvex objective functions,
- $O(\kappa \log(1/\epsilon))$ iterations for optimizing PL objectives, where κ is the condition number,

while only requires bounded local stochastic gradient variance assumption and global L -smoothness assumption.

Table 1.4 summarizes the iteration complexities, per-agent IFO complexities and extra assumptions for BEER and baseline algorithms for optimizing nonconvex objectives using minibatch stochastic gradients,

where the IFO complexities are results after using appropriate batch sizes, extra assumptions are those needed in addition to BEER’s assumptions.

Let σ^2 denote the local gradient variance. When all algorithms use full gradient updates, which is equal to $\sigma = 0$, in terms of the variance-independent terms in the iteration complexities, BEER improves the dependency on ϵ from $O(\epsilon^{-3})$ to $O(\epsilon^{-2})$ and matches the centralized SGD.

When the local gradient variance is present, using appropriate batch sizes for all algorithms, BEER still attains the best iteration complexity. Furthermore, as shown in Table 1.4, SQuARM-SGD and CHOCO-SGD all require bounded gradient assumption, and DeepSqueeze requires bounded dissimilarity assumption, which results in a significant performance degeneration when those assumptions are not satisfied. On the contrary, BEER maintains the same convergence guarantees under arbitrary data heterogeneity, which highlights BEER’s adaptivity to more optimization problems.

1.2.4 Decentralized private stochastic algorithm

This subsection highlights our contributions of developing PORTER (cf. Chapter 5), which is a resource-efficient and differentially private decentralized stochastic optimization algorithm for nonconvex ERM problems. PORTER is the first private decentralized optimization algorithm to incorporate communication compression and gradient clipping, to improve communication efficiency and converge for arbitrary objective functions without the bounded gradient assumption, respectively. We establish convergence analysis both with and without the bounded gradient assumption, and show explicit dependency on the mixing rate and compression parameter.

Algorithm development The focus of PORTER is preserving the privacy of each agent while being resource-efficient. To address the privacy concern, PORTER applies gradient clipping to ensure gradients are bounded for any objective function, then adds privacy perturbation to clipped gradients. To improve communication efficiency, PORTER incorporates gradient tracking [ZM10] without extra mixing to construct an estimate of the global gradient at each agent, and utilizes communication compression together with error feedback [RSF21] to reduce the amount of data to be communicated without hurting the convergence rate. Lastly, to improve computation efficiency, PORTER queries a mini-batch of samples to compute stochastic gradient per agent at each iteration.

Privacy constraints PORTER achieves (ϵ, δ) -local differential privacy (LDP) (cf. Theorem 9), which protects privacy leakage between any two agents, hence is a stronger guarantee than conventional differential privacy that only protects privacy leakage to external adversaries. Privacy is protected by adding Gaussian

Algorithm	Privacy	Compression operator	Bounded gradient	Utility	Communication rounds
DP-SGD [ACG ⁺ 16]	DP	-	✓	ϕ_m	-
DDP-SRM [WJEG19]	DP	-	✓	$\frac{1}{n}\phi_m$	$n^2 d\phi_m^{-1}$
Soteria-SGD ⁽¹⁾ [LZLC22]	LDP	Unbiased	✓	$(1 + \theta^{1/2})\left(\frac{1+\omega}{n}\right)^{1/2}\phi_m$	$(1 + \theta^{1/2})\left(\frac{n}{1+\omega}\right)^{2/3}d\phi_m^{-1}$
PORTER (Algorithm 8)	LDP	General	✓	$\frac{1}{(1-\alpha)^{8/3}\rho^{4/3}}\phi_m$	ϕ_m^{-2}
PORTER (Algorithm 7)	-	General	✗	$\frac{1}{(1-\alpha)^{8/3}\rho^{4/3}}\phi_m^{1/2}$	$(1-\alpha)^{8/3}\rho^{4/3}\phi_m^{-1}$

Table 1.5: Comparison of final utility upper bounds communication complexities of different stochastic optimization algorithms that achieves (ϵ, δ) -DP/LDP. $\phi_m = \frac{\sqrt{d \log(1/\delta)}}{m\epsilon}$ is the baseline utility. Big-O notation (defined in Section 1.4) is omitted for simplicity. DP-SGD is a single-machine optimization algorithm that serves as a baseline, to show the overhead brought in by the distributed setting. DDP-SRM and Soteria-SGD are server/client distributed algorithms, but DDP-SRM doesn't use communication compression.

⁽¹⁾ Here $\theta = (1 + \omega)^{3/2}n^{-1/2}$.

perturbations to clipped stochastic gradients at each agent [DR14], where the variance of the perturbation is decided according to the target final utility and the norm of gradients (which is bounded by the clipping parameter τ).

Efficiency analysis Table 1.5 presents final utilities and corresponding communication complexities for PORTER and baseline algorithms under (ϵ, δ) -DP/LDP, where we show explicit convergence rates of PORTER for general compression operators (cf. Definition 5). When assuming the bounded gradient assumption (cf. Assumption 8), PORTER reaches $O((1 - \alpha)^{-8/3}\rho^{-4/3}\phi_m)$ utility in $O(\phi_m^{-2})$ iterations. Otherwise, PORTER reaches $O((1 - \alpha)^{-4/3}\rho^{-2/3}\phi_m^{1/4})$ ℓ_2 utility in no more than $O((1 - \alpha)^{4/3}\rho^{2/3}\phi_m^{-1})$ iterations. This convergence analysis is the first that specifies explicate dependency on mixing rate and compression parameter for decentralized private optimization algorithms.

To begin with, we first assume Assumption 8 holds, Using general compression operators (cf. Definition 5), PORTER reaches $O((1 - \alpha)^{-8/3}\rho^{-4/3}\phi_m)$ squared ℓ_2 utility in ϕ_m^{-2} iterations. The iteration complexity is relatively higher compared to the $O(\phi_m^{-1})$ complexity of baseline algorithms, due to extra errors induced by the decentralized setting.

However, the bounded gradient assumption is rarely met in practice, which makes it necessary for private optimization algorithms to apply gradient clipping to converge on more general objective functions. We propose a novel framework to analyze algorithms that uses gradient clipping without bounded gradient assumption, and show that using general compression operators, Algorithm 7 can reach $O((1 - \alpha)^{-4/3}\rho^{-2/3}\phi_m^{1/4})$ ℓ_2 utility in no more than $O((1 - \alpha)^{4/3}\rho^{2/3}\phi_m^{-1})$ iterations.

To conclude, PORTER is a promising approach to efficiently and privately optimizing decentralized nonconvex optimization problems by leveraging gradient tracking, communication compression, error feedback, gradient clipping and privacy perturbation in a refined manner, which is backed by novel convergence analysis that shows explicit rates under various assumptions.

1.3 Related works

Gradient tracking Gradient tracking [ZM10, QL18] provides a systematic approach to estimate the global gradient at each agent, which allows one to easily design decentralized optimization algorithms based on existing centralized algorithms. This idea is first incorporated to adjust distributed gradient descent (DGD) to ensure its linear convergence using a constant step size [NOS17, QL18, LSY19, XXK17, YYZS18, SS19, XSKK19]. Using the same idea, [LCCC20] extends Newton-style algorithms as well as stochastic variance-reduced algorithms with performance guarantees for optimizing strongly convex functions, and [ZY19, SLH20, XKK22b, XKK22a, LLC22, HSZ⁺22, LY22] design stochastic optimization algorithms for nonconvex problems.

Newton-type methods for distributed optimization Distributed Approximate NEwton-type Method (DANE) [SSZ14] exhibits appealing performance in both theory and practice. AIDE [RKR⁺16] relaxes the local optimization problem to an inexact solver and applies acceleration techniques, which improves the communication complexity to match the lower bound. CEASE [FGW21] further extends DANE to optimizing objective functions with proximal terms and improves the analysis.

Stochastic variance-reduced methods Many variants of stochastic variance-reduced gradient based methods have been proposed for finite-sum optimization for finding first-order stationary points, including but not limited to SVRG [JZ13, AZH16, RHS⁺16], SCSG [LJCJ17], SVRG+ [LL18], SAGA [DBLJ14], SARA [NLST17, NvP⁺22], SPIDER [FLLZ18], SpiderBoost [WJZ⁺19], SSRGD [Li19], ZeroSARA [LR21b] and PAGE [LBZR21]. SVRG/SVRG+/SCSG/SAGA utilize stochastic variance-reduced gradients as a corrected estimator of the full gradient, but can only achieve a sub-optimal IFO complexity of $O(N + N^{2/3}L/\epsilon)$. Other algorithms such as SARA, SPIDER, SpiderBoost, SSRGD and PAGE adopt stochastic recursive gradients to improve the IFO complexity to $O(N + N^{1/2}L/\epsilon)$, which is optimal indicated by the lower bound [FLLZ18, LBZR21].

Decentralized stochastic nonconvex optimization There has been a flurry of recent activities in decentralized nonconvex optimization. D-PSGD [LZZ⁺17] and SGP [ALBR19] extend stochastic gradient descent

(SGD) to solve the nonconvex decentralized expectation minimization problems with sub-optimal rates. However, due to the noisy stochastic gradients, D-PSGD can only use diminishing step size to ensure convergence, and SGP uses a small step size on the order of $1/K$, where K denotes the total iterations, [KDG03, XB04, Sha07, BJ13, LZZ⁺17, WJ21] also propose decentralized algorithms with similar structures. D^2 [TLY⁺18] introduces a variance-reduced correction term to D-PSGD, which allows a constant step size and hence reaches a better convergence rate.

GT-SAGA [XKK22b] further uses SAGA-style updates and reaches a convergence rate that matches SAGA [DBLJ14, RSPS16]. However, GT-SAGA requires to store a variable table, which leads to a high memory complexity. D-GET [SLH20] and GT-SARAH [XKK22a] adopt equivalent recursive local gradient estimators to enable the use of constant step sizes without extra memory usage. The IFO complexity of GT-SARAH is optimal in the restrictive range $m \gtrsim \frac{n}{(1-\alpha)^6}$, while DESTRESS achieves the optimal IFO over all parameter regimes.

In addition to variance reduction techniques, performing multiple mixing steps between local updates can greatly improve the dependence of the network in convergence rates, which is equivalent of communicating over a better-connected communication graph for the agents, which in turn leads to a faster convergence (and a better overall efficiency) due to better information mixing. This technique is applied by a number of recent literatures [BBKW19, PLW20, BBW21, LCCC20, HAD⁺21, IW22, SBB⁺17, SBB⁺18, LFYL20, YZLZ20, GF20, LDS21].

Communication compression Communication efficiency is critical to decentralized optimization algorithms, as communication can quickly become bottleneck of the system as the number of agents and the size of the model increase. This has led to the development of communication compression (or quantization) technique, which can significantly reduce the communication burden by transferring compressed information without hurting the convergence too much. [DSZOR15, AGL⁺17] adopt gradient compression to create a server/client distributed stochastic gradient descent, however, the large variance of compressed gradients leads to a sub-optimal convergence rate. [SFD⁺14] first proposes using error feedback to compensate for the variance induced by compression. [SCJ18, AHJ⁺18, MGTR19, LKQR20, GBLR21, LR21a] all equip similar mechanism to improve convergence for server/client distributed optimization algorithms, and [RSF21] proposes EF21 that formalizes the error feedback mechanism and reaches a $O(1/T)$ convergence rate for smooth nonconvex objective functions. [TGZ⁺18, KSJ19, SDGD21, TMHP20, ZLL⁺22, YCC⁺23, LLP23] further extend communication compression schemes to the decentralized setting.

On the other hand, many compressed methods are proposed recently such as [AGL⁺17, KFJ18, TGZ⁺18, SCJ18, KSJ19, LKQR20, GBLR21, LR21a, RSF21, FSG⁺21, ZBLR21].

Private optimization algorithms The concern of leaking sensitive data has been increasing with the rapid development of large-scale machine learning systems. To address this concern, the concept of differential privacy is widely adopted [DMNS06, Dwo08, DR14], which is the possibility of a system to leak information under an adversarial attack. A popular approach to protect privacy is adding a noise to the model or gradients, so that the algorithm will converge to an “inexact” solution. [ACG⁺16, WYX17, INS⁺19, FKT20, CWH20] apply this method to design differentially private optimization algorithms for the single server setting, while [HDJ⁺20, ACCÖ21, NBD22, DLC⁺22] consider differential privacy for server/client distributed setting. However, merely protecting the privacy of all agents against external adversary is insufficient in the decentralized setting. It is also necessary to protect each agent’s privacy from leaking to other agents, which leads to the development of locally differentially private server/client and decentralized optimization algorithms [DJW13, ABCP13, CABP13, XYD19, CSU⁺19, WXDX20, ZZY⁺21, LZLC22, ZCH⁺22, ZKL18, ZKL20, CEBM22, ZP23, MS23].

Gradient clipping Gradient clipping has gained significant attention in recent years. Earlier works e.g. [PMB13, BBP13, KLL16, KFI17, YGG17], use gradient clipping as a pure heuristic to solve gradient vanishing/exploding problems without theoretical understandings. Then, [ZHSJ20, ZKV⁺20, ZJFW20, RLDJ23] introduce theoretical analysis trying to understand the impact on the convergence rate when an algorithm adopts gradient clipping. With the recent advancements in differentially private optimization algorithms, a series of works, e.g. [CWH20, ZCH⁺22, DKX⁺22, ?, ?], apply this technique to limit the magnitude of gradients, so that the variance of privacy perturbation can be decided without the bounded gradient assumption.

1.4 Notation

Throughout this thesis, we use lowercase and uppercase boldface letters to represent vectors and matrices, respectively. We use $\|\cdot\|_{\text{op}}$ for matrix operator norm, $\|\cdot\|_{\text{F}}$ for Frobenius norm, $\|\cdot\|_2$ for 2-norm, \otimes for the Kronecker product, \odot for the element-wise multiplication, \mathbf{I}_n for the n -dimensional identity matrix, $\langle \mathbf{x}, \mathbf{y} \rangle$ for the inner product of two vectors \mathbf{x} and \mathbf{y} , $\mathbf{1}_n$ for the n -dimensional all-one vector and $\mathbf{0}_{d \times n}$ for the $(d \times n)$ -dimensional zero matrix. For two real functions $f(\cdot)$ and $g(\cdot)$ defined on \mathbb{R}^+ , we say $f(x) = O(g(x))$ or $f(x) \lesssim g(x)$ if there exists some universal constant $M > 0$ such that $f(x) \leq Mg(x)$. The notation $f(x) = \Omega(g(x))$ or $f(x) \gtrsim g(x)$ means $g(x) = O(f(x))$.

Let $\mathbf{x}_i \in \mathbb{R}^d$ be the optimization variable at agent i . We define the matrix of all optimization variables

and the average vector as

$$\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n} \quad \bar{\mathbf{x}} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad (1.5)$$

In addition, for a matrix of variables, we introduce the distributed gradient $\nabla F(\mathbf{X}) \in \mathbb{R}^{d \times n}$ as

$$\nabla F(\mathbf{X}) := [\nabla f_1(\mathbf{x}_1), \nabla f_2(\mathbf{x}_2), \dots, \nabla f_n(\mathbf{x}_n)] \in \mathbb{R}^{d \times n}, \quad (1.6)$$

and the global gradient of the matrix $\nabla f(\mathbf{X}) \in \mathbb{R}^{d \times n}$ as

$$\nabla f(\mathbf{x}) := [\nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_n)]. \quad (1.7)$$

1.5 Thesis organization

The rest of this thesis is organized as following: Chapters 2 to 5 develop and analyze Network-DANE, DESTRESS, BEER and PORTER, respectively, with extensive numerical experiments; Chapter 6 concludes this thesis and proposes future directions.

Chapter 2

Decentralized Newton-style algorithm

There is growing interest in large-scale machine learning and optimization over decentralized networks, e.g., in the context of multi-agent learning and federated learning. Due to the imminent need to alleviate the communication burden, the investigation of communication-efficient distributed optimization algorithms — particularly for empirical risk minimization — has flourished in recent years. A large fraction of these algorithms has been developed for the server/client setting, relying on the presence of a central parameter server that can communicate with all agents.

This chapter decentralized optimization problems defined by (1.2), where objective functions don't have finite-sum structure, and each agent is only allowed to aggregate information from its neighbors over a network (namely, no centralized coordination is present). By properly adjusting the global gradient estimate via local averaging in conjunction with proper correction, we develop a communication-efficient approximate Newton-type method, called *Network-DANE*, which generalizes DANE to accommodate decentralized scenarios. Our key ideas can be applied, in a systematic manner, to obtain decentralized versions of other server/client distributed algorithms. A notable example is the development of *Network-SVRG*, which employs variance reduction at each agent to further accelerate local computation. We establish linear convergence of *Network-DANE* and *Network-SVRG* for strongly convex losses, which shed light on the impacts of data homogeneity, network connectivity, and local averaging upon the rate of convergence. We further extend *Network-DANE* to composite optimization by allowing a nonsmooth penalty term. Numerical evidence is provided to demonstrate the appealing performance of our algorithms over competitive baselines, in terms of both communication and computation efficiency. Our work suggests that by performing a judiciously chosen amount of local communication and computation per iteration, the overall efficiency can be substantially improved.

This chapter is based on our previous publication [[LCCC20](#)].

2.1 Preliminaries

This section introduces two important concepts that are crucial to `Network-DANE` and `Network-SVRG`, and the DANE algorithm that inspires the development of `Network-DANE`.

2.1.1 Dynamic average consensus

Assume that each agent generates some *time-varying* quantity $r_j^{(t)}$ (e.g., the current local parameter or gradient estimates). Let $\mathbf{r}^{(t)} = [r_1^{(t)}, \dots, r_n^{(t)}]^\top$. To track the dynamic average $\frac{1}{n} \sum_{j=1}^n r_j^{(t)} = \frac{1}{n} \mathbf{1}_n^\top \mathbf{r}^{(t)}$ at each agent, [ZM10] proposes a simple tracking algorithm: suppose each agent maintains an estimate $q_j^{(t)}$ in the t -th iteration, and the network collectively adopts the following update rule

$$\mathbf{q}^{(t)} = \mathbf{W}\mathbf{q}^{(t-1)} + \mathbf{r}^{(t)} - \mathbf{r}^{(t-1)}, \quad (2.1)$$

where $\mathbf{q}^{(t)} = [q_1^{(t)}, \dots, q_n^{(t)}]^\top$. The first term $\mathbf{W}\mathbf{q}^{(t-1)}$ represents the standard local information mixing operation (meaning that each agent updates its own estimate by a weighted average of its neighbors' estimates), the second term $\mathbf{r}^{(t)} - \mathbf{r}^{(t-1)}$ tracks the temporal difference. A crucial property of (2.1) is

$$\mathbf{1}_n^\top \mathbf{q}^{(t)} = \mathbf{1}_n^\top \mathbf{r}^{(t)}, \quad (2.2)$$

which indicates that the average of $\{q_i^{(t)}\}_{1 \leq i \leq n}$ dynamically tracks the average of $\{r_i^{(t)}\}_{1 \leq i \leq n}$. We shall adapt this procedure in our algorithmic development, in the hope of reliably tracking the global gradients (i.e., the average of the local, and often time-varying, gradients at all agents).

2.1.2 Chebyshev's acceleration

Performing one round of mixing for $\mathbf{x} \in \mathbb{R}^{nd}$ using mixing matrix \mathbf{W} , the resulting vector at agent i is $\tilde{x}_i = \sum_j w_{ij} x_j$. If the communication network is not well-connected, we may need to perform $K > 1$ rounds of communications to improve the communication quality. If implemented as repeatedly mixing the mixed variable, the result will be equivalent as mixing using \mathbf{W}^K , and the resulting vector at agent i is $\tilde{x}_i = \sum_j (\mathbf{W}^K)_{ij} x_j$.

However, the simple implementation is not optimal, as we can construct a polynomial of \mathbf{W} to further minimize the mixing rate of the resulting mixing matrix. Let $P_k(\mathbf{W}) = c_k \mathbf{W}^k + c_{k-1} \mathbf{W}^{k-1} + \dots + c_1 \mathbf{W} + c_0$. Optimize the mixing rate of $P_k(\mathbf{W})$ leads to the solution

$$P_k(x) = 1 - \frac{T_k(x/\alpha_0)}{T_k(1/\alpha_0)},$$

where $T_k(x)$ is the Chebyshev polynomial defined by

$$T_0(x) = 1$$

$$\begin{aligned} T_1(x) &= x \\ T_{k+1}(x) &= 2xT_k(x) - T_{k-1}(x). \end{aligned}$$

When α is close to 1, the mixing rate of $P_K(W)$ is approximately $\alpha_{\text{cheby}} \sim 1 - \sqrt{2(1 - \alpha_0)}$.

2.1.3 Additional notation

We define the following (nd) -dimensional stacked vectors

$$\mathbf{x}^{(t)} := [\mathbf{x}_1^{(t)\top}, \dots, \mathbf{x}_n^{(t)\top}]^\top, \quad \mathbf{y}^{(t)} := [\mathbf{y}_1^{(t)\top}, \dots, \mathbf{y}_n^{(t)\top}]^\top, \quad \mathbf{s}^{(t)} := [\mathbf{s}_1^{(t)\top}, \dots, \mathbf{s}_n^{(t)\top}]^\top. \quad (2.3)$$

With a slight abuse of notation, we introduce the stacked distributed gradient $\nabla F(\mathbf{x}) \in \mathbb{R}^{nd}$ and the stacked global gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^{nd}$ of an (nd) -dimensional vector \mathbf{x} as follows:

$$\nabla F(\mathbf{x}) := [\nabla f_1(\mathbf{x}_1)^\top, \dots, \nabla f_n(\mathbf{x}_n)^\top]^\top, \quad \nabla f(\mathbf{x}) := [\nabla f(\mathbf{x}_1)^\top, \dots, \nabla f(\mathbf{x}_n)^\top]^\top. \quad (2.4)$$

2.1.4 The DANE algorithm

The DANE algorithm is a popular communication-efficient approximate Newton method developed for the server/client model [SSZ14]. Here, we review some key features of DANE. (i) Each agent performs an update using both the local loss function $f_j(\cdot)$ and the gradient $\nabla f(\cdot)$ of the global loss function (obtained via the parameter server). (ii) In the t -th iteration, the j -th agent solves the following problem to update its local estimate $\mathbf{x}_j^{(t)}$:

$$\mathbf{x}_j^{(t)} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f_j(\mathbf{x}) - \left\langle \nabla f_j(\bar{\mathbf{x}}^{(t)}) - \nabla f(\bar{\mathbf{x}}^{(t)}), \mathbf{x} \right\rangle + \frac{\mu}{2} \|\mathbf{x} - \bar{\mathbf{x}}^{(t)}\|_2^2 \right\}, \quad (2.5)$$

where $\mu \geq 0$ is the regularization parameter.¹ Implementing this algorithm requires two rounds of communications per iteration.

- (a) The parameter server first collects all local estimates $\{\mathbf{x}_j^{(t-1)}\}_{1 \leq j \leq n}$ and computes the average global parameter estimate $\bar{\mathbf{x}}^{(t)} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^{(t-1)}$; this is then sent back to all agents.
- (b) The parameter server collects all local gradients evaluated at the point $\bar{\mathbf{x}}^{(t)}$, computes the global gradient $\nabla f(\bar{\mathbf{x}}^{(t)}) = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\bar{\mathbf{x}}^{(t)})$, and shares it with all agents.

The DANE algorithm has been demonstrated as a competitive baseline whose communication efficiency improves, in some sense, with the increase of data size [SSZ14]; see [FGW21] for its proximal variation and

¹In [SSZ14], the second term in (2.5) takes the form $\nabla f_j(\bar{\mathbf{x}}^{(t)}) - \bar{\eta} \nabla f(\bar{\mathbf{x}}^{(t)})$. We set $\bar{\eta} = 1$ without loss of generality following the analysis in [FGW21].

improved theoretical analysis. To see the reason why DANE is an approximate Newton-type algorithm, consider the case when the local loss functions in all agents are quadratic and takes the form (2.10). The local optimization subproblem (2.5) in DANE can be solved in closed form, with $\mathbf{x}_j^{(t)}$ given by²

$$\mathbf{x}_j^{(t)} = \bar{\mathbf{x}}^{(t)} - \underbrace{(\mathbf{H}_j + \mu \mathbf{I}_d)}_{\text{local Hessian}}^{-1} \nabla f(\bar{\mathbf{x}}^{(t)}). \quad (2.6)$$

Clearly, this can be interpreted as

$$\mathbf{x}_j^{(t)} = \text{local parameter estimate} - (\text{local Hessian})^{-1} (\text{global gradient}),$$

which is an approximate Newton-type update rule (since we invoke the local Hessian to approximate the true global Hessian). It is worth noting that the algorithm proceeds without communicating the local Hessians.

2.2 The Network-DANE algorithm

The DANE algorithm was originally developed for the server/client setting. In the network setting, however, agents can no longer compute (2.5) locally, due to the absence of centralization enabled by the parameter server; more specifically, agents have access to neither $\bar{\mathbf{x}}^{(t)}$ nor $\nabla f(\bar{\mathbf{x}}^{(t)})$, both of which are required when solving (2.5). To address this lack of global information, one might naturally wonder whether we can simply replace global averaging by local averaging; that is, replacing $\bar{\mathbf{x}}^{(t)}$ and $\nabla f(\bar{\mathbf{x}}^{(t)})$ by $\frac{1}{|\mathcal{N}_j|} \sum_{i \in \mathcal{N}_j} \mathbf{x}_i^{(t-1)}$ and $\frac{1}{|\mathcal{N}_j|} \sum_{i \in \mathcal{N}_j} \nabla f_i(\mathbf{x}_i^{(t-1)})$, respectively, in the j -th agent. However, this simple idea fails to guarantee convergence in local agents. For instance, the local estimation errors may stay flat (but nonvanishing) — as opposed to converging to zero — as the iterations progress, primarily due to imperfect information sharing.

With this convergence issue in mind, our key idea is composed of the following components.

- The first ingredient is to maintain an additional estimate of the global gradient in each agent — denoted by $\mathbf{s}_j^{(t)}$ in the j -th agent. This additional gradient estimate is updated via dynamic average consensus (2.8), in the hope of tracking the global gradient evaluated at $\mathbf{y}_j^{(t)}$ in the j -th agent ($1 \leq j \leq n$), i.e., $\mathbf{s}_j^{(t)}$ attempts to track $\nabla f(\mathbf{y}_j^{(t)})$. Here, $\mathbf{y}_j^{(t)}$ stands for the parameter estimate obtained by local neighborly averaging in the t -th iteration (see Algorithm 1 for details). As the algorithm converges, $\{\mathbf{y}_j^{(t)}\}_{1 \leq j \leq n}$ is expected to reach consensus, allowing $\mathbf{s}_j^{(t)}$ ($1 \leq j \leq n$) to converge to the true global gradient as well.

²See [SSZ14] for a short derivation.

Algorithm 1 Network-DANE

```

1 input: initial parameter estimates  $\mathbf{x}_j^{(0)} \in \mathbb{R}^d$  ( $1 \leq j \leq n$ ), regularization parameter  $\mu$ .
2 initialization: set  $\mathbf{y}_j^{(0)} = \mathbf{x}_j^{(0)}$ ,  $\mathbf{s}_j^{(0)} = \nabla f_j(\mathbf{y}_j^{(0)})$  for all agents  $1 \leq j \leq n$ .
3 for  $t = 1, 2, \dots$  do
4   for Agents  $1 \leq j \leq n$  in parallel do
5     Set  $\mathbf{y}_j^{(t),0} = \mathbf{x}_j^{(t-1)}$  and  $\mathbf{s}_j^{(t),0} = \mathbf{s}_j^{(t-1)}$ .
6     for  $k = 1, 2, \dots, K$  do
7       Receive information  $\mathbf{y}_i^{(t),k-1}$  and  $\mathbf{s}_i^{(t),k-1}$  from its neighbors  $i \in \mathcal{N}_j$ .
8       Aggregate parameter estimates from neighbors:
9       
$$\mathbf{y}_j^{(t),k} = \sum_{i \in \mathcal{N}_j} w_{ji} \mathbf{y}_i^{(t),k-1}, \quad \mathbf{s}_j^{(t),k} = \sum_{i \in \mathcal{N}_j} w_{ji} \mathbf{s}_i^{(t),k-1} \quad (2.7)$$

10      Set the local parameter estimate to  $\mathbf{y}_j^{(t)} = \mathbf{y}_j^{(t),K}$ .
11      Update the global gradient estimate by aggregated local information and gradient tracking:
12      
$$\mathbf{s}_j^{(t)} = \mathbf{s}_j^{(t),K} + \underbrace{\nabla f_j(\mathbf{y}_j^{(t)}) - \nabla f_j(\mathbf{y}_j^{(t-1)})}_{\text{gradient tracking}}. \quad (2.8)$$

13      Update the parameter estimate by solving:
14      
$$\mathbf{x}_j^{(t)} = \underset{z \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f_j(z) - \langle \nabla f_j(\mathbf{y}_j^{(t)}) - \mathbf{s}_j^{(t)}, z \rangle + \frac{\mu}{2} \|z - \mathbf{y}_j^{(t)}\|_2^2 \right\}. \quad (2.9)$$


```

- In addition, we also allow multiple rounds of mixing within each iteration, i.e., (2.7), which is helpful in accelerating convergence when the network exhibits a high degree of locality. In essence, by applying K rounds of mixing, we improve the mixing rate from α to α^K . As we shall see later, choosing a proper (but not too large) K suffices to achieve the desired trade-off between the rate of information sharing and iteration complexity, which helps reduce the overall communication and computation cost. This step of extra averaging can be implemented in an efficient manner via the Chebyshev acceleration scheme [AS14, SBB⁺17].

Armed with such improved global gradient estimates, we propose to solve a modified local optimization subproblem (2.9) in Network-DANE, which approximates the original Newton-type problem (2.5) by replacing $\nabla f(\bar{\mathbf{x}}^{(t)})$ with the local surrogate $\mathbf{s}_j^{(t)}$. The proposed local subproblem (2.9) is convex and can be solved efficiently via, say, Nesterov's accelerated gradient methods. The whole algorithm is presented in

Algorithm 1.

Remark 1. *It is certainly possible to employ more general mixing matrices in (2.7). For instance, in mobile computing scenarios with moving agents, one might prefer using time-varying mixing matrices in order to accommodate the topology changes over time. We omit such extensions for brevity.*

2.3 Convergence guarantees

2.3.1 Assumptions, metrics and parameters

This section formally introduces additional assumptions, key parameters, and error metrics required for convergence analysis of Network-DANE.

To begin with, we introduce Assumptions 3 and 4 that characterize local objective functions.

Assumption 3 (strongly convex and smooth local objective function). *The local objective function $f_i(\mathbf{x})$ at each agent is strongly convex and smooth, namely, $\sigma \mathbf{I} \preceq \nabla^2 f_j(\mathbf{x}) \preceq L \mathbf{I}$ ($1 \leq j \leq n$) for some quantities $0 < \sigma \leq L$. Define $\kappa = L/\sigma$ is the condition number.*

Assumption 4 (quadratic local objective function). *The local objective function $f_i(\mathbf{x})$ at each agent is quadratic w.r.t. \mathbf{x} , i.e., taking the form of*

$$f_i(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{H}_i \mathbf{x} + \mathbf{b}_i^\top \mathbf{x} + c_i, \quad (2.10)$$

where $\mathbf{b}_i \in \mathbb{R}^d$, $c_i \in \mathbb{R}$, and $\mathbf{H}_i = \nabla^2 f_i(\mathbf{x}) \in \mathbb{R}^{d \times d}$ is a fixed symmetric and positive semidefinite matrix.

Next, we define the homogeneity parameter following [CZC⁺20, FGW21].

Definition 4 (Homogeneity parameter). *Let $f(\mathbf{x})$ and $f_j(\mathbf{x})$ be as defined in (1.2). The homogeneity parameter β is defined as*

$$\beta := \max_{1 \leq j \leq n} \beta_j \quad \text{with} \quad \beta_j := \sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla^2 f_j(\mathbf{x}) - \nabla^2 f(\mathbf{x})\|. \quad (2.11)$$

As it turns out, β is bounded by the smoothness parameter of $f(\mathbf{x})$, i.e., $\beta \leq L$.³ On the other end, as the local loss functions f_j 's become similar with each other, β will become smaller. Therefore, β is a key quantity measuring the similarity of data across agents.

³To prove this inequality, we note from the minimax theorem of eigenvalues and the triangle inequality that

$$\beta \leq \max_j \left\{ \sup_{\substack{\mathbf{x} \in \mathbb{R}^d \\ \|\mathbf{v}\|_2=1}} \mathbf{v}^\top \left(\frac{n-1}{n} \nabla^2 f_j(\mathbf{x}) \right) \mathbf{v} - \inf_{\substack{\mathbf{x} \in \mathbb{R}^d \\ \|\mathbf{v}\|_2=1}} \mathbf{v}^\top \left(\frac{1}{n} \sum_{i:i \neq j} \nabla^2 f_i(\mathbf{x}) \right) \mathbf{v} \right\} = \left(1 - \frac{1}{n}\right)(L - \sigma) \leq L.$$

Remark 2. *If the local data follow certain statistical models, it is possible to show that β decreases as the local data size m grows. For example, [SSZ14] shows that if the data samples at all agents are i.i.d. and Assumption 2 holds, then with probability at least $1 - \delta$ over the samples, we have $\beta < \sqrt{\frac{32L^2}{m} \log \frac{nd}{\delta}}$ – implying β decreases at the rate of $1/\sqrt{m}$.*

In addition to notation defined in Section 1.4, we define an extra (nd) -dimensional vector \mathbf{y} analogous to (1.5).

To characterize the convergence behavior of our algorithm, we need to simultaneously track several interrelated error metrics as follows:

- (1) the convergence error: $\|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2$;
- (2) the parameter consensus error: $\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2$;
- (3) the gradient estimation error: $\|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \nabla f(\mathbf{y}^{(t)})\|_2$.

In this thesis, an algorithm is said to converge linearly at a rate $\rho \in (0, 1)$ if there exists some constant $C > 0$ such that the following holds for all $t \geq 1$:

$$\max \left\{ \sqrt{n} \|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2, \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2, L^{-1} \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 \right\} \leq C\rho^t.$$

In addition, an algorithm is said to reach ϵ -accuracy if the left-hand side of the above expression is bounded by ϵ .

2.3.2 Convergence guarantees of Network-DANE for quadratic loss

This section establishes linear convergence of Network-DANE when the objective functions are quadratic.

Theorem 1 (Network-DANE under quadratic loss, arbitrary K). *Suppose that Assumptions 1, 3 and 4 hold. Set $K > 0$, the effective mixing rate becomes α^K . Set μ large enough so that $\sigma + \mu \geq \frac{140L}{(1-\alpha^K)^2} \left(\frac{\beta}{\sigma} + 1\right)$. Then Network-DANE converges linearly at a rate ρ_1 obeying*

$$\rho_1 := \max \left\{ \frac{1 + \theta_1}{2}, \alpha^K + \frac{140\kappa}{1 - \alpha^K} \left(\frac{\sigma + \beta}{\sigma + \mu} \right), \frac{1 + \alpha^K}{2} + \frac{2\beta}{\sigma + \mu} \right\}, \quad (2.12)$$

where θ_1 is defined by

$$\theta_1 := 1 - \frac{\sigma}{\sigma + \mu} + \frac{L}{L + \mu} \frac{\beta^2}{(\sigma + \mu)(\sigma + \mu - \beta)}. \quad (2.13)$$

Remark 3. *It turns out that $\theta_1 \in (0, 1)$ is the convergence rate of DANE in the server/client setting under quadratic losses [SSZ14, Theorem 1].*

It is worth noting that we do not optimized constants in the above theorem, as our primary focus is the order of convergence rate. If the regularization parameter μ is sufficiently large, one can guarantee that $\theta_1 < 1$ and hence DANE converges at a linear rate when optimizing quadratic losses [SSZ14]. We can clearly see that (2.12) is always greater than θ_1 , which is the price we pay for consensus under the network setting. Fortunately, by properly setting μ , we can still guarantee that $\rho_1 < 1$, which in turn enables linear convergence of Network-DANE.

In view of (2.12), if the network is sufficiently connected (i.e., α is small), or if the data are sufficiently homogeneous (i.e., β is small), we can use a smaller parameter μ , which makes θ_1 (defined in (2.13)) smaller and results in faster convergence. In summary, Network-DANE takes fewer iterations to converge when α and β are both small. After some basic calculations, the complexity of Network-DANE for quadratic losses is formalized in the following corollary.

Corollary 1. Set $\mu + \sigma = \frac{180L}{(1-\alpha^K)^2} \left(\frac{\beta}{\sigma} + 1\right)$. Under the assumptions of Theorem 1, one has

$$\rho_1 \leq 1 - \left(\frac{1-\alpha^K}{20}\right)^2 \frac{1}{\kappa} \frac{1}{(\beta/\sigma + 1)}. \quad (2.14)$$

To reach ϵ -accuracy, Network-DANE takes at most $O\left(\frac{\kappa(\beta/\sigma+1)\log(1/\epsilon)}{(1-\alpha^K)^2}\right)$ iterations, and $O\left(K \cdot \frac{\kappa(\beta/\sigma+1)\log(1/\epsilon)}{(1-\alpha^K)^2}\right)$ communication rounds.

If we set the number of local averaging rounds to be $K = 1$, then the iteration complexity can be directly compared with other existing results. If the homogeneous parameter β obeys $\beta = O(\sigma)$, then the convergence rate can be improved to $O(\kappa \log(1/\epsilon)/(1-\alpha)^2)$; this is much faster than the corrected DGD [QL18] with gradient tracking, which converges in $O(\kappa^2 \log(1/\epsilon)/(1-\alpha)^2)$ iterations. The convergence rate of Network-DANE degenerates to that of DGD [QL18] with gradient tracking under the worst condition $\beta = \Theta(L)$. This observation highlights the communication efficiency of Network-DANE by harnessing the homogeneity of data across different agents. We emphasize that this is an important feature of our analysis, where the convergence rate adapts with respect to the data homogeneity.

Benefits of extra local averaging (i.e., $K > 1$). The careful reader might have noticed that the rate established above scales poorly with respect to the network parameter, namely, $(1-\alpha)^{-1}$, when $K = 1$. One remedy is to consider the case with $K > 1$, where Network-DANE performs K rounds of communications per iteration. On the one hand, the effective network parameter becomes α^K that can be made arbitrarily small by setting K sufficiently large, thus leading to faster convergence; on the other hand, the total number of communications is K times larger than the number of iterations, meaning that we might end up with a higher communication complexity. As an example, invoking Corollary 1, the total communication cost to

reach ϵ -accuracy is given by

$$O(K \cdot \kappa(1 + \beta/\sigma) \log(1/\epsilon)/(1 - \alpha^K)^2).$$

Therefore, by judiciously choosing K , it is possible to significantly improve the overall communication complexity, especially when α is close to 1. For example, by setting $K \asymp 1/\log(1/\alpha) = O(1/(1 - \alpha))$, we can ensure $\alpha^K \asymp 1/2$ and reduce the communication complexity to $O(\kappa \cdot (\beta/\sigma + 1) \log(1/\epsilon)/(1 - \alpha))$, thus improving the dependence with the graph topology.

The following theorem shows an improved result following a refined analysis, which improves the dependence simultaneously with respect to both κ and $(1 - \alpha)^{-1}$.

Theorem 2 (Network-DANE under quadratic loss, optimized K). *Instate the assumptions of Theorem 1. Set K and μ large enough so that $\alpha^K \leq 1/(2\kappa)$ and $\sigma + \mu \geq 360\sigma \left(\frac{\beta^2}{\sigma^2} + 1\right)$. To reach ϵ -accuracy, Network-DANE takes at most $O((\beta^2/\sigma^2 + 1) \log(1/\epsilon))$ iterations, and $O\left(\log \kappa \cdot \frac{(\beta^2/\sigma^2 + 1) \log(1/\epsilon)}{1 - \alpha}\right)$ communications rounds.*

When we set K as suggested in Theorem 2, the iteration complexity becomes independent of the network topology. Moreover, it matches the rate of DANE in the server/client setting [SSZ14] when $\beta = O(\sigma)$, which is $O(\log(1/\epsilon))$ and further independent of the condition number κ .

In terms of network dependence, the communication complexity improves from $O(1/(1 - \alpha)^2)$ to $O(1/(1 - \alpha))$. By implementing the extra averaging step in an efficient manner via the well-known Chebyshev acceleration scheme [AS14, SBB⁺17], the dependence of the communication complexity with respect to $(1 - \alpha)^{-1}$ can be further improved to $O((1 - \alpha)^{-1/2})$. The final communication complexity of Network-DANE for quadratic losses thus becomes

$$O\left(\log \kappa \cdot \frac{(\beta^2/\sigma^2 + 1) \log(1/\epsilon)}{(1 - \alpha)^{1/2}}\right).$$

Therefore, the total amount of communication is significantly reduced using extra averaging, where it scales only logarithmically with respect to κ .

2.3.3 Convergence guarantees of Network-DANE for strongly convex loss

This section establishes the linear convergence of Network-DANE for general smooth and strongly convex loss functions, where the rate is worse than that for quadratic losses.

Theorem 3. *Assume Assumptions 1 and 3 hold. Set $K > 0$, and μ large enough so that $\sigma + \mu \geq \frac{170\kappa L}{(1 - \alpha^K)^2}$. Then Network-DANE converges linearly at a rate ρ_2 obeying*

$$\rho_2 := \max \left\{ \frac{1 + \theta_2}{2}, \alpha^K + \frac{170\kappa}{1 - \alpha^K} \left(\frac{L}{\sigma + \mu} \right), \frac{1 + \alpha^K}{2} + \frac{2\beta}{\sigma + \mu} \right\}, \quad (2.15)$$

where θ_2 is given by

$$\theta_2 := 1 - \frac{\sigma}{\sigma + \mu} + \frac{\beta}{\sigma + \mu} \sqrt{1 - \left(\frac{\mu}{\sigma + \mu}\right)^2}. \quad (2.16)$$

Remark 4. Note that $\theta_2 \in (0, 1)$ is precisely the convergence rate of DANE in the server/client setting [FGW21, Theorem 3.1].

Similar to Theorem 1, one can guarantee $\theta_2 < 1$ and $\rho_2 < 1$ by setting the regularization parameter μ sufficiently large. Therefore, Network-DANE can converge at a linear rate for a general class of smooth and strongly convex problems. Comparing the convergence rates of Network-DANE derived for the above two different losses (i.e., comparing (2.13) with (2.16)), we see that: when the loss functions are non-quadratic, θ_2 is generally greater than θ_1 ⁴. This happens since the Hessian matrices associated with the non-quadratic loss functions may vary across different points, which is also the reason why the convergence rate of Network-DANE derived for the general case degenerates to the worst-case rate. After some basic calculations, the complexity of Network-DANE under strongly convex losses is formalized by the following corollary.

Corollary 2. Set $\sigma + \mu = \frac{180\kappa L}{(1-\alpha^K)^2}$. Under the assumptions of Theorem 3, one has

$$\rho_2 \leq 1 - \left(\frac{1-\alpha^K}{20}\right)^2 \frac{1}{\kappa^2}. \quad (2.17)$$

To reach ϵ -accuracy, Network-DANE takes at most $O\left(\frac{\kappa^2 \log(1/\epsilon)}{(1-\alpha^K)^2}\right)$ iterations and $O\left(K \cdot \frac{\kappa^2 \log(1/\epsilon)}{(1-\alpha^K)^2}\right)$ communication rounds.

When $K = 1$, the communication complexity of Network-DANE is $O\left(\frac{\kappa^2 \log(1/\epsilon)}{(1-\alpha)^2}\right)$, which is rather pessimistic and does not improve with data homogeneity. Similar to Theorem 2, we can improve this by optimizing K properly. We have the following theorem, which is parallel to Theorem 2.

Theorem 4 (Network-DANE under strongly convex loss, optimized K). *Instate the assumptions of Theorem 3. Set K and μ large enough so that $\alpha^K \leq 1/(2\kappa)$ and $\sigma + \mu \geq 360L\left(\frac{\beta}{\sigma} + 1\right)$. To reach ϵ -accuracy, Network-DANE takes at most $O\left(\kappa(\beta/\sigma + 1) \log(1/\epsilon)\right)$ iterations and $O\left(\log \kappa \cdot \frac{\kappa(\beta/\sigma + 1) \log(1/\epsilon)}{1-\alpha}\right)$ communication rounds.*

The improved rate in Theorem 4 improves as the local data become more homogeneous, recovering a feature that has been highlighted previously. Similar to earlier discussions in Section 2.1.2, by using the Chebyshev acceleration scheme [AS14, SBB⁺17], the final communication complexity of Network-DANE for strongly convex losses becomes

$$O\left(\log \kappa \cdot \frac{\kappa(\beta/\sigma + 1) \log(1/\epsilon)}{(1-\alpha)^{1/2}}\right).$$

⁴This is because $\sqrt{\frac{\sigma^2 + 2\sigma\mu}{(\sigma + \mu)^2}} \geq \frac{\sigma}{\sigma + \mu}$.

Remark 5. The homogeneity parameter β defined in Definition 4 measures the largest deviation of local Hessians from the global Hessian. A refined analysis using local deviation β_j is possible by permitting different regularization parameters μ_j in (2.9) for different agents.

2.4 Extension to nonsmooth composite optimization

Network-DANE can be extended for nonsmooth composite optimization, by properly adjusting the local optimization step, leveraging proximal variants of DANE [FGW21] and SVRG [XZ14]. For simplicity, we present the proximal variant of Network-DANE and leave its theoretical analysis to future work.

Consider the following regularized empirical risk minimization problem:

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) + g(x) \triangleq \frac{1}{N} \sum_{i=1}^N \ell(x; z_i) + g(x), \quad (2.18)$$

where $f(\cdot)$ and $f_j(\cdot)$ are defined as in (1.2), and $g(\cdot)$ is a deterministic convex regularizer that can be nonsmooth. This type of problem has wide applications, where it is desirable to promote additional structures or incorporate prior knowledge about the solution through adding a deterministic regularization term $g(x)$. We can extend Network-DANE to solve (2.18) by adding the proximal term into the local optimization step, as detailed in Algorithm 2, which is a direct extension of Algorithm 1. Section 2.6 numerically verifies the effectiveness of Algorithm 2.

Algorithm 2 Network-DANE for nonsmooth composite optimization

- 1 Replace the local optimization sub-problem (2.9) of Network-DANE by the following:
- 2 **Input:** $\mathbf{y}_j^{(t)}$, $\mathbf{s}_j^{(t)}$, regularization parameter μ .
- 3 Update the parameter estimate by solving:

$$\mathbf{x}_j^{(t)} = \underset{z \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f_j(z) + g(z) - \langle \nabla f_j(\mathbf{y}_j^{(t)}) - \mathbf{s}_j^{(t)}, z \rangle + \frac{\mu}{2} \|z - \mathbf{y}_j^{(t)}\|_2^2 \right\}. \quad (2.19)$$

2.5 Extension with variance reduction

The design of Network-DANE suggests a systematic approach to obtain decentralized versions of other algorithms. We illustrate this by reducing local computation of Network-DANE using variance reduction. Stochastic variance reduction methods are a popular class of stochastic optimization algorithms, developed to allow for constant step sizes and faster convergence in finite-sum optimization [JZ13, XZ14, NLST17]. It is therefore natural to ask whether such variance reduction techniques can be leveraged in a network setting to further save local computation without compromising communication.

Algorithm 3 Network-SVRG

1 Replace the local optimization subproblem (2.9) of Network-DANE by the following:

2 **Input:** $\mathbf{y}_j^{(t)}$, $\mathbf{s}_j^{(t)}$, step size δ , number of local iterations S .

3 **Initialization:** set $\mathbf{u}_j^{(t),0} = \mathbf{y}_j^{(t)}$, $\mathbf{v}_j^{(t),0} = \mathbf{s}_j^{(t)}$.

4 **for** $s = 1, \dots, S$ **do**

5 $\mathbf{u}_j^{(t),s} = \mathbf{u}_j^{(t),s-1} - \delta \mathbf{v}_j^{(t),s-1}$.

6 Sample \mathbf{z} from \mathcal{M}_j uniformly at random, then,

$$\mathbf{v}_j^{(t),s} = \nabla \ell(\mathbf{u}_j^{(t),s}; \mathbf{z}) - \nabla \ell(\mathbf{u}_j^{(t),0}; \mathbf{z}) + \mathbf{v}_j^{(t),0}. \quad (2.20)$$

7 Choose the new parameter estimate $\mathbf{x}_j^{(t)}$ from $\{\mathbf{u}_j^{(t),1}, \dots, \mathbf{u}_j^{(t),S}\}$ uniformly at random.

Inspired by the connection between DANE and SVRG [KMR15], we introduce Network-SVRG in Algorithm 3, a decentralized version of SVRG [JZ13] tailored to the network setting, with the assistance of gradient tracking. In particular, the inner loops of SVRG [JZ13] are adopted to replace the local computation subproblem (2.9) of Network-DANE, where the reference to the global gradient is replaced by $\mathbf{s}_j^{(t)}$ to calculate the variance-reduced stochastic gradient.

The convergence analysis of Algorithm 3 is more challenging due to the biased stochastic gradient involved in each local iteration. Theorem 5 establishes the linear convergence of Network-SVRG for strongly convex losses, as long as β is sufficiently small and the number of mixing rounds K is sufficiently large. Again, we have not strived to improve the pre-constants specified in the theorem.

Theorem 5. *Assume Assumptions 1 and 3 hold and $\beta/\sigma \leq 1/200$. Set K large enough such that $\alpha^K \simeq 1/\kappa$ and S large enough, Network-SVRG converges linearly. To reach ε -accuracy, Network-SVRG takes at most $O(\log(1/\varepsilon))$ iterations and $O\left(\log \kappa \cdot \frac{\log(1/\varepsilon)}{1-\alpha}\right)$ communication rounds.*

The proof of Theorem 5 can be found in Appendix A.4. Theorem 5 implies that: as long as the local data are sufficiently similar (so that β does not exceed the order of σ), by performing $O(\log \kappa / (1 - \alpha))$ rounds of local communication per iteration, Network-SVRG converges in $O(\log(1/\varepsilon))$ iterations independent of κ . This performance guarantee matches its counterpart in the server/client setting [CZC+20]. Altogether, Network-SVRG achieves appealing computation and communication complexities simultaneously. By further adopting the Chebyshev acceleration scheme [AS14, SBB⁺17], the final communication complexity of Network-SVRG is at most

$$O\left(\log \kappa \cdot \frac{\log(1/\varepsilon)}{(1-\alpha)^{1/2}}\right).$$

It is straightforward to extend this idea to obtain decentralized variants of other stochastic variance

reduced algorithms such as Katyusha [AZ17], by replacing the local computation step (2.9) by the inner loop update rules of the stochastic methods of interest. For the sake of brevity, this paper does not pursue such “plug-and-play” extensions.

Remark 6. *Our convergence theory of Network-SVRG requires $\beta \lesssim \sigma$, which is consistent with its counterpart in the server/client setting [CZC⁺20]. In contrast, Network-DANE is guaranteed to converge linearly in the entire range of β by setting μ sufficiently large. One scheme to relax this requirement, as analyzed in [CZC⁺20], is to add a regularization term, similar to the last term in (2.9), that penalizes the distance to the previous estimate. However, this might come at a price of slower convergence. We leave this to future investigation.*

2.6 Numerical experiments

We evaluate the performance of the proposed algorithms⁵ for solving both strongly convex and nonconvex problems, in order to demonstrate the appealing performance in terms of communication-computation trade-offs.

Throughout this section, we set the number of agents $n = 20$. We use symmetric fastest distributed linear averaging (FDLA) matrices [XB04] generated according to the communication graph as the mixing matrix W for aggregating $x_j^{(t)}$ in (2.7). For aggregating $s_j^{(t)}$ in (2.7), we use a convex combination of I and W such that its diagonal elements are greater than 0.1, which makes the algorithm more stable in practice. The same regularization parameter μ is used for DANE and Network-DANE. We generate connected random communication graphs using an Erdős-Rényi graph with the probability of connectivity $p = 0.3$ (if not specified). For each experiment, we use the same random starting point $x^{(0)}$ and mixing matrix W for all algorithms. To solve the local optimization subproblems, we use Nesterov’s accelerated gradient descent for at most 100 iterations for DANE and Network-DANE.

2.6.1 Experiments on synthetic data

We conduct five synthetic numerical experiments based on linear regression to investigate the performance of our algorithms. The same data generation method is used for all synthetic experiments. We generate $m = 1000$ samples of dimension $d = 40$, denoted by A_i , randomly from $\mathcal{N}(\mathbf{0}, \Sigma)$ i.i.d. for each agent, where Σ is a diagonal matrix with $\Sigma_{ii} = i^{-\varrho}$. By changing ϱ , we can change the condition number κ . Data samples are generated according to the linear model $\mathbf{b}_i = A_i x_0 + \xi_i$, with a random signal x_0 and i.i.d. noise $\xi_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For DANE and Network-DANE, we set $\mu = 5 \times 10^{-10}$ when $\kappa = 10$ and $\mu = 5 \times 10^{-4}$ when

⁵In our experiments of Network-SVRG, we use the last iterate $u_j^{(t),S}$ as the new parameter estimate locally, which is more practical; though our analysis only handles the case where the new parameter estimate is selected uniformly at random from previous iterates.

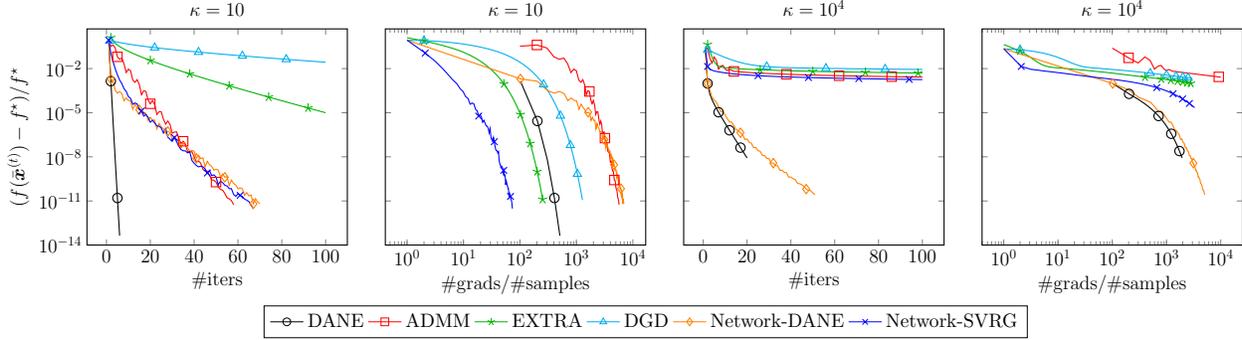


Figure 2.1: The relative optimality gap with respect to the number of iterations and gradient evaluations under different conditioning $\kappa = 10$ (left two panels) and $\kappa = 10^4$ (right two panels) for linear regression.

$\kappa = 10^4$. For *Network-SVRG*, we set the step size $\delta = 0.1/(L + \sigma + 2\mu)$ and the number of local iterations $S = 0.05m$.

Comparison with existing algorithms. To make a fair comparison with other algorithms, no extra local averaging is adopted in this experiment, i.e., the number of mixing rounds is set to $K = 1$. The loss function at each agent is given as $f_i(x) = \frac{1}{2m} \|A_i x - b_i\|_2^2$. We plot the relative optimality gap, given as $(f(\bar{x}^{(t)}) - f^*)/f^*$, where $\bar{x}^{(t)}$ is the average parameter of all agents at the t -th iteration, and f^* is the optimal value. We compare the proposed *Network-DANE* (cf. Algorithm 1) and *Network-SVRG* (cf. Algorithm 3) with the server/client algorithms *DANE* [SSZ14] and *ADMM* [BPC+11],⁶ and two popular network-distributed gradient descent algorithms, referred to as *DGD* [QL18] and *EXTRA* [SLWY15a].

Figure 2.1 shows the relative optimality gap with respect to the number of iterations as well as the number of gradient evaluations under different condition numbers $\kappa = 10$ and $\kappa = 10^4$ for linear regression. In both experiments, *Network-DANE* and *Network-SVRG* outperform *DGD* and *EXTRA* in terms of the numbers of communication rounds. *Network-SVRG* has similar communication rounds with *ADMM* but only communicates locally. *Network-DANE* is quite insensitive to the condition number, performing almost as well as the *DANE* algorithm in the ill-conditioned case, but operates in a fully decentralized setting. *Network-SVRG* further outperforms other algorithms in terms of gradient evaluations in most settings, especially for well-conditioned cases.

Benefits of extra local mixing (communication) per iteration. We conduct synthetic experiments to investigate the communication-computation trade-off observed in Theorem 4 when employing multiple rounds of mixing within every iteration. Following the suggestion of the theory, we use a poorly connected

⁶We apply *ADMM* to the constrained optimization problem, which amounts to the centrally-distributed setting, $\min_{x_i} \frac{1}{n} \sum f_i(x_i)$ s.t. $x_i = x$. Note that *ADMM* can also be applied to the network-distributed setting, which is not shown here since our network algorithms already outperform *ADMM* in the centrally-distributed setting.

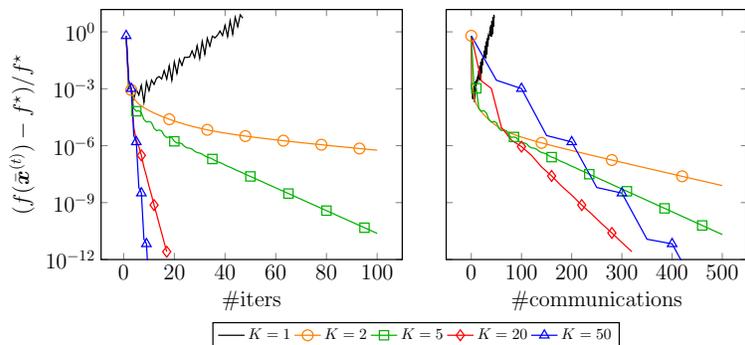


Figure 2.2: The relative optimality gap with respect to the number of iterations and communication rounds under different rounds of mixing K for Network-DANE over a poorly connected graph.

network with mixing rate $\alpha_0 = 0.944$ for communication, which is generated by an Erdős-Rényi graph with $p = 0.2$. For illustration, we consider the relative optimality gap for a linear regression problem with $\kappa = 10$, with respect to the number of iterations and communication rounds for Network-DANE and Network-SVRG, under different values of K (no Chebyshev acceleration is employed), shown in Figure 2.2. Due to poor connectivity, Network-DANE and Network-SVRG fail to converge when using moderate parameters. However, by using a larger K , due to improvement in consensus, Network-DANE converge faster in terms of the number of iterations. Notice that after certain threshold, further increasing K will not improve the convergence rate in terms of communication rounds.

Effects of local computation for Network-SVRG. We conduct an experiment to analyze the effect of different numbers of local stochastic iterations for Network-SVRG. Throughout this experiment, we run our algorithms on a linear regression problem with $\kappa = 10$ and Erdős-Rényi graph ($p = 0.2$) as the communication graph. Figure 2.3 shows the number of communication rounds and the number of gradient evaluations till converge for different numbers of local iterations. It is clear that with too few local iterations, Network-SVRG converges very slow and requires more communication. As soon as S is above a threshold, which is around $0.05m$ local iterations, the communication rounds no longer decrease. Therefore, in our experiments, we set the number of local iterations as $S = 0.05m$ to ensure satisfactory convergence rate while using an economical amount of local computation.

Effects of network topology. We conduct another experiment to compare the effect of network topology on linear regression problem with $\kappa = 10$. We generate communication graphs with different topology settings. Figure 2.4 shows the relative optimality gap with respect to the number of iterations and gradient evaluations for Network-DANE for Erdős-Rényi graph ($p = 0.3$), a 4×5 grid graph, a star graph, and a ring graph. The performance degrades as the network becomes less connected (where $1 - \alpha_0$ becomes small)

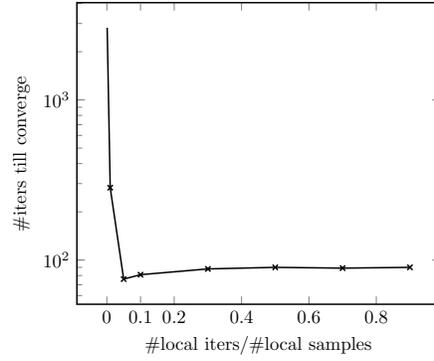


Figure 2.3: Number of communication rounds till converge with respect to different numbers of local iterations for Network-SVRG.

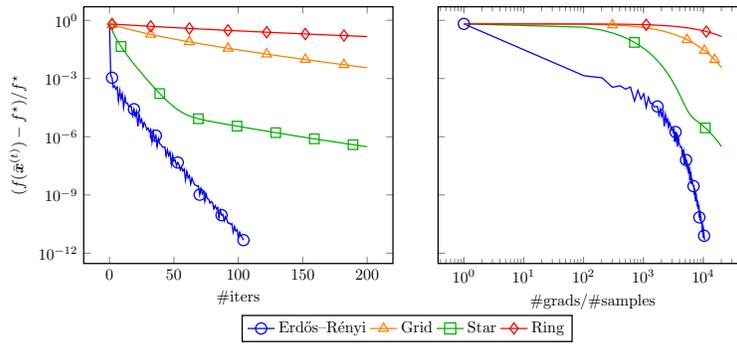


Figure 2.4: Performance of Network-DANE under different network topologies.

[NOR18].

Experiments for nonsmooth composite optimization We consider the ℓ_1 -norm regularized linear regression, where the loss function of each agent is given as $\tilde{f}_i(\mathbf{x}) = f_i(\mathbf{x}) + g(\mathbf{x}) = \frac{1}{2m} \|\mathbf{A}_i \mathbf{x} - \mathbf{b}_i\|_2^2 + 0.01 \|\mathbf{x}\|_1$, and the communication graph are generated in the same way as Figure 2.1. The condition number κ is also defined in the same way as earlier. We compare the performance of Network-DANE with CEASE [FGW21], which is the proximal version of DANE in the server/client setting, ADMM, and PG-EXTRA, which is the proximal version of EXTRA [SLWY15b]. For CEASE and Network-DANE, we set $\mu = 10^{-4}$ when $\kappa = 10$ and $\mu = 10^{-1}$ when $\kappa = 10^4$, and use FISTA [BT09] to solve the ℓ_1 -norm regularized local problems for computation efficiency. Figure 2.5 plots the relative optimality gap $\|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^*\|_2 / \|\mathbf{x}^*\|_2$ with respect to the number of iterations and the number of gradient evaluations for different algorithms under different condition numbers. In both experiments, Network-DANE outperformed ADMM and PG-EXTRA in both metrics, and achieves similar convergence behavior as CEASE, though at a slower rate due to optimizing over a decentralized topology.

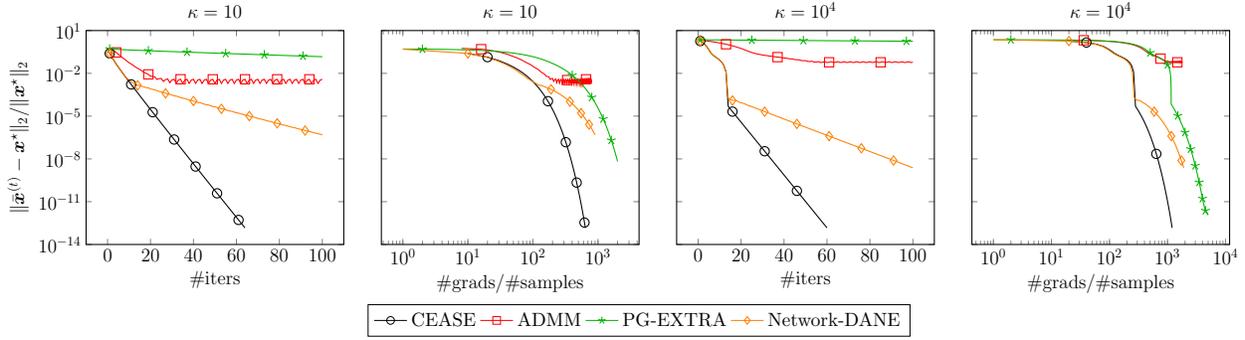


Figure 2.5: The relative optimality gap with respect to the number of iterations and gradient evaluations under different conditioning $\kappa = 10$ (left two panels) and $\kappa = 10^4$ (right two panels) for linear regression with ℓ_1 -norm regularization.

2.6.2 Experiments on real data

We perform two experiments on real data to further evaluate the performance of the proposed algorithms for both convex and nonconvex problems.

Binary classification using logistic regression. We use regularized logistic regression to solve a binary classification problem using the Gisette dataset.⁷ We split the Gisette dataset to $n = 20$ agents, where each agent receives $m = 300$ training samples of dimension $d = 5000$. The loss function at each agent is given as

$$f_i(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \log \left(1 + (2b_i^{(j)} - 1) \exp(\mathbf{x}^\top \mathbf{a}_i^{(j)}) \right) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2,$$

where $\mathbf{a}_i^{(j)} \in \mathbb{R}^d$ and $b_i^{(j)} \in \{0, 1\}$ are samples stored at agent i . For DANE and Network-DANE, we set $\mu = 5 \times 10^{-9}$ when $\kappa = 2$ and $\mu = 5 \times 10^{-1}$ when $\kappa = 100$. The condition number is controlled by changing the regularization λ . Figure 2.6 shows the results. In both cases, our algorithms exhibit compelling performance over other decentralized optimization algorithms especially in terms of communication efficiency.

1-hidden-layer neural network training. Though our theory only applies to the strongly convex case, we examine Network-SVRG in the nonconvex case, by training a one-hidden-layer neural network with 64 hidden neurons and sigmoid activations for a classification task using the MNIST dataset. We split 60,000 training samples to 20 agents and use an Erdős-Rényi graph with $p = 0.3$ for communications. Figure 2.7 plots the training loss and testing accuracy against the number of iterations and gradient evaluations for different algorithms, where centralized ADMM and decentralized stochastic algorithm (DSGD) are plotted as baselines. Being more communication-efficient than DSGD, and more computation-efficient than ADMM, Network-SVRG reach a desirable balance between computation and communication efficacies.

⁷The dataset can be found at <https://archive.ics.uci.edu/ml/datasets/Gisette>.

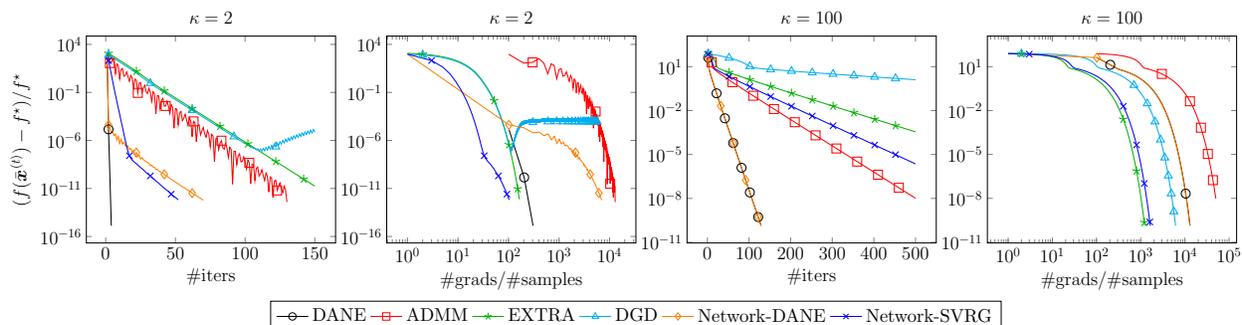


Figure 2.6: The relative optimality gap with respect to the number of iterations and gradient evaluations under different conditioning $\kappa = 2$ (left two panels) and $\kappa = 100$ (right two panels) for logistic regression using the Gisette dataset.

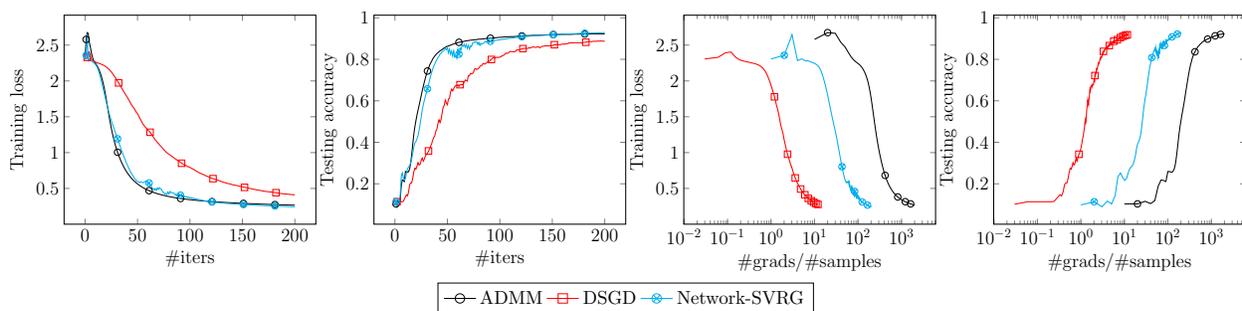


Figure 2.7: The training loss and testing accuracy with respect to the number of iterations (left two panels) and gradient evaluations (right two panels) for different algorithms on the MNIST dataset.

Chapter 3

Decentralized stochastic recursive gradient algorithm

Emerging applications in multi-agent environments such as internet-of-things, networked sensing, autonomous systems and federated learning, call for decentralized algorithms for finite-sum optimizations that are resource-efficient in terms of both computation and communication. In this chapter, we consider the prototypical setting where the agents work collaboratively to minimize the sum of local loss functions (cf. (1.3)) by only communicating with their neighbors over a predetermined network topology. We develop a new algorithm, called DEcentralized STochastic REcurSive gradient methodS (DESTRESS) for nonconvex finite-sum optimization, which matches the optimal incremental first-order oracle (IFO) complexity of centralized algorithms for finding first-order stationary points, while maintaining communication efficiency. Detailed theoretical and numerical comparisons corroborate that the resource efficiencies of DESTRESS improve upon prior decentralized algorithms over a wide range of parameter regimes. DESTRESS leverages several key algorithm design ideas including randomly activated stochastic recursive gradient updates with mini-batches for local computation, gradient tracking with extra mixing (i.e., multiple gossiping rounds) for per-iteration communication, together with careful choices of hyper-parameters and new analysis frameworks to provably achieve a desirable computation-communication trade-off.

This chapter is based on our previous publication [LLC22].

3.1 The DESTRESS algorithm

We propose DESTRESS for finding first-order order stationary points of nonconvex finite-sum problems. Throughout, we define (nd) -dimensional stacked vectors $\mathbf{u}^{(t)}$, $\mathbf{v}^{(t)}$, $\mathbf{g}^{(t)}$ and $\mathbf{s}^{(t)}$ analogously to Section 2.1.3. Motivated by stochastic recursive gradient methods in the centralized setting, DESTRESS has a nested loop

structure:

1. The inner loop refines the parameter estimate $\mathbf{u}^{(t),0} = \mathbf{x}^{(t-1)}$ by performing randomly activated stochastic recursive gradient updates (3.1), where the stochastic recursive gradient $\mathbf{v}^{(t),s}$ is updated in (3.1b) and (3.1c) via mixing mini-batch stochastic gradients from activated agents' local datasets.
2. The outer loop adopts dynamic average consensus to estimate and track the global gradient $\nabla F(\mathbf{x}^{(t)})$ at each agent by $\mathbf{s}^{(t)}$ in (3.2), which allows the next inner loop to start from a less noisy starting gradient $\mathbf{v}^{(t+1),0} = \mathbf{s}^{(t)}$. A key property of (3.2)—which is a direct consequence of dynamic average consensus—is that the average of $\mathbf{s}^{(t)}$ equals to the dynamic average of local gradients, i.e., $\bar{\mathbf{s}}^{(t)} = \frac{1}{n} \sum_{i \in [n]} \mathbf{s}_i^{(t)} = \frac{1}{n} \sum_{i \in [n]} \nabla f_i(\mathbf{x}_i^{(t)})$.

To enable better information sharing and faster convergence, inspired by [LCCC20], we allow DESTRESS to perform a few rounds of mixing or gossiping whenever communication takes place. Specifically, DESTRESS performs K_{out} and K_{in} mixing steps for the outer and inner loops respectively per iteration, which is equivalent to using

$$\mathbf{W}_{\text{out}} = \mathbf{W}^{K_{\text{out}}} \quad \text{and} \quad \mathbf{W}_{\text{in}} = \mathbf{W}^{K_{\text{in}}}$$

as mixing matrices, and correspondingly a network with better connectivity; see (3.2), (3.1a) and (3.1c). Note that Algorithm 4 is written in matrix notation, where the mixing steps are described by $\mathbf{W}_{\text{in}} \otimes \mathbf{I}_n$ or $\mathbf{W}_{\text{out}} \otimes \mathbf{I}_n$ and applied to all agents simultaneously. The extra mixing steps can be implemented by Chebyshev acceleration [AS14] with improved communication efficiency.

Compared with existing decentralized algorithms based on stochastic variance-reduced algorithms such as D-GET [SLH20] and GT-SARAH [XKK22a], DESTRESS utilizes different gradient estimators and communication protocols: First, DESTRESS produces a sequence of reference points $\{\mathbf{x}^{(t)}\}$ that converge to a global first-order stationary point and corresponding global gradient estimates $\{\mathbf{s}^{(t)}\}$ that are updated by full gradient computations, so that inner loops can refine $\mathbf{x}^{(t)}$ using stochastic recursive gradients based on accurate gradient estimates; second, the communication and computation in DESTRESS are paced differently due to the introduction of extra mixing, which allow a more flexible trade-off schemes between different types of resources; last but not least, the random activation of stochastic recursive gradient updates further saves local computation, especially when the local sample size is small compared to the number of agents.

¹The stochastic gradients will not be computed if $\lambda_i^{(t),s} = 0$.

Algorithm 4 DESTRESS for decentralized nonconvex finite-sum optimization

-
- 1 **input:** initial parameter $\bar{\mathbf{x}}^{(0)}$, step size η , activation probability p , batch size b , number of outer loops T , number of inner loops S , and number of communication (extra mixing) steps K_{in} and K_{out} .
 - 2 **initialization:** set $\mathbf{x}_i^{(0)} = \bar{\mathbf{x}}^{(0)}$ and $\mathbf{s}_i^{(0)} = \nabla f(\bar{\mathbf{x}}^{(0)})$ for all agents $1 \leq i \leq n$.
 - 3 **for** $t = 1, \dots, T$ **do**
 - 4 Set inner loop initial parameters $\mathbf{u}^{(t),0} = \mathbf{x}^{(t-1)}$ and $\mathbf{v}^{(t),0} = \mathbf{s}^{(t-1)}$.
 - 5 **for** $s = 1, \dots, S$ **do**
 - 6 Each agent i samples a mini-batch $\mathcal{Z}_i^{(t),s}$ of size b from \mathcal{M}_i uniformly at random, $\lambda_i^{(t),s} \sim \mathcal{B}(p)$ where $\mathcal{B}(p)$ denotes the Bernoulli distribution with parameter p ,¹ and then performs the following updates:

$$\mathbf{u}^{(t),s} = (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d)(\mathbf{u}^{(t),s-1} - \eta \mathbf{v}^{(t),s-1}), \quad (3.1a)$$

$$\mathbf{g}_i^{(t),s} = \frac{\lambda_i^{(t),s}}{pb} \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t),s}} \left(\nabla \ell(\mathbf{u}_i^{(t),s}; \mathbf{z}_i) - \nabla \ell(\mathbf{u}_i^{(t),s-1}; \mathbf{z}_i) \right) + \mathbf{v}_i^{(t),s-1}, \quad (3.1b)$$

$$\mathbf{v}^{(t),s} = (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) \mathbf{g}^{(t),s}. \quad (3.1c)$$
 - 7 Set the new parameter estimate $\mathbf{x}^{(t)} = \mathbf{u}^{(t),S}$.
 - 8 Update the global gradient estimate by aggregated local information and gradient tracking:

$$\mathbf{s}^{(t)} = (\mathbf{W}_{\text{out}} \otimes \mathbf{I}_d) \left(\mathbf{s}^{(t-1)} + \nabla F(\mathbf{x}^{(t)}) - \nabla F(\mathbf{x}^{(t-1)}) \right) \quad (3.2)$$
-
- 9 **output:** $\mathbf{x}^{\text{out}} \sim \text{Uniform}(\{\mathbf{u}_i^{(t),s-1} | i \in [n], t \in [T], s \in [S]\})$.
-

3.2 Convergence guarantees

This section presents the performance guarantees of DESTRESS. Due to the nonconvexity, first-order algorithms are generally guaranteed to converge to only first-order stationary points of the global loss function $f(\cdot)$, which is defined in Definition 3.

Theorem 6 shows that DESTRESS converges in expectation to an approximate first-order stationary point, under suitable parameter choices.

Theorem 6 (First-order optimality). *Assume Assumptions 1 and 2 hold. Set $p \in (0, 1]$, K_{in} , K_{out} , S , b and η to be positive and satisfy*

$$\alpha^{K_{\text{in}}} \leq p \quad \text{and} \quad \eta L \leq \frac{(1 - \alpha^{K_{\text{in}}})^3 (1 - \alpha^{K_{\text{out}}})}{10(1 + \alpha^{K_{\text{in}}} \alpha^{K_{\text{out}}} \sqrt{npb}) (\sqrt{S/(npb)} + 1)}. \quad (3.3)$$

The output produced by Algorithm 4 satisfies

$$\mathbb{E} \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 < \frac{4}{\eta TS} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^* \right). \quad (3.4)$$

If there is only one agent, i.e., $n = 1$, the mixing rate will be $\alpha = 0$, we can choose $K_{\text{in}} = K_{\text{out}} = p = 1$, and Theorem 6 reduces to [NvP⁺22, Theorem 1], its counterpart in the centralized setting. For

general decentralized settings with arbitrary mixing schedules, Theorem 6 provides a comprehensive characterization of the convergence rate, where an ϵ -first-order stationary point can be found in expectation in a total of

$$TS = O\left(\frac{\mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^*}{\eta\epsilon^2}\right)$$

iterations, where T is the number of outer iterations and S is the number of inner iterations. Clearly, a larger step size η , as allowable by (3.3), hints on a smaller iteration complexity, and hence a smaller IFO complexity.

There are two conditions in (3.3). On one end, K_{in} needs to be large enough (i.e., perform more rounds of extra mixing) to counter the effect when p is small (i.e., we compute less stochastic gradients every iteration), or when α is close to 1 (i.e., the network is poorly connected). On the other end, the step size η needs to be small enough to account for the requirement of the step size in the centralized setting, as well as the effect of imperfect communication due to decentralization. For well connected networks where $\alpha \ll 1$, the terms introduced by the decentralized setting will diminish—indicating the iteration complexity is close to that of the centralized setting. For poorly connected networks, carefully designing the mixing matrix and other parameters can ensure a desirable trade-off between convergence speed and communication cost. The following corollary provides specific parameter choices for DESTRESS to achieve the optimal per-agent IFO complexity.

Corollary 3 (Complexity for finding first-order stationary points). *Under conditions of Theorem 6, set $S = \lceil \sqrt{mn} \rceil$, $b = \lceil \sqrt{m/n} \rceil$, $p = \frac{\sqrt{m/n}}{\lceil \sqrt{m/n} \rceil}$, $K_{\text{out}} = \left\lceil \frac{\log(\sqrt{npb+1})}{(1-\alpha)^{1/2}} \right\rceil$, $K_{\text{in}} = \left\lceil \frac{\log(2/p)}{(1-\alpha)^{1/2}} \right\rceil$, $\eta = \frac{1}{640L}$, and implement the mixing steps using Chebyshev’s acceleration [AS14]. To reach an ϵ -first-order stationary point, in expectation, DESTRESS takes $O\left(m + \frac{(m/n)^{1/2}L}{\epsilon^2}\right)$ IFO calls per agent, and $O\left(\frac{\log\left(\frac{(n/m)^{1/2}+2}{(1-\alpha)^{1/2}}\right) \cdot ((mn)^{1/2} + \frac{L}{\epsilon^2})}{(1-\alpha)^{1/2}}\right)$ rounds of communication.*

DESTRESS achieves a network-independent IFO complexity that matches the optimal complexity in the centralized setting. In addition, when the accuracy $\epsilon^2 \lesssim L/(mn)^{1/2}$, DESTRESS reaches a communication complexity of $O\left(\frac{1}{(1-\alpha)^{1/2}} \cdot \frac{L}{\epsilon^2}\right)$, which is independent of the sample size.

It is worthwhile to further highlight the role of the random activation probability p in achieving the optimal IFO by allowing “fractional” batch size. Note that the batch size is set as $b = \lceil \sqrt{m/n} \rceil$, where m is the local sample size, and n is the number of agents.

1. When the local sample size is large, i.e., $m \geq n$, we can approximate $b \approx \sqrt{m/n}$ and $p \approx 1$. In fact, Corollary 3 continues to hold with $p = 1$ in this regime.

	Erdős-Rényi graph	2-D grid graph	Path graph
$1 - \alpha$ (spectral gap)	1	$\frac{1}{n \log n}$	$\frac{1}{n^2}$
D-GET [SLH20]	$m + \frac{m^{1/2}L}{\epsilon^2}$	$m + \frac{m^{1/2}n^2L}{\epsilon^2}$	$m + \frac{m^{1/2}n^4L}{\epsilon^2}$
GT-SARAH [XKK22a]	$m + \max \left\{ 1, \left(\frac{m}{n}\right)^{1/3}, \left(\frac{m}{n}\right)^{1/2} \right\} \cdot \frac{L}{\epsilon^2}$	$m + \max \left\{ n^2, m^{1/3}n^{2/3}, \left(\frac{m}{n}\right)^{1/2} \right\} \cdot \frac{L}{\epsilon^2}$	$m + \max \left\{ n^4, m^{1/3}n^{5/3}, \left(\frac{m}{n}\right)^{1/2} \right\} \cdot \frac{L}{\epsilon^2}$
DESTRESS (Algorithm 4)	$(mn)^{1/2} + \frac{L}{\epsilon^2}$	$m^{1/2}n + \frac{n^{1/2}L}{\epsilon^2}$	$(mn^3)^{1/2} + \frac{nL}{\epsilon^2}$
Improvement factors for ϵ -independent term	$\left(\frac{m}{n}\right)^{1/2}$	$\frac{m^{1/2}}{n}$	$\frac{m^{1/2}}{n^{3/2}}$
Improvement factors for ϵ -dependent term	$\max \left\{ 1, \left(\frac{m}{n}\right)^{1/3}, \left(\frac{m}{n}\right)^{1/2} \right\}$	$\max \left\{ n^{3/2}, m^{1/3}n^{1/6}, \frac{m^{1/2}}{n} \right\}$	$\max \left\{ n^3, m^{1/3}n^{2/3}, \frac{m^{1/2}}{n^{3/2}} \right\}$

Table 3.1: Detailed comparisons of the communication complexities of D-GET, GT-SARAH and DESTRESS under three graph topologies, where the last two rows delineate the improve factors of DESTRESS over existing algorithms. The communication savings become significant especially when $m = \Omega\left(\frac{n}{1-\alpha}\right)$. The complexities are simplified by plugging the bound on the spectral gap $1 - \alpha$ from [NOR18, Proposition 5]. m, n, L are defined in Section 1.1 and α is the mixing rate defined in (1.4). The big- O (defined in Section 1.4) notation and logarithmic terms are omitted for simplicity.

2. However, when the number of agents is large, i.e., $n > m$, the batch size $b = 1$ and $p = \sqrt{m/n} < 1$, which mitigates the potential computation waste by only selecting a subset of agents to perform local computation, compared to the case when we naively set $p = 1$.

Therefore, by introducing random activation, we can view $pb = \sqrt{m/n}$ as the effective batch size at each agent, which allows fractional values and leads to the optimal IFO complexity in all scenarios.

To gain further insights in terms of the communication savings of DESTRESS, Table 3.1 further compares the communication complexities of decentralized algorithms for finding first-order stationary points under three common network settings, which highlights the communication improvement over existing works.

3.3 Numerical experiments

This section provides numerical experiments on real datasets to evaluate our proposed algorithm DESTRESS with comparisons against two existing baselines: DSGD [NO09, LZZ⁺17] and GT-SARAH [XKK22a].

For all experiments, we shuffle the datasets and normalize the samples by subtracting the mean and dividing the standard deviation. We use the same number of agents, FDLA matrices, communication graphs and Chebyshev’s acceleration scheme as in Section 2.6. In addition, since $m \gg n$ in all experiments, we set $p = 1$ for simplicity. To ensure convergence, DSGD adopts a diminishing step size schedule. All parameters are tuned manually for the best performance. We defer detailed descriptions of baseline algorithms to Appendix B.1.

3.3.1 Logistic regression with nonconvex regularization

To begin with, we employ logistic regression with nonconvex regularization to solve a binary classification problem using the Gisette dataset.² We split the Gisette dataset to $n = 20$ agents, where each agent receives $m = 300$ training samples. The sample loss function is given as

$$\ell(\mathbf{x}; \{\mathbf{f}, l\}) = \log \left(1 + (2l - 1) \exp(\mathbf{x}^\top \mathbf{f}) \right) + \lambda \sum_{i=1}^d \frac{x_i^2}{1 + x_i^2},$$

where $\{\mathbf{f}, l\}$ represents a training tuple, $\mathbf{f} \in \mathbb{R}^d$ is the feature vector and $l \in \{0, 1\}$ is the label, and λ is the regularization parameter. For this experiment, we set $\lambda = 0.1$. Table 3.2 specifies parameter setting we use for each graph.

Algorithms	DSGD		DESTRESS						GT-SARAH		
Parameters	η_0	b	η	p	K_{in}	K_{out}	b	S	η	b	S
Erdős-Rényi	1	10	0.01	1	2	2	10	10	0.001	10	10
Grid	1	10	0.01	1	2	3	10	10	0.001	10	10
Path	0.1	10	0.01	1	8	8	10	10	0.0001	10	10

Table 3.2: Parameter settings for the experiments on regularized logistic regression in Figure 3.1.

Figure 3.1 shows the train gradient norm and testing accuracy for all algorithms. DESTRESS significantly outperforms other algorithms both in terms of communication and computation. It is worth noting that, DSGD converges very fast at the beginning of training, but cannot sustain the progress due to the diminishing schedule of step sizes. On the contrary, the variance-reduced algorithms can converge with a constant step size, and hence converge better overall. Moreover, due to the refined gradient estimation and information mixing designs, DESTRESS can bear a larger step size than GT-SARAH, which leads to the fastest convergence and best overall performance. In addition, a larger number of extra mixing steps leads to a better performance when the communication graph is less connected.

3.3.2 One-hidden-layer neural network training

Next, we compare the performance of DESTRESS to DSGD and GT-SARAH for training a one-hidden-layer neural network with 64 hidden neurons and sigmoid activations for classifying the MNIST dataset [Den12]. We evenly split the MNIST dataset to $n = 20$ agents, where each agent receives $m = 3,000$ training samples. Table 3.3 specifies parameter setting we use for each graph.

Figure 3.2 plots the training gradient norm and testing accuracy against the number of communication rounds and gradient evaluations for all algorithms. DESTRESS significantly outperforms GT-SARAH in terms of computation and communication costs due to the larger step size and extra mixing. Different

²The dataset can be accessed at <https://archive.ics.uci.edu/ml/datasets/Gisette>.

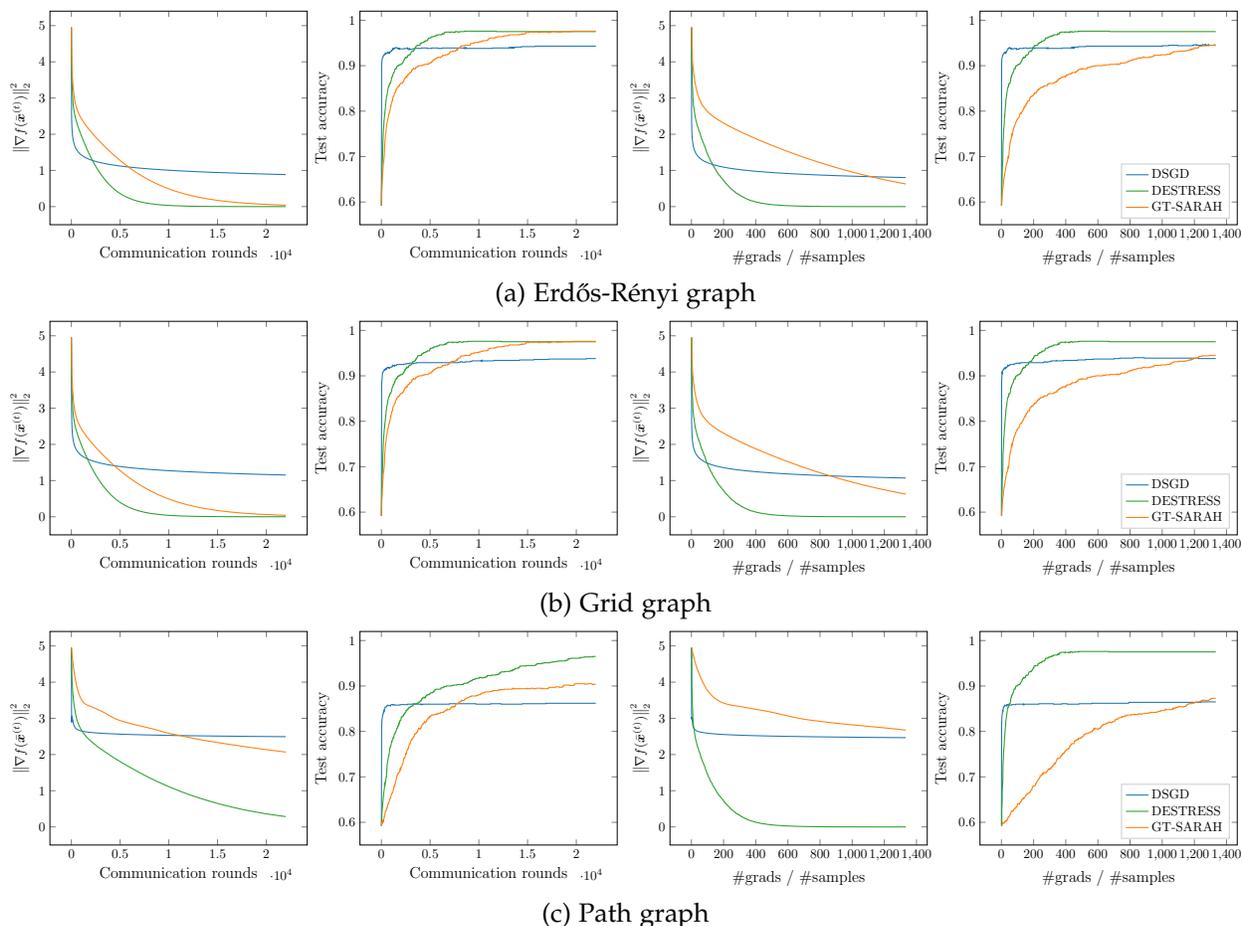


Figure 3.1: The train gradient norm and testing accuracy with respect to the number of communication rounds (left two panels) and gradient evaluations (right two panels) for DSGD, GT-SARAH and DESTRESS when training logistic regression model with nonconvex regularization on the Gisetite dataset. Due to the initial full-gradient computation, the gradient evaluations of DESTRESS and GT-SARAH do not start from 0.

Algorithms	DSGD		DESTRESS					GT-SARAH			
Parameters	η_0	b	η	p	K_{in}	K_{out}	b	S	η	b	S
Erdős-Rényi	1	100	1	1	2	2	100	10	0.1	100	10
Grid	1	100	1	1	3	4	100	10	0.1	100	10
Path	0.1	100	1	1	8	10	100	10	0.0001	100	10

Table 3.3: Parameter settings for the experiments on neural network training in Figure 3.2.

from the previous experiment, DSGD performs the best for Erdős-Rényi graph and grid graph that are well connected, while converges slower than DESTRESS on path graph.

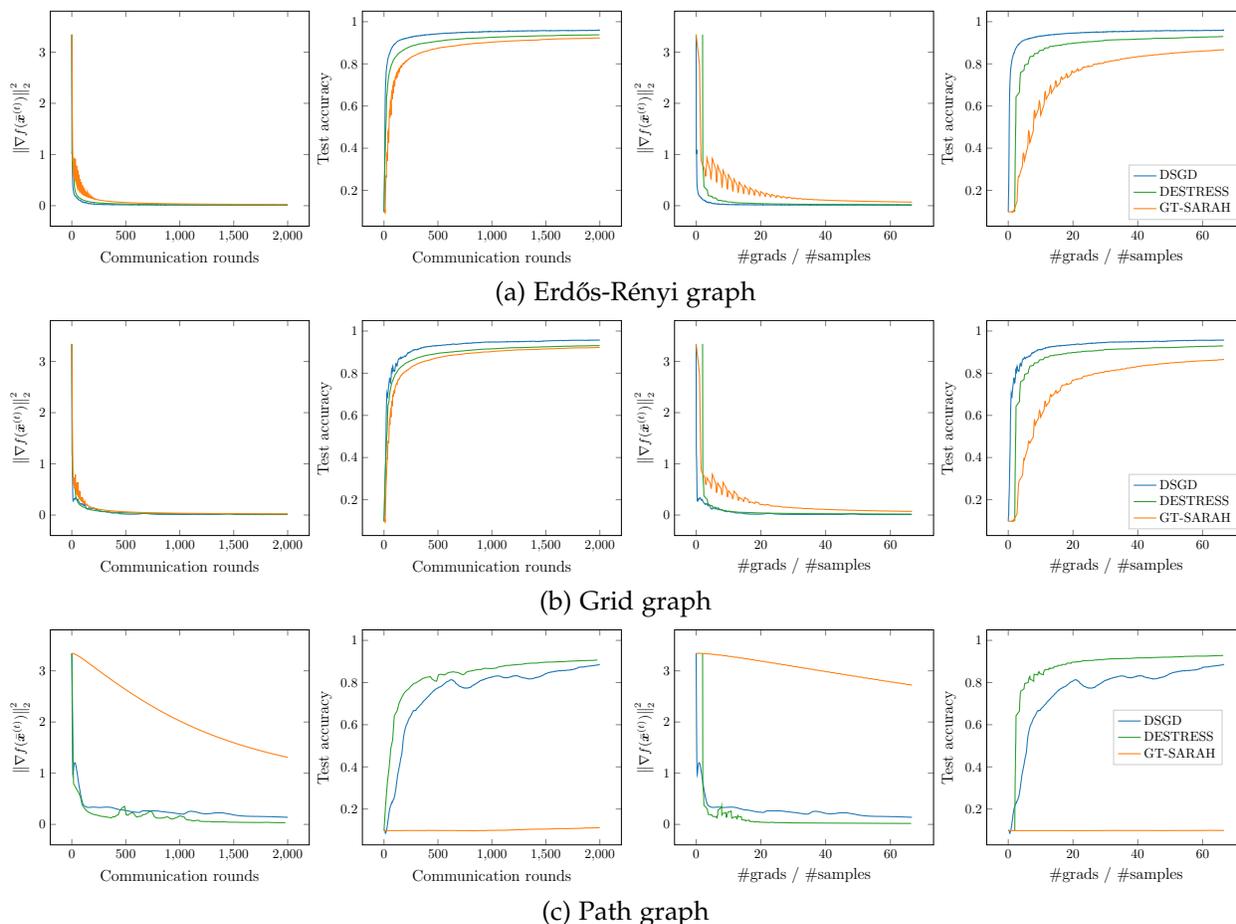


Figure 3.2: The train gradient norm and testing accuracy with respect to the number of communication rounds (left two panels) and gradient evaluations (right two panels) for DSGD, GT-SARAH and DESTRESS when training a one-hidden-layer neural network on the MNIST dataset. Due to the initial full-gradient computation, the gradient evaluations of DESTRESS and GT-SARAH do not start from 0.

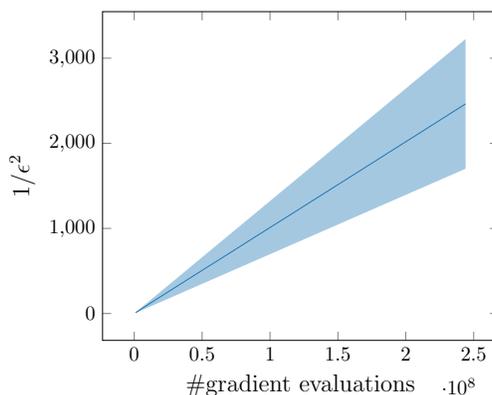


Figure 3.3: The convergence precision $1/\epsilon^2$ with respect to the number of total gradient evaluations for neural network training averaged over 64 experiments. The shade shows the variance.

3.3.3 Convergence and gradient computations

The last experiment investigates the convergence precision $1/\epsilon^2$ of DESTRESS with respect to the number of gradient evaluations. Under the same experimental setup, we conduct 64 different runs where each run starts from a different initial point. The convergence precision is computed by the inverse of the running average of the squared gradient norms. The results, including mean and variance, are shown in Figure 3.3, which numerically validate the linear relation indicated by Corollary 3.

Chapter 4

Decentralized stochastic algorithm with communication compression

Communication efficiency has been widely recognized as the bottleneck for large-scale decentralized machine learning applications in multi-agent or federated environments. To tackle the communication bottleneck, there have been many efforts to design algorithms with communication compression for decentralized nonconvex optimization (cf. (1.3)), where the clients are only allowed to communicate a small amount (commonly measured in bits) of quantized information with their neighbors over a predefined graph topology. Despite significant efforts, the state-of-the-art algorithm in the nonconvex setting still suffers from a slower rate of convergence $O(G/\epsilon^3)$ compared with their uncompressed counterpart to reach an ϵ -first-order stationary point for nonconvex objectives, where G measures the data heterogeneity across different clients. We propose BEER, which adopts communication compression with gradient tracking, and show it converges at a *faster rate* of $O(1/\epsilon^2)$. This significantly improves over the state-of-the-art rate, by matching the rate without compression even under arbitrary data heterogeneity. Numerical experiments are also provided to corroborate our theory and confirm the practical superiority of BEER in the data heterogeneous regime.

This chapter is based on our previous publication [ZLL⁺22].

4.1 Preliminaries

4.1.1 Assumptions

We first introduce Assumptions 5 and 6, which are imposed on the global loss function. Assumption 5 only assumes global loss function is L -smooth, which makes it weaker than Assumption 2. The Polyak-

Łojasiewicz (PL) condition [Pol63] described in Assumption 6 can lead to linear convergence even when the function is nonconvex. It's worth noting that Assumption 6 is weaker than strong convexity.

Assumption 5 (*L-smooth objective function*). A function $f(\mathbf{x})$ is *L-smooth* if $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the following inequality holds for some $L \geq 0$:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2.$$

Assumption 6 (*PL condition*). A function $f(\mathbf{x})$ satisfies the *Polyak-Łojasiewicz condition* if $\forall \mathbf{x} \in \mathbb{R}^d$, the following inequality holds for some $\mu > 0$:

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2\mu(f(\mathbf{x}) - f^*),$$

where $f^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$.

Assumption 7 is a standard assumption for stochastic optimization algorithms, which implies a mini-batch stochastic gradient of batch size b has a variance of σ^2/b .

Assumption 7 (*Bounded local variance*). The variance of a local stochastic gradient at a uniformly randomly selected data sample \mathbf{z}_i is bounded if $\forall i \in [n]$ and $\forall \mathbf{x} \in \mathbb{R}^d$, the following inequality holds:

$$\mathbb{E}_{\mathbf{z}_i \sim \mathcal{M}_i} \|\nabla \ell(\mathbf{x}; \mathbf{z}_i) - \nabla f_i(\mathbf{x})\|_2^2 \leq \sigma^2.$$

4.1.2 Compression operators

Using compression operators to compress gradients or any data needed to be communicated can substantially improve communication efficiency [TGZ⁺18, SCJ18, KSJ19, RSF21, FSG⁺21].

Definition 5 defines a randomized general compression operator that only guarantees the expected compression error $\mathbb{E}\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|_2^2$ is less than the magnitude of original message $\|\mathbf{x}\|_2^2$. If the compression operator is unbiased, it falls into the unbiased compression operator category [AGL⁺17, KFJ18, MGTR19, LR20] defined in Definition 5.

Definition 5 (*General compression operator*). A randomized map $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a ρ -compression operator if $\forall \mathbf{x} \in \mathbb{R}^d$ and some $\rho \in [0, 1]$, the following inequality holds:

$$\mathbb{E}\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|_2^2 \leq (1 - \rho)\|\mathbf{x}\|_2^2.$$

Definition 6 (*Unbiased compression operator*). A randomized map $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is an ω -unbiased compression operator if $\forall \mathbf{x} \in \mathbb{R}^d$, and there exists some $\omega > 0$, the followings hold:

$$\mathbb{E}[\mathcal{C}(\mathbf{x})] = \mathbf{x} \quad \text{and} \quad \mathbb{E}\|\mathcal{C}(\mathbf{x}) - \mathbf{x}\|_2^2 \leq \omega\|\mathbf{x}\|_2^2.$$

It is possible to convert a general (biased) compression operator to an unbiased one. Example 1 shows an example that random sparsification can be converted between these two categories by applying a constant scaling. Examples 2 and 3 are two examples of general compression operators.

Example 1 (Random sparsification). *Random sparsification keeps an element from a d -dimensional vector with probability $\frac{k}{d}$. Let $\mathbf{u} \in \mathbb{R}^d$ where $u_i \sim B(\frac{k}{d})$.*

- Biased random sparsification is defined as $\text{random}_{k,\text{biased}}(\mathbf{x}) = \mathbf{u} \odot \mathbf{x}$, which satisfies Definition 5 with $\rho = \frac{k}{d}$.
- Unbiased random sparsification is defined as $\text{random}_{k,\text{unbiased}}(\mathbf{x}) = \frac{d}{k} \mathbf{u} \odot \mathbf{x}$, which satisfies Definition 6 with $\omega = \frac{d}{k} - 1$.

Example 2 (Random dithering). *Random dithering [AGL⁺17] quantizes the message to b bits after adding random noise, which is defined as*

$$\text{gsgd}_b(\mathbf{x}) := \frac{\|\mathbf{x}\|_2}{\tau} \cdot \text{sign}(\mathbf{x}) \cdot 2^{-(b-1)} \cdot \left\lfloor \frac{2^{(b-1)}|\mathbf{x}|}{\|\mathbf{x}\|_2} + \mathbf{u} \right\rfloor$$

where $\tau = 1 + \min \left\{ \frac{d}{2^{2(b-1)}}, \frac{\sqrt{d}}{2^{(b-1)}} \right\}$, and \mathbf{u} is the random dithering vector uniformly sampled from $[0, 1]^d$. gsgd_b satisfies Definition 5 with $\rho = 1/\tau$.

Example 3 (top_k). top_k [AHJ⁺18, SCJ18] keeps k elements that have the largest absolute values and sets other elements to 0, which is defined as

$$\text{top}_k(\mathbf{x}) := \mathbf{x} \odot \mathbf{u}(\mathbf{x}),$$

where $[\mathbf{u}(\mathbf{x})]_i = 1$ if the absolute value of the i -th element is one of the k -largest absolute values, otherwise $[\mathbf{u}(\mathbf{x})]_i = 0$. top_k satisfies Definition 5 with $\rho = k/d$.

4.1.3 Additional notation

We denote a gradient of local objective function as $\tilde{\nabla} f_i(\mathbf{x}) := \nabla \ell(\mathbf{x}; \mathbf{z}_i)$, where $\mathbf{z}_i \in \mathcal{Z}_i$ is uniformly randomly sampled from local dataset \mathcal{Z}_i . In addition, for a batch of b i.i.d. uniform random samples $\{\mathbf{z}_{i,k} \in \mathcal{Z}_i\}_{k \in [b]}$, the mini-batch stochastic gradient is defined as $\tilde{\nabla}_b f_i(\mathbf{x}) := \frac{1}{b} \sum_{k=1}^b \nabla \ell(\mathbf{x}; \mathbf{z}_{i,k})$.

We define the distributed stochastic gradient and mini-batch stochastic gradient $\tilde{\nabla} F(\mathbf{X})$, $\tilde{\nabla}_b F(\mathbf{X})$, and distributed compression operator $\mathcal{C}(\mathbf{X})$ analogously to (1.7).

4.2 The BEER algorithm

This section presents BEER (cf. Algorithm 5) for decentralized nonconvex optimization with compressed communication.

At the t -th iteration, BEER maintains the current model estimates $\mathbf{X}^{(t)}$ and the global gradient estimates $\mathbf{V}^{(t)}$ across the clients. At the crux of its design, BEER also tracks and maintains two control sequences $\mathbf{H}^{(t)}$ and $\mathbf{G}^{(t)}$ that serve as compressed surrogates of $\mathbf{X}^{(t)}$ and $\mathbf{V}^{(t)}$, respectively. In particular, these two control sequences are updated by aggregating the received compressed messages alone (cf. Line 5 and Line 8).

It then boils down to how to carefully update these quantities in each iteration with communication compression. To begin, note that for each client i , BEER not only maintains its own parameters $\{\mathbf{x}_i^{(t)}, \mathbf{v}_i^{(t)}, \mathbf{h}_i^{(t)}, \mathbf{g}_i^{(t)}\}$, but also the control variables from its neighbors, namely, $\{\mathbf{h}_j^{(t)}\}_{j \in \mathcal{N}(i)}$ and $\{\mathbf{g}_j^{(t)}\}_{j \in \mathcal{N}(i)}$.

Each iteration can be roughly broken into three steps.

- **Update model estimate:** Each agent i first updates its model $\mathbf{x}_i^{(t+1)}$ according to Line 3, by a gradient-style update with a correction term using compressed surrogates of models, i.e., $\{\mathbf{h}_j^{(t)}\}_{j \in \mathcal{N}(i)}$. This update step incorporates aggregated information and compensates for the compression errors, thus leads to better consensus among agents and improved communication efficiency.
- **Update global gradient estimate:** Each client i updates the global gradient estimate $\mathbf{v}_i^{(t+1)}$ according to Line 6, where the last correction term—based on the difference of the gradients at consecutive models—is known as a trick called *gradient tracking* [QL18, DLS16, NOS17]. The use of gradient tracking is critical: as shall be seen momentarily, it contributes to the key difference from CHOCO-SGD that enables the fast rate of $O(1/\epsilon^2)$ without any bounded dissimilarity or bounded gradient assumptions. Indeed, if we remove the control sequence $\mathbf{G}^{(t)}$ and substitute Lines 6-8 by $\mathbf{V}^{(t+1)} = \tilde{\nabla}_b F(\mathbf{X}^{(t+1)})$, we recover CHOCO-SGD from BEER.
- **Update compressed surrogates with communication:** To update $\{\mathbf{h}_j^{(t)}\}_{j \in \mathcal{N}(i)}$, each client i first computes a compressed message $\mathbf{q}_{h,i}^{(t+1)}$ that encodes the difference $\mathbf{x}_i^{(t+1)} - \mathbf{h}_i^{(t)}$, and broadcasts to its neighbors (cf. Line 4). Then, each client i updates $\{\mathbf{h}_j^{(t)}\}_{j \in \mathcal{N}(i)}$ by aggregating the received compressed messages $\{\mathbf{q}_{h,j}^{(t+1)}\}_{j \in \mathcal{N}(i)}$ following Line 5. The updates of $\{\mathbf{g}_j^{(t)}\}_{j \in \mathcal{N}(i)}$ can be performed similarly. Moreover, all the compressed messages can be sent in a single communication round at one iteration, i.e., the communications in Lines 4 and 7 can be performed at once. This leverages EF21 [RSF21] for communication compression, which is a *better and simpler* algorithm that deals with biased compression operators compared with the error feedback (or error compensation, EF/EC) framework [KRSJ19, SK20]. Using the control sequence $\mathbf{G}^{(t)}$, BEER does not need to apply EF/EC explicitly and can deal with the error implicitly.

Algorithm 5 BEER: Better comprESSION for dEcentRALized optimization

-
- 1 **Input:** Initial point $\mathbf{X}^{(0)} = \mathbf{x}^{(0)}\mathbf{1}^\top$, $\mathbf{G}^{(0)} = \mathbf{H}^{(0)} = \mathbf{0}_{d \times n}$, $\mathbf{V}^{(0)} = \nabla F(\mathbf{X}^{(0)})$, step size η , mixing step size γ , minibatch size b .
 - 2 **for** $t = 0, 1, \dots$ **do**
 - 3 $\mathbf{X}^{(t+1)} = \mathbf{X}^{(t)} + \gamma\mathbf{H}^{(t)}(\mathbf{W} - \mathbf{I}) - \eta\mathbf{V}^{(t)}$
 - 4 $\mathbf{Q}_h^{(t+1)} = \mathcal{C}(\mathbf{X}^{(t+1)} - \mathbf{H}^{(t)})$ \triangleright agent i sends $q_{h,i}^{(t+1)}$ to all its neighbors
 - 5 $\mathbf{H}^{(t+1)} = \mathbf{H}^{(t)} + \mathbf{Q}_h^{(t+1)}$
 - 6 $\mathbf{V}^{(t+1)} = \mathbf{V}^{(t)} + \gamma\mathbf{G}^{(t)}(\mathbf{W} - \mathbf{I}) + \tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)})$
 - 7 $\mathbf{Q}_g^{(t+1)} = \mathcal{C}(\mathbf{V}^{(t+1)} - \mathbf{G}^{(t)})$ \triangleright agent i sends $q_{g,i}^{(t+1)}$ to all its neighbors
 - 8 $\mathbf{G}^{(t+1)} = \mathbf{G}^{(t)} + \mathbf{Q}_g^{(t+1)}$
-

4.3 Convergence guarantees

This section presents convergence guarantees of BEER for nonconvex objective functions (cf. Section 4.3.1) and PL objective functions (cf. Section 4.3.2). The convergence analysis is based on a Lyapunov function tailored for BEER, given by

$$\Phi^{(t)} = \mathbb{E}[f(\bar{\mathbf{x}}^{(t)})] - f^* + \frac{c_1 L}{n} \Omega_1^{(t)} + \frac{c_2(1-\alpha)^2}{nL} \Omega_2^{(t)} + \frac{c_3 L}{n} \Omega_3^{(t)} + \frac{c_4(1-\alpha)^4}{nL} \Omega_4^{(t)}, \quad (4.1)$$

where $\{c_i\}$ are some constants depend on specific settings, $\mathbb{E}[f(\bar{\mathbf{x}}^{(t)})] - f^*$ represents the sub-optimality gap, and the errors $\{\Omega_i^{(t)}\}$ are defined by

$$\begin{aligned} \Omega_1^{(t)} &:= \mathbb{E} \|\mathbf{H}^{(t)} - \mathbf{X}^{(t)}\|_{\mathbb{F}}^2, & \Omega_2^{(t)} &:= \mathbb{E} \|\mathbf{G}^{(t)} - \mathbf{V}^{(t)}\|_{\mathbb{F}}^2, \\ \Omega_3^{(t)} &:= \mathbb{E} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)}\mathbf{1}^\top\|_{\mathbb{F}}^2, & \Omega_4^{(t)} &:= \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)}\mathbf{1}^\top\|_{\mathbb{F}}^2. \end{aligned} \quad (4.2)$$

Here, $\Omega_1^{(t)}$ and $\Omega_2^{(t)}$ denote the compression errors for $\mathbf{X}^{(t)}$ and $\mathbf{V}^{(t)}$ when approximated using the compressed surrogates $\mathbf{H}^{(t)}$ and $\mathbf{G}^{(t)}$, respectively, and $\Omega_3^{(t)}$ and $\Omega_4^{(t)}$ denote the consensus errors of $\mathbf{X}^{(t)}$ and $\mathbf{V}^{(t)}$.

4.3.1 Convergence for nonconvex objective functions

Theorem 7 presents the convergence results of BEER for nonconvex objective functions

Theorem 7 (Convergence in the nonconvex setting). *Suppose Assumptions 1, 5 and 7 hold. There exist absolute constants $c_1, c_2, c_3, c_4, c_\gamma, c_\eta > 0$, such that if we set $\gamma = c_\gamma \rho(1-\alpha)$, $\eta = c_\eta \gamma(1-\alpha)^2/L$, for the Lyapunov function $\Phi^{(t)}$ defined in (4.1), it holds*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 \leq \frac{2(\Phi^{(0)} - \Phi^{(T)})}{\eta T} + \frac{36c_4\sigma^2}{c_\gamma b \rho L}.$$

Theorem 7 shows the convergence results for BEER.

Full gradient case When we can access the full gradient, which is equivalent to $\sigma^2 = 0$, BEER converges at a rate of $O(\epsilon^{-2})$. This rate is faster than the $O(\epsilon^{-3})$ rate by CHOCO-SGD [KSJ19] and DeepSqueeze [TLQ⁺19], and the $O(\epsilon^{-4})$ rate by SQuARM-SGD [SDGD21].

More specifically, to achieve an ϵ -first-order stationary point, BEER needs

$$O\left(\frac{1}{(1-\alpha)^3\rho\epsilon^2}\right)$$

communication rounds, where α and ρ are the mixing rate (cf. Definition 1) and the compression parameter (cf. Definition 5), respectively. In comparison, the state-of-the-art algorithm CHOCO-SGD [KSJ19] converges in $O\left(\frac{G}{(1-\alpha)^2\rho\epsilon^3}\right)$ communication rounds, with G being the bounded gradient parameter, namely, $\mathbb{E}_{z_i \sim \mathcal{M}_i} \|\nabla \ell(\mathbf{x}; z_i)\|^2 \leq G^2$. Therefore, BEER improves over CHOCO-SGD not only in terms of a better dependency on ϵ , but also removing the bounded gradient assumption, which is significant since in practice, G can be excessively large due to data heterogeneity across the clients.

The dependency on ρ of BEER is consistent with other compression schemes, such as CHOCO-SGD, DeepSqueeze and SQuARM-SGD for the nonconvex setting, as well as LEAD [KKJ⁺21] and EF-C-GT [LLHP22] for the strongly convex setting.

As for the dependency on $(1-\alpha)^{-1}$, BEER is slightly worse than CHOCO-SGD, where CHOCO-SGD has a dependency of $O((1-\alpha)^{-2})$ whereas BEER has a dependency of $O((1-\alpha)^{-3})$. This degeneration is also seen in the analysis of uncompressed decentralized algorithms using gradient tracking [SLH20, XKK22b], where the rate $O((1-\alpha)^{-2})$ is worse than the rate of $O((1-\alpha)^{-1})$ for basic decentralized SGD algorithms [KDG03, LZZ⁺17] by a factor of $1-\alpha$. In addition, both BEER and CHOCO-SGD use small mixing step size γ to guarantee convergence, which makes the dependency on $(1-\alpha)^{-1}$ worse than their uncompressed counterparts.

Stochastic gradient case In the presence of local variance, the squared gradient norm of BEER has an additional term that scales on the order of $O\left(\frac{\sigma^2}{\rho b}\right)$ (ignoring other parameters). By choosing a sufficiently large minibatch size b , i.e., $b \geq O\left(\frac{\sigma^2}{\rho\epsilon^2}\right)$, BEER maintains the iteration complexity

$$O\left(\frac{1}{(1-\alpha)^3\rho\epsilon^2}\right)$$

to reach an ϵ -first-order stationary point, without the bounded gradient assumption, thus inheriting similar advantages over CHOCO-SGD as discussed earlier. In terms of local computation, the gradient oracle complexity on a single client of BEER is

$$O\left(\frac{1}{(1-\alpha)^3\rho\epsilon^2} + \frac{\sigma^2}{(1-\alpha)^3\rho^2\epsilon^4}\right).$$

While our focus is on communication efficiency, to gain more insights, Table 1.4 summarizes both the communication rounds and the gradient complexity for different decentralized stochastic methods. While BEER does not require the bounded gradient assumption, it may lead to a worse gradient complexity in the data homogeneous setting due to the use of large minibatch size. Fortunately, this only impacts the local computation cost, and does not exacerbate the communication complexity, which is often the bottleneck. It is of great interest to further refine the design and analysis of BEER in terms of the gradient complexity.

4.3.2 Convergence for PL objectives functions

Theorem 8 (Convergence under PL condition). *Suppose Assumptions 1 and 5 to 7 hold. There exist absolute constants $c_1, c_2, c_3, c_4, c_\gamma, c_\eta > 0$, such that if we set $\gamma = c_\gamma \rho(1 - \alpha)$, $\eta = c_\eta \gamma(1 - \alpha)^2/L$, for the Lyapunov function $\Phi^{(t)}$ in (4.1), it holds*

$$\Phi^{(T)} \leq (1 - \mu\eta)^T \Phi^{(0)} + \frac{36c_4\sigma^2}{c_\gamma L b \rho}.$$

Theorem 8 shows the convergence guarantees for BEER under the PL condition. If we can access full local gradients ($\sigma = 0$), BEER converges linearly to the global optimum f^* at a rate of $O(\kappa \log(1/\epsilon))$. When using stochastic gradients, BEER converges linearly to a neighborhood of size $O(\frac{\sigma^2}{\rho b})$ around the global optimum.

4.4 Numerical experiments

This section presents numerical experiments on real-world datasets to showcase BEER’s superior ability to handle data heterogeneity over agents, by running each experiment on unshuffled datasets and comparing the performances with the state-of-the-art baseline algorithms both with and without communication compression.

For all experiments, we split unshuffled datasets evenly to 10 agents that are connected by a ring topology, so that we can simulate the scenario with high data heterogeneity across agents. Approximately, for the a9a dataset, 5 agents receive data with label 1 and others receive data with label 0; for the MNIST dataset, agent i receives data with label i . We use the FDLA matrix [XB04] as the mixing matrix to perform weighted information aggregation to accelerate convergence. For each experiment, all algorithms are initialized to the same starting point, and uses best-tuned learning rates and batch sizes.

Sections 4.4.1 and 4.4.2 compares BEER with 1) CHOCO-SGD [KSJ19], which is the state-of-the-art nonconvex decentralized optimizing algorithm using communication compression, and 2) DSGD [LZZ⁺17] and D^2 [TLY⁺18], which are decentralized optimization algorithms without compression. Section 4.4.3 further evaluates the impact of communication network and compression operators.

4.4.1 Logistic regression with nonconvex regularization

We first run experiments on logistic regression with a nonconvex regularizer [WJZ⁺19] on the a9a dataset [CL11]. Similar to Section 3.3.1, following [WJZ⁺19], the objective function over a datum (a, b) is defined as

$$\ell(x; \{f, l\}) = \log \left(1 + l \exp(-x^\top f) \right) + \lambda \sum_{i=1}^d \frac{x_i^2}{1 + x_i^2},$$

where $\{f, l\}$ represents a training tuple, $f \in \mathbb{R}^d$ is the feature vector and $l \in \{0, 1\}$ is the label, and λ is the regularization parameter which is set to $\lambda = 0.05$.

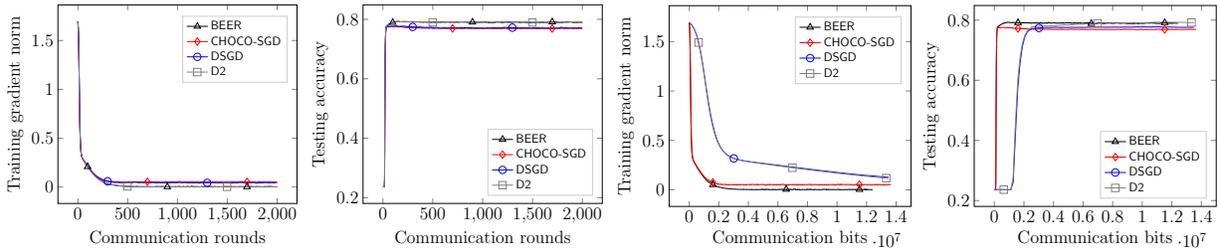


Figure 4.1: The training gradient norm and testing accuracy against communication rounds (left two panels) and communication bits (right two panels) for logistic regression with nonconvex regularization on unshuffled a9a dataset. Both BEER and CHOCO-SGD employ the biased gsgd_5 compression [AGL⁺17].

Figure 4.1 plots the training gradient norm and testing accuracy against communication rounds and communication bits for logistic regression with nonconvex regularization. The algorithms with communication compression (BEER and CHOCO-SGD [KSJ19]) converge faster than the uncompressed algorithms (DSGD [LZZ⁺17] and D^2 [TLY⁺18]) in terms of the communication bits. However, CHOCO-SGD fails to converge to a small gradient norm due to its inability to tolerate a high level of data dissimilarity across different agents. In contrast, BEER and D^2 converge to the smallest gradient norm, while BEER outperforms D^2 in terms of communication bits. The results for testing accuracy are similar, where BEER achieves the best testing accuracy and is the fastest.

4.4.2 One-hidden-layer neural network training

Similar to Section 3.3.2, we evaluate BEER by training a 1-hidden layer neural network on the MNIST dataset [LJB⁺95]. The network uses 64 hidden neurons, sigmoid activation functions and cross-entropy loss, where the loss function over a training pair $\{f, l\}$ is defined as

$$\ell(x; (f, l)) = \text{CrossEntropy}(\text{softmax}(\mathbf{W}_2 \text{sigmoid}(\mathbf{W}_1 f + c_1) + c_2), l).$$

Here the model parameter is defined by $x = \text{vec}(\mathbf{W}_1, c_1, \mathbf{W}_2, c_2)$, where the dimensions of the network parameters $\mathbf{W}_1, c_1, \mathbf{W}_2, c_2$ are $64 \times 784, 64 \times 1, 10 \times 64$, and 10×1 , respectively.

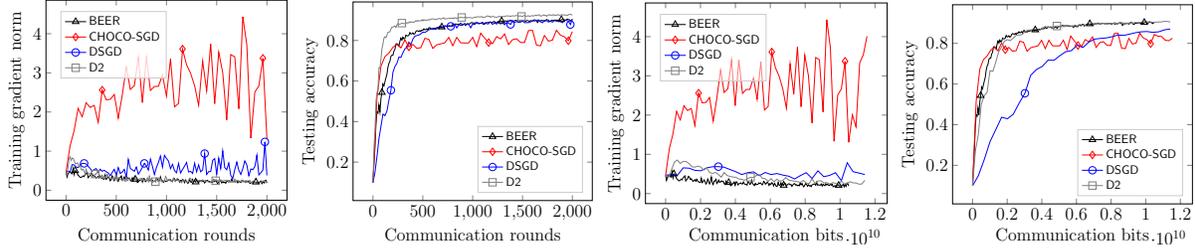


Figure 4.2: The training gradient norm and testing accuracy against communication rounds (left two panels) and communication bits (right two panels) for classification on unshuffled MNIST dataset using a 1-hidden-layer neural network. Both BEER and CHOCO-SGD employ the biased gsgd_{20} compression [AGL⁺17].

Figure 4.2 plots the training gradient norm and testing accuracy against communication rounds and communication bits for 1-hidden-layer neural network training. In terms of the final training gradient norm, BEER converges to a solution comparable to D^2 but at a lower communication cost, while CHOCO-SGD and DSGD cannot converge due to the high data heterogeneity. In terms of testing accuracy, BEER and D^2 have very similar performance, and outperform CHOCO-SGD and DSGD.

4.4.3 Network topology and compression operators

We further investigate the impact of communication network topology and compression operators on the performance of BEER. We follow the same setup as Section 4.4.1 to run logistic regression with nonconvex regularization ($\lambda = 0.05$) on the unshuffled a9a dataset, by splitting it evenly to 40 agents. All experiments use the same best-tuned step size $\eta = 0.5$, batch size $b = 100$ and $\gamma = 0.7$.

Impacts of network topology Figure 4.3 shows the training gradient norm and testing accuracy of BEER with respect to the communication rounds over different network topologies using the gsgd_5 compression [AGL⁺17]. Experimented topologies are ring topology ($\alpha = 0.978$), star topology ($\alpha = 0.951$), grid topology ($\alpha = 0.937$), and Erdős-Rényi topology with connectivity probability $p = 0.5$ and $p = 0.2$ ($\alpha = 0.49$ and $\alpha = 0.23$, respectively). Despite the huge differences in mixing rates, BEER can use nearly the same hyper-parameters to obtain similar performance. The experiments complement our theoretical analysis and show that BEER may converge way better in practice despite its cubic dependency of $(1 - \alpha)^{-1}$ in Theorem 7.

Impacts of compression schemes Figure 4.4 shows the training gradient norm and testing accuracy of BEER with respect to the communication rounds and communication bits on a ring topology using different compression schemes: no compression, gsgd_5 (cf. Example 2) and top_{10} (cf. Example 3). The

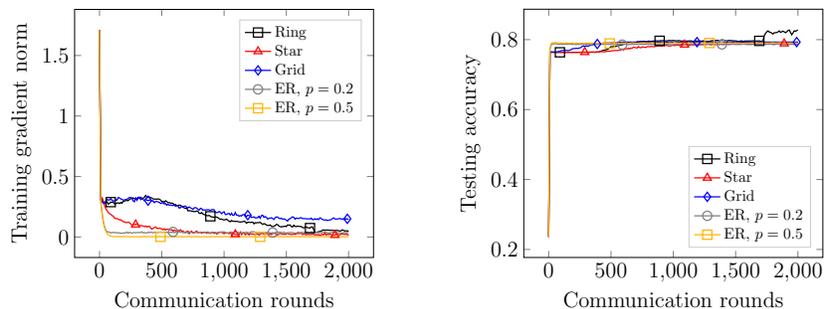


Figure 4.3: The training gradient norm and testing accuracy against communication rounds for BEER using the biased gsgd_5 compression [AGL⁺17] for logistic regression with nonconvex regularization on unshuffled a9a dataset.

compression parameters are chosen such that each compression operator transfers similar number of bits per communication round. All experiments use the same best-tuned step size $\eta = 0.5$, batch size $b = 100$ and $\gamma = 0.7$, except that we use $\eta = 0.005$ and $\gamma = 0.8$ for top_{10} compression.

In terms of communication bits, using compression operators improves over the uncompressed baseline, because all algorithms with compression converge to a solution with lower gradient norm and higher testing accuracy at a lower communication cost. In terms of communication rounds and testing accuracy, different compression operators can lead to significantly behaviors. For example, using gsgd_5 compression operator leads to faster converges than without compression, but using the top_{10} compression operator leads to slower converges than without compression. In sum, BEER with gsgd_5 reaches the highest final testing accuracy while behaves similar to BEER without compression in terms of communication rounds, which indicates the benefit of using communication compression.

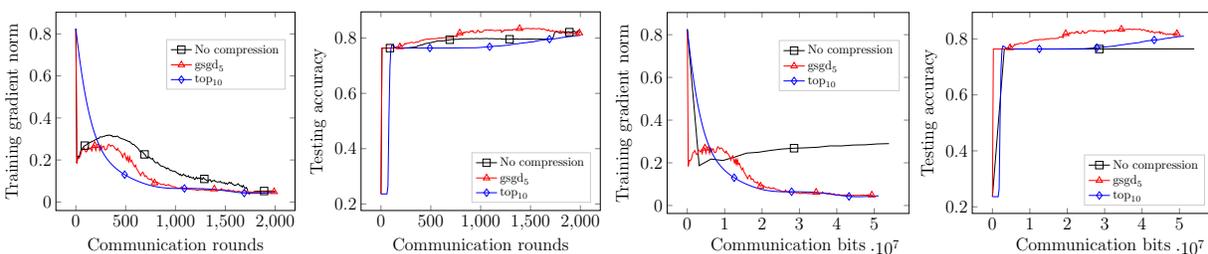


Figure 4.4: The training gradient norm and testing accuracy against communication rounds (top two panels) and communication bits (bottom two panels) for BEER using different compression schemes for logistic regression with nonconvex regularization on unshuffled a9a dataset.

Chapter 5

Decentralized private stochastic algorithm with communication compression

To enable large-scale machine learning in bandwidth-hungry environments such as wireless networks, significant progress has been made recently in designing communication-efficient federated learning algorithms with the aid of communication compression. On the other end, privacy-preserving, especially at the client level, is another important desideratum that has not been addressed simultaneously in the presence of advanced communication compression techniques yet. In this chapter, based on BEER, we propose PORTER, a decentralized stochastic optimization algorithm for decentralized nonconvex ERM problems (cf. (1.3)) which enhances the communication efficiency of private decentralized learning with communication compression, by exploiting general compression operators, gradient clipping and local differential privacy.

We provide the *first* theoretical analysis of private decentralized optimization algorithms using communication compression and gradient clipping, and show explicit dependency on the mixing rate and compression parameter. Furthermore, we provide numerical evidence that shows PORTER converges in similar communication rounds without sacrificing privacy nor utility compared to server/client private optimization algorithm SoteriaFL-SGD [LZLC22].

5.1 Preliminaries

5.1.1 Local differential privacy

In decentralized learning systems, all agents share potentially sensitive information with their neighbors. If some agents are exploited by adversaries, the system will face a risk of privacy leaking even though

the system-level privacy is protected. Therefore, we introduce local differential privacy (LDP) defined Definition 7 following [DJW13, ABCP13, CABP13, XYD19, WXDX20, ZZY⁺21, LZLC22], which protects each agent’s privacy from leaking to other agents.

Definition 7 (Local differential privacy (LDP)). *A randomized mechanism $\mathcal{M} : \mathcal{Z} \rightarrow \mathcal{R}$ with domain \mathcal{Z} and range \mathcal{R} satisfies (ϵ, δ) -local differential privacy for client i , if for any two neighboring dataset $\mathbf{Z}_i, \mathbf{Z}'_i \in \mathcal{Z}$ at client i and for any subset of outputs $\mathbf{R} \subseteq \mathcal{R}$, it holds that*

$$\mathbb{P}(\mathcal{M}(\mathbf{Z}_i) \in \mathbf{R}) \leq e^\epsilon \mathbb{P}(\mathcal{M}(\mathbf{Z}'_i) \in \mathbf{R}) + \delta. \quad (5.1)$$

The two datasets \mathbf{Z}_i and \mathbf{Z}'_i are called neighboring if they are only different by one data point at client i .

Definition 7 is a more stricter privacy guarantee because it can imply general differential privacy (DP). Consequently, LDP requires a larger perturbation variance than general DP. To identify the impact of the decentralized LDP setting compared to centralized DP setting, we define the baseline utility

$$\phi_m = \frac{\sqrt{d \log(1/\delta)}}{m\epsilon}, \quad (5.2)$$

which can be understood as the final utility of a centralized system with m data samples that guarantees (ϵ, δ) -DP. For typical private problems, the local sample size m has to be large enough for the privacy perturbation to work, we impose a mild assumption that $\phi_m < 1$. For example, the problem defined in (1.3) has in total mn data samples, running an (ϵ, δ) -DP algorithm on one server that can access all data will achieve $\frac{1}{n}\phi_m$ utility in $n\phi_m^{-1}$ iterations.

5.1.2 The SoteriaFL algorithm framework

This section reviews SoteriaFL [LZLC22] (cf. Algorithm 6) as it has inspired the development of PORTER. SoteriaFL is a unified framework for differentially private optimization algorithms with unbiased communication compression in the server/client setting, which employs shifted communication to compensate for the error induced by compression, and can use various local sub-routines to produce an estimate of global gradients, e.g., stochastic gradient descent (SGD) or stochastic variance reduced gradient (SVRG).

During each iteration, each agent 1) computes an estimate of the global gradient using a local sub-routine, 2) adds a Gaussian perturbation to the gradient estimate to preserve privacy, and 3) communicates the perturbed gradient with its neighbors using the shifted communication scheme. Instead of directly compressing the perturbed gradients, SoteriaFL maintains a reference $\mathbf{s}_i^{(t)}$ and compresses the shifted message $\mathbf{g}_i^{(t)} - \mathbf{s}_i^{(t)}$ at each agent. This extra shift operation allows SoteriaFL achieve much better convergence behavior (fewer communication rounds) than existing algorithms.

Algorithm 6 SoteriaFL (a unified framework for compressed private FL)

```

1  initial point  $\mathbf{x}^{(0)}$ , step size  $\eta_t$ , shift step size  $\gamma_t$ , variance  $\sigma_p^2$ , initial reference  $\mathbf{s}_i^{(0)} = \mathbf{0}_d$ 
2  for  $t = 0, 1, 2, \dots, T$  do
3      for each node  $i \in [n]$  do in parallel
4          Compute local gradient estimator  $\tilde{\mathbf{g}}_i^{(t)}$   $\triangleright$  It allows many methods, e.g., SGD, SVRG, and SAGA.
5          Privacy:  $\mathbf{g}_i^{(t)} = \tilde{\mathbf{g}}_i^{(t)} + \boldsymbol{\zeta}_i^{(t)}$ , where  $\boldsymbol{\zeta}_i^{(t)} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ 
6          Compression: let  $\mathbf{v}_i^{(t)} = \mathcal{C}_i^t(\mathbf{g}_i^{(t)} - \mathbf{s}_i^{(t)})$  and send to the server  $\triangleright$  Shifted compression.
7          Update shift  $\mathbf{s}_i^{t+1} = \mathbf{s}_i^{(t)} + \gamma_t \mathcal{C}_i^t(\mathbf{g}_i^{(t)} - \mathbf{s}_i^{(t)})$ 
8      end each node Server aggregates compressed information  $\mathbf{v}^t = \mathbf{s}^t + \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 
9       $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{v}^t$ 
10      $\mathbf{s}^{t+1} = \mathbf{s}^t + \gamma_t \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t$ 

```

SoteriaFL-SGD is obtained by using mini-batch stochastic gradients as the gradient estimator for Line 4. The theoretical analysis shows explicit dependency on the compression parameter, which helps to better understand the communication-utility trade-off for private server/client optimization algorithms, as show in Table 1.5. Compare the centralized DP setting, the server/client algorithm differentially private algorithm DDP-SRM [WJEG19] reaches the same $\frac{1}{n}\phi_m$ utility but at a worse $n^2 d \phi_m^{-1}$ iteration complexity. SoteriaFL-SGD is also a server/client algorithm but with LDP guarantee, which reaches an even worse utility but with potentially less iterations due to the stronger privacy constraint.

5.1.3 Gradient clipping

In practice, gradient clipping is frequently adopted to ensure the gradients are within a predetermined region, so that the variance of privacy perturbation can be decided accordingly. The clipping operator we adopt is a smooth clipping operator defined in Definition 8, which scales a vector into a ball of radius τ centered at the origin. Another widely used clipping operator is the piece-wise linear clipping operator defined in Definition 9, which scales inputs whenever its gradient norm is larger than τ and does nothing otherwise.

Definition 8 (Smooth clipping operator). For $\mathbf{x} \in \mathbb{R}^d$, the clipping operator is defined as

$$\text{Clip}_\tau(\mathbf{x}) = \frac{\tau}{\tau + \|\mathbf{x}\|_2} \mathbf{x}.$$

For $\mathbf{X} \in \mathbb{R}^{d \times n}$, the distributed clipping operator is defined as

$$\text{Clip}_\tau(\mathbf{X}) = [\text{Clip}_\tau(\mathbf{x}_1), \dots, \text{Clip}_\tau(\mathbf{x}_n)].$$

Definition 9 (Piece-wise linear clipping operator). For $\mathbf{x} \in \mathbb{R}^d$, the clipping operator is defined as

$$\text{Clip}_\tau(\mathbf{x}) = \mathbf{x} \min \{1, \tau / \|\mathbf{x}\|_2\}.$$

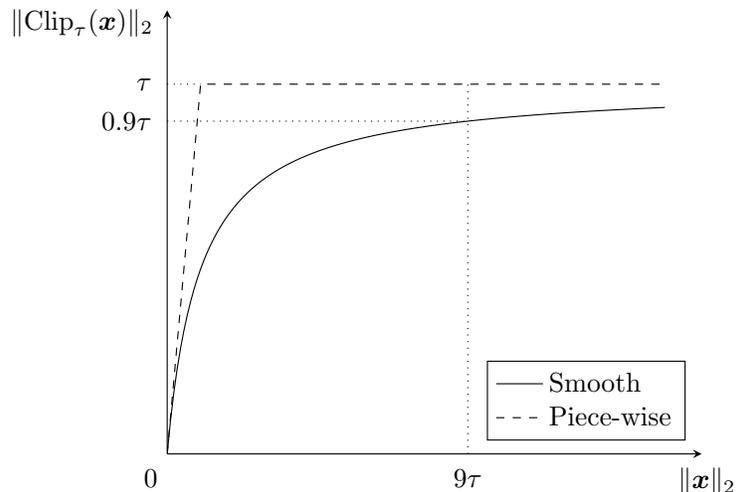


Figure 5.1: Illustration of input norm and clipped norm for the smooth clipping operator (Definition 8) and piece-wise linear clipping operator (Definition 9), where τ is the clipping parameter.

Figure 5.1 plots the norm of a vector before and after clipping for these two clipping operators, which shows most fundamental difference is when $\|\mathbf{x}\|_2 = \tau$, the piece-wise linear compression operator is not smooth, which leads to more difficulties in theoretical analysis. When $\|\mathbf{x}\|_2$ is small, these clipping operators behave like identity transformation, which keeps the stationary point property of objective functions.

5.2 The PORTER algorithm

We propose PORTER (Algorithm 7), a novel stochastic private decentralized optimization algorithm for finding ν -stationary points of nonconvex finite-sum problems based on BEER, where gradient estimates V , stochastic gradients G , perturbation noise E , compressed surrogate Q_x and $Q_{v'}$, and their corresponding agent-wise values are defined analogously to (1.5).

Key ingredients of PORTER are: 1) gradient clipping, which ensures the norm of clipped gradients are bounded by τ ; 2) privacy perturbation, which adds a Gaussian noise to clipped gradients to achieve privacy constraints; 3) compression with error feedback [RSF21], which accelerates the convergence with biased compression operators; and 4) stochastic gradient tracking, tracks the global gradient locally at each agent. The algorithm details explained in Section 4.2 also applies to PORTER, except PORTER doesn't explicitly show compressed messages as Lines 4 and 7, and PORTER adds perturbation to the stochastic

Algorithm 7 PORTER

```

1 input:  $\bar{\mathbf{x}}^{(0)}, \eta, \gamma, \tau, b, \sigma_p, T$ 
2 initialize:  $\mathbf{V}^{(0)} = \mathbf{Q}_v^{(0)} = \mathbf{G}_p^{(0)} = \mathbf{0}_{d \times n}, \mathbf{Q}_x^{(0)} = \mathbf{X}^{(0)} = \bar{\mathbf{x}}^{(0)} \mathbf{1}_n^\top$ 
3 for  $t = 1, \dots, T$  do
4      $\mathbf{G}_\tau^{(t)} = \frac{1}{b} \sum_{\mathbf{Z} \in \mathcal{Z}^{(t)}} \text{Clip}_\tau(\nabla \ell(\mathbf{X}^{(t-1)}; \mathbf{Z}))$ 
5      $\mathbf{G}_p^{(t)} = \mathbf{G}_\tau^{(t)} + \mathbf{E}^{(t)}$ , where  $e_i^{(t)} \sim \mathcal{N}(\mathbf{0}_d, \sigma_p^2 \mathbf{I}_d)$ 
6      $\mathbf{Q}_v^{(t)} = \mathbf{Q}_v^{(t-1)} + \mathcal{C}(\mathbf{V}^{(t-1)} - \mathbf{Q}_v^{(t-1)})$  ▷ Communication
7      $\mathbf{V}^{(t)} = \mathbf{V}^{(t-1)} + \gamma \mathbf{Q}_v^{(t)} (\mathbf{W} - \mathbf{I}_n) + \mathbf{G}_p^{(t)} - \mathbf{G}_p^{(t-1)}$ 
8      $\mathbf{Q}_x^{(t)} = \mathbf{Q}_x^{(t-1)} + \mathcal{C}(\mathbf{X}^{(t-1)} - \mathbf{Q}_x^{(t-1)})$  ▷ Communication
9      $\mathbf{X}^{(t)} = \mathbf{X}^{(t-1)} + \gamma \mathbf{Q}_x^{(t)} (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t)}$ 
10 output:  $\mathbf{x}_{out} \sim \text{Uniform}(\{\mathbf{x}_i^{(t)} | i \in [n], t \in [T]\})$ .
    
```

gradients. PORTER initializes gradient-related variables to $\mathbf{0}_d$ and other variables to the same random value $\bar{\mathbf{x}}^{(0)}$, which improves the algorithm's stability in early iterations and simplifies analysis but has no impact on convergence rates.

5.3 Local differential privacy guarantee

Algorithm 7 adds a Gaussian noise $e_i^{(t)} \sim \mathcal{N}(\mathbf{0}_d, \sigma_p^2 \mathbf{I}_d)$ to ensure privacy. Theorem 9 proves our algorithm is (ϵ, δ) -private when setting the variance of Gaussian perturbation to be $\sigma_p^2 = T\tau^2 b^2 \phi_m^2 / d$. The proof is deferred to Appendix D.1.

Theorem 9 (Local differential privacy). *Let $\phi_m = \frac{\sqrt{d \log(1/\delta)}}{m\epsilon}$. For any $\epsilon \leq T/m^2$ and $\delta \in (0, 1)$, Algorithm 7 is (ϵ, δ) -LDP after T iterations if we set*

$$\sigma_p^2 = \frac{T\tau^2 b^2 \log(1/\delta)}{m^2 \epsilon^2} = T\tau^2 b^2 \phi_m^2 / d. \quad (5.3)$$

Using clipping operators guarantees all gradients' ℓ_2 norms are bounded by τ , i.e., $\forall i, t, \|\mathbf{g}_{\tau, i}^{(t)}\|_2 \leq \tau$. Therefore, Theorem 9 holds for PORTER throughout.

5.4 Convergence with bounded gradient assumption

This section theoretically analyzes the convergence properties of PORTER under bounded gradient assumption without gradient clipping. Section 5.3 shows PORTER is (ϵ, δ) -local differential private when using a specific perturbation variance $\sigma_p^2 = T\tau^2 b^2 \phi_m^2 / d$, then presents convergence analysis.

Assumption 8 (Bounded sample gradient). For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and any datum \mathbf{z} in dataset \mathcal{Z} , it holds

$$\|\nabla \ell(\mathbf{x}; \mathbf{z})\|_2 \leq \tau.$$

With Assumption 8, PORTER can omit the clipping operator and reduce to Algorithm 8. Here, we assume $b = 1$.

Algorithm 8 PORTER with bounded gradient assumption

```

1 input:  $\bar{\mathbf{x}}^{(0)}, \eta, \gamma, \tau, \sigma_p, T$ 
2 initialize:  $\mathbf{V}^{(0)} = \mathbf{Q}_v^{(0)} = \mathbf{G}_p^{(0)} = \mathbf{0}_{d \times n}, \mathbf{Q}_x^{(0)} = \mathbf{X}^{(0)} = \bar{\mathbf{x}}^{(0)} \mathbf{1}_n^\top$ 
3 for  $t = 1, \dots, T$  do
4    $\mathbf{G}_\tau^{(t)} = \nabla \ell(\mathbf{X}^{(t-1)}; \mathbf{Z}^{(t)})$ 
5    $\mathbf{G}_p^{(t)} = \mathbf{G}_\tau^{(t)} + \mathbf{E}^{(t)}$ , where  $\mathbf{e}_i^{(t)} \sim \mathcal{N}(\mathbf{0}_d, \sigma_p^2 \mathbf{I}_d)$ 
6    $\mathbf{Q}_v^{(t)} = \mathbf{Q}_v^{(t-1)} + \mathcal{C}(\mathbf{V}^{(t-1)} - \mathbf{Q}_v^{(t-1)})$  ▷ Communicate
7    $\mathbf{V}^{(t)} = \mathbf{V}^{(t-1)} + \gamma \mathbf{Q}_v^{(t)} (\mathbf{W} - \mathbf{I}_n) + \mathbf{G}_p^{(t)} - \mathbf{G}_p^{(t-1)}$ 
8    $\mathbf{Q}_x^{(t)} = \mathbf{Q}_x^{(t-1)} + \mathcal{C}(\mathbf{X}^{(t-1)} - \mathbf{Q}_x^{(t-1)})$  ▷ Communicate
9    $\mathbf{X}^{(t)} = \mathbf{X}^{(t-1)} + \gamma \mathbf{Q}_x^{(t)} (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t)}$ 
10 output:  $\mathbf{x}_{out} \sim \text{Uniform}(\{\mathbf{x}_i^{(t)} | i \in [n], t \in [T]\})$ .
```

Theorem 10 presents the convergence result of PORTER using general compression operators (Definition 5). The proof is deferred to Appendix D.2.

Theorem 10 (Convergence using general compression operators). Assume Assumptions 1, 2 and 8 hold. Use general compression operators (Definition 5). Set $\gamma = O((1 - \alpha)\rho)$, $\eta = O(\gamma^{4/3}(1 - \alpha)^{4/3}\phi_m/L)$, $T = \phi_m^{-2}$, and $\sigma_p^2 = T\tau^2\phi_m^2/d$. Let $\Delta = \mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^*$, Algorithm 8 reaches $O((1 - \alpha)^{-8/3}\rho^{-4/3}\phi_m \cdot \max\{\tau^2, L\Delta\})$ utility in T iterations.

We can compare each term of the final utility in Theorem 10 with the results of the server/client algorithm SoteriaFL-SGD, due to a lack of theoretical analysis of comparable algorithms. 1) PORTER will always reach the same final utility, whose dependency on the compression operator is $O(\rho^{-4/3})$. We can convert SoteriaFL-SGD's rate (cf. Table 1.5) using the approximation $\rho = 1/(1 + \omega)$. When n is not large enough compared to ω , or using aggressive compression (ω is large), SoteriaFL-SGD reaches $O(n^{-1}\rho^{-2}\phi_m)$ utility, which is worse than PORTER, even though PORTER is a decentralized optimization algorithm. 2) In terms of dependency on the baseline utility ϕ_m , due to extra complexities induced by the decentralized setting, PORTER takes $O(\phi_m^{-2})$ iterations to reach $O(\phi_m)$ utility, which is worse than other algorithms, e.g., SoteriaFL-SGD only takes $O(\phi_m^{-1})$ iterations.

5.5 Convergence without bounded gradient assumption

When assuming the bounded gradient assumption (cf. Assumption 8), PORTER can skip the clipping operator, and $G_{\tau,i}^{(t)}$ becomes an unbiased estimator of the local gradient $\nabla f_i(x_i^{(t)})$. However, Assumption 8 is rarely met in reality. For example, the gradient of a quadratic loss function is not bounded. Therefore, it is of interest to examine convergence without the strong bounded gradient assumption, where we utilize the gradient clipping operator $\text{Clip}_\tau(\cdot)$ to ensure gradients are bounded.

However, it is necessary to introduce milder assumptions, Assumptions 9 and 10, which limit the deviation of local objective functions to the global objective function. Because stochastic gradients at different agents will lose correct scaling after clipping, which will break the stationary point property at local minima. For example, consider optimizing a one-dimensional objective function over a 3-agent fully connected communication graph. Initialize every agent with $\bar{x}^{(0)} = x^*$, where x^* is one of global local minima. Assume local gradients are $g_1 = 1$, $g_2 = 1$ and $g_3 = -2$, which leads to a global gradient of $\mathbf{0}_3$. If we apply gradient clipping with $\tau = 1$, $\text{Clip}_1(g_1) = 0.5$, $\text{Clip}_1(g_2) = 0.5$ and $\text{Clip}_1(g_3) = -0.66$, and the global gradient is $g' = 0.33 \neq 0$.

Assumption 9 (Bounded gradient dissimilarity). $\forall x \in \mathbb{R}^d$ and $i \in [n]$, $\|\nabla f(x) - \nabla f_i(x)\|_2 \leq \frac{1}{12} \|\nabla f(x)\|_2$.

Assumption 10 (Bounded local variance). $\forall x \in \mathbb{R}^d$ and $i \in [n]$, $\mathbb{E}_{z \sim \mathcal{Z}_i} \|\nabla f_i(x) - \nabla \ell(x; z)\|_2^2 \leq \sigma_g^2$.

As an initial step to understand the clipping operation, we study a variant of PORTER, which applies the clipping to the mini-batch gradient, i.e. Line 4 is modified to

$$\mathbf{G}^{(t)} = \frac{1}{b} \sum_{\mathbf{Z} \in \mathcal{Z}^{(t)}} \nabla \ell(\mathbf{X}^{(t-1)}; \mathbf{Z}) \quad (5.4a)$$

$$\mathbf{G}_\tau^{(t)} = \text{Clip}_\tau(\mathbf{G}^{(t)}), \quad (5.4b)$$

which is also a widely used operation in deep learning training. However, the privacy guarantee does not apply to this variation anymore. Theorem 11 describes the convergence behavior of the modified version of Algorithm 7 in this case. The proof is deferred to Appendix D.3.

Theorem 11 (Convergence without bounded gradient assumptions). *Assume Assumptions 1, 2, 9 and 10 hold and use general compression operators (Definition 5). Let $\Delta = \mathbb{E}[f(\bar{x}^{(0)})] - f^*$ and $v = O(\rho^{-\frac{2}{3}}(1 - \alpha)^{-\frac{4}{3}} \sigma_g^{1/2} \phi_m^{1/4})$. Set $\gamma = O((1 - \alpha)\rho)$, $\eta = O(L^{-1} \gamma^{\frac{4}{3}} (1 - \alpha)^{\frac{4}{3}} \sqrt{bL\Delta\phi_m}/\tau)$, $b = O(\sigma_g^2 v^{-2})$, $\tau = v$, $\sigma_p^2 = T\tau^2 \phi_m^2/d$, and $T = (b\phi_m)^{-1} = O(\rho^{\frac{4}{3}}(1 - \alpha)^{\frac{8}{3}} \phi_m^{-1})$, Algorithm 7 with mini-batch clipping (5.4) reaches a v -first-order stationary point in no more than T iterations, i.e.,*

$$\min_{t \in [T]} \mathbb{E} \|\nabla f(\bar{x}^{(t)})\|_2 \leq v. \quad (5.5)$$

Theorem 11 establishes a simple framework for analyzing algorithms that use gradient clipping without strong assumptions, which results in a clean deterministic bound of the iterations to reach a ν -stationary point. In contrast, most works that analyze gradient clipping require stronger assumptions. We believe our technique can be applied to analyze the original version of PORTER too, which will be left to future work.

5.6 Numerical experiments

This section presents numerical experiments on the a9a dataset to numerically examine the performance of PORTER, with comparison to the state-of-the-art server/client private optimization algorithm SoteriaFL-SGD, which also utilizes communication compression and guarantees local differential privacy. More specifically, we evaluate the convergence of utility and accuracy in terms of communication rounds and communication bits, to analyze the privacy-utility-communication trade-offs of different optimization algorithms.

For all experiments, we split shuffled datasets evenly to 10 agents that are connected by an Erdős-Rényi random graph with connecting probability $p = 0.8$. We use the FDLA matrix [XB04] as the mixing matrix to perform weighted information aggregation to accelerate convergence. We use biased random sparsification (cf. Example 1) for all algorithms where $k = \lfloor \frac{d}{20} \rfloor$, i.e., the compressor randomly selects 5% elements from each vector. We also apply gradient clipping with $\tau = 1$ to all algorithms for simplicity. For each experiment, all algorithms are initialized to the same starting point, and use best-tuned learning rates, batch size 1 and $\sigma_p = \frac{\tau \sqrt{T \log(1/\delta)}}{mc}$.

5.6.1 Logistic regression with nonconvex regularization

We run experiments on logistic regression with nonconvex regularization [WJZ⁺19] on the a9a dataset [CL11] for different privacy settings. Similar to Section 3.3.1, following [WJZ⁺19], the objective function is defined as

$$\ell(\mathbf{x}; \{f, l\}) = \log \left(1 + l \exp(-\mathbf{x}^\top \mathbf{f}) \right) + \lambda \sum_{i=1}^d \frac{x_i^2}{1 + x_i^2},$$

where $\{f, l\}$ represents a training tuple, $\mathbf{f} \in \mathbb{R}^d$ is the feature vector and $l \in \{0, 1\}$ is the label, and λ is the regularization parameter which is set to $\lambda = 0.2$.

Figure 5.2 and Figure 5.3 show the convergence results of PORTER and SoteriaFL-SGD for logistic regression with nonconvex regularization on the a9a dataset to reach $(10^{-2}, 10^{-3})$ -LDP and $(10^{-1}, 10^{-3})$ -LDP, respectively. Under $(10^{-2}, 10^{-3})$ -LDP, which is a more stricter privacy setting, PORTER converges faster than SoteriaFL-SGD in test accuracy, while PORTER converges slightly slower in train utility. Under $(10^{-1}, 10^{-3})$ -LDP, PORTER performs SoteriaFL-SGD slightly worse than SoteriaFL-SGD. The results highlight PORTER's

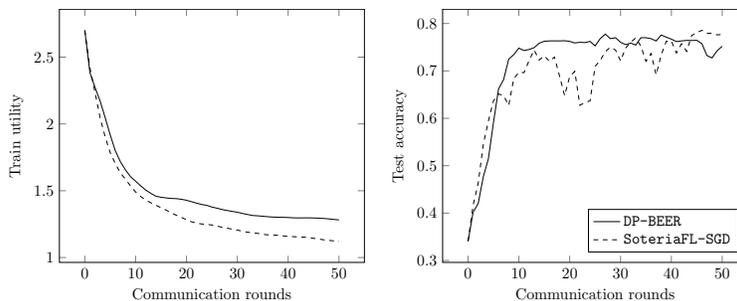


Figure 5.2: The train utility and test accuracy vs. communication rounds for logistic regression with nonconvex regularization on the a9a dataset when guaranteeing $(10^{-2}, 10^{-3})$ -LDP. Both PORTER and SoteriaFL-SGD employ $\text{random}_{162, \text{biased}}$ compression (cf. Example 1).

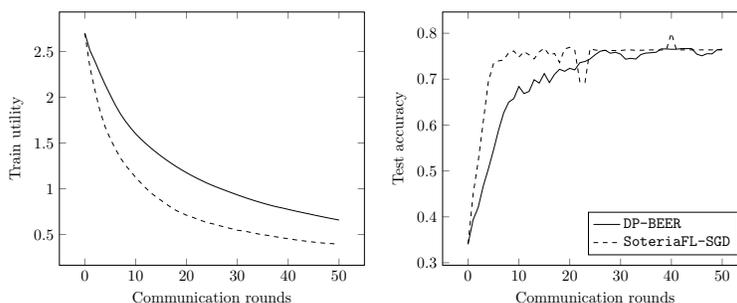


Figure 5.3: The train utility and test accuracy vs. communication rounds for logistic regression with nonconvex regularization on the a9a dataset when guaranteeing $(10^{-1}, 10^{-3})$ -LDP. Both PORTER and SoteriaFL-SGD employ $\text{random}_{162, \text{biased}}$ compression (cf. Example 1).

communication efficiency by showing it can achieve similar performance as its server/client counterpart, i.e. SoteriaFL-SGD, especially under strict privacy constraints.

5.6.2 One-hidden-layer neural network training

Similar to Section 4.4.2, we evaluate PORTER by training a one-hidden layer neural network on the MNIST dataset [LJB⁺95]. The network uses 64 hidden neurons, sigmoid activation functions and cross-entropy loss, where the loss function over a training pair $\{f, l\}$ is defined as

$$\ell(x; (f, l)) = \text{CrossEntropy}(\text{softmax}(\mathbf{W}_2 \text{sigmoid}(\mathbf{W}_1 f + \mathbf{c}_1) + \mathbf{c}_2), l).$$

Here the model parameter is defined by $x = \text{vec}(\mathbf{W}_1, \mathbf{c}_1, \mathbf{W}_2, \mathbf{c}_2)$, where the dimensions of the network parameters $\mathbf{W}_1, \mathbf{c}_1, \mathbf{W}_2, \mathbf{c}_2$ are $64 \times 784, 64 \times 1, 10 \times 64$, and 10×1 , respectively.

Figures 5.4 and 5.5 show the convergence results of PORTER and SoteriaFL-SGD for training a one-hidden-layer neural network on the MNIST dataset to reach $(10^{-2}, 10^{-3})$ -LDP and $(10^{-1}, 10^{-3})$ -LDP, respectively. Under both privacy settings, PORTER converges at a similar rate as SoteriaFL-SGD in train utility. However, in terms of convergence in test accuracy, PORTER outperforms SoteriaFL-SGD

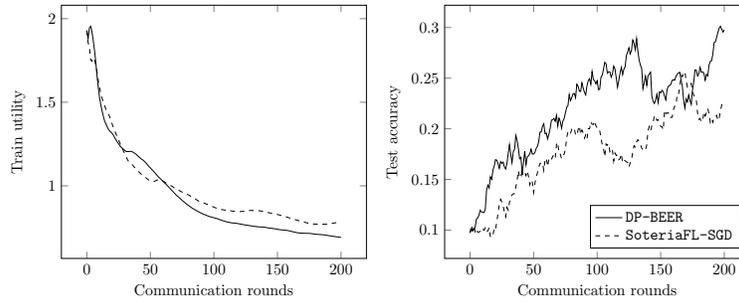


Figure 5.4: The train utility and test accuracy vs. communication rounds for training a one-hidden-layer neural network on the MNIST dataset when guaranteeing $(10^{-2}, 10^{-3})$ -LDP. Both PORTER and SoteriaFL-SGD employ $\text{random}_{2583, \text{biased}}$ compression (cf. Example 1).

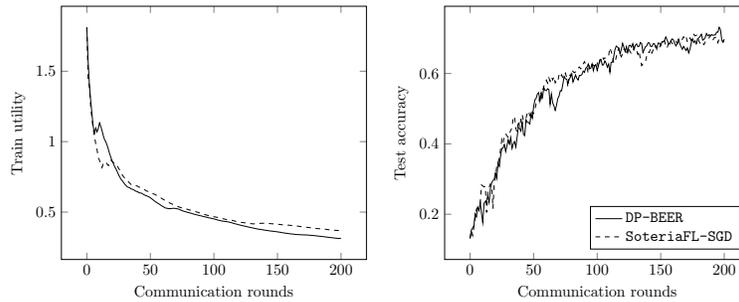


Figure 5.5: The train utility and test accuracy vs. communication rounds for training a one-hidden-layer neural network on the MNIST dataset when guaranteeing $(10^{-1}, 10^{-3})$ -LDP. Both PORTER and SoteriaFL-SGD employ $\text{random}_{2583, \text{biased}}$ compression (cf. Example 1).

under the more stricter $(10^{-2}, 10^{-3})$ -LDP, while the two algorithms have similar performance under the other setting. This experiment again empathizes PORTER’s communication efficiency by comparing to the server/client algorithm SoteriaFL-SGD.

Chapter 6

Conclusions

The main focuses of this thesis are theoretical and practical communication and computation efficiency of decentralized optimization algorithms, where these algorithms leverage techniques such as gradient tracking, Newton-type updates, stochastic recursive gradients, communication compression and differential privacy.

We start from the *Network-DANE* algorithm in Chapter 2, which adapts Newton-type updates to the decentralized setting using gradient tracking and extra mixing, and achieves a condition number free (up to a log factor) communication complexity for quadratic objectives. *Network-DANE* indicates that it is possible to perform more local computation, i.e., solving a strongly convex optimization local problem at each agent, and more communication per iteration, i.e., extra mixing, to reduce overall communication complexity.

Next, we propose *DESTRESS*, a stochastic recursive gradient algorithm, in Chapter 3. *DESTRESS* employs recursive stochastic gradients, stochastic gradient tracking and extra mixing, which results in an optimal IFO-complexity for arbitrary nonconvex ERM problems and improved communication complexity upon existing works. *DESTRESS* reaches both communication and computation efficiency at the same time, which emphasizes the efficacy of gradient tracking, extra mixing and stochastic algorithms.

Then, we propose *BEER* and *PORTER* in Chapters 4 and 5, respectively. While these two algorithms share a similar structure that uses gradient tracking, communication compression and error feedback, they focus on different perspectives. *BEER* achieves an improved communication complexity but at the cost of worse IFO complexity, while being able to converge under high data heterogeneity. On top of *BEER*, *PORTER* adds gradient clipping and privacy perturbation to gradients to protect the privacy of each agent, with explicit utility and communication complexity. The development of *BEER* and *PORTER* shows that by using gradient tracking and communication compression, the efficiency of decentralized optimization algorithms can be significantly improved, and these algorithms can be easily extended to suit new problems, e.g.,

privacy-preserving decentralized optimization.

Appendix A

Appendix for Chapter 2

A.1 Derivation of Equation (2.6)

We observe that

$$f_j(\mathbf{x}) - \langle \nabla f_j(\bar{\mathbf{x}}^{(t)}), \mathbf{x} \rangle = \frac{1}{2} \mathbf{x}^\top \mathbf{H}_j \mathbf{x} - \mathbf{x}^\top \mathbf{H}_j \bar{\mathbf{x}}^{(t)} + \text{constant} = \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}^{(t)})^\top \mathbf{H}_j (\mathbf{x} - \bar{\mathbf{x}}^{(t)}) + \text{constant},$$

which allows us to derive a closed-form expression for $\mathbf{x}_j^{(t)}$ as follows

$$\begin{aligned} \mathbf{x}_j^{(t)} &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}^{(t)})^\top \mathbf{H}_j (\mathbf{x} - \bar{\mathbf{x}}^{(t)}) + \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \mathbf{x} - \bar{\mathbf{x}}^{(t)} \rangle + \frac{\mu}{2} \|\mathbf{x} - \bar{\mathbf{x}}^{(t)}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}^{(t)})^\top (\mathbf{H}_j + \mu \mathbf{I}) (\mathbf{x} - \bar{\mathbf{x}}^{(t)}) + \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \mathbf{x} - \bar{\mathbf{x}}^{(t)} \rangle \right\} \\ &= \bar{\mathbf{x}}^{(t)} - (\mathbf{H}_j + \mu \mathbf{I}_d)^{-1} \nabla f(\bar{\mathbf{x}}^{(t)}). \end{aligned}$$

A.2 Proof of Theorem 1 and Theorem 2

This sections proves the convergence rate of Network-DANE for quadratic losses. When local and global loss functions are quadratic, we can solve (2.9) explicitly. Specifically, Algorithm 1 can be alternatively written as Algorithm 9 below.

For simplicity, we first proceed with the proof as if $K = 1$, then replace α with α^K in the last step, and we let $\bar{\mathbf{H}} = \nabla^2 f(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n \mathbf{H}_j$ be the Hessian of the global loss function. From the definition of the homogeneity parameter β , we have $\|\bar{\mathbf{H}} - \mathbf{H}_j\|_2 \leq \beta$ for all $j = 1, \dots, n$. In addition, we recall the notation in Definition 2, (2.3) and (2.4), and define the error vector as follows

$$\mathbf{e}^{(t)} = \begin{bmatrix} \sqrt{n} \|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2 \\ \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2 \\ L^{-1} \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 \end{bmatrix}. \quad (\text{A.2})$$

Algorithm 9 Network-DANE for quadratic losses (2.10)1 **for** $t = 1, 2, \dots$ **do**

2

$$\mathbf{y}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d) \mathbf{x}^{(t-1)}, \quad (\text{A.1a})$$

$$\mathbf{s}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d) \mathbf{s}^{(t-1)} + \mathbf{H}(\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}), \quad (\text{A.1b})$$

$$\mathbf{x}^{(t)} = \mathbf{y}^{(t-1)} - (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)}, \quad (\text{A.1c})$$

where $\mathbf{y}^{(t)}$ and $\mathbf{s}^{(t)}$ are defined in (2.3), $\mathbf{H} := \text{diag}(\mathbf{H}_1, \dots, \mathbf{H}_n) \in \mathbb{R}^{nd \times nd}$, and \mathbf{H}_i is defined in (2.10).

Establishing the convergence of Network-DANE relies on characterization of the per-iteration dynamics of $\mathbf{e}^{(t)}$ for quadratic losses. Towards this end, we state the following key lemma — which is established in Appendix A.5 — that plays a crucial role in the analysis.

Lemma 1. Let $\eta = \frac{1}{\sigma + \mu}$ and $\gamma = \frac{L}{L + \mu}$. Suppose that Assumptions 3 and 4 hold. Then one has

$$\mathbf{e}^{(t)} \leq \underbrace{\begin{bmatrix} \theta_1 & \gamma\eta\beta + \eta\beta & \eta^2 L\beta \\ \alpha\gamma\eta\beta & \alpha + \alpha\eta L & \alpha\eta L \\ \frac{\beta}{L} + \theta_1 \frac{\beta}{L} + \alpha\gamma\eta\beta \frac{\beta}{L} & \alpha \frac{\beta}{L} + \alpha + 1 + \gamma\eta\beta \frac{\beta}{L} + \eta\beta \frac{\beta}{L} + \alpha \frac{\beta}{L} + \alpha\eta\beta & \alpha + \gamma\eta\beta \frac{\beta}{L} + \alpha\eta\beta \end{bmatrix}}_{=: \mathbf{G}} \mathbf{e}^{(t-1)}. \quad (\text{A.3})$$

Here, $\mathbf{a} \leq \mathbf{b}$ indicates that $a_i \leq b_i$ for all entries i .

In what follows, we invoke this result to establish Theorems 1 and 2 separately.

A.2.1 Proof of Theorem 1

By the choice of μ stated in Theorem 1, we can show that

$$\gamma < 1 \quad \text{and} \quad \eta\beta \leq \eta L < 1. \quad (\text{A.4})$$

In view of Lemma 1, we can obtain

$$\mathbf{e}^{(t)} \leq \mathbf{G}_1 \mathbf{e}^{(t-1)}$$

with a simplified matrix

$$\mathbf{G}_1 := \begin{bmatrix} \theta_1 & 2\eta\beta & \eta^2 L\beta \\ \alpha\gamma\eta\beta & \alpha + \alpha\eta L & \alpha\eta L \\ 3\frac{\beta}{L} & 7 & \alpha + 2\eta\beta \end{bmatrix}, \quad (\text{A.5})$$

where $\mathbf{e}^{(t)}$ is defined in (A.2). We first invoke an argument from [WYWH18] to show that $\mathbf{e}^{(t)}$ converges linearly at a rate not exceeding $\rho(\mathbf{G}_1)$. Given that \mathbf{G}_1 is a positive matrix (i.e. all of its entries are strictly

greater than zero), one can invoke the Perron-Frobenius Theorem to show that: there exists a real-valued positive number $\rho(\mathbf{G}_1) \in \mathbb{R}$ — the spectral radius of \mathbf{G}_1 — such that (i) $\rho(\mathbf{G}_1)$ is an algebraically simple eigenvalue of \mathbf{G}_1 associated with a strictly positive eigenvector χ , (ii) all other eigenvalues of \mathbf{G}_1 are strictly smaller in magnitude than $\rho(\mathbf{G}_1)$. Therefore, there exists some constant $C > 0$ such that $e_0 \leq C\chi$, and consequently,

$$e^{(1)} \leq \mathbf{G}_1 e^{(0)} \leq C\mathbf{G}_1 \chi = C\rho(\mathbf{G}_1)\chi. \quad (\text{A.6})$$

Invoking this argument recursively for all t , we arrive at

$$e^{(t)} \leq C(\rho(\mathbf{G}_1))^t \chi. \quad (\text{A.7})$$

Therefore, the rest of this proof boils down to upper bounding $\rho(\mathbf{G}_1)$. Rearrange the characteristic polynomial of \mathbf{G}_1 , given by

$$\begin{aligned} f_1(\lambda) &= \det(\lambda \mathbf{I} - \mathbf{G}_1) \\ &= (\lambda - \theta_1)p_1(\lambda) + \alpha\gamma\eta^2\beta^2(2\alpha + 4\eta\beta - 2\theta_1 - 7\eta L) - 3\eta^2\beta^2(\alpha - \alpha\eta L + \theta_1), \end{aligned} \quad (\text{A.8})$$

where $p_1(\lambda)$ is the following function obtained by direct computation

$$p_1(\lambda) = (\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) - 7\alpha\eta L - 2\alpha\gamma\eta^2\beta^2 - 3\eta^2\beta^2. \quad (\text{A.9})$$

From the Perron-Frobenius Theorem, we know that $\rho(\mathbf{G}_1)$ is a simple positive root of $f_1(\lambda)$ (so that $f_1(\rho(\mathbf{G}_1)) = 0$). However, it is difficult to compute it directly. In what follows, we seek to first upper bound $\rho(\mathbf{G}_1)$ by

$$\rho_1 := \lambda_0 = \max \left\{ \frac{1 + \theta_1}{2}, \alpha + \frac{140\eta L}{1 - \alpha} \left(\frac{\beta}{\sigma} + 1 \right), \frac{1 + \alpha}{2} + 2\eta\beta \right\}, \quad (\text{A.10})$$

then demonstrate that $\lambda_0 < 1$, which ensures linear convergence, and finally replace α with α^K for the final result.

Step 1: bounding $\rho(\mathbf{G}_1)$ by λ_0 . The following calculation aims to verify the fact that: for all $\lambda \geq \lambda_0$, one has $f_1(\lambda) > 0$, and hence $\rho(\mathbf{G}_1) \leq \lambda_0$. Recall the definition of θ_1 in (2.13). When $\lambda \geq \lambda_0 \geq \frac{1 + \theta_1}{2}$, one has

$$\begin{aligned} \lambda - \theta_1 &\geq \frac{1 - \theta_1}{2} \\ &= \frac{1}{2} \frac{\sigma}{\sigma + \mu} \left(1 - \frac{L}{L + \mu} \frac{\beta}{\sigma + \mu - \beta} \frac{\beta}{\sigma} \right) \\ &\geq \frac{1}{4} \frac{\sigma}{\sigma + \mu}. \end{aligned} \quad (\text{A.11})$$

In order for the last inequality to hold, we must make sure that

$$\begin{cases} \sigma + \mu \geq \frac{3\beta^2}{\sigma}, & \text{if } \beta \geq \sigma; \\ \sigma + \mu \geq 3\sigma, & \text{otherwise.} \end{cases} \quad (\text{A.12})$$

Note that the above relationship is guaranteed by the condition $\sigma + \mu \geq \frac{140L}{(1-\alpha)^2} \left(\frac{\beta}{\sigma} + 1\right)$. When $\lambda \geq \lambda_0$, using (A.4), we can lower bound the first term of $p_1(\lambda)$ by

$$\begin{aligned} (\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) &\geq \frac{1-\alpha}{2} \left(\frac{140\eta L}{1-\alpha} \left(\frac{\beta}{\sigma} + 1\right) - \alpha\eta L\right) \\ &> 69\eta L \left(\frac{\beta}{\sigma} + 1\right). \end{aligned}$$

We can lower bound $p_1(\lambda)$ by incorporating (A.4) as

$$\begin{aligned} p_1(\lambda) &= (\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) - 7\alpha\eta L - 2\alpha\gamma\eta^2\beta^2 - 3\eta^2\beta^2 \\ &> 69\eta L \left(\frac{\beta}{\sigma} + 1\right) - 12\eta L \\ &> 68\kappa\eta\beta. \end{aligned} \quad (\text{A.13})$$

As a result of (A.11) and (A.13), when $\lambda \geq \lambda_0$, the characteristic polynomial (A.8) satisfies

$$\begin{aligned} f_1(\lambda) &\geq (\lambda - \theta_1)p_1(\lambda) + \alpha\gamma\eta^2\beta^2(2\alpha + 4\eta\beta - 2\theta_1 - 7\eta L) - 3\eta^2\beta^2(\alpha - \alpha\eta L + \theta_1) \\ &> \frac{1}{4}\eta\sigma \cdot 68\kappa\eta\beta - 9\alpha\gamma\eta^2\beta^2 - 3\eta^2\beta^2(\alpha + \theta_1) \\ &> 17\eta\beta\eta L - 9\alpha\gamma\eta^2\beta^2 - 6\eta^2\beta^2 > 0. \end{aligned}$$

Therefore, any λ that exceeds λ_0 cannot be a root of $f_1(\cdot)$. This implies that the spectral radius $\rho(\mathbf{G}_1)$, of necessity, obeys $\rho(\mathbf{G}_1) < \lambda_0$.

Step 2: bounding λ_0 . This step verifies that all three terms in (A.10) are smaller than 1, thus leading to the conclusion $\lambda_0 < 1$.

- First, observe that if (A.12) is satisfied, we have $\frac{1+\theta_1}{2} \leq 1 - \frac{1}{4}\eta\sigma < 1$.
- When $\sigma + \mu \geq \frac{140L}{(1-\alpha)^2} \left(\frac{\beta}{\sigma} + 1\right)$, the second term in (A.10) obeys $\alpha + \frac{140\eta L}{1-\alpha} \left(\frac{\beta}{\sigma} + 1\right) \leq 1$.
- Finally, the third term in (A.10) is also less than 1, since

$$\frac{1+\alpha}{2} + 2\eta\beta \leq \frac{1+\alpha}{2} + \frac{(1-\alpha)^2}{70} \frac{\beta}{\frac{\beta}{\sigma} + 1} \frac{1}{L} \leq \frac{1+\alpha}{2} + \frac{(1-\alpha)^2}{70} \leq 1 - \frac{1-\alpha}{2} + \frac{1-\alpha}{70} < 1.$$

Step 3: replacing α with α^K . This step gives the final results as

$$\rho_1 = \max \left\{ \frac{1+\theta_1}{2}, \alpha^K + \frac{140\eta L}{1-\alpha^K} \left(\frac{\beta}{\sigma} + 1\right), \frac{1+\alpha^K}{2} + 2\eta\beta \right\}.$$

A.2.2 Proof of Theorem 2

By the assumption $\sigma + \mu \geq 360\sigma \left(\frac{\beta^2}{\sigma^2} + 1 \right)$ and $\alpha^K \leq \frac{1}{2\kappa}$ as we can prove that $\eta\beta < 1$ and $\alpha^K\eta L \leq \frac{1}{2}$. The characteristic polynomial (A.8) in Appendix A.2.1 can then be lower bounded by

$$\begin{aligned}
f_1(\lambda) &= \det(\lambda I - \mathbf{G}_1) \\
&= (\lambda - \theta_1) \left((\lambda - \alpha^K - \alpha^K\eta L)(\lambda - \alpha^K - 2\eta\beta) - 7\alpha^K\eta L - 2\alpha^K\gamma\eta^2\beta^2 - 3\eta^2\beta^2 \right) \\
&\quad + \alpha^K\gamma\eta^2\beta^2(2\alpha^K + 4\eta\beta - 2\theta_1 - 7\eta L) - 3\eta^2\beta^2(\alpha^K - \alpha^K\eta L + \theta_1) \\
&\geq (\lambda - \theta_1) \left((\lambda - \alpha^K - \frac{1}{2}\eta\sigma)(\lambda - \alpha^K - 2\eta\beta) - \frac{7}{2}\eta\sigma - \eta\sigma\eta^2\beta^2 - 3\eta^2\beta^2 \right) \\
&\quad + \alpha^K\gamma\eta^2\beta^2(2\alpha^K + 4\eta\beta - 2\theta_1 - 7\eta L) - 3\eta^2\beta^2(\alpha^K - \alpha^K\eta L + \theta_1), \tag{A.14}
\end{aligned}$$

provided that λ obeys

$$\lambda \geq \max \left\{ \frac{1 + \theta_1}{2}, \alpha^K + 180\eta\sigma \left(\frac{\beta^2}{\sigma^2} + 1 \right), \frac{1 + \alpha^K}{2} + 2\eta\beta \right\}.$$

Given that all conditions in (A.12) are satisfied, we can show $\eta^2\beta^2 \leq \eta\sigma \cdot \frac{\beta^2}{360\sigma^2(\beta^2/\sigma^2 + 1)} < \eta\sigma < 1$. One can thus continue to lower bound (A.14) by

$$\begin{aligned}
f_1(\lambda) &> (\lambda - \theta_1) \left((\lambda - \alpha^K - \frac{1}{2}\eta\sigma)(\lambda - \alpha^K - 2\eta\beta) - 8\eta\sigma \right) - 11\eta^2\beta^2 \\
&> \frac{1}{4}\eta\sigma \left\{ \frac{1}{4} \left[180\eta\sigma \left(\frac{\beta^2}{\sigma^2} + 1 \right) - \frac{1}{2}\eta\sigma \right] - 8\eta\sigma \right\} - 11\eta^2\beta^2 \\
&> \frac{1}{4}\eta\sigma \left\{ 45\eta\beta \frac{\beta}{\sigma} + 44\eta\sigma - 8\eta\sigma \right\} - 11\eta^2\beta^2 \\
&> \frac{45}{4}\eta\beta - 11\eta^2\sigma^2 \\
&> 0.
\end{aligned}$$

Consequently, following similar arguments as in Appendix A.2.1, we can show that: under the conditions of Theorem 2, the spectral radius of \mathbf{G}_1 can be upper bounded by

$$\rho(\mathbf{G}_1) \leq 1 - \frac{C}{\frac{\beta^2}{\sigma^2} + 1},$$

where C is some sufficiently small positive constant. This immediately tells us that: to reach ϵ -accuracy, Network-DANE takes at most $O\left(\left(\frac{\beta^2}{\sigma^2} + 1\right) \log(1/\epsilon)\right)$ iterations. For each iteration, Network-DANE needs

$$K \asymp \frac{\log(1/2\kappa)}{\log \alpha} \lesssim \frac{\log \kappa}{1 - \alpha}$$

rounds of communication, where we have used the elementary inequality $1 - \alpha < \log(1/\alpha)$. Putting all this together leads to a communication complexity at most $O\left(\log \kappa \cdot \frac{(\beta^2/\sigma^2 + 1) \log(1/\epsilon)}{1 - \alpha}\right)$.

A.3 Proofs of Theorem 3 and Theorem 4

This sections establishes the convergence rate of Network-DANE for smooth and strongly convex loss functions, following the analysis approach adopted in the proof of Theorem 1. Similarly, we first proceed with the proof as if $K = 1$, then replace α with α^K in the last step. The following key lemma plays a crucial role, which characterizes the per-iteration dynamics of the proposed Network-DANE for general smooth strongly convex losses. The proof of this lemma is deferred to Appendix A.6.

Lemma 2. *Recall the notation in Lemma 1. Suppose that Assumption 3 holds, and $(\frac{\beta}{\sigma+\mu})^2 \leq \frac{\sigma}{\sigma+2\mu}$. One has*

$$e^{(t)} \leq \underbrace{\begin{bmatrix} \theta_2 & \eta L & \gamma \eta L \\ \alpha \gamma \eta L & \alpha + \alpha \eta L & \alpha \eta L \\ \frac{\beta}{L} + \theta_2 \frac{\beta}{L} + \alpha \gamma \eta \beta & \alpha + 1 + \alpha \frac{\beta}{L} + \eta \beta + \alpha \frac{\beta}{L} + \alpha \eta \beta & \alpha + \gamma \eta \beta + \alpha \eta \beta \end{bmatrix}}_{=: G'} e^{(t-1)}. \quad (\text{A.15})$$

Here, $e^{(t)}$ is the error vector defined in (A.2), and the notation $\mathbf{a} \leq \mathbf{b}$ indicates that $a_i \leq b_i$ for all entries i .

A.3.1 Proof of Theorem 3

Under the conditions of Theorem 3, the inequalities stated in (A.4) remain valid. In addition, when $\sigma + \mu = \frac{170\kappa L}{(1-\alpha)^2}$, we can verify that

$$\left(\frac{\beta}{\sigma + \mu}\right)^2 = \frac{(1-\alpha)^4 \beta^2}{170^2 \kappa^2 L^2} \leq \frac{(1-\alpha)^2}{170^2 \kappa^2} < \frac{1}{2} \cdot \frac{(1-\alpha)^2}{170 \kappa^2} = \frac{1}{2} \cdot \frac{\sigma}{\sigma + \mu} < \frac{\sigma}{\sigma + 2\mu}.$$

When $\sigma + \mu \geq \frac{170\kappa L}{(1-\alpha)^2}$, the LHS decreases faster than the RHS, thus the requirement of Lemma 2 is met. In view of Lemma 2 as well as the fact $\theta_2 \leq 1$, we can replace G' by a simplified matrix that dominates G' :

$$\mathbf{G}_2 := \begin{bmatrix} \theta_2 & 2\eta L & \gamma \eta L \\ \alpha \gamma \eta L & \alpha + \alpha \eta L & \alpha \eta L \\ 3\frac{\beta}{L} & 7 & \alpha + 2\eta \beta \end{bmatrix}. \quad (\text{A.16})$$

The above matrix \mathbf{G}_2 is similar to \mathbf{G}_1 in (A.5) in the quadratic case, except that the quantity β in the first two rows of \mathbf{G}_1 is replaced by L (thus leading to a worse convergence rate).

Similar to the proof of Theorem 1, we shall upper bound $\rho(\mathbf{G}_2)$ — the spectral radius of \mathbf{G}_2 . To locate the eigenvalues of \mathbf{G}_2 , we rearrange the characteristic polynomial of \mathbf{G}_2 as follows

$$\begin{aligned} f_2(\lambda) &= \det(\lambda \mathbf{I} - \mathbf{G}_2) \\ &= (\lambda - \theta_2) p_2(\lambda) + \alpha \gamma \eta^2 L^2 (2\alpha + 4\eta \beta - 2\theta_2 - 7\gamma) - 3\eta \beta (2\alpha \eta L - \gamma(\alpha + \alpha \eta L - \theta_2)), \end{aligned} \quad (\text{A.17})$$

where $p_2(\lambda)$ is the following function obtained by direct computation

$$p_2(\lambda) = (\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) - 7\alpha\eta L - 2\alpha\gamma\eta^2 L^2 - 3\gamma\eta\beta.$$

From the Perron-Frobenius Theorem, $\rho(\mathbf{G}_2)$ is a simple positive root of the equation $f_2(\lambda) = 0$. However, it is hard to calculate it directly. In what follows, we seek to first upper bound $\rho(\mathbf{G}_2)$ by

$$\rho_2 := \lambda_0 = \max \left\{ \frac{1 + \theta_2}{2}, \alpha + \frac{170\kappa\eta L}{1 - \alpha}, \frac{1 + \alpha}{2} + 2\eta\beta \right\}, \quad (\text{A.18})$$

and then demonstrate that $\lambda_0 < 1$, which in turn ensures linear convergence.

Step 1: bounding $\rho(\mathbf{G}_2)$ by λ_0 . The following calculation aims to verify the fact that $f_2(\lambda) > 0$ holds for all $\lambda \geq \lambda_0$, so that $\rho(\mathbf{G}_2) \leq \lambda_0$. Recalling the definition of θ_2 in Lemma 2, we see that when $\lambda \geq \lambda_0 \geq \frac{1 + \theta_2}{2}$,

$$\begin{aligned} \lambda - \theta_2 &\geq \frac{1 - \theta_2}{2} \\ &= \frac{1}{2}\eta \left(\sigma - \beta \sqrt{(1 - \eta\mu)(1 + \eta\mu)} \right) \\ &\geq \frac{1}{2}\eta \left(\sigma - \beta \sqrt{2(1 - \eta\mu)} \right) > \frac{1}{4}\eta\sigma, \end{aligned} \quad (\text{A.19})$$

where we have used the fact $\eta\mu < 1$ to reach the second inequality. For the last inequality to hold, we need to make sure

$$\begin{cases} \sigma + \mu \geq \frac{10\beta^2}{\sigma}, & \beta \geq \sigma \\ \sigma + \mu \geq 10\sigma, & \text{otherwise} \end{cases} \quad (\text{A.20})$$

which is guaranteed by the assumption $\sigma + \mu \geq \frac{170\kappa L}{(1 - \alpha)^2}$.

Similarly, when $\lambda \geq \lambda_0$, the first term of $p_2(\lambda)$ can be lower bounded by

$$(\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) \geq \frac{1 - \alpha}{2} \left(\frac{170\kappa\eta L}{1 - \alpha} - \alpha\eta L \right) > 80\kappa\eta L.$$

Then, using (A.4) we can bound $p_2(\lambda)$ by

$$\begin{aligned} p_2(\lambda) &= (\lambda - \alpha - \alpha\eta L)(\lambda - \alpha - 2\eta\beta) - 7\alpha\eta L - 2\alpha\gamma\eta^2 L^2 - 3\gamma\eta\beta \\ &> 80\kappa\eta L - 12\eta L \geq 68\kappa\eta L. \end{aligned} \quad (\text{A.21})$$

By virtue of (A.19) and (A.21), it is seen that when $\lambda \geq \lambda_0$, the characteristic polynomial $f_2(\lambda)$ in (A.17) satisfies

$$f_2(\lambda) > \frac{1}{4}\eta\sigma \cdot 68\kappa\eta L - 8\alpha\gamma\eta^2 L^2 - 9\eta\beta\eta L > 0.$$

Therefore, any λ that exceeds λ_0 cannot possibly be a root of $f_2(\cdot)$. This implies that the spectral radius necessarily obeys $\rho(\mathbf{G}_2) < \lambda_0$.

Step 2: bounding λ_0 . This step verifies that the three terms in the expression of λ_0 in (A.18) is smaller than 1, allowing us to conclude that $\lambda_0 < 1$.

- First, observe that if (A.20) is satisfied, then we have $\frac{1+\theta_2}{2} \leq 1 - \frac{1}{4}\eta\sigma < 1$.
- When $\sigma + \mu \geq \frac{170\kappa L}{(1-\alpha)^2}$, the second term is $\alpha + \frac{170\kappa\eta L}{1-\alpha} \leq 1$.
- We conclude the proof by checking that the third term is also less than 1, namely,

$$\frac{1+\alpha}{2} + 2\eta\beta \leq \frac{1+\alpha}{2} + \frac{(1-\alpha)^2}{85} \frac{1}{\kappa L} \beta \leq \frac{1+\alpha}{2} + \frac{(1-\alpha)^2}{85} \leq 1 - \frac{1-\alpha}{2} + \frac{1-\alpha}{85}.$$

Step 3: replacing α with α^K . This step gives the final results as

$$\rho_2 = \max \left\{ \frac{1+\theta_2}{2}, \alpha^K + \frac{170\kappa\eta L}{1-\alpha^K}, \frac{1+\alpha^K}{2} + 2\eta\beta \right\}.$$

A.3.2 Proof of Theorem 4

We first verify the assumption of Lemma 2. When $\sigma + \mu = 360L \left(\frac{\beta}{\sigma} + 1 \right)$,

$$\left(\frac{\beta}{\sigma + \mu} \right)^2 = \frac{\beta^2}{360^2 L^2 \left(\frac{\beta}{\sigma} + 1 \right)^2} \leq \frac{\beta}{360^2 \kappa L \left(\frac{\beta}{\sigma} + 1 \right)} < \frac{1}{2} \cdot \frac{1}{360\kappa \left(\frac{\beta}{\sigma} + 1 \right)} = \frac{1}{2} \cdot \frac{\sigma}{\sigma + \mu} < \frac{\sigma}{\sigma + 2\mu}.$$

Therefore, Lemma 2 still holds.

By the assumption $\alpha^K \leq \frac{1}{2\kappa}$, we can further lower bound the characteristic polynomial (A.17) in Appendix A.3.1 as follows:

$$\begin{aligned} f_2(\lambda) &= \det(\lambda I - \mathbf{G}_2) \\ &= (\lambda - \theta_2) \left((\lambda - \alpha^K - \alpha^K \eta L)(\lambda - \alpha^K - 2\eta\beta) - 7\alpha^K \eta L - 2\alpha^K \gamma \eta^2 L^2 - 3\gamma \eta \beta \right) \\ &\quad + \alpha^K \gamma \eta^2 L^2 \left(2\alpha^K + 4\eta\beta - 2\theta_2 - 7\gamma \right) - 3\eta\beta \left(2\alpha^K \eta L - \gamma(\alpha^K + \alpha^K \eta L - \theta_2) \right) \\ &\geq (\lambda - \theta_2) \left((\lambda - \alpha^K - \frac{1}{2}\eta\sigma)(\lambda - \alpha^K - 2\eta\beta) - \frac{7}{2}\eta\sigma - \eta\sigma \eta^2 L^2 - 3\gamma \eta \beta \right) \\ &\quad - \eta\sigma \eta^2 L^2 \left(\theta_2 + \frac{7}{2}\gamma \right) - 3\eta\beta \left(\eta\sigma + \gamma\theta_2 \right) \\ &> (\lambda - \theta_2) \left((\lambda - \alpha^K - \frac{1}{2}\eta\sigma)(\lambda - \alpha^K - 2\eta\beta) - 8\eta\sigma \right) - 5\eta\sigma \eta^2 L^2 - 6\eta\beta \eta L, \end{aligned} \tag{A.22}$$

providing λ obeys

$$\lambda \geq \max \left\{ \frac{1+\theta_2}{2}, \alpha^K + 180\eta L \left(\frac{\beta}{\sigma} + 1 \right), \frac{1+\alpha^K}{2} + 2\eta\beta \right\}.$$

We can further lower bound (A.22) by

$$f_2(\lambda) \geq \frac{1}{4}\eta\sigma \left\{ \frac{1}{4} \left[180\eta L \left(\frac{\beta}{\sigma} + 1 \right) - \frac{1}{2}\eta\sigma \right] - 8\eta\sigma \right\} - 5\eta\sigma \eta^2 L^2 - 6\eta\beta \eta L > 0,$$

as long as μ satisfies $\sigma + \mu \geq 360L(\frac{\beta}{\sigma} + 1)$. Therefore, following similar arguments as adopted in Appendix A.3.1, the spectral radius of G_2 can be upper bounded by

$$\rho(G_2) \leq 1 - \frac{C}{\kappa(\frac{\beta}{\sigma} + 1)},$$

where C is a small positive constant. Consequently, to reach ϵ -accuracy, Network-DANE takes at most $O\left(\kappa(\frac{\beta}{\sigma} + 1) \log(1/\epsilon)\right)$ iterations and $O\left(\log \kappa \cdot \frac{\kappa(\beta/\sigma+1) \log(1/\epsilon)}{1-\alpha}\right)$ communication rounds.

A.4 Proof of Theorem 5

The proof strategy of Theorem 5 is similar in spirit to the convergence proof of Network-DANE, where we will carefully build a linear system that tracks the coupling of the consensus error and the optimization error. Under the assumptions in Theorem 5, we can assume that $1 - 3\alpha\kappa - 3\beta/\sigma > 0$. Let

$$\zeta = 1/(1 - 3\alpha\kappa - 3\beta/\sigma).$$

In what follows, we first introduce two key lemmas that connect the convergence behavior of Network-SVRG in the network setting to their server/client counterparts (namely, D-SVRG) studied in [CZC⁺20]. Lemma 3, proved in Appendix A.7, creates the linear system characterizing the iteration dynamics of Network-SVRG.

Lemma 3. *Under the assumptions in Theorem 5, Network-SVRG satisfies*

$$\mathbb{E}[e^{(t)}] \leq \underbrace{\begin{bmatrix} \left(\nu(1 + 3\alpha\kappa + 4\frac{\beta}{\sigma}) + \frac{\beta}{\sigma}\right) \zeta & 8\frac{\beta}{\sigma}\zeta & \alpha\zeta/\kappa & \zeta/16 \\ 1/2 & 0 & 0 & 0 \\ 8(\frac{\beta}{\sigma})^2 & 64(\frac{\beta}{\sigma})^2 & 4\alpha^2 & \alpha\kappa/2 \\ 64\alpha\kappa & 0 & 0 & 0 \end{bmatrix}}_{:=G_3} \mathbb{E}[e^{(t-1)}], \quad (\text{A.23})$$

where the error vector is defined as

$$e^{(t)} = \begin{bmatrix} \sum_{j=1}^n (f(x_j^{(t)}) - f(\mathbf{y}^*)) \\ \sum_{j=1}^n (f(y_j^{(t)}) - f(\mathbf{y}^*)) / 2 \\ \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2^2 / \sigma \\ 32L \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2 / \alpha \end{bmatrix}.$$

Here, $\nu \leq \frac{1}{2} \frac{\sigma - 2\beta}{\sigma - 3\beta}$ is the convergence rate of D-SVRG in the server/client setting under the same assumptions [CZC⁺20, Theorem 1].

Since every term in the matrices of linear systems of Lemma 3 is non-negative, all eigenvalues of G_3 are bounded by the maximum of the sum of rows according to the Gershgorin circle theorem. For

Network-SVRG, by setting $\alpha = \frac{1}{70\kappa}$, which needs $K \asymp O(\log_{\alpha} 1/\kappa) = O(\log \kappa / (1 - \alpha))$, we can ensure that the sum of the first row is bounded by $5/6$, and the sums of other rows are also bounded by a constant smaller than 1, under the assumption $\beta \leq \sigma/200$. Therefore, invoking the Gershgorin circle theorem, the spectral radius is bounded by a constant smaller than 1. To achieve ε -accuracy, the total number of iterations needed is $O(\log(1/\varepsilon))$ and thus the communication complexity is $O\left(\log \kappa \cdot \frac{\log(1/\varepsilon)}{1-\alpha}\right)$.

A.5 Proof of Lemma 1

The proof is divided into several steps. (i) In Appendix A.5.1, we bound the convergence error $\sqrt{n}\|\bar{\mathbf{x}}^{(t)} - \mathbf{y}^*\|_2$; (ii) in Appendix A.5.2, we bound the parameter consensus error $\|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2$; (iii) in Appendix A.5.3, we bound the gradient estimation error $\|\mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2$; (iv) finally, we create induction inequalities of $\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2$, $\sqrt{n}\|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2$ and $\|\mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2$ in Appendix A.5.4 to conclude the proof.

A.5.1 Convergence error

We begin by defining an auxiliary variable \mathbf{x}_j^+ , which can be seen as the result of one local iterate (2.9) of the original DANE algorithm initialized at $\bar{\mathbf{y}}^{(t-1)}$:

$$\mathbf{x}_j^+ = \operatorname{argmin}_{\mathbf{x}} \left\{ f_j(\mathbf{x}) - \left\langle \nabla f_j(\bar{\mathbf{y}}^{(t-1)}) - \nabla f(\bar{\mathbf{y}}^{(t-1)}), \mathbf{x} \right\rangle + \frac{\mu}{2} \|\mathbf{x} - \bar{\mathbf{y}}^{(t-1)}\|_2^2 \right\}. \quad (\text{A.24})$$

Following the same convention as in previous definitions, we also define

$$\bar{\mathbf{x}}^+ = \frac{1}{n} \sum_j \mathbf{x}_j^+. \quad (\text{A.25})$$

Given that the function we optimize at each agent is strongly convex, the local optimality conditions of (A.24) and (2.9) are as follows:

$$\nabla f_j(\mathbf{x}_j^+) + \mu(\mathbf{x}_j^+ - \mathbf{y}^*) = \nabla(f_j - f)(\bar{\mathbf{y}}^{(t-1)}) + \mu(\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*), \quad (\text{A.26a})$$

$$\nabla f_j(\mathbf{x}_j^{(t-1)}) + \mu(\mathbf{x}_j^{(t-1)} - \mathbf{y}^*) = \nabla f_j(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)} + \mu(\mathbf{y}_j^{(t-1)} - \mathbf{y}^*). \quad (\text{A.26b})$$

Taking the average of (A.26) over $j = 1, \dots, n$, we obtain another set of optimality conditions:

$$\frac{1}{n} \sum_j \nabla f_j(\mathbf{x}_j^+) + \mu(\bar{\mathbf{x}}^+ - \mathbf{y}^*) = \mu(\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*), \quad (\text{A.27a})$$

$$\frac{1}{n} \sum_j \nabla f_j(\mathbf{x}_j^{(t-1)}) + \mu(\bar{\mathbf{x}}^{(t-1)} - \mathbf{y}^*) = \mu(\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*), \quad (\text{A.27b})$$

where we use the fact $\sum_j \mathbf{s}_j^{(t-1)} = \sum_j \nabla f_j(\mathbf{y}_j^{(t-1)})$ due to the property of gradient tracking (2.2).

In view of the triangle inequality, the convergence error can be decomposed as

$$\|\bar{\mathbf{x}}^{(t-1)} - \mathbf{y}^*\|_2 \leq \|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2 + \|\bar{\mathbf{x}}^+ - \mathbf{y}^*\|_2, \quad (\text{A.28})$$

where the first term is the error caused by inaccurate gradient estimate, and the second term is the progress of DANE initialized at $\bar{\mathbf{y}}^{(t-1)}$.

1. For the first term $\|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2$, we first plug in the Hessian of the quadratic losses to solve for $\mathbf{x}_j^{(t-1)}$ and \mathbf{x}_j^+ explicitly as

$$\mathbf{x}_j^{(t-1)} = \mathbf{y}_j^{(t-1)} - (\mathbf{H}_j + \mu \mathbf{I}_d)^{-1} \mathbf{s}_j^{(t-1)}, \quad (\text{A.29a})$$

$$\mathbf{x}_j^+ = \bar{\mathbf{y}}^{(t-1)} - (\mathbf{H}_j + \mu \mathbf{I}_d)^{-1} \nabla f(\bar{\mathbf{y}}^{(t-1)}). \quad (\text{A.29b})$$

The first error term $\|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2$ can be written as

$$\begin{aligned} & \|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2 \\ = & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{x}^{(t-1)} - \mathbf{x}^+) \right\|_2 \\ = & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} - (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) + (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right) \right\|_2 \\ = & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)})) \right\|_2, \end{aligned}$$

where the last line follows from the definition of $\bar{\mathbf{y}}^{(t-1)}$. Then, we add and subtract $(\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1}$ and rearrange terms, obtaining

$$\begin{aligned} & \|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2 \\ = & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} - (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} \right) (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)})) \right. \\ & \left. + \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)})) \right\|_2 \\ = & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}} - \mathbf{H}) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)})) \right. \\ & \left. + \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}} - \mathbf{H}) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)})) \right. \\ & \left. + \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{H} - \mathbf{I}_n \otimes \bar{\mathbf{H}}) (\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \quad (\text{A.30}) \\ \leq & \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \right\|_2 \left\| (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}} - \mathbf{H}) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} \right\|_2 \|\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)})\|_2 \\ & + \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \right\|_2 \left\| (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}} - \mathbf{H}) (\mathbf{I}_{nd} + \mu \mathbf{I}_n \otimes \bar{\mathbf{H}}^{-1})^{-1} \right\|_2 \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 \\ & + \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \right\|_2 \left\| (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{H} - \mathbf{I}_n \otimes \bar{\mathbf{H}}) \right\|_2 \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2. \end{aligned}$$

The last term in (A.30) follows from the identity

$$\left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}))$$

$$\begin{aligned}
&= (\bar{\mathbf{H}} + \mu \mathbf{I}_d)^{-1} \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{s}^{(t-1)} - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)})) \\
&= (\bar{\mathbf{H}} + \mu \mathbf{I}_d)^{-1} \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{H}} \bar{\mathbf{y}}^{(t-1)}) \\
&= (\bar{\mathbf{H}} + \mu \mathbf{I}_d)^{-1} \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{H} \bar{\mathbf{y}}^{(t-1)}) \\
&= \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} \mathbf{H} (\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \\
&= \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{H} - \mathbf{I}_n \otimes \bar{\mathbf{H}}) (\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}).
\end{aligned}$$

Taken together with the identity $\|\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d\|_2 = \frac{1}{\sqrt{n}}$, the assumption $\|\mathbf{H}_j - \bar{\mathbf{H}}\|_2 \leq \beta$, and the bound $\|(\mathbf{H} + \mu \mathbf{I}_{nd})^{-1}\|_2 \leq \frac{1}{\sigma + \mu}$ and $\|(\mathbf{I}_{nd} + \mu \mathbf{I}_n \otimes \bar{\mathbf{H}}^{-1})^{-1}\|_2 \leq \frac{L}{L + \mu}$, we can further bound (A.30) by

$$\begin{aligned}
\sqrt{n} \|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2 &\leq \frac{1}{\sigma + \mu} \frac{\beta}{\sigma + \mu} \|\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)})\|_2 \\
&\quad + \left(\frac{L}{L + \mu} \frac{\beta}{\sigma + \mu} + \frac{\beta}{\sigma + \mu} \right) \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2.
\end{aligned} \tag{A.31}$$

2. Regarding the second term $\|\bar{\mathbf{x}}^+ - \mathbf{y}^*\|_2$, we provide a slightly improved bound compared to [SSZ14].

In view of (A.29b),

$$\begin{aligned}
\|\bar{\mathbf{x}}^+ - \mathbf{y}^*\|_2 &= \left\| \bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^* - \frac{1}{n} \sum_j (\mathbf{H}_j + \mu \mathbf{I}_d)^{-1} \nabla f(\bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\
&= \left\| \left(\mathbf{I} - \frac{1}{n} \sum_{i=1}^n (\mathbf{H}_i + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right) (\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*) \right\|_2 \\
&\leq \left\| \mathbf{I} - \frac{1}{n} \sum_{i=1}^n (\mathbf{H}_i + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right\|_2 \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2.
\end{aligned} \tag{A.32}$$

Then, we use the triangle inequality to break the convergence rate in (A.32) into two parts:

$$\begin{aligned}
&\left\| \mathbf{I} - \frac{1}{n} \sum_{i=1}^n (\mathbf{H}_i + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right\|_2 \\
&\leq \left\| \mathbf{I} - (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right\|_2 + \left\| \frac{1}{n} \sum_{i=1}^n \left((\mathbf{H}_i + \mu \mathbf{I})^{-1} - (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \right) \bar{\mathbf{H}} \right\|_2.
\end{aligned} \tag{A.33}$$

When $\bar{\mathbf{H}} \succeq \sigma \mathbf{I}_d$, it is straightforward to check that the first term of (A.33) is upper bounded by

$$\left\| \mathbf{I} - (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right\|_2 \leq 1 - \frac{\sigma}{\sigma + \mu}.$$

Regarding the second term of (A.33), let $\Delta_i := \mathbf{H}_i - \bar{\mathbf{H}}$ and use the definition of β , one derives

$$\left\| (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \Delta_i \right\|_2 \leq \left\| (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \right\|_2 \cdot \|\Delta_i\|_2 \leq \frac{\beta}{\sigma + \mu} < 1 \tag{A.34}$$

under our hypothesis $\beta < \mu + \sigma$. In addition,

$$\left\| \frac{1}{n} \sum_{i=1}^n \left((\mathbf{H}_i + \mu \mathbf{I})^{-1} - (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \right) \bar{\mathbf{H}} \right\|_2$$

$$= \left\| \frac{1}{n} \sum_{i=1}^n \left(\sum_{m=0}^{\infty} (-1)^m [(\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \Delta_i]^m (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} - (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \right) \bar{\mathbf{H}} \right\|_2 \quad (\text{A.35})$$

$$= \left\| \frac{1}{n} \sum_{i=1}^n \left(\sum_{m=2}^{\infty} (-1)^m [(\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \Delta_i]^m (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \right) \bar{\mathbf{H}} \right\|_2 \quad (\text{A.36})$$

$$\leq \frac{1}{n} \sum_{i=1}^n \sum_{m=2}^{\infty} \|(\bar{\mathbf{H}} + \mu \mathbf{I})^{-1}\|_2^m \cdot \|\Delta_i\|_2^m \cdot \|(I + \mu \bar{\mathbf{H}}^{-1})^{-1}\|_2$$

$$\leq \sum_{m=2}^{\infty} (\sigma + \mu)^{-m} \beta^m \frac{L}{L + \mu} = \frac{L}{L + \mu} \frac{\beta^2}{(\sigma + \mu)(\sigma + \mu - \beta)}.$$

Here, the line (A.35) is an expansion based on the Neumann series (whose convergence is guaranteed by (A.34))

$$(\mathbf{H}_i + \mu \mathbf{I})^{-1} = (\bar{\mathbf{H}} + \mu \mathbf{I} + \Delta_i)^{-1} = (\mathbf{I} + (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \Delta_i)^{-1} (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1}$$

$$= \left\{ \sum_{m=0}^{\infty} (-1)^m [(\bar{\mathbf{H}} + \mu \mathbf{I})^{-1} \Delta_i]^m \right\} (\bar{\mathbf{H}} + \mu \mathbf{I})^{-1}.$$

The identity (A.36) holds since $\sum_{i=1}^n \Delta_i = \mathbf{0}$, and hence the summation in (A.36) effectively starts at $m = 2$.

Putting the above two bounds together back in (A.33), we arrive at

$$\left\| \mathbf{I} - \frac{1}{n} \sum_{i=1}^n (\mathbf{H}_i + \mu \mathbf{I})^{-1} \bar{\mathbf{H}} \right\|_2 \leq \theta_1 = 1 - \frac{\sigma}{\sigma + \mu} + \frac{L}{L + \mu} \frac{\beta^2}{(\sigma + \mu)(\sigma + \mu - \beta)}. \quad (\text{A.37})$$

Putting together (A.31) and (A.37), and plugging back into (A.28), we can bound the convergence error by:

$$\sqrt{n} \|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2 = \sqrt{n} \|\bar{\mathbf{x}}^{(t-1)} - \mathbf{y}^*\|_2$$

$$\leq \theta_1 \sqrt{n} \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2 + \frac{1}{\sigma + \mu} \frac{\beta}{\sigma + \mu} \|\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)})\|_2$$

$$+ \left(\frac{L}{L + \mu} \frac{\beta}{\sigma + \mu} + \frac{\beta}{\sigma + \mu} \right) \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2. \quad (\text{A.38})$$

A.5.2 Consensus error

Using the identity $\bar{\mathbf{y}}^{(t)} = \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \mathbf{y}^{(t)}$ and the update rule (A.1c), we can demonstrate that

$$\left\| \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)} \right\|_2$$

$$= \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \mathbf{y}^{(t)} \right\|_2$$

$$= \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{W} \otimes \mathbf{I}_d) \left(\mathbf{y}^{(t-1)} - (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right) \right\|_2$$

$$\leq \left\| \left(\mathbf{W}^K - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right\|_2 \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} - \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right) \right\|_2 \quad (\text{A.39})$$

$$\leq \alpha \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 + \alpha \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right\|_2, \quad (\text{A.40})$$

where (A.39) is due to the following equality:

$$\left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{W} \otimes \mathbf{I}_d) = \left[\left(\mathbf{W}^K - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right),$$

which holds because the property of the averaging operator $\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right)$,

$$\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) = \left[\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \left(\mathbf{I}_d - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \right] \otimes \mathbf{I}_n = \mathbf{0},$$

and the fact that $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$.

We rearrange the second term in (A.40) as

$$\begin{aligned} & \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right\|_2 \\ &= \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) \right. \\ & \quad + \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \left(\nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right) \\ & \quad \left. + \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \left(\nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \mathbf{y}^*) \right) \right\|_2 \\ &= \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) \right. \\ & \quad + \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}}) (\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \\ & \quad \left. + \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} - (\mathbf{I}_n \otimes \bar{\mathbf{H}} + \mu \mathbf{I}_{nd})^{-1} \right) (\mathbf{I}_n \otimes \bar{\mathbf{H}}) (\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{y}^*) \right\|_2. \end{aligned}$$

Using similar trick as in (A.30), the above quantity can be further upper bounded as

$$\begin{aligned} & \left\| \left(\mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right\|_2 \\ & \leq \left\| \mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right\|_2 \left\| (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \right\|_2 \left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\ & \quad + \left\| \mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right\|_2 \left\| (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} \right\|_2 \left\| \mathbf{I}_n \otimes \bar{\mathbf{H}} \right\|_2 \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} \right\|_2 \\ & \quad + \sqrt{n} \left\| \mathbf{I}_{nd} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right\|_2 \left\| (\mathbf{H} + \mu \mathbf{I}_{nd})^{-1} (\mathbf{I}_n \otimes \bar{\mathbf{H}} - \mathbf{H}) (\mathbf{I}_{nd} + \mu \mathbf{I}_n \otimes \bar{\mathbf{H}}^{-1})^{-1} \right\|_2 \left\| \bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^* \right\|_2. \quad (\text{A.41}) \end{aligned}$$

Combine (A.40) and (A.41), we conclude that

$$\begin{aligned} \left\| \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)} \right\|_2 & \leq \left(\alpha + \frac{\alpha L}{\sigma + \mu} \right) \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} \right\|_2 + \frac{\alpha}{\sigma + \mu} \left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\ & \quad + \frac{\alpha L}{L + \mu} \frac{\beta}{\sigma + \mu} \sqrt{n} \left\| \bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^* \right\|_2. \quad (\text{A.42}) \end{aligned}$$

A.5.3 Gradient estimation error

In view of the fundamental theorem of calculus and the definition of β , it holds that

$$\left\| \nabla(f - f_j)(\mathbf{x}) - \nabla(f - f_j)(\mathbf{y}) \right\|_2 = \left\| \left[\int_0^1 \nabla^2(f - f_j)(c\mathbf{x} + (1-c)\mathbf{y}) dc \right] (\mathbf{x} - \mathbf{y}) \right\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2.$$

To begin, the update formulas (2.7) and (2.8) are equivalent to

$$\mathbf{y}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d) \mathbf{x}^{(t-1)}, \quad (\text{A.43})$$

$$\mathbf{s}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d) \mathbf{s}^{(t-1)} + \nabla F(\mathbf{y}^{(t)}) - \nabla F(\mathbf{y}^{(t-1)}). \quad (\text{A.44})$$

Note that, since

$$\begin{aligned} \left(\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right)^K &= \left(\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \left(\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \cdots \left(\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \\ &= \left(\mathbf{W}^2 - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \cdots \left(\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) = \mathbf{W}^K - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top, \end{aligned}$$

we have the mixing rate of \mathbf{W}^K is

$$\alpha := \|\mathbf{W}^K - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top\| = \|\mathbf{W} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top\|^K = \alpha_0^K.$$

In view of the equivalent update rule (A.44),

$$\begin{aligned} \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 &= \left\| (\mathbf{W} \otimes \mathbf{I}_d) \mathbf{s}^{(t-1)} + \nabla F(\mathbf{y}^{(t)}) - \nabla F(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t)}) \right\|_2 \\ &= \left\| (\mathbf{W} \otimes \mathbf{I}_d) \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) + (\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) \right. \\ &\quad \left. + \nabla F(\mathbf{y}^{(t)}) - \nabla F(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t)}) \right\|_2 \\ &= \left\| (\mathbf{W} \otimes \mathbf{I}_d) \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) + \nabla(F - f)(\mathbf{y}^{(t)}) \right. \\ &\quad \left. + (\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) - \nabla F(\mathbf{y}^{(t-1)}) \right\|_2 \end{aligned}$$

Subtract and add $\left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right)$, $\nabla(f - F)(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)})$ and $\nabla(f - F)(\mathbf{1}_n \otimes \mathbf{y}^*)$ to the previous equation, and rearrange terms,

$$\begin{aligned} \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 &= \left\| \left[(\mathbf{W} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) \right. \\ &\quad \left. + \nabla(F - f)(\mathbf{y}^{(t)}) - \nabla(F - f)(\mathbf{1}_n \otimes \mathbf{y}^*) \right. \\ &\quad \left. + (\mathbf{W} \otimes \mathbf{I}_d) \left(\nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \mathbf{y}^*) \right) - \left[\nabla F(\mathbf{y}^{(t-1)}) - \nabla F(\mathbf{1}_n \otimes \mathbf{y}^*) \right] \right. \\ &\quad \left. + \left[\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) \right\|_2 \\ &\leq \alpha \|\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + \beta \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^*\|_2 \\ &\quad + \left\| (\mathbf{W} \otimes \mathbf{I}_d) \left(\nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \mathbf{y}^*) \right) - \left[\nabla F(\mathbf{y}^{(t-1)}) - \nabla F(\mathbf{1}_n \otimes \mathbf{y}^*) \right] \right. \\ &\quad \left. + \left[\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right) \right\|_2. \quad (\text{A.45}) \end{aligned}$$

Using the facts $\left[\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \mathbf{s}^{(t-1)} = \left[\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \nabla F(\mathbf{y}^{(t-1)})$ and $\left[\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \nabla(F - f)(\mathbf{1}_n \otimes \mathbf{y}^*) = \mathbf{0}$, the last term of (A.45) becomes

$$\left\| \left[(\mathbf{W} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right] \left(\nabla(f - F)(\mathbf{y}^{(t-1)}) - \nabla(f - F)(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right) \right\|_2$$

$$\begin{aligned}
& + \left(\nabla(f - F)(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) - \nabla(f - F)(\mathbf{1}_n \otimes \mathbf{y}^*) \right) \\
& + \left[(\mathbf{W} \otimes \mathbf{I}_d) - \mathbf{I}_{nd} \right] \left(\nabla F(\mathbf{y}^{(t-1)}) - \nabla F(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right) \Big\|_2 \\
\leq & \left\| (\mathbf{W} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right\|_2 \left\| \nabla(f - F)(\mathbf{y}^{(t-1)}) - \nabla(f - F)(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\
& + \left\| \nabla(f - F)(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) - \nabla(f - F)(\mathbf{1}_n \otimes \mathbf{y}^*) \right\|_2 \\
& + \left\| (\mathbf{W} \otimes \mathbf{I}_d) - \mathbf{I}_{nd} \right\|_2 \left\| \nabla F(\mathbf{y}^{(t-1)}) - \nabla F(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\
\leq & \alpha \beta \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} \right\|_2 + \beta \sqrt{n} \left\| \bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^* \right\|_2 + (\alpha + 1)L \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} \right\|_2. \tag{A.46}
\end{aligned}$$

We used $\left\| (\mathbf{W} \otimes \mathbf{I}_d) - \mathbf{I}_{nd} \right\|_2 = \left\| (\mathbf{W} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) + \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) - \mathbf{I}_{nd} \right\|_2 \leq \left\| (\mathbf{W} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \right\|_2 + \left\| \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) - \mathbf{I}_{nd} \right\|_2 \leq \alpha + 1$ to obtain the last inequality.

Combining (A.45) and (A.46), we obtain the bound

$$\begin{aligned}
\left\| \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}) \right\|_2 \leq & \alpha \left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 + \beta \left\| \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)} \right\|_2 + \beta \sqrt{n} \left\| \bar{\mathbf{y}}^{(t)} - \mathbf{y}^* \right\|_2 \\
& + (\alpha \beta + (\alpha + 1)L) \left\| \mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)} \right\|_2 + \beta \sqrt{n} \left\| \bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^* \right\|_2. \tag{A.47}
\end{aligned}$$

A.5.4 Linear system

Recall the definitions $\eta = \frac{1}{\sigma + \mu}$, $\gamma = \frac{L}{L + \sigma}$ and the error vector (2.13). Combining (A.38), (A.42) and (A.47) leads to the matrix \mathbf{G} defined in (A.3).

A.6 Proof of Lemma 2

The proof follows the same procedures as the proof of Lemma 1. (i) In Appendix A.6.1, we bound the convergence error $\sqrt{n} \left\| \bar{\mathbf{y}}^{(t)} - \mathbf{y}^* \right\|_2$; (ii) in Appendix A.6.2, we bound the parameter consensus error $\left\| \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)} \right\|_2$; (iii) finally, using the bound we obtained in Appendix A.5.3 of the gradient estimation error, we create induction inequalities of $\left\| \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)} \right\|_2$, $\sqrt{n} \left\| \bar{\mathbf{y}}^{(t)} - \mathbf{y}^* \right\|_2$ and $L^{-1} \left\| \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}^{(t)}) \right\|_2$ in Appendix A.6.3 to conclude the proof. For consistency and simplicity, we use the same definitions of \mathbf{x}^+ in (A.25), $\eta = \frac{1}{\sigma + \mu}$, and $\gamma = \frac{L}{L + \sigma}$ as in the proof of Lemma 1.

A.6.1 Convergence error

We continue to decompose the convergence error as (A.28), and bound the two terms respectively.

1. For the term $\left\| \bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+ \right\|_2$, we first subtract (A.26a) from (A.26b), which gives

$$\nabla f_j(\mathbf{x}_j^{(t-1)}) - \nabla f_j(\mathbf{x}_j^+) + \mu(\mathbf{x}_j^{(t-1)} - \mathbf{x}_j^+) = \nabla f(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)}$$

$$+ \nabla(f - f_j)(\bar{\mathbf{y}}^{(t-1)}) - \nabla(f - f_j)(\mathbf{y}_j^{(t-1)}) + \mu(\mathbf{y}_j^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}),$$

then use the strong convexity of $f_j(\cdot)$ and the definition of β to bound both sides,

$$\begin{aligned} & \|\nabla f_j(\mathbf{x}_j^{(t-1)}) - \nabla f_j(\mathbf{x}_j^+) + \mu(\mathbf{x}_j^{(t-1)} - \mathbf{x}_j^+)\|_2 \geq (\sigma + \mu)\|\mathbf{x}_j^{(t-1)} - \mathbf{x}_j^+\|_2, \\ & \|\nabla f(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)} + \nabla(f - f_j)(\bar{\mathbf{y}}^{(t-1)}) - \nabla(f - f_j)(\mathbf{y}_j^{(t-1)}) + \mu(\mathbf{y}_j^{(t-1)} - \bar{\mathbf{y}}^{(t-1)})\|_2 \\ & \leq (\beta + \mu)\|\mathbf{y}_j^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2 + \|\nabla f(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)}\|_2. \end{aligned}$$

Therefore, combining the above two inequalities, we have

$$\|\mathbf{x}_j^{(t-1)} - \mathbf{x}_j^+\|_2 \leq \frac{1}{\sigma + \mu}\|\nabla f(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)}\|_2 + \frac{\beta + \mu}{\sigma + \mu}\|\mathbf{y}_j^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2. \quad (\text{A.48})$$

Subtracting the optimality conditions in (A.27),

$$\begin{aligned} \mathbf{0} & \in \frac{1}{n} \sum_j \nabla f_j(\mathbf{x}_j^{(t-1)}) - \frac{1}{n} \sum_j \nabla f_j(\mathbf{x}_j^+) + \mu(\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+) \\ & = \frac{1}{n} \sum_j (\nabla f_j(\mathbf{x}_j^{(t-1)}) - L\mathbf{x}_j^{(t-1)}) - \frac{1}{n} \sum_j (\nabla f_j(\mathbf{x}_j^+) - L\mathbf{x}_j^+) + (L + \mu)(\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+). \end{aligned}$$

Note the gradient of the function $Lx - \nabla f_j(x)$ is a $(L - \sigma)$ -Lipschitz function. Taking the ℓ_2 norm and plugging in (A.48), we have

$$\begin{aligned} \|\bar{\mathbf{x}}^{(t-1)} - \bar{\mathbf{x}}^+\|_2 & \leq \frac{1}{L + \mu} \left\| \frac{1}{n} \sum_j \left([L\mathbf{x}_j^{(t-1)} - \nabla f_j(\mathbf{x}_j^{(t-1)})] - [L\mathbf{x}_j^+ - \nabla f_j(\mathbf{x}_j^+)] \right) \right\|_2 \\ & \leq \frac{1}{L + \mu} \frac{1}{n} \sum_j \left\| [L\mathbf{x}_j^{(t-1)} - \nabla f_j(\mathbf{x}_j^{(t-1)})] - [L\mathbf{x}_j^+ - \nabla f_j(\mathbf{x}_j^+)] \right\|_2 \\ & \leq \frac{L - \sigma}{L + \mu} \frac{1}{n} \sum_j \|\mathbf{x}_j^{(t-1)} - \mathbf{x}_j^+\|_2 \\ & \leq \frac{L - \sigma}{L + \mu} \frac{1}{\sigma + \mu} \frac{1}{n} \sum_j \|\nabla f(\mathbf{y}_j^{(t-1)}) - \mathbf{s}_j^{(t-1)}\|_2 + \frac{L - \sigma}{L + \mu} \frac{\beta + \mu}{\sigma + \mu} \frac{1}{n} \sum_j \|\mathbf{y}_j^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2, \quad (\text{A.49}) \end{aligned}$$

where the last line follows (A.48).

2. For the second term $\|\bar{\mathbf{x}}^+ - \mathbf{y}^*\|_2$, because of the assumption $(\frac{\beta}{\sigma + \mu})^2 \leq \frac{\sigma}{\sigma + 2\mu}$, we can invoke [FGW21, Theorem 3.1], which is a careful analysis of the error of DANE, and bound the error as

$$\|\bar{\mathbf{x}}^+ - \mathbf{y}^*\|_2 \leq \frac{\frac{\beta}{\sigma + \mu} \sqrt{\sigma^2 + 2\sigma\mu} + \mu}{\sigma + \mu} \|\bar{\mathbf{y}} - \mathbf{y}^*\|_2 := \theta_2 \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2. \quad (\text{A.50})$$

Putting together (A.49) and (A.50), and plugging back into (A.28), we can bound the convergence error by:

$$\begin{aligned} \sqrt{n} \|\bar{\mathbf{y}}^{(t)} - \mathbf{y}^*\|_2 & = \sqrt{n} \|\bar{\mathbf{x}}^{(t-1)} - \mathbf{y}^*\|_2 \\ & \leq \theta_2 \sqrt{n} \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2 + \frac{1}{L + \mu} \frac{L}{\sigma + \mu} \|\nabla f(\mathbf{y}^{(t-1)}) - \mathbf{s}^{(t-1)}\|_2 \\ & \quad + \frac{\beta + \mu}{L + \mu} \frac{L}{\sigma + \mu} \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2. \quad (\text{A.51}) \end{aligned}$$

A.6.2 Consensus error

Let $\mathbf{H}_j^{(t)} = \int_0^1 \nabla^2 f_j(c\mathbf{x}_j^{(t)} + (1-c)\mathbf{y}_j^{(t)}) dc$ and $\mathbf{H}^{(t)} = \text{diag}(\mathbf{H}_1^{(t)}, \mathbf{H}_2^{(t)}, \dots, \mathbf{H}_n^{(t)})$. Via the fundamental theorem of calculus, we can solve for $\mathbf{x}_j^{(t-1)}$ from the optimality condition (A.26b) as

$$\mathbf{x}_j^{(t-1)} = \mathbf{y}_j^{(t-1)} - (\mathbf{H}_j^{(t-1)} + \mu \mathbf{I}_d)^{-1} \mathbf{s}_j^{(t-1)}. \quad (\text{A.52})$$

Similar to (A.40), we decompose the consensus error as

$$\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2 \leq \alpha \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 + \alpha \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right\|_2 \quad (\text{A.53})$$

Then, we bound (A.53). Adding and subtracting terms and using the triangle inequality,

$$\begin{aligned} & \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \mathbf{s}^{(t-1)} \right\|_2 \\ & \leq \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) + \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right) \right\|_2 \\ & \quad + \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \end{aligned} \quad (\text{A.54})$$

We can bound the first term in (A.54) as

$$\begin{aligned} & \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \left(\mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) + \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right) \right\|_2 \\ & \leq \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \right\|_2 \left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) + \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\ & \leq \frac{1}{\sigma + \mu} \left(\left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 + \left\| \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \right) \\ & \leq \frac{1}{\sigma + \mu} \left(\left\| \mathbf{s}^{(t-1)} - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 + L \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 \right) \end{aligned} \quad (\text{A.55})$$

Then, for the second term in (A.54),

$$\begin{aligned} & \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\ & = \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left((\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} - ((L + \mu) \mathbf{I}_{nd})^{-1} \right) \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\ & \leq \left\| (\mathbf{H}^{(t-1)} + \mu \mathbf{I}_{nd})^{-1} (L \mathbf{I}_{nd} - \mathbf{H}^{(t-1)}) ((L + \mu) \mathbf{I}_{nd})^{-1} \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}) \right\|_2 \\ & \leq \frac{L - \sigma}{L + \mu} \frac{L}{\sigma + \mu} \sqrt{n} \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2 \end{aligned} \quad (\text{A.56})$$

Therefore, by combing (A.53), (A.54), (A.55) and (A.56), we can bound the consensus error by:

$$\begin{aligned} \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2 & \leq \left(\alpha + \frac{\alpha L}{\sigma + \mu} \right) \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 \\ & \quad + \frac{\alpha}{\sigma + \mu} \left\| \nabla f(\mathbf{y}^{(t-1)}) - \mathbf{s}^{(t-1)} \right\|_2 + \frac{\alpha L}{L + \mu} \frac{L}{\sigma + \mu} \sqrt{n} \|\bar{\mathbf{y}}^{(t-1)} - \mathbf{y}^*\|_2. \end{aligned} \quad (\text{A.57})$$

A.6.3 Linear system

Combining (A.47), (A.57), (A.51), we reach the matrix claimed in (A.15).

A.7 Proof of Lemma 3

The proof follows similar procedures as the proof of Lemma 1. (i) In Appendix A.7.1, we bound the expected function value convergence errors $\mathbb{E}[\sum_{j=1}^n (f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*))]$ and $\mathbb{E}[\sum_{j=1}^n (f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*))]$; (ii) in Appendix A.7.2, we bound the expected parameter consensus error $\mathbb{E}\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2$; (iii) in Appendix A.7.3, we bound the expected parameter consensus error $\mathbb{E}\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2$; (iv) finally, we create induction inequalities of $\mathbb{E}[\sum_{j=1}^n (f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*))]$, $\mathbb{E}[\sum_{j=1}^n (f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*))]$, $\mathbb{E}\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2$ and $\mathbb{E}\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2$ to conclude the proof. Expectations in this section are conditioned on $\mathbf{x}^{(t-1)}$, $\mathbf{y}^{(t-1)}$ and $\mathbf{s}^{(t-1)}$, if not specified.

A.7.1 Function value convergence error

First, we bound the function value convergence error of $\mathbf{y}^{(t)}$ using the previous estimate $\mathbf{x}^{(t-1)}$. By the strong convexity of $f(\cdot)$ and the assumption of $\alpha \leq 1/\kappa$,

$$\begin{aligned}
\sum_{j=1}^n f(\mathbf{y}_j^{(t)}) &\leq n f(\bar{\mathbf{y}}^{(t-1)}) + \frac{L}{2} \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2 \\
&\leq n f(\bar{\mathbf{x}}^{(t-1)}) + \frac{\alpha^2 L}{2} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 \\
&\leq n f(\bar{\mathbf{x}}^{(t-1)}) + \frac{\sigma}{2} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 \\
&= \sum_{j=1}^n \left(f(\bar{\mathbf{x}}^{(t-1)}) + \left\langle \nabla f(\bar{\mathbf{x}}^{(t-1)}), \mathbf{x}_j^{(t-1)} - \bar{\mathbf{x}}^{(t-1)} \right\rangle + \frac{\sigma}{2} \|\mathbf{x}_j^{(t-1)} - \bar{\mathbf{x}}^{(t)}\|_2^2 \right) \\
&\leq \sum_{j=1}^n f(\mathbf{x}_j^{(t-1)}). \tag{A.58}
\end{aligned}$$

Next, we bound the function value convergence error after local update, $\sum_{j=1}^n (f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*))$. By constructing the following helper function, we can connect local updates of Network-SVRG to that of D-SVRG [CZC⁺20], which is the counterpart of SVRG in the server/client setting. For agent j at the t th time, we define the corrected sample loss function as

$$\tilde{\ell}^{(j)}(\mathbf{x}; \mathbf{z}) = \ell(\mathbf{x}; \mathbf{z}) + \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x} - \mathbf{y}_j^{(t)} \right\rangle.$$

Then, define the corrected local and global loss functions as

$$h_i^{(t,j)}(\mathbf{x}) = \frac{1}{m} \sum_{\mathbf{z} \in \mathcal{M}_i} \tilde{\ell}^{(j)}(\mathbf{x}; \mathbf{z}) = f_i(\mathbf{x}) + \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x} - \mathbf{y}_j^{(t)} \right\rangle,$$

$$h^{(t,j)}(\mathbf{x}) = \frac{1}{n} \sum_i h_i^{(t,j)}(\mathbf{x}) = f(\mathbf{x}) + \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x} - \mathbf{y}_j^{(t)} \right\rangle. \quad (\text{A.59})$$

Here, $h^{(t,j)}(\cdot)$ and $h_i^{(t,j)}(\cdot)$ are σ -strongly convex and L -smooth functions, and $\|h_i^{(t,j)}(\mathbf{x}) - h^{(t,j)}(\mathbf{x})\|_2 \leq \beta$ by the definition of β . Let $h_*^{(t,j)}$ denote the optimum value of $h^{(t,j)}(\cdot)$.

The key observation is that the local update (2.20) at agent j is the same as the update at agent j when applying D-SVRG to optimize $h^{(t,j)}$ initialized with $\mathbf{y}_j^{(t)}$. This is true because $\forall \mathbf{z} \in \mathcal{M}_j$, the sample gradient and global gradient used in D-SVRG updates at $\mathbf{y}_j^{(t)}$ satisfy

$$\nabla \tilde{\ell}^{(j)}(\mathbf{u}; \mathbf{z}) - \nabla \tilde{\ell}^{(j)}(\mathbf{u}'; \mathbf{z}) = \nabla \ell(\mathbf{u}'; \mathbf{z}) - \nabla \ell(\mathbf{u}; \mathbf{z}), \quad \text{and} \quad \nabla h^{(t,j)}(\mathbf{y}_j^{(t)}) = \mathbf{s}_j^{(t)},$$

which agree with (2.20). Therefore, we can apply [CZC⁺20, Theorem 1] to bound the optimization error of optimizing $h^{(t,j)}$

$$\mathbb{E} \left[h^{(t,j)}(\mathbf{x}_j^{(t)}) - h_*^{(t,j)} \right] < \nu \left(h^{(t,j)}(\mathbf{y}_j^{(t)}) - h_*^{(t,j)} \right), \quad (\text{A.60})$$

where $\mathbf{x}_j^{(t)}$ is the output at agent j produced by running one iteration of Alg. 3, which is also the output of running one iteration of D-SVRG at the same agent, ν is the convergence rate of D-SVRG, which can be bounded by $\nu \leq 1 - \frac{1}{2} \frac{\sigma - 2\beta}{\sigma - 3\beta}$ when choosing step size $\delta = \frac{1}{40L} \left(1 - \frac{4\beta}{\sigma}\right)$ and the number of local updates $S = 160 \frac{L}{\sigma} \left(1 - \frac{4\beta}{\sigma}\right)^{-2}$.

Next, we relate function value descent of $h^{(t,j)}$ to the function value descent of f . Plug in (A.59) and rearrange terms,

$$\begin{aligned} f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) &= h^{(t,j)}(\mathbf{x}_j^{(t)}) - (1 - \nu)f(\mathbf{y}^*) - \nu f(\mathbf{y}^*) - \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}_j^{(t)} \right\rangle \\ &= h^{(t,j)}(\mathbf{x}_j^{(t)}) - (1 - \nu)h^{(t,j)}(\mathbf{y}^{\text{opt}}) - \nu f(\mathbf{y}^*) \\ &\quad - \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}_j^{(t)} - (1 - \nu)(\mathbf{y}^{\text{opt}} - \mathbf{y}_j^{(t)}) \right\rangle \\ &\leq h^{(t,j)}(\mathbf{x}_j^{(t)}) - (1 - \nu)h_*^{(t,j)} - \nu f(\mathbf{y}^*) \\ &\quad - \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}_j^{(t)} - (1 - \nu)(\mathbf{y}^{\text{opt}} - \mathbf{y}_j^{(t)}) \right\rangle \\ &= h^{(t,j)}(\mathbf{x}_j^{(t)}) - h_*^{(t,j)} + \nu \left(h_*^{(t,j)} - f(\mathbf{y}^*) \right) \\ &\quad - \left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}_j^{(t)} - (1 - \nu)(\mathbf{y}^{\text{opt}} - \mathbf{y}_j^{(t)}) \right\rangle, \end{aligned}$$

where we used $h^{(t,j)}(\mathbf{y}^*) \geq h_*^{(t,j)}$ and $\nu \leq 1$ to reach the last inequality.

Taking expectation on both sides and combining with (A.60), we reach the following function value descent of $f(\cdot)$:

$$\mathbb{E} \left[f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) \right] \leq \nu \left(h^{(t,j)}(\mathbf{y}_j^{(t)}) - h_*^{(t,j)} \right) + \nu \left(h_*^{(t,j)} - f(\mathbf{y}^*) \right)$$

$$\begin{aligned}
& - \mathbb{E} \left[\left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}_j^{(t)} - (1 - \nu) (\mathbf{y}^{\text{opt}} - \mathbf{y}_j^{(t)}) \right\rangle \right] \\
& = \nu \left(f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*) \right) - \mathbb{E} \left[\left\langle \mathbf{s}_j^{(t)} - \nabla f(\mathbf{y}_j^{(t)}), \mathbf{x}_j^{(t)} - \mathbf{y}^{\text{opt}} - \nu(\mathbf{y}_j^{(t)} - \mathbf{y}^{\text{opt}}) \right\rangle \right],
\end{aligned}$$

where the last line follows from (A.59). Summing the previous inequality over all agents and using matrix notations, we obtain the following inequality

$$\begin{aligned}
\mathbb{E} \left[\sum_{j=1}^n f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) \right] & \leq \nu \left[\sum_{j=1}^n f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*) \right] - \mathbb{E} \left[\left\langle \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}), \mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}} \right\rangle \right] \\
& \quad + \nu \mathbb{E} \left[\left\langle \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}), \mathbf{y}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}} \right\rangle \right].
\end{aligned} \tag{A.61}$$

Our next step is to carefully bound the last two error terms in (A.61).

$$\begin{aligned}
& \left| \left\langle \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}), \mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}} \right\rangle \right| \\
& \leq \|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2 \\
& \leq \left(\alpha \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + 2L \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2 \right. \\
& \quad \left. + 2\beta \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2 + \beta \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2 \right) \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2 \\
& \leq \frac{1}{2} \alpha L^{-1} \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2^2 + \alpha^{-1} L \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t-1)}\|_2^2 + \frac{3}{2} \alpha L \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2 \\
& \quad + \beta \|\mathbf{y}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2 + \frac{\beta}{2} \|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2 + \frac{3\beta}{2} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2,
\end{aligned} \tag{A.62}$$

where the first inequality is due to (A.72), and the last inequality is obtained by Cauchy-Schwarz inequality.

Similar to (A.61), because of the strong convexity of loss functions, we have

$$\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \mathbf{y}^*\|_2^2 \leq \frac{2}{\sigma} \sum_j \left(f(\mathbf{y}_j^{(t)}) - f(\mathbf{y}^*) \right).$$

Then, we can further bound (A.62) as

$$\begin{aligned}
\left| \left\langle \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}), \mathbf{x}^{(t)} - \mathbf{y}^* \right\rangle \right| & \leq \frac{1}{2} \alpha L^{-1} \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2^2 + \alpha^{-1} L \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2^2 \\
& \quad + \frac{2\beta}{\sigma} \sum_{j=1}^n \left(f(\mathbf{y}_j^{(t-1)}) - f(\mathbf{y}^*) \right) + \frac{\beta}{\sigma} \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*) \right) \\
& \quad + \left(\frac{3\beta}{\sigma} + 3\kappa\alpha \right) \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) \right).
\end{aligned} \tag{A.63}$$

Similarly, we have the same bound applicable for the last term of (A.61):

$$\begin{aligned}
\left| \left\langle \mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)}), \mathbf{y}^{(t)} - \mathbf{y}^{\text{opt}} \right\rangle \right| & \leq \frac{1}{2} \alpha L^{-1} \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2^2 + \alpha^{-1} L \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2^2 \\
& \quad + \frac{2\beta}{\sigma} \sum_{j=1}^n \left(f(\mathbf{y}_j^{(t-1)}) - f(\mathbf{y}^*) \right) + \frac{\beta}{\sigma} \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*) \right) \\
& \quad + \left(\frac{3\beta}{\sigma} + 3\kappa\alpha \right) \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*) \right),
\end{aligned} \tag{A.64}$$

where the last term is due to (A.58).

Put together (A.62), (A.63) and (A.64) and taking expectation, we reach the following bound

$$\begin{aligned}
\mathbb{E} \left[\sum_{j=1}^n \left(f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) \right) \right] &\leq \left(\nu(1 + 3\alpha\kappa + \frac{4\beta}{\sigma}) + \frac{\beta}{\sigma} \right) \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*) \right) \\
&\quad + \alpha L^{-1} \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2^2 + 2\alpha^{-1} L \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2^2 \\
&\quad + \frac{4\beta}{\sigma} \sum_{j=1}^n \left(f(\mathbf{y}_j^{(t-1)}) - f(\mathbf{y}^*) \right) + \left(\frac{3\beta}{\sigma} + 3\kappa\alpha \right) \mathbb{E} \left[\sum_{j=1}^n \left(f(\mathbf{x}_j^{(t)}) - f(\mathbf{y}^*) \right) \right].
\end{aligned} \tag{A.65}$$

Rearranging terms, we proved the advertised bound.

A.7.2 Consensus error

We first bound the consensus error $\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2 / (\alpha L)$. Similar to (A.40),

$$\begin{aligned}
\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2 &\leq \alpha^2 \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 \\
&= \alpha^2 \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2 - n\alpha^2 \|\mathbf{y}^* - \bar{\mathbf{x}}^{(t-1)}\|_2^2 \\
&\leq \alpha^2 \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \mathbf{y}^{\text{opt}}\|_2^2.
\end{aligned} \tag{A.66}$$

Then, using the strong convexity of $f(\cdot)$,

$$\begin{aligned}
\|\mathbf{y}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{y}}^{(t)}\|_2^2 &\leq \alpha^2 \sum_{j=1}^n \|\mathbf{x}_j^{(t-1)} - \mathbf{y}^{\text{opt}}\|_2^2 \\
&\leq \frac{2\alpha^2}{\sigma} \sum_{j=1}^n \left(f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*) \right).
\end{aligned} \tag{A.67}$$

A.7.3 Gradient estimation error

To bound the gradient estimation error, we note that

$$\begin{aligned}
\|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 &= \|(\mathbf{W} \otimes \mathbf{I}_d) \mathbf{s}^{t-1} + \nabla F(\mathbf{y}^{(t)}) - \nabla F(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t)})\|_2 \\
&= \|(\mathbf{W} \otimes \mathbf{I}_d) (\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})) + (\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)}) \\
&\quad + \nabla F(\mathbf{y}^{(t)}) - \nabla F(\mathbf{y}^{(t-1)}) + \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t)})\|_2 \\
&\leq \|(\mathbf{W} \otimes \mathbf{I}_d) (\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)}))\|_2 + \|(\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)})\|_2 \\
&\quad + \|\nabla(F - f)(\mathbf{y}^{(t)}) + \nabla(F - f)(\mathbf{y}^{(t-1)})\|_2.
\end{aligned} \tag{A.68}$$

We then bound the three terms in (A.68) respectively.

1. The first term can be bounded as

$$\begin{aligned}
& \|(\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)}))\|_2 \\
&= \left\| (\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})) - \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})) \right\|_2 \\
&\quad + \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})) \right\|_2 \\
&\leq \alpha \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})) \right\|_2 \\
&= \alpha \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\nabla(F-f)(\mathbf{y}^{(t-1)}) - \nabla(F-f)(\mathbf{y}^{\text{opt}})) \right\|_2 \\
&\leq \alpha \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + \beta \|\mathbf{y}^{(t-1)} - \mathbf{y}^{\text{opt}}\|_2, \tag{A.69}
\end{aligned}$$

where we used the fact $\left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \right\|_2 = 1$ and the definition of β to reach the last inequality.

2. As for the second term in (A.68), we have

$$\begin{aligned}
& \left\| (\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\
&\leq \left\| (\mathbf{W} \otimes \mathbf{I}_d) \nabla f(\mathbf{y}^{(t-1)}) - \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\
&\quad + \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\
&\leq 2 \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \nabla f(\mathbf{y}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)}) \right\|_2 \\
&\leq 2 \|\nabla f(\bar{\mathbf{y}}^{(t-1)}) - \nabla f(\mathbf{y}^{(t-1)})\|_2 \\
&\leq 2L \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2, \tag{A.70}
\end{aligned}$$

where the third inequality follows from the similar trick we used to obtain (A.66).

3. Using the triangle inequality and the definition of β , the last term in (A.68) can be bounded by

$$\|\nabla(F-f)(\mathbf{y}^{(t)}) + \nabla(F-f)(\mathbf{y}^{(t-1)})\|_2 \leq \beta \|\mathbf{y}^{(t)} - \mathbf{y}^{\text{opt}}\|_2 + \beta \|\mathbf{y}^{(t-1)} - \mathbf{y}^{\text{opt}}\|_2. \tag{A.71}$$

Combining (A.68), (A.69), (A.70) and (A.71), the gradient estimation error can be bounded by

$$\begin{aligned}
\|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2 &\leq \alpha \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2 + 2\beta \|\mathbf{y}^{(t-1)} - \mathbf{y}^{\text{opt}}\|_2 \\
&\quad + \beta \|\mathbf{y}^{(t)} - \mathbf{y}^{\text{opt}}\|_2 + 2L \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2. \tag{A.72}
\end{aligned}$$

Because of the strong convexity, $\|\mathbf{y} - \mathbf{y}^{\text{opt}}\|_2^2 \leq \frac{2}{\sigma} \sum_{j=1}^n (f(\mathbf{y}_j) - f(\mathbf{y}^*))$. Combining with (A.58), we reached the following bound

$$\begin{aligned}
\|\mathbf{s}^{(t)} - \nabla f(\mathbf{y}^{(t)})\|_2^2 &\leq 4\alpha^2 \|\mathbf{s}^{t-1} - \nabla f(\mathbf{y}^{(t-1)})\|_2^2 + \frac{32\beta^2}{\sigma} \sum_{j=1}^n (f(\mathbf{y}_j^{(t-1)}) - f(\mathbf{y}^*)) \\
&\quad + \frac{8\beta^2}{\sigma} \sum_{j=1}^n (f(\mathbf{x}_j^{(t-1)}) - f(\mathbf{y}^*)) + 16L^2 \|\mathbf{y}^{(t-1)} - \bar{\mathbf{y}}^{(t-1)}\|_2^2. \tag{A.73}
\end{aligned}$$

A.7.4 Linear System

Combining (A.58), (A.67), (A.65), and (A.73), we obtain the claimed linear system.

Appendix B

Appendix for Chapter 3

B.1 Experiment details

For completeness, we list two baseline algorithms, DSGD [NO09, LZZ⁺17] (cf. Algorithm 10) and GT-SARAH [XKK22a] (cf. Algorithm 11), which are compared numerically against the proposed DESTRESS algorithm in Section 3.3.

Algorithm 10 Decentralized stochastic gradient descent (DSGD)

- 1 **input:** initial parameter $\bar{x}^{(0)}$, initial step size η_0 , number of iterations T .
 - 2 **initialization:** set $x_i^{(0)} = \bar{x}^{(0)}$.
 - 3 **for** $t = 1, \dots, T$ **do**
 - 4 Each agent i samples a mini-batch $\mathcal{Z}_i^{(t)}$ from \mathcal{M}_i uniformly at random, and then performs the following updates:
$$\mathbf{g}_i^{(t)} = \frac{1}{b} \sum_{z_i \in \mathcal{Z}_i^{(t)}} \nabla \ell(\mathbf{u}_i^{(t)}; z_i).$$
 - 5 Update via local communication: $\mathbf{x}^{(t+1)} = (\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{x}^{(t)} - \frac{\eta_0}{\sqrt{t}} \mathbf{g}^{(t)})$.
 - 6 **output:** $\mathbf{x}^{\text{out}} = \bar{\mathbf{x}}^{(T)}$.
-

B.2 Proof of Theorem 6

For notation simplicity, let

$$\alpha_{\text{in}} = \alpha^{K_{\text{in}}}, \quad \alpha_{\text{out}} = \alpha^{K_{\text{out}}}$$

throughout the proof. We define the global gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^{nd}$ of an (nd) -dimensional vector $\mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top$ analogously to Section 2.1.3.

Algorithm 11 GT-SARAH

```

1 input: initial parameter  $\bar{\mathbf{x}}^{(0)}$ , step size  $\eta$ , number of outer loops  $T$ , number of inner loops  $q$ .
2 initialization: set  $\mathbf{v}^{(0)} = \mathbf{y}^{(0)} = \nabla F(\mathbf{x}^{(0)})$ .
3 for  $t = 1, \dots, T$  do
4   Update via local communication  $\mathbf{x}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d)\mathbf{x}^{(t-1)} - \eta\mathbf{y}^{(t-1)}$ .
5   if  $\text{mod}(t, q) = 0$  then
6      $\mathbf{v}^{(t)} = \nabla F(\mathbf{x}^{(t)})$ .
7   else
8     Each agent  $i$  samples a mini-batch  $\mathcal{Z}_i^{(t)}$  from  $\mathcal{M}_i$  uniformly at random, and then performs the
      following updates:
          
$$\mathbf{v}_i^{(t)} = \frac{1}{b} \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t)}} (\nabla \ell(\mathbf{x}_i^{(t)}; \mathbf{z}_i) - \nabla \ell(\mathbf{x}_i^{(t-1)}; \mathbf{z}_i)) + \mathbf{v}_i^{(t-1)}.$$

9   Update via local communication  $\mathbf{y}^{(t)} = (\mathbf{W} \otimes \mathbf{I}_d)\mathbf{y}^{(t-1)} + \mathbf{v}^{(t)} - \mathbf{v}^{(t-1)}$ .
10 output:  $\mathbf{x}^{\text{out}} = \bar{\mathbf{x}}^{(T)}$ .
```

The following fact is a straightforward consequence of our assumption on the mixing matrix \mathbf{W} in Definition 1.

Fact 1. Let $\mathbf{x} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top$, and $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$, where $\mathbf{x}_i \in \mathbb{R}^d$. For a mixing matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ satisfying Definition 1, we have

1. $(\frac{1}{n}\mathbf{1}_n^\top \otimes \mathbf{I}_d)(\mathbf{W} \otimes \mathbf{I}_d)\mathbf{x} = (\frac{1}{n}\mathbf{1}_n^\top \otimes \mathbf{I}_d)\mathbf{x} = \bar{\mathbf{x}}$;
2. $(\mathbf{I}_{nd} - (\frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top) \otimes \mathbf{I}_d)(\mathbf{W} \otimes \mathbf{I}_d) = (\mathbf{W} \otimes \mathbf{I}_d - (\frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top) \otimes \mathbf{I}_d)(\mathbf{I}_{nd} - (\frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top) \otimes \mathbf{I}_d)$.

To begin with, we introduce a key lemma that upper bounds the norm of the gradient of the global loss function evaluated at the average local estimates over n agents, in terms of the function value difference at the beginning and the end of the inner loop, the gradient estimation error, and the norm of gradient estimates.

Lemma 4 (Inner loop induction). Assume Assumption 2 holds. After $S \geq 1$ inner loops, one has

$$\begin{aligned} \sum_{s=0}^{S-1} \|\nabla f(\bar{\mathbf{u}}^{(t),s})\|_2^2 &\leq \frac{2}{\eta} \left(f(\bar{\mathbf{u}}^{(t),0}) - f(\bar{\mathbf{u}}^{(t),S}) \right) \\ &\quad + \sum_{s=0}^{S-1} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 - (1 - \eta L) \sum_{s=0}^{S-1} \|\bar{\mathbf{v}}^{(t),s}\|_2^2. \end{aligned}$$

Proof of Lemma 4. The local update rule (3.1a), combined with Lemma 1, yields

$$\bar{\mathbf{u}}^{(t),s+1} = \bar{\mathbf{u}}^{(t),s} - \eta \bar{\mathbf{v}}^{(t),s}.$$

By Assumption 2, we have

$$\begin{aligned}
f(\bar{\mathbf{u}}^{(t),s+1}) &= f(\bar{\mathbf{u}}^{(t),s} - \eta \bar{\mathbf{v}}^{(t),s}) \\
&\leq f(\bar{\mathbf{u}}^{(t),s}) - \langle \nabla f(\bar{\mathbf{u}}^{(t),s}), \eta \bar{\mathbf{v}}^{(t),s} \rangle + \frac{L}{2} \|\eta \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
&= f(\bar{\mathbf{u}}^{(t),s}) - \frac{\eta}{2} \|\nabla f(\bar{\mathbf{u}}^{(t),s})\|_2^2 + \frac{\eta}{2} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \|\bar{\mathbf{v}}^{(t),s}\|_2^2, \tag{B.1}
\end{aligned}$$

where the last equality is obtained by applying $-\langle \mathbf{a}, \mathbf{b} \rangle = \frac{1}{2}(\|\mathbf{a} - \mathbf{b}\|_2^2 - \|\mathbf{a}\|_2^2 - \|\mathbf{b}\|_2^2)$. Summing over $s = 0, \dots, S-1$ finishes the proof. \square

Because the output \mathbf{x}^{out} is chosen from $\{\mathbf{u}_i^{(t),s-1} | i \in [n], t \in [T], s \in [S]\}$ uniformly at random, we can compute the expectation of the output's gradient as follows:

$$\begin{aligned}
nTSE \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 &= \sum_{i=1}^n \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\mathbf{u}_i^{(t),s})\|_2^2 \\
&\stackrel{(i)}{=} \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\mathbf{u}^{(t),s})\|_2^2 \\
&= \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\mathbf{u}^{(t),s}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}) + \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s})\|_2^2 \\
&\stackrel{(ii)}{\leq} 2 \sum_{t=1}^T \sum_{s=0}^{S-1} \left(\mathbb{E} \|\nabla f(\mathbf{u}^{(t),s}) - \nabla f(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s})\|_2^2 + \mathbb{E} \|\nabla f(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s})\|_2^2 \right) \\
&\stackrel{(iii)}{\leq} 2 \sum_{t=1}^T \sum_{s=0}^{S-1} \left(L^2 \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + n \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s})\|_2^2 \right), \tag{B.2}
\end{aligned}$$

where (i) follows from the change of notation using the stacked gradient, (ii) follows from the Cauchy-Schwartz inequality, and (iii) follows from Assumption 2. Then, in view of Lemma 4, (B.2) can be further bounded by

$$\begin{aligned}
nTSE \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 &\leq \frac{4n}{\eta} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^* \right) + 2L^2 \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 \\
&\quad + 2n \sum_{t=1}^T \sum_{s=0}^{S-1} \left(\mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 - (1 - \eta L) \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \right), \tag{B.3}
\end{aligned}$$

where we use $\bar{\mathbf{u}}^{(t),0} = \bar{\mathbf{x}}^{(t)}$ and $f(\bar{\mathbf{u}}^{(t),S}) \geq f^*$.

Next, we present Lemmas 5 and 6 to bound the double sum in (B.3), whose proofs can be found in Appendix B.4 and Appendix B.5, respectively.

Lemma 5 (Sum of inner loop errors). *Assuming all conditions in Theorem 6 hold. For all $t > 0$, we can bound the summation of inner loop errors as*

$$\begin{aligned}
&2L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + 2n \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
&\leq \frac{64L^2}{1 - \alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1 \right) \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 + 2\alpha_{\text{in}}^2 \mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 + \frac{2n}{25} \sum_{s=1}^S \mathbb{E} \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2.
\end{aligned}$$

Lemma 6 (Sum of outer loop gradient estimation error and consensus error). *Assuming all conditions in Theorem 6 hold. We have*

$$\frac{64L^2}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1\right) \sum_{t=1}^T \mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 + 2\alpha_{\text{in}}^2 \sum_{t=1}^T \mathbb{E} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2 \leq \frac{11n}{25} \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2.$$

Using Lemma 5, (B.3) can be bounded as follows:

$$\begin{aligned} nTSE \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 &< \frac{4n}{\eta} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(t),0})] - f^* \right) - 2n \left(\frac{24}{25} - \eta L \right) \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \\ &+ \frac{64L^2}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1\right) \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 + 2\alpha_{\text{in}}^2 \sum_{t=0}^{T-1} \mathbb{E} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2, \end{aligned} \quad (\text{B.4})$$

where we bound the sum of inner loop errors $L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2$ and $n \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2$ by the initial value of each inner loop $\mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2$ and $\mathbb{E} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2$, and the summation of the norm of average inner loop gradient estimator $n \sum_{s=1}^S \mathbb{E} \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2$.

By Lemma 6, (B.4) can be further bounded as

$$\begin{aligned} nTSE \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 &\leq \frac{4n}{\eta} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(t),0})] - f^* \right) - 2n \left(\frac{37}{50} - \eta L \right) \sum_{t=1}^T \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \\ &< \frac{4n}{\eta} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(t),0})] - f^* \right), \end{aligned}$$

which concludes the proof.

B.3 Proof of Corollary 3

Without loss of generality, we assume $n \geq 2$. Otherwise, the problem reduces to the centralized setting with a single agent $n = 1$, and the bound holds trivially. We will confirm the choice of parameters in Corollary 3 in the following paragraphs, and finally obtain the IFO complexity and communication complexity.

Step size η We first assume $\alpha_{\text{in}} \leq \frac{p}{2} \leq \frac{1}{2}$ and $\alpha_{\text{out}} \leq \frac{1}{\sqrt{npb+1}} \leq \frac{1}{2}$, which will be proved to hold shortly, then we can verify the step size choice meets the requirement in (3.3) as:

$$\frac{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})}{1+\alpha_{\text{in}}^{K_{\text{in}}}\alpha_{\text{out}}^{K_{\text{out}}}\sqrt{pnb}} \cdot \frac{1}{10L(\sqrt{S/(npb)}+1)} \geq \frac{(1/2)^4}{2} \cdot \frac{1}{20L} = \frac{1}{640L}.$$

Mixing steps K_{in} and K_{out} Using Chebyshev's acceleration [AS14] to implement the mixing steps, it amounts to an improved mixing rate of $\alpha_{\text{cheb}} \asymp 1 - \sqrt{2(1-\alpha)}$, when the original mixing rate α is close to 1. Set $K_{\text{in}} = \left\lceil \frac{\log(2/p)}{\sqrt{1-\alpha}} \right\rceil$ and $K_{\text{out}} = \left\lceil \frac{\log(\sqrt{npb+1})}{\sqrt{1-\alpha}} \right\rceil$. We are now positioned to examine the effective mixing rate $\alpha_{\text{in}} = \alpha_{\text{cheb}}^{K_{\text{in}}}$ and $\alpha_{\text{out}} = \alpha_{\text{cheb}}^{K_{\text{out}}}$, as follows

$$\alpha_{\text{out}} = \alpha_{\text{cheb}}^{K_{\text{out}}} \stackrel{\text{(i)}}{\leq} \alpha_{\text{cheb}}^{\frac{\log(\sqrt{npb+1})}{\sqrt{1-\alpha}}} \asymp \alpha_{\text{cheb}}^{\frac{\sqrt{2}\log(\sqrt{npb+1})}{1-\alpha_{\text{cheb}}}} \stackrel{\text{(ii)}}{\leq} \alpha_{\text{cheb}}^{\frac{\sqrt{2}\log(\sqrt{npb+1})}{-\log \alpha_{\text{cheb}}}} < \frac{1}{\sqrt{npb+1}} \stackrel{\text{(iii)}}{\leq} \frac{1}{2},$$

where (i) follows from $K_{\text{out}} = \left\lceil \frac{\log(\sqrt{npb}+1)}{\sqrt{1-\alpha}} \right\rceil$, (ii) follows from $\log x \leq x - 1, \forall x > 0$, and (iii) follows from $n \geq 1$ and $b \geq 1$. By a similar argument, we have $\alpha_{\text{in}} = \alpha_{\text{cheb}}^{K_{\text{in}}} \leq \frac{p}{2}$.

Complexity Plugging in the selected parameters into (3.4) in Theorem 6, We have

$$\mathbb{E} \|\nabla f(\mathbf{x}^{\text{out}})\|_2^2 \leq \frac{4}{\eta TS} \left(\mathbb{E}[f(\bar{\mathbf{x}}^{(t),0})] - f^* \right) = O\left(\frac{L}{T\sqrt{mn}}\right).$$

Consequently, the outer iteration complexity is $T = O\left(1 + \frac{L}{(mn)^{1/2}\epsilon^2}\right)$. With this in place, we summarize the communication and IFO complexities as follows:

- The communication complexity is

$$\begin{aligned} T \cdot (SK_{\text{in}} + K_{\text{out}}) &= O\left(\frac{(mn)^{1/2} \log(2(n/m)^{1/2} + 2) + \log((mn)^{1/4} + 1)}{\sqrt{1-\alpha}} \cdot \left(1 + \frac{L}{(mn)^{1/2}\epsilon^2}\right)\right) \\ &= O\left(\frac{\log((n/m)^{1/2} + 2)}{\sqrt{1-\alpha}} \cdot ((mn)^{1/2} + \frac{L}{\epsilon^2})\right), \end{aligned}$$

where we use $2/p = \frac{2\lceil\sqrt{m/n}\rceil}{\sqrt{m/n}} \leq \frac{2(\sqrt{m/n}+1)}{\sqrt{m/n}} = 2(\sqrt{n/m} + 1)$ to bound K_{in} .

- The IFO complexity is $T \cdot (Spb + 2m) = O\left(m + \frac{(m/n)^{1/2}L}{\epsilon^2}\right)$.

B.4 Proof of Lemma 5

This section proves Lemma 5. Appendices B.4.1 and B.4.2 bounds the expected inner loop gradient estimation error and consensus errors by their previous values and the sum of inner loop gradient estimator's norms, Appendix B.4.3 then creates a linear system to compute the summation of inner loop errors using their initial values of each inner loop, which concludes the proof.

B.4.1 Sum of inner loop gradient estimation errors

To begin with, note that the gradient estimation error at the s -th inner loop iteration can be written as

$$\begin{aligned} &\mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 \\ &= \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2 \\ &= \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}) - \nabla F(\mathbf{u}^{(t),s}) \right) + \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2 \\ &\leq 2\mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}) - \nabla F(\mathbf{u}^{(t),s}) \right) \right\|_2^2 + 2\mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2 \\ &\leq \frac{2L^2}{n} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + 2\mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2, \end{aligned} \tag{B.5}$$

where the first equality follows from (2.4), and the last inequality is due to Assumption 2. To continue, the expectation of the second term in (B.5) can be bounded as

$$\begin{aligned}
& \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2 \\
&= \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s}) - (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}) + (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}) \right) \right\|_2^2 \\
&\stackrel{(i)}{=} \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s}) - (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}) \right) \right\|_2^2 \\
&\quad + \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}) \right\|_2^2 \\
&\stackrel{(ii)}{=} \sum_{k=1}^s \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\nabla F(\mathbf{u}^{(t),k}) - \mathbf{v}^{(t),k}) - (\nabla F(\mathbf{u}^{(t),k-1}) - \mathbf{v}^{(t),k-1}) \right) \right\|_2^2 \\
&\quad + \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),0}) - \mathbf{v}^{(t),0}) \right\|_2^2 \\
&\stackrel{(iii)}{=} \sum_{k=1}^s \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\nabla F(\mathbf{u}^{(t),k}) - \mathbf{v}^{(t),k}) - (\nabla F(\mathbf{u}^{(t),k-1}) - \mathbf{v}^{(t),k-1}) \right) \right\|_2^2. \tag{B.6}
\end{aligned}$$

Here, (i) follows from the expectation with respect to the activating indicator $\lambda_i^{(t),s}$ and random samples $\mathcal{Z}^{(t),s}$, conditioned on $\mathbf{u}^{(t),s-1}$ and $\mathbf{v}^{(t),s-1}$:

$$\begin{aligned}
& \mathbb{E} \left[\left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s}) \middle| \mathbf{u}^{(t),s-1}, \mathbf{v}^{(t),s-1} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{u}_i^{(t),s}) - \bar{\mathbf{v}}^{(t),s-1} \\
&\quad - \mathbb{E} \left[\frac{1}{npb} \sum_i \lambda_i^{(t),s} \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t),s}} \left(\nabla \ell(\mathbf{u}_i^{(t),s}; \mathbf{z}_i) - \nabla \ell(\mathbf{u}_i^{(t),s-1}; \mathbf{z}_i) \right) \middle| \mathbf{u}^{(t),s-1}, \mathbf{v}^{(t),s-1} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{u}_i^{(t),s}) - \frac{1}{np} \sum_i \mathbb{E}[\lambda_i^{(t),s}] \left(\nabla f_i(\mathbf{u}_i^{(t),s}) - \nabla f_i(\mathbf{u}_i^{(t),s-1}) \right) - \bar{\mathbf{v}}^{(t),s-1} \\
&= \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{u}_i^{(t),s-1}) - \bar{\mathbf{v}}^{(t),s-1} \\
&= \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}), \tag{B.7}
\end{aligned}$$

(ii) follows by recursively applying the relation obtained from (i); and (iii) follows from the property of gradient tracking, i.e.

$$\left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \nabla F(\mathbf{u}^{(t),0}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{u}_i^{(t),0}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^{(t-1)}) = \bar{\mathbf{s}}^{(t-1)} = \bar{\mathbf{v}}^{(t),0}, \tag{B.8}$$

which leads to $\left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),0}) - \mathbf{v}^{(t),0}) = \mathbf{0}$.

We now continue to bound each term in (B.6), which can be viewed as the variance of the stochastic gradient, as

$$\mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) \left((\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s}) - (\nabla F(\mathbf{u}^{(t),s-1}) - \mathbf{v}^{(t),s-1}) \right) \right\|_2^2$$

$$\begin{aligned}
&\stackrel{(i)}{=} \mathbb{E} \left\| \frac{1}{nb} \sum_{i=1}^n \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t),s}} \left((\nabla f_i(\mathbf{u}_i^{(t),s}) - \nabla f_i(\mathbf{u}_i^{(t),s-1})) - \frac{\lambda_i^{(t),s}}{p} (\nabla \ell(\mathbf{u}_i^{(t),s}; \mathbf{z}_i) - \nabla \ell(\mathbf{u}_i^{(t),s-1}; \mathbf{z}_i)) \right) \right\|_2^2 \\
&\stackrel{(ii)}{=} \frac{1}{n^2 b^2} \sum_{i=1}^n \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t),s}} \mathbb{E} \left\| (\nabla f_i(\mathbf{u}_i^{(t),s}) - \nabla f_i(\mathbf{u}_i^{(t),s-1})) - \frac{\lambda_i^{(t),s}}{p} (\nabla \ell(\mathbf{u}_i^{(t),s}; \mathbf{z}_i) - \nabla \ell(\mathbf{u}_i^{(t),s-1}; \mathbf{z}_i)) \right\|_2^2 \\
&\stackrel{(iii)}{=} \frac{1}{n^2 p^2 b^2} \sum_{i=1}^n \sum_{\mathbf{z}_i \in \mathcal{Z}_i^{(t),s}} \mathbb{E} \left[(\lambda_i^{(t),s})^2 \right] \mathbb{E} \left\| \nabla \ell(\mathbf{u}_i^{(t),s}; \mathbf{z}_i) - \nabla \ell(\mathbf{u}_i^{(t),s-1}; \mathbf{z}_i) \right\|_2^2 \\
&\quad - \frac{1}{n^2 b} \mathbb{E} \left\| \nabla F(\mathbf{u}^{(t),s}) - \nabla F(\mathbf{u}^{(t),s-1}) \right\|_2^2 \\
&\leq \frac{L^2}{n^2 p b} \mathbb{E} \left\| \mathbf{u}^{(t),s} - \mathbf{u}^{(t),s-1} \right\|_2^2, \tag{B.9}
\end{aligned}$$

where (i) follows from the update rules (3.1b) and (3.1c), (ii) follows from the independence of samples and $\mathbb{E}[\lambda_i^{(t),s}] = p$, (iii) follows from similar argument with (B.7), and the last inequality follows from Assumption 2 and $\mathbb{E}[(\lambda_i^{(t),s})^2] = p$.

In view of (3.1a), the difference between inner loop variables in (B.9) can be bounded deterministically as

$$\begin{aligned}
&\left\| \mathbf{u}^{(t),s} - \mathbf{u}^{(t),s-1} \right\|_2^2 \\
&= \left\| (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) (\mathbf{u}^{(t),s-1} - \eta \mathbf{v}^{(t),s-1}) - \mathbf{u}^{(t),s-1} \right\|_2^2 \\
&\stackrel{(i)}{=} \left\| \left((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \mathbf{I}_{nd} \right) (\mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1}) \right. \\
&\quad \left. - \eta \left((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}) - \eta \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1} \right\|_2^2 \\
&\stackrel{(ii)}{=} \left\| \left((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \mathbf{I}_{nd} \right) (\mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1}) \right. \\
&\quad \left. - \eta \left((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}) \right\|_2^2 + \eta^2 n \left\| \bar{\mathbf{v}}^{(t),s-1} \right\|_2^2 \\
&\leq 8 \left\| \mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1} \right\|_2^2 + 2\alpha_{\text{in}}^2 \eta^2 \left\| \mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1} \right\|_2^2 + \eta^2 n \left\| \bar{\mathbf{v}}^{(t),s-1} \right\|_2^2, \tag{B.10}
\end{aligned}$$

where (i) and (ii) follow from $((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \mathbf{I}_{nd})(\mathbf{1}_n \otimes \bar{\mathbf{x}}) = \mathbf{0}$ and $((\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - (\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes \mathbf{I}_d)(\mathbf{1}_n \otimes \bar{\mathbf{x}}) = \mathbf{0}$ for any mean vector $\bar{\mathbf{x}}$; and the last inequality follows from the property of the mixing matrix $\|(\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - \mathbf{I}_{nd}\|_{\text{op}} \leq 2$ and $\|(\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) - (\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes \mathbf{I}_d\|_{\text{op}} \leq \alpha_{\text{in}}$.

Plugging (B.9) and (B.10) into (B.6), we can further obtain

$$\begin{aligned}
&\mathbb{E} \left\| \left(\left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s} \right) \right\|_2^2 \\
&\leq \frac{L^2}{n^2 p b} \sum_{k=1}^s \mathbb{E} \left\| \mathbf{u}^{(t),k} - \mathbf{u}^{(t),k-1} \right\|_2^2 \\
&\leq \frac{8L^2}{n^2 p b} \sum_{k=0}^{s-1} \mathbb{E} \left\| \mathbf{u}^{(t),k} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),k} \right\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{n^2 p b} \sum_{k=0}^{s-1} \mathbb{E} \left\| \mathbf{v}^{(t),k} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),k} \right\|_2^2 + \frac{\eta^2 L^2}{n p b} \sum_{k=0}^{s-1} \mathbb{E} \left\| \bar{\mathbf{v}}^{(t),k} \right\|_2^2.
\end{aligned}$$

Using (B.5) and the previous inequality, we can bound the summation of inner loop gradient estimation errors as

$$\begin{aligned}
& \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
& \leq \frac{2L^2}{n} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + 2 \sum_{s=0}^{S-1} \mathbb{E} \left\| \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d \right) (\nabla F(\mathbf{u}^{(t),s}) - \mathbf{v}^{(t),s}) \right\|_2^2 \\
& \leq \frac{2L^2}{n} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + \frac{16L^2}{n^2 pb} \sum_{s=0}^{S-1} \sum_{k=0}^{s-1} \mathbb{E} \|\mathbf{u}^{(t),k} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),k}\|_2^2 \\
& \quad + \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{n^2 pb} \sum_{s=0}^{S-1} \sum_{k=0}^{s-1} \mathbb{E} \|\mathbf{v}^{(t),k} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),k}\|_2^2 + \frac{2\eta^2 L^2}{n pb} \sum_{s=0}^{S-1} \sum_{k=0}^{s-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),k}\|_2^2 \\
& \leq \left(\frac{8S}{n pb} + 1 \right) \cdot \frac{2L^2}{n} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 \\
& \quad + \frac{4S\alpha_{\text{in}}^2 \eta^2 L^2}{n^2 pb} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 + \frac{2S\eta^2 L^2}{n pb} \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2,
\end{aligned}$$

where the last inequality is obtained by relaxing the upper bound of the summation w.r.t. k from $s-1$ to $S-1$.

The quantity of interest can be now bounded as

$$\begin{aligned}
& 2L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + 2n \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
& \leq \left(\frac{4S}{n pb} + 1 \right) \cdot 8L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 \\
& \quad + \frac{8S\alpha_{\text{in}}^2 \eta^2 L^2}{n pb} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 + \frac{4S\eta^2 L^2}{pb} \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \tag{B.11}
\end{aligned}$$

B.4.2 Sum of inner loop consensus errors

Using the update rule (3.1a), the variable consensus error can be expanded deterministically as follows:

$$\begin{aligned}
\|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 &= \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) (\mathbf{u}^{(t),s-1} - \eta \mathbf{v}^{(t),s-1}) \right\|_2^2 \\
&\stackrel{(i)}{\leq} \alpha_{\text{in}}^2 \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{u}^{(t),s-1} - \eta \mathbf{v}^{(t),s-1}) \right\|_2^2 \\
&\leq \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2}{1 - \alpha_{\text{in}}^2} \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2, \tag{B.12}
\end{aligned}$$

where (i) follows from the fact

$$\left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{W} \otimes \mathbf{I}_d) = \left(\mathbf{W} \otimes \mathbf{I}_d - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right)$$

and the definition of the mixing rate. The last inequality follows from the elementary inequality $2\langle \mathbf{a}, \mathbf{b} \rangle \leq$

$$\frac{1 - \alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{a}\|_2^2 + \frac{1 + \alpha_{\text{in}}^2}{1 - \alpha_{\text{in}}^2} \|\mathbf{b}\|_2^2, \text{ so that } \|\mathbf{a} + \mathbf{b}\|_2^2 \leq \frac{2}{1 + \alpha_{\text{in}}^2} \|\mathbf{a}\|_2^2 + \frac{2}{1 - \alpha_{\text{in}}^2} \|\mathbf{b}\|_2^2.$$

Furthermore, using the update rules (3.1b) and (3.1c) and defining $\Lambda^{(t),s} = \frac{1}{p} \text{diag}(\lambda_1^{(t),s}, \lambda_2^{(t),s}, \dots, \lambda_n^{(t),s}) \otimes \mathbf{I}_d$, the gradient consensus error can be similarly expanded as follows:

$$\begin{aligned}
& \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
&= \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) \mathbf{g}^{(t),s} \right\|_2^2 \\
&\leq \alpha_{\text{in}}^2 \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \mathbf{g}^{(t),s} \right\|_2^2 \\
&\stackrel{(i)}{\leq} \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2 \\
&\quad + \frac{2\alpha_{\text{in}}^2}{1 - \alpha_{\text{in}}^2} \cdot \frac{1}{b} \sum_{\mathbf{z} \in \mathcal{Z}^{(t),s}} \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \Lambda^{(t),s} (\nabla \ell(\mathbf{u}^{(t),s}; \mathbf{z}) - \nabla \ell(\mathbf{u}^{(t),s-1}; \mathbf{z})) \right\|_2^2 \\
&\stackrel{(ii)}{\leq} \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2 + \frac{2\alpha_{\text{in}}^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \|\mathbf{u}^{(t),s} - \mathbf{u}^{(t),s-1}\|_2^2 \\
&\stackrel{(iii)}{\leq} \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2 + \frac{2\alpha_{\text{in}}^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \left(8 \|\mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1}\|_2^2 \right. \\
&\quad \left. + 2\alpha_{\text{in}}^2 \eta^2 \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2 + \eta^2 n \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2 \right) \\
&= \left(\frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} + \frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \right) \|\mathbf{v}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s-1}\|_2^2 \\
&\quad + \frac{16\alpha_{\text{in}}^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \|\mathbf{u}^{(t),s-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s-1}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \cdot n \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2, \tag{B.13}
\end{aligned}$$

where the second term in (i) is obtained by Jensen's inequality, (ii) follows from Assumption 2 and $\|\Lambda^{(t),s}\|_{\text{op}} \leq \frac{1}{p}$, and (iii) follows from (B.10).

B.4.3 Linear system

Let $\mathbf{e}^{(t),s} = \begin{bmatrix} L^2 \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 \\ \mathbb{E} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 \end{bmatrix}$, and $\mathbf{b}^{(t),s} = \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \begin{bmatrix} 0 \\ n \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \end{bmatrix}$. By taking expectation of (B.12) and (B.13), we can construct the following linear system

$$\begin{aligned}
\mathbf{e}^{(t),s} &\leq \begin{bmatrix} \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} & \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{1 - \alpha_{\text{in}}^2} \\ \frac{16\alpha_{\text{in}}^2}{(1 - \alpha_{\text{in}}^2) p^2} & \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} + \frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1 - \alpha_{\text{in}}^2) p^2} \end{bmatrix} \mathbf{e}^{(t),s-1} + \mathbf{b}^{(t),s-1} \\
&\leq \underbrace{\begin{bmatrix} \alpha_{\text{in}} & \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{1 - \alpha_{\text{in}}} \\ \frac{16\alpha_{\text{in}}^2}{(1 - \alpha_{\text{in}}) p^2} & \alpha_{\text{in}} + \frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1 - \alpha_{\text{in}}) p^2} \end{bmatrix}}_{=: \mathbf{G}_{\text{in}}} \mathbf{e}^{(t),s-1} + \mathbf{b}^{(t),s-1} = \mathbf{G}_{\text{in}} \mathbf{e}^{(t),s-1} + \mathbf{b}^{(t),s-1}, \tag{B.14}
\end{aligned}$$

where the second inequality is due to $2\alpha_{\text{in}} < 1 + \alpha_{\text{in}}^2$ and $1 + \alpha_{\text{in}} \geq 1$. Telescope the above inequality to obtain

$$\mathbf{e}^{(t),s} \leq \mathbf{G}_{\text{in}}^s \mathbf{e}^{(t),0} + \sum_{k=1}^s \mathbf{G}_{\text{in}}^{s-k} \mathbf{b}^{(t),k-1}. \tag{B.15}$$

Thus, the sum of the consensus errors can be bounded by

$$\begin{aligned}
\sum_{s=0}^{S-1} \mathbf{e}^{(t),s} &\leq \mathbf{e}^{(t),0} + \sum_{s=1}^{S-1} \left(\mathbf{G}_{\text{in}}^s \mathbf{e}^{(t),0} + \sum_{k=1}^s \mathbf{G}_{\text{in}}^{s-k} \mathbf{b}_i^{(t),k-1} \right) \\
&= \sum_{s=0}^{S-1} \mathbf{G}_{\text{in}}^s \mathbf{e}^{(t),0} + \sum_{s=1}^{S-1} \sum_{k=1}^s \mathbf{G}_{\text{in}}^{s-k} \mathbf{b}^{(t),k-1} \\
&\stackrel{\text{(i)}}{=} \sum_{s=0}^{S-1} \mathbf{G}_{\text{in}}^s \mathbf{e}^{(t),0} + \sum_{k=1}^{S-1} \sum_{s=0}^{S-1-k} \mathbf{G}_{\text{in}}^s \mathbf{b}^{(t),k-1} \\
&\stackrel{\text{(ii)}}{\leq} \sum_{s=0}^{S-1} \mathbf{G}_{\text{in}}^s \mathbf{e}^{(t),0} + \sum_{k=1}^{S-1} \sum_{s=0}^{S-1} \mathbf{G}_{\text{in}}^s \mathbf{b}^{(t),k-1} \\
&\stackrel{\text{(iii)}}{\leq} \sum_{s=0}^{\infty} \mathbf{G}_{\text{in}}^s \left(\mathbf{e}^{(t),0} + \sum_{s=0}^{S-1} \mathbf{b}^{(t),s} \right)
\end{aligned} \tag{B.16}$$

where (i) follows by changing the order of summation, (ii) and (iii) follows from the nonnegativity of \mathbf{G}_{in} and $\mathbf{b}^{(t),s}$ respectively. To continue, we begin with the following claim about \mathbf{G}_{in} which will be proved momentarily.

Claim 1. *Under the choice of η in Theorem 6, the eigenvalues of \mathbf{G}_{in} are in $(-1, 1)$, and the Neumann series converges,*

$$\sum_{s=0}^{\infty} \mathbf{G}_{\text{in}}^s = (\mathbf{I}_2 - \mathbf{G}_{\text{in}})^{-1} \leq \begin{bmatrix} \frac{2}{1-\alpha_{\text{in}}} & \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^3} \\ \frac{32\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}})^3 p^2} & \frac{2}{1-\alpha_{\text{in}}} \end{bmatrix}. \tag{B.17}$$

Let $\mathfrak{G}_{\text{in}}^\top = \left[8 \left(\frac{4S}{npb} + 1 \right) \quad \frac{8S\alpha_{\text{in}}^2 \eta^2 L^2}{pnb} \right]$, in view of Claim 1, the summation of consensus errors in (B.11) can be bounded as

$$\begin{aligned}
&\left(\frac{4S}{npb} + 1 \right) \cdot 8L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + \frac{8S\alpha_{\text{in}}^2 \eta^2 L^2}{npb} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
&= \mathfrak{G}_{\text{in}}^\top \sum_{s=0}^{S-1} \mathbf{e}^{(t),s} \\
&\leq \mathfrak{G}_{\text{in}}^\top \left(\sum_{s=0}^{\infty} \mathbf{G}_{\text{in}}^s \right) \left(\mathbf{e}^{(t),0} + \sum_{k=0}^{S-1} \mathbf{b}^{(t),k} \right) \\
&\leq \mathfrak{G}_{\text{in}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{in}})^{-1} \left(\mathbf{e}^{(t),0} + \sum_{s=0}^{S-1} \mathbf{b}^{(t),s} \right),
\end{aligned}$$

and

$$\begin{aligned}
\mathfrak{G}_{\text{in}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{in}})^{-1} &\leq \left[\frac{16}{1-\alpha_{\text{in}}} \left(\frac{4S}{npb} + 1 \right) + \frac{32\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}})^3 p^2} \cdot \frac{8S\alpha_{\text{in}}^2 \eta^2 L^2}{pnb} \quad \frac{32\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^3} \cdot \left(\frac{4S}{npb} + 1 \right) + \frac{2}{1-\alpha_{\text{in}}} \cdot \frac{8S\alpha_{\text{in}}^2 \eta^2 L^2}{pnb} \right] \\
&\leq \left[\frac{16}{1-\alpha_{\text{in}}} \left(\frac{4S}{npb} + 1 \right) + \frac{3\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}}) p^2} \cdot \frac{128\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^3} \cdot \left(\frac{S}{npb} + 1 \right) + \frac{16\alpha_{\text{in}}^2 (1-\alpha_{\text{in}})}{100} \right] \\
&\leq \left[\frac{64}{1-\alpha_{\text{in}}} \left(\frac{S}{npb} + 1 \right) \quad 2\alpha_{\text{in}}^2 \right],
\end{aligned}$$

where we use (3.3), $\eta L \leq \frac{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})}{10(1+\alpha_{\text{in}}\alpha_{\text{out}}\sqrt{npb})(\sqrt{S/(npb)}+1)} \leq \frac{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})}{10(\sqrt{S/(npb)}+1)}$, to prove the last two inequalities.

Therefore, (B.11) can be bounded as

$$\begin{aligned}
& 2L^2 \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{u}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),s}\|_2^2 + 2n \sum_{s=0}^{S-1} \mathbb{E} \|\nabla f(\bar{\mathbf{u}}^{(t),s}) - \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
& \leq \boldsymbol{\zeta}_{\text{in}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{in}})^{-1} \left(\mathbf{e}^{(t),0} + \sum_{s=0}^{S-1} \mathbf{b}^{(t),s} \right) + \frac{4S\eta^2 L^2}{pb} \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \\
& \leq \frac{64L^2}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1 \right) \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 + 2\alpha_{\text{in}}^2 \mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 \\
& \quad + \left(\frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2} + \frac{4S\eta^2 L^2}{npb} \right) \cdot n \sum_{s=1}^S \mathbb{E} \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2 \\
& < \frac{64L^2}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1 \right) \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 + 2\alpha_{\text{in}}^2 \mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 \\
& \quad + \frac{2n}{25} \sum_{s=1}^S \mathbb{E} \|\bar{\mathbf{v}}^{(t),s-1}\|_2^2,
\end{aligned}$$

where the last inequality is proved by incorporating (3.3) as $\frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2} \leq \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^2} < \frac{4\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}})^2} \cdot \frac{(1-\alpha_{\text{in}})^6}{100} \leq \frac{1}{25}$ and $\frac{4S\eta^2 L^2}{npb} \leq \frac{S}{npb} \cdot \frac{4}{100(\sqrt{S/(npb)}+1)^2} < \frac{1}{25}$.

Proof of Claim 1. By the definition of \mathbf{G}_{in} in (B.14), the characteristic polynomial of \mathbf{G}_{in} is

$$f(\lambda) = (\alpha_{\text{in}} - \lambda) \left(\alpha_{\text{in}} + \frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}}) p^2} - \lambda \right) - \frac{32\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2}.$$

By (3.3), $\eta L \leq \frac{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})}{10(1+\alpha_{\text{in}}\alpha_{\text{out}}\sqrt{npb})(\sqrt{S/(npb)}+1)} \leq \frac{(1-\alpha_{\text{in}})^3}{10}$ and $\alpha_{\text{in}} \leq p$, we have $\frac{32\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2} \leq \frac{32\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^2} \leq \frac{32}{100} \alpha_{\text{in}}^2 (1-\alpha_{\text{in}})^4 < 1$, so that $f(-1) \geq 1 - \frac{32\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2} > 0$, and

$$\begin{aligned}
f(1) &= (1-\alpha_{\text{in}})^2 - \frac{4\alpha_{\text{in}}^4 \eta^2 L^2}{p^2} - \frac{32\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2} \\
&\geq (1-\alpha_{\text{in}})^2 - \frac{36\alpha_{\text{in}}^4 \eta^2 L^2}{(1-\alpha_{\text{in}})^2 p^2} \\
&> (1-\alpha_{\text{in}})^2 - \frac{36}{100} (1-\alpha_{\text{in}})^4 > 0.
\end{aligned}$$

Because $f(\alpha_{\text{in}}) \leq 0$, all eigenvalues of \mathbf{G}_{in} are in $(-1, 1)$, then the Neumann series converges, yielding

$$\begin{aligned}
& \sum_{s=0}^{\infty} \mathbf{G}_{\text{in}}^s = (\mathbf{I}_2 - \mathbf{G}_{\text{in}})^{-1} \\
&= \frac{1-\alpha_{\text{in}}}{(1-\alpha_{\text{in}})^4 p^2 - 4((1-\alpha_{\text{in}})^2 + 8)\alpha_{\text{in}}^4 \eta^2 L^2} \begin{bmatrix} (1-\alpha_{\text{in}})^2 p^2 - 4\alpha_{\text{in}}^4 \eta^2 L^2 & 2\alpha_{\text{in}}^2 \eta^2 L^2 p^2 \\ 16\alpha_{\text{in}}^2 & (1-\alpha_{\text{in}})^2 p^2 \end{bmatrix} \\
&\leq \frac{1-\alpha_{\text{in}}}{(1-\alpha_{\text{in}})^4 p^2 - 4((1-\alpha_{\text{in}})^2 + 8)\alpha_{\text{in}}^4 \eta^2 L^2} \begin{bmatrix} (1-\alpha_{\text{in}})^2 p^2 & 2\alpha_{\text{in}}^2 \eta^2 L^2 p^2 \\ 16\alpha_{\text{in}}^2 & (1-\alpha_{\text{in}})^2 p^2 \end{bmatrix}
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(i)}{\leq} \frac{1 - \alpha_{\text{in}}}{(1 - \alpha_{\text{in}})^4 p^2 - 36\alpha_{\text{in}}^4 \eta^2 L^2} \begin{bmatrix} (1 - \alpha_{\text{in}})^2 p^2 & 2\alpha_{\text{in}}^2 \eta^2 L^2 p^2 \\ 16\alpha_{\text{in}}^2 & (1 - \alpha_{\text{in}})^2 p^2 \end{bmatrix} \\
&\stackrel{(ii)}{\leq} \begin{bmatrix} \frac{2}{1 - \alpha_{\text{in}}} & \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}})^3} \\ \frac{32\alpha_{\text{in}}^2}{(1 - \alpha_{\text{in}})^3 p^2} & \frac{2}{1 - \alpha_{\text{in}}} \end{bmatrix},
\end{aligned}$$

where (i) and (ii) follow the fact $(1 - \alpha_{\text{in}})^2 \leq 1$, and $(1 - \alpha_{\text{in}})^4 p^2 - 36\alpha_{\text{in}}^4 \eta^2 L^2 \geq (1 - \alpha_{\text{in}})^4 p^2 - \frac{36}{100}\alpha_{\text{in}}^4 (1 - \alpha_{\text{in}})^6 \geq (1 - \alpha_{\text{in}})^4 p^2 - \frac{36}{100}\alpha_{\text{in}}^2 (1 - \alpha_{\text{in}})^6 p^2 > \frac{1}{2}(1 - \alpha_{\text{in}})^4 p^2$ due to (3.3). \square

B.5 Proof of Lemma 6

This section proves Lemma 6. In the following subsections, Appendices B.5.1 and B.5.2 derive induction inequalities for the consensus errors and Appendix B.5.3 creates a linear system of consensus errors to compute the summation.

B.5.1 Sum of outer loop variable consensus errors

The variable consensus error can be bounded deterministically as following,

$$\begin{aligned}
&\|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 \\
&= \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \mathbf{x}^{(t)} \right\|_2^2 \\
&\stackrel{(i)}{=} \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \mathbf{u}^{(t),S} \right\|_2^2 \\
&\stackrel{(ii)}{=} \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{W}_{\text{in}} \otimes \mathbf{I}_d) (\mathbf{u}^{(t),S-1} - \eta \mathbf{v}^{(t),S-1}) \right\|_2^2 \\
&\leq \alpha_{\text{in}}^2 \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{u}^{(t),S-1} - \eta \mathbf{v}^{(t),S-1}) \right\|_2^2 \\
&\leq \frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{u}^{(t),S-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),S-1}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2}{1 - \alpha_{\text{in}}^2} \|\mathbf{v}^{(t),S-1} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),S-1}\|_2^2,
\end{aligned}$$

where (i) uses $\mathbf{x}^{(t)} = \mathbf{u}^{(t),S}$, (ii) uses the update rule (3.1a), and the last two inequalities follow from similar reasoning as (B.12). Apply the same reasoning to $\frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \|\mathbf{u}^{(t),S-1} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),S-1}\|_2^2$ and use $\frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \leq 1$, we can prove

$$\begin{aligned}
\|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 &\leq \left(\frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \right)^S \|\mathbf{u}^{(t),0} - \mathbf{1}_n \otimes \bar{\mathbf{u}}^{(t),0}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2}{1 - \alpha_{\text{in}}^2} \sum_{s=0}^{S-1} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 \\
&= \left(\frac{2\alpha_{\text{in}}^2}{1 + \alpha_{\text{in}}^2} \right)^S \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2}{1 - \alpha_{\text{in}}^2} \sum_{s=0}^{S-1} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2, \quad (\text{B.18})
\end{aligned}$$

where the last equality follows from $\mathbf{x}^{(t-1)} = \mathbf{u}^{(t),0}$.

Take expectation of the previous inequality, by (B.16), we can further compute the summation in (B.18) as follows

$$\begin{aligned} \sum_{s=0}^{S-1} \mathbb{E} \|\mathbf{v}^{(t),s} - \mathbf{1}_n \otimes \bar{\mathbf{v}}^{(t),s}\|_2^2 &\leq \frac{32\alpha_{\text{in}}^2 L^2}{(1-\alpha_{\text{in}})^3 p^2} \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 \\ &\quad + \frac{2}{1-\alpha_{\text{in}}} \left(\mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}}^2) p^2} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\| \right). \end{aligned}$$

Together with $\mathbf{x}^{(t)} = \mathbf{u}^{(t),0}$ and $\mathbf{s}^{(t)} = \mathbf{v}^{(t),0}$, (B.18) can be further bounded as

$$\begin{aligned} \mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 &\leq \left(\left(\frac{2\alpha_{\text{in}}^2}{1+\alpha_{\text{in}}^2} \right)^S + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{1-\alpha_{\text{in}}^2} \cdot \frac{32\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}})^3 p^2} \right) \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 \\ &\quad + \frac{2\alpha_{\text{in}}^2 \eta^2}{1-\alpha_{\text{in}}^2} \cdot \frac{2}{1-\alpha_{\text{in}}} \left(\mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}}^2) p^2} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\| \right) \\ &< \alpha_{\text{in}} \mathbb{E} \|\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}\|_2^2 \\ &\quad + \frac{4\alpha_{\text{in}}^2 \eta^2}{(1-\alpha_{\text{in}})^2} \left(\mathbb{E} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}}) p^2} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\| \right). \quad (\text{B.19}) \end{aligned}$$

The last inequality is obtained by using (3.3) and the fact that $0 \leq \alpha_{\text{in}} < 1$ as follows

$$\begin{aligned} \left(\frac{2\alpha_{\text{in}}^2}{1+\alpha_{\text{in}}^2} \right)^S + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{1-\alpha_{\text{in}}^2} \cdot \frac{32\alpha_{\text{in}}^2}{(1-\alpha_{\text{in}})^3 p^2} &= \left(\frac{2\alpha_{\text{in}}^2}{1+\alpha_{\text{in}}^2} \right)^S + \frac{\alpha_{\text{in}}^2 (1-\alpha_{\text{in}})^2}{1+\alpha_{\text{in}}} \cdot \frac{64\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^6 p^2} \\ &< \frac{2\alpha_{\text{in}}^2}{1+\alpha_{\text{in}}^2} + \frac{64}{100} \cdot \frac{\alpha_{\text{in}}^2 (1-\alpha_{\text{in}})^2}{1+\alpha_{\text{in}}} \\ &\leq \frac{2\alpha_{\text{in}}^2}{1+\alpha_{\text{in}}^2} + \frac{\alpha_{\text{in}} (1-\alpha_{\text{in}})^2}{1+\alpha_{\text{in}}^2} = \alpha_{\text{in}}. \end{aligned}$$

B.5.2 Sum of outer loop gradient estimation consensus errors

In view of the update rule for the gradient tracking term (3.2) and reorganize terms,

$$\begin{aligned} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2 &= \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \mathbf{s}^{(t)} \right\|_2^2 \\ &= \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) (\mathbf{W}_{\text{out}} \otimes \mathbf{I}_d) \left(\mathbf{s}^{(t-1)} + \nabla F(\mathbf{x}^{(t)}) - \nabla F(\mathbf{x}^{(t-1)}) \right) \right\|_2^2 \\ &\leq \frac{2\alpha_{\text{out}}^2}{1+\alpha_{\text{out}}^2} \|\mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)}\|_2^2 \\ &\quad + \frac{2\alpha_{\text{out}}^2}{1-\alpha_{\text{out}}^2} \left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{x}^{(t)}) - \nabla F(\mathbf{x}^{(t-1)}) \right) \right\|_2^2, \quad (\text{B.20}) \end{aligned}$$

which follows from similar reasonings as (B.12). The second term can be further decomposed as

$$\begin{aligned} &\left\| \left(\mathbf{I}_{nd} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \otimes \mathbf{I}_d \right) \left(\nabla F(\mathbf{x}^{(t)}) - \nabla F(\mathbf{x}^{(t-1)}) \right) \right\|_2^2 \\ &\leq \|\nabla F(\mathbf{x}^{(t)}) - \nabla F(\mathbf{x}^{(t-1)})\|_2^2 \\ &\leq L^2 \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)} - (\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}) + (\mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)})\|_2^2 \end{aligned}$$

$$\begin{aligned}
&= L^2 \left\| (\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}) - (\mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)}) \right\|_2^2 + nL^2 \left\| \bar{\mathbf{x}}^{(t)} - \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 \\
&\leq 2L^2 \left\| \mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)} \right\|_2^2 + 2L^2 \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 + S\eta^2 L^2 \cdot n \sum_{s=0}^{S-1} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2, \tag{B.21}
\end{aligned}$$

where the last line follows from the update rule (3.1a) by identifying $\bar{\mathbf{x}}^{(t)} - \bar{\mathbf{x}}^{(t-1)} = \eta \sum_{s=0}^{S-1} \bar{\mathbf{v}}^{(t),s}$ and Cauchy-Schwartz inequality.

With (B.21), (B.20) can be further bounded as follows

$$\begin{aligned}
\left\| \mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)} \right\|_2^2 &\leq \frac{2\alpha_{\text{out}}^2}{1 + \alpha_{\text{out}}^2} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 + \frac{2\alpha_{\text{out}}^2}{1 - \alpha_{\text{out}}^2} \left(2L^2 \left\| \mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)} \right\|_2^2 \right. \\
&\quad \left. + 2L^2 \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 + S\eta^2 L^2 \cdot n \sum_{s=0}^{S-1} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2 \right) \\
&\leq \alpha_{\text{out}} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 + \frac{2\alpha_{\text{out}}^2}{1 - \alpha_{\text{out}}} \left(2L^2 \left\| \mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)} \right\|_2^2 \right. \\
&\quad \left. + 2L^2 \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 + S\eta^2 L^2 \cdot n \sum_{s=0}^{S-1} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2 \right). \tag{B.22}
\end{aligned}$$

Combine with (B.19), after taking expectations, (B.22) can be further bounded as

$$\begin{aligned}
\mathbb{E} \left\| \mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)} \right\|_2^2 &< \alpha_{\text{out}} \mathbb{E} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \mathbb{E} \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 \\
&\quad + \frac{2\alpha_{\text{out}}^2 S\eta^2 L^2}{1 - \alpha_{\text{out}}} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2 + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \left(\alpha_{\text{in}} \mathbb{E} \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 \right. \\
&\quad \left. + \frac{4\alpha_{\text{in}}^2 \eta^2}{(1 - \alpha_{\text{in}})^2} \left(\mathbb{E} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 + \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}}) p^2} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2 \right) \right) \\
&= \left(\alpha_{\text{out}} + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2 \eta^2}{(1 - \alpha_{\text{in}})^2} \right) \mathbb{E} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 \\
&\quad + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} (1 + \alpha_{\text{in}}) \mathbb{E} \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 \\
&\quad + \left(\frac{2\alpha_{\text{out}}^2 S\eta^2 L^2}{1 - \alpha_{\text{out}}} + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2 \eta^2}{(1 - \alpha_{\text{in}})^2} \cdot \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}}) p^2} \right) \cdot n \sum_{s=0}^{S-1} \mathbb{E} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2 \\
&\stackrel{(i)}{<} \left(\alpha_{\text{out}} + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2 \eta^2}{(1 - \alpha_{\text{in}})^2} \right) \mathbb{E} \left\| \mathbf{s}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t-1)} \right\|_2^2 \\
&\quad + \frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} (1 + \alpha_{\text{in}}) \mathbb{E} \left\| \mathbf{x}^{(t-1)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t-1)} \right\|_2^2 \\
&\quad + \frac{3\alpha_{\text{out}}^2 S\eta^2 L^2}{1 - \alpha_{\text{out}}} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \left\| \bar{\mathbf{v}}^{(t),s} \right\|_2^2, \tag{B.23}
\end{aligned}$$

where (i) is obtained by applying the condition in (3.3) as follows

$$\begin{aligned}
\frac{4\alpha_{\text{out}}^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2 \eta^2}{(1 - \alpha_{\text{in}})^2} \cdot \frac{2\alpha_{\text{in}}^2 \eta^2 L^2}{(1 - \alpha_{\text{in}}) p^2} &= \frac{\alpha_{\text{out}}^2 \eta^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{32\alpha_{\text{in}}^4 \eta^2 L^2}{(1 - \alpha_{\text{in}})^3 p^2} \\
&\leq \frac{\alpha_{\text{out}}^2 S\eta^2 L^2}{1 - \alpha_{\text{out}}} \cdot \frac{32\alpha_{\text{in}}^2 (1 - \alpha_{\text{in}})^6}{100(1 - \alpha_{\text{in}})^3} \\
&\leq \frac{\alpha_{\text{out}}^2 S\eta^2 L^2}{1 - \alpha_{\text{out}}},
\end{aligned}$$

where the inequalities are obtained by using $S \geq 1$ and $0 \leq \alpha_{\text{in}} < 1$.

B.5.3 Linear system

Defining $\mathbf{e}^{(t)} := \mathbf{e}^{(t),0} = \begin{bmatrix} L^2 \mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 \\ \mathbb{E} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2 \end{bmatrix}$ and $\mathbf{b}'^{(t)} = \begin{bmatrix} \frac{8\alpha_{\text{in}}^4 \eta^4 L^4}{(1-\alpha_{\text{in}})^3 p^2} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \\ \frac{3\alpha_{\text{out}}^2 S \eta^2 L^2}{1-\alpha_{\text{out}}} \cdot n \sum_{s=0}^{S-1} \mathbb{E} \|\bar{\mathbf{v}}^{(t),s}\|_2^2 \end{bmatrix}$, we construct a linear system by putting together (B.19) and (B.23) as

$$\mathbf{e}^{(t)} \leq \underbrace{\begin{bmatrix} \alpha_{\text{in}} & \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^2} \\ \frac{4\alpha_{\text{out}}^2}{1-\alpha_{\text{out}}} (1 + \alpha_{\text{in}}) & \alpha_{\text{out}} + \frac{4\alpha_{\text{out}}^2}{1-\alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^2} \end{bmatrix}}_{=: \mathbf{G}_{\text{out}}} \mathbf{e}^{(t-1)} + \mathbf{b}'^{(t)} = \mathbf{G}_{\text{out}} \mathbf{e}^{(t-1)} + \mathbf{b}'^{(t)}. \quad (\text{B.24})$$

Then, following the same argument as (B.16), we obtain

$$\sum_{t=1}^T \mathbf{e}^{(t)} \leq \sum_{t=1}^{\infty} \mathbf{G}_{\text{out}}^t \left(\mathbf{e}^{(0)} + \sum_{t=1}^T \mathbf{b}'^{(t)} \right). \quad (\text{B.25})$$

Before continuing, we state the following claim about \mathbf{G}_{out} which will be proven momentarily.

Claim 2. *Under the choice of η in Theorem 6, the eigenvalues of \mathbf{G}_{out} are in $(-1, 1)$, and the Neumann series converges,*

$$\sum_{t=0}^{\infty} \mathbf{G}_{\text{out}}^t = (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} \leq \begin{bmatrix} \frac{2}{1-\alpha_{\text{in}}} & \frac{8\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^3 (1-\alpha_{\text{out}})} \\ \frac{16\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} & \frac{2}{1-\alpha_{\text{out}}} \end{bmatrix}.$$

With Claim 2 in hand, and the fact that $\mathbf{e}^{(0)} = \mathbf{0}$, we can bound the summation of outer loop consensus errors by

$$\begin{aligned} & \frac{64L^2}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1 \right) \sum_{t=1}^T \mathbb{E} \|\mathbf{x}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{x}}^{(t)}\|_2^2 + \frac{2\alpha_{\text{in}}^2}{1-\alpha_{\text{in}}} \sum_{t=1}^T \mathbb{E} \|\mathbf{s}^{(t)} - \mathbf{1}_n \otimes \bar{\mathbf{s}}^{(t)}\|_2^2 \\ &= \boldsymbol{\zeta}_{\text{out}}^\top \sum_{t=1}^T \mathbf{e}^{(t)} \\ &\leq \boldsymbol{\zeta}_{\text{out}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} \left(\mathbf{e}^{(0)} + \sum_{t=1}^T \mathbf{b}'^{(t)} \right) \\ &= \boldsymbol{\zeta}_{\text{out}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} \sum_{t=1}^T \mathbf{b}'^{(t)}, \end{aligned} \quad (\text{B.26})$$

where $\boldsymbol{\zeta}_{\text{out}}^\top = \begin{bmatrix} \frac{64}{1-\alpha_{\text{in}}} \cdot \left(\frac{S}{npb} + 1 \right) & \frac{2\alpha_{\text{in}}^2}{1-\alpha_{\text{in}}} \end{bmatrix}$.

Note that by elementary calculations,

$$\begin{aligned} & \boldsymbol{\zeta}_{\text{out}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} \\ &\leq \begin{bmatrix} \frac{64}{1-\alpha_{\text{in}}} \left(\frac{S}{npb} + 1 \right) & 2\alpha_{\text{in}}^2 \end{bmatrix} \begin{bmatrix} \frac{2}{1-\alpha_{\text{in}}} & \frac{8\alpha_{\text{in}}^2 \eta^2 L^2}{(1-\alpha_{\text{in}})^3 (1-\alpha_{\text{out}})} \\ \frac{16\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} & \frac{2}{1-\alpha_{\text{out}}} \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \left[\frac{64}{1-\alpha_{\text{in}}} \left(\frac{S}{npb} + 1 \right) \cdot \frac{2}{1-\alpha_{\text{in}}} + \frac{32\alpha_{\text{in}}^2\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} \quad \frac{64}{1-\alpha_{\text{in}}} \left(\frac{S}{npb} + 1 \right) \cdot \frac{8\alpha_{\text{in}}^2\eta^2L^2}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})} + \frac{4\alpha_{\text{in}}^2}{1-\alpha_{\text{out}}} \right] \\
&\stackrel{(i)}{<} \left[\frac{128}{(1-\alpha_{\text{in}})^2} \left(\frac{S}{npb} + 1 \right) + \frac{32\alpha_{\text{in}}^2\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} \quad 6\alpha_{\text{in}}^2 + \frac{4\alpha_{\text{in}}^2}{1-\alpha_{\text{out}}} \right] \\
&\stackrel{(ii)}{<} \left[\frac{128}{(1-\alpha_{\text{in}})^2} \left(\frac{S}{npb} + 1 \right) + \frac{32\alpha_{\text{in}}^2\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} \quad \frac{10\alpha_{\text{in}}^2}{1-\alpha_{\text{out}}} \right],
\end{aligned}$$

where we use (3.3) to prove (i), and $1/(1-\alpha_{\text{in}}) \geq 1$ and $1/(1-\alpha_{\text{out}}) \geq 1$ to prove (ii).

Thus, (B.26) can be bounded using (3.3) as

$$\begin{aligned}
&\mathfrak{G}_{\text{out}}^\top (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} \begin{bmatrix} \frac{8\alpha_{\text{in}}^4\eta^4L^4}{(1-\alpha_{\text{in}})^3p^2} \\ \frac{3\alpha_{\text{out}}^2S\eta^2L^2}{1-\alpha_{\text{out}}} \end{bmatrix} \\
&\leq \left(\frac{128}{(1-\alpha_{\text{in}})^2} \left(\frac{S}{npb} + 1 \right) + \frac{32\alpha_{\text{in}}^2\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} \right) \frac{8\alpha_{\text{in}}^4\eta^4L^4}{(1-\alpha_{\text{in}})^3p^2} + \frac{10\alpha_{\text{in}}^2}{1-\alpha_{\text{out}}} \cdot \frac{3\alpha_{\text{out}}^2S\eta^2L^2}{1-\alpha_{\text{out}}} \\
&= \frac{1024\alpha_{\text{in}}^4\eta^4L^4}{(1-\alpha_{\text{in}})^5p^2} \left(\frac{S}{npb} + 1 \right) + \frac{256\alpha_{\text{in}}^6\alpha_{\text{out}}^2\eta^4L^4}{(1-\alpha_{\text{in}})^4(1-\alpha_{\text{out}})^2p^2} + \frac{30\alpha_{\text{in}}^2\alpha_{\text{out}}^2npb \cdot S/(npb)}{(1-\alpha_{\text{out}})^2} \cdot \eta^2L^2 \\
&\leq 11\alpha_{\text{in}}^4\eta^2L^2 + 3\alpha_{\text{in}}^6\alpha_{\text{out}}^2\eta^2L^2 + \frac{30}{100} \\
&< \frac{11}{25},
\end{aligned}$$

which concludes the proof.

Proof of Claim 2. For simplicity, denote $c = \frac{4\alpha_{\text{in}}^2\eta^2L^2}{(1-\alpha_{\text{in}})^2}$ and $d = \frac{4\alpha_{\text{out}}^2}{1-\alpha_{\text{out}}}$. Then \mathbf{G}_{out} can be written as

$$\mathbf{G}_{\text{out}} = \begin{bmatrix} \alpha_{\text{in}} & c \\ d(1+\alpha_{\text{in}}) & \alpha_{\text{out}} + cd \end{bmatrix},$$

whose characteristic polynomial is

$$f(\lambda) = (\alpha_{\text{in}} - \lambda)(\alpha_{\text{out}} + cd - \lambda) - (1 + \alpha_{\text{in}})cd.$$

First, note that $f(1)$ can be bounded by

$$\begin{aligned}
f(1) &= (\alpha_{\text{in}} - 1)(\alpha_{\text{out}} + cd - 1) - (1 + \alpha_{\text{in}})cd \\
&= (1 - \alpha_{\text{in}})(1 - \alpha_{\text{out}}) - 2cd > 0,
\end{aligned}$$

where the last inequality is due to the choice of η , namely,

$$cd = \frac{4\alpha_{\text{out}}^2}{1-\alpha_{\text{out}}} \cdot \frac{4\alpha_{\text{in}}^2\eta^2L^2}{(1-\alpha_{\text{in}})^2} \leq \frac{1}{6}(1-\alpha_{\text{in}})(1-\alpha_{\text{out}}).$$

Combined with the trivial fact that $f(-1) > 0$ and $f(\alpha_{\text{in}}) \leq 0$, all eigenvalues of \mathbf{G}_{out} are in $(-1, 1)$.

Consequently, the Neumann series converges, leading to

$$\sum_{t=0}^{\infty} \mathbf{G}_{\text{out}}^t = (\mathbf{I}_2 - \mathbf{G}_{\text{out}})^{-1} = \begin{bmatrix} \frac{(1-\alpha_{\text{in}})^2(1-\alpha_{\text{out}})^2 - 16\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})^2 - 32\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2} & \frac{4\alpha_{\text{in}}^2(1-\alpha_{\text{out}})\eta^2L^2}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})^2 - 32\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2} \\ \frac{4(1-\alpha_{\text{in}})^2(1+\alpha_{\text{in}})\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})^2 - 32\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2} & \frac{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})^2 - 32\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2} \end{bmatrix}$$

$$\leq \begin{bmatrix} \frac{2}{1-\alpha_{\text{in}}} & \frac{8\alpha_{\text{in}}^2\eta^2L^2}{(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})} \\ \frac{16\alpha_{\text{out}}^2}{(1-\alpha_{\text{in}})(1-\alpha_{\text{out}})^2} & \frac{2}{1-\alpha_{\text{out}}} \end{bmatrix},$$

where we use the condition in (3.3) to prove $32\alpha_{\text{in}}^2\alpha_{\text{out}}^2\eta^2L^2 \leq \frac{32}{100}(1-\alpha_{\text{in}})^6(1-\alpha_{\text{out}})^2 < \frac{1}{2}(1-\alpha_{\text{in}})^3(1-\alpha_{\text{out}})^2$ to bound the denominator.

□

Appendix C

Appendix for Chapter 4

C.1 Technical lemmas

We first recall some classical inequalities that helps our derivation.

Proposition 1. *Let $\{\mathbf{v}_1, \dots, \mathbf{v}_\tau\}$ be a set of τ vectors in \mathbb{R}^d . Then, $\forall \beta > 0$, we have*

$$\langle \mathbf{v}_i, \mathbf{v}_j \rangle \leq \frac{\beta}{2} \|\mathbf{u}\|^2 + \frac{1}{2\beta} \|\mathbf{v}\|^2, \quad (\text{C.1})$$

$$\|\mathbf{v}_i + \mathbf{v}_j\|^2 \leq (1 + \beta) \|\mathbf{v}_i\|^2 + \left(1 + \frac{1}{\beta}\right) \|\mathbf{v}_j\|^2, \quad (\text{C.2})$$

$$\left\| \sum_{i=1}^{\tau} \mathbf{v}_i \right\|^2 \leq \tau \sum_{i=1}^{\tau} \|\mathbf{v}_i\|^2. \quad (\text{C.3})$$

Here, (C.1) is referred as the Cauchy-Schwarz inequality, (C.2) and (C.3) are referred as Young's inequality.

Additional notation The following notation will be used throughout our proof:

$$\nabla F(\bar{\mathbf{x}}) := [\nabla f_1(\bar{\mathbf{x}}), \nabla f_2(\bar{\mathbf{x}}), \dots, \nabla f_n(\bar{\mathbf{x}})], \quad \tilde{\nabla}_b F(\bar{\mathbf{x}}) := [\tilde{\nabla}_b f_1(\bar{\mathbf{x}}), \tilde{\nabla}_b f_2(\bar{\mathbf{x}}), \dots, \tilde{\nabla}_b f_n(\bar{\mathbf{x}})].$$

Properties of the mixing matrix We make note of several useful properties of the mixing matrix in the following lemma.

Lemma 7. *Let W be a mixing matrix satisfying Definition 1. For any matrix $M \in \mathbb{R}^{d \times n}$ and $\bar{\mathbf{m}} = \frac{1}{n} M \mathbf{1}_n$, we have*

$$\left\| MW - \bar{\mathbf{m}} \mathbf{1}_n^\top \right\|_F^2 = \left\| MW - \bar{\mathbf{m}} \mathbf{1}_n^\top W \right\|_F^2 \leq \alpha \left\| M - \bar{\mathbf{m}} \mathbf{1}_n^\top \right\|_F^2. \quad (\text{C.4})$$

In addition, for any $\gamma \in (0, 1]$, the matrix $\tilde{W} = I + \gamma(W - I)$ satisfies Definition 1 with a spectral gap at least $\gamma(1 - \alpha)$.

Proof. The first claim follows from the spectral decomposition of \mathbf{W} . Since \mathbf{W} is a doubly stochastic matrix, the largest absolute eigenvalue of \mathbf{W} is 1 and the corresponding eigenvector is $\mathbf{1}_n$. Let $\mathbf{v}_2, \dots, \mathbf{v}_n$ be the eigenvectors of \mathbf{W} corresponding to the remaining eigenvalues. Then, we have

$$\left\| \mathbf{M}\mathbf{W} - \overline{\mathbf{m}}\mathbf{1}_n^\top \right\|_F^2 = \left\| \mathbf{M}\mathbf{W} - \overline{\mathbf{m}}\mathbf{1}_n^\top \mathbf{W} \right\|_F^2 = \sum_{i=1}^r \left\| \mathbf{W}(\mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n) \right\|_2^2,$$

where the first equality follows from $\mathbf{1}_n^\top \mathbf{W} = \mathbf{1}_n^\top$, \mathbf{m}_i denotes the transpose of i -th row of matrix \mathbf{M} , and $\overline{\mathbf{m}}_i$ denotes the average of \mathbf{m}_i . Now we decompose $\mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n$ using the eigenvectors of \mathbf{W} . Noting that

$$\mathbf{1}_n^\top (\mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n) = \mathbf{1}_n^\top \mathbf{m}_i - \mathbf{1}_n^\top \mathbf{1}_n \frac{1}{n} \mathbf{1}_n^\top \mathbf{m}_i = 0,$$

and thus we can write

$$\mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n = \sum_{j=2}^n c_j \mathbf{v}_j$$

for some $\{c_j\}_{j=2}^n$. Then, we have

$$\left\| \mathbf{W}(\mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n) \right\|_2^2 = \left\| \mathbf{W} \sum_{j=2}^n c_j \mathbf{v}_j \right\|_2^2 \leq (1 - (1 - \alpha))^2 \sum_{j=2}^n c_j^2 \leq (1 - (1 - \alpha)) \sum_{j=2}^n c_j^2 = (1 - (1 - \alpha)) \left\| \mathbf{m}_i - \overline{\mathbf{m}}_i\mathbf{1}_n \right\|_2^2,$$

and we conclude the proof of this claim.

For the second claim, recall the fact that if \mathbf{v} is an eigenvector of \mathbf{W} corresponding to the eigenvalue λ , then \mathbf{v} is also an eigenvector of $\widetilde{\mathbf{W}}$ with the corresponding eigenvalue $(1 - \gamma) + \gamma\lambda$. This claim follows from simple computation based on this relation. □

A key consequence of gradient tracking Before diving in the proofs of the main theorems, we record a key property of gradient tracking. Specifically, we have the following lemma.

Lemma 8. *If $\overline{\mathbf{v}}^{(0)} = \frac{1}{n} \widetilde{\nabla}_b F(\mathbf{X}^{(0)})\mathbf{1}_n$, then for any $t \geq 1$, we have*

$$\overline{\mathbf{v}}^{(t)} = \frac{1}{n} \widetilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n, \tag{C.5}$$

and

$$\overline{\mathbf{x}}^{(t+1)} = \overline{\mathbf{x}}^{(t)} - \frac{\eta}{n} \widetilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n. \tag{C.6}$$

Proof. We first prove (C.5) by induction. For the base case ($t = 0$), the relation (C.5) is obviously true by the means of initialization. Now suppose that at the t -th iteration, the relation (C.5) is true, i.e.,

$$\overline{\mathbf{v}}^{(t)} = \frac{1}{n} \widetilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n,$$

then at the $(t + 1)$ -th iteration, we have

$$\begin{aligned}
\bar{\mathbf{v}}^{(t+1)} &= \frac{1}{n} \mathbf{V}^{(t+1)} \mathbf{1}_n \\
&= \frac{1}{n} \mathbf{V}^{(t)} \mathbf{1}_n + \frac{1}{n} \gamma \mathbf{G}^{(t)} (\mathbf{W} - \mathbf{I}) \mathbf{1}_n + \frac{1}{n} \left(\tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)}) \right) \mathbf{1}_n \\
&= \frac{1}{n} \mathbf{V}^{(t)} \mathbf{1}_n + \frac{1}{n} \left(\tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)}) \right) \mathbf{1}_n \\
&= \frac{1}{n} \tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) \mathbf{1}_n.
\end{aligned} \tag{C.7}$$

where (C.7) follows from the update rule of BEER (cf. Line 6), the penultimate line follows from $\mathbf{W}\mathbf{1}_n = \mathbf{1}_n$, and the last line follows from the induction hypothesis at the t -th iteration. Thus the induction hypothesis is also true at the $(t + 1)$ -th iteration, and we complete the proof of (C.5).

For (C.6), it follows from the update rule of BEER (cf. Line 3) that

$$\begin{aligned}
\bar{\mathbf{x}}^{(t+1)} &= \bar{\mathbf{x}}^{(t)} + \frac{\gamma}{n} \mathbf{H}^{(t)} (\mathbf{W} - \mathbf{I}) \mathbf{1} - \frac{\eta}{n} \mathbf{V}^{(t)} \mathbf{1} \\
&= \bar{\mathbf{x}}^{(t)} - \eta \bar{\mathbf{v}}^{(t)} = \bar{\mathbf{x}}^{(t)} - \frac{\eta}{n} \tilde{\nabla}_b F(\mathbf{X}^{(t)}) \mathbf{1}_n,
\end{aligned}$$

where the second line uses $\mathbf{W}\mathbf{1}_n = \mathbf{1}_n$ and (C.5). □

C.2 Recursive relations of main errors

For convenience, we repeat the definitions in (4.2) below.

$$\begin{aligned}
(\text{compression approximation error:}) \quad \Omega_1^{(t)} &= \mathbb{E} \left\| \mathbf{H}^{(t)} - \mathbf{X}^{(t)} \right\|_F^2, & \Omega_2^{(t)} &= \mathbb{E} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_F^2, \\
(\text{consensus error:}) \quad \Omega_3^{(t)} &= \mathbb{E} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2, & \Omega_4^{(t)} &= \mathbb{E} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2, \\
(\text{gradient norm:}) \quad \Omega_5^{(t)} &= \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2.
\end{aligned}$$

Lemma 9 creates induction inequalities for $\{\Omega_i^{(t)}\}$ that captures the approximation errors induced by compression and the consensus errors due to the decentralized setting. Then, we can show that the Lyapunov function defined in (4.1) descends, which leads to the claimed convergence results in Theorem 7 and Theorem 8.

Lemma 9. *Suppose Assumptions 5 and 7 hold, then for any $t \geq 0$, we have*

$$\Omega_1^{(t+1)} \leq \left(1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho} \right) \Omega_1^{(t)} + 0 \cdot \Omega_2^{(t)} + \frac{6\gamma^2 C}{\rho} \Omega_3^{(t)} + \frac{6\eta^2}{\rho} \Omega_4^{(t)} + \frac{6n\eta^2}{\rho} \Omega_5^{(t)}, \tag{C.8a}$$

$$\begin{aligned}
\Omega_2^{(t+1)} &\leq \frac{18\gamma^2 C L^2}{\rho} \Omega_1^{(t)} + \left(1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho} \right) \Omega_2^{(t)} + \frac{18\gamma^2 C L^2}{\rho} \Omega_3^{(t)} \\
&\quad + \frac{6\gamma^2 C + 18L^2 \eta^2}{\rho} \Omega_4^{(t)} + \frac{18L^2 \eta^2 n}{\rho} \Omega_5^{(t)} + \frac{12n\sigma^2}{b\rho},
\end{aligned} \tag{C.8b}$$

$$\Omega_3^{(t+1)} \leq \frac{6\gamma C}{(1-\alpha)} \Omega_1^{(t)} + 0 \cdot \Omega_2^{(t)} + \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \Omega_3^{(t)} + \frac{6\eta^2}{\gamma(1-\alpha)} \Omega_4^{(t)} + 0 \cdot \Omega_5^{(t)}, \quad (\text{C.8c})$$

$$\begin{aligned} \Omega_4^{(t+1)} &\leq \frac{18\gamma CL^2}{(1-\alpha)} \Omega_1^{(t)} + \frac{6\gamma C}{(1-\alpha)} \Omega_2^{(t)} + \frac{18\gamma CL^2}{(1-\alpha)} \Omega_3^{(t)} + \left(1 - \frac{\gamma(1-\alpha)}{2} + \frac{18L^2\eta^2}{\gamma(1-\alpha)}\right) \Omega_4^{(t)} \\ &\quad + \frac{18n\eta^2 L^2}{\gamma(1-\alpha)} \cdot \Omega_5^{(t)} + \frac{12n\sigma^2}{b\gamma(1-\alpha)}, \end{aligned} \quad (\text{C.8d})$$

where

$$C = \|\mathbf{W} - \mathbf{I}\|_{\text{op}}^2 = \sigma_{\max}(\mathbf{W} - \mathbf{I})^2 \quad (\text{C.9})$$

is the square of the maximum singular value of the matrix $\mathbf{W} - \mathbf{I}$.

Note that the eigenvalues of \mathbf{W} and \mathbf{I} all lies in $[-1, 1]$, and thus clearly $C \leq 4$.

Proof. We will establish the inequalities in (C.8) one by one.

Bounding $\Omega_1^{(t)}$ in (C.8a) First from the update rule of BEER (cf. Line 5), we have

$$\begin{aligned} \left\| \mathbf{H}^{(t+1)} - \mathbf{X}^{(t+1)} \right\|_{\text{F}}^2 &= \left\| \mathbf{H}^{(t)} + \mathcal{C}(\mathbf{X}^{(t+1)} - \mathbf{H}^{(t)}) - \mathbf{X}^{(t+1)} \right\|_{\text{F}}^2 \\ &\leq (1-\rho) \left\| \mathbf{X}^{(t+1)} - \mathbf{H}^{(t)} \right\|_{\text{F}}^2 \\ &\leq \left(1 - \frac{\rho}{2}\right) \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\text{F}}^2 + \frac{2}{\rho} \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_{\text{F}}^2, \end{aligned} \quad (\text{C.10})$$

where the first inequality comes from the definition of compression operators (Definition 5) and the second inequality comes from Young's inequality. It then boils down to bound $\left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_{\text{F}}^2$, for which we have

$$\begin{aligned} \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_{\text{F}}^2 &= \left\| \gamma \mathbf{H}^{(t)} (\mathbf{W} - \mathbf{I}) - \eta \mathbf{V}^{(t)} \right\|_{\text{F}}^2 \\ &= \left\| \gamma (\mathbf{H}^{(t)} - \mathbf{X}^{(t)}) (\mathbf{W} - \mathbf{I}) + \gamma (\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top) (\mathbf{W} - \mathbf{I}) - \eta \mathbf{V}^{(t)} \right\|_{\text{F}}^2 \\ &\leq 3\gamma^2 C \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\text{F}}^2 + 3\gamma^2 C \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_{\text{F}}^2 + 3\eta^2 \left\| \mathbf{V}^{(t)} \right\|_{\text{F}}^2 \\ &= 3\gamma^2 C \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\text{F}}^2 + 3\gamma^2 C \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_{\text{F}}^2 + 3\eta^2 \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\text{F}}^2 + 3\eta^2 n \|\bar{\mathbf{v}}^{(t)}\|_2^2, \end{aligned} \quad (\text{C.12})$$

where in the first line we use the update rule of BEER (cf. Line 3), in the second line we use the property of the mixing matrix $\mathbf{1}_n^\top \mathbf{W} = \mathbf{1}_n^\top$, and in the third line, we apply Young's inequality (cf. (C.3)). In the fourth line, we use $\|\mathbf{v}\|_2^2 = \|\mathbf{v} - \bar{v} \mathbf{1}_n\|_2^2 + n\bar{v}^2$ for any vector \mathbf{v} with an average \bar{v} . Plugging this back into (C.10), we get

$$\begin{aligned} \left\| \mathbf{H}^{(t+1)} - \mathbf{X}^{(t+1)} \right\|_{\text{F}}^2 &\leq \left(1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho}\right) \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\text{F}}^2 + \frac{6\gamma^2 C}{\rho} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_{\text{F}}^2 \\ &\quad + \frac{6\eta^2}{\rho} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\text{F}}^2 + \frac{6n\eta^2}{\rho} \|\bar{\mathbf{v}}^{(t)}\|_2^2. \end{aligned}$$

Plugging in the definitions of $\Omega_i^{(t)}$, we obtain (C.8a).

Bounding $\Omega_2^{(t)}$ in (C.8b) Similar to the derivation of (C.8a), by applying the update rule of $\mathbf{G}^{(t)}$ in BEER (Line 8), and Young's inequality, we have

$$\begin{aligned} \left\| \mathbf{V}^{(t+1)} - \mathbf{G}^{(t+1)} \right\|_{\mathbb{F}}^2 &= \left\| \mathbf{G}^{(t)} + \mathcal{C}(\mathbf{V}^{(t+1)} - \mathbf{G}^{(t)}) - \mathbf{V}^{(t+1)} \right\|_{\mathbb{F}}^2 \\ &\leq (1 - \rho) \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t+1)} \right\|_{\mathbb{F}}^2 \\ &\leq \left(1 - \frac{\rho}{2}\right) \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + \frac{2}{\rho} \left\| \mathbf{V}^{(t+1)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2. \end{aligned} \quad (\text{C.13})$$

It then boils down to bound $\left\| \mathbf{V}^{(t+1)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2$. By the update rule of BEER (cf. Line 6), we have

$$\begin{aligned} \left\| \mathbf{V}^{(t+1)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 &= \left\| \gamma \mathbf{G}^{(t)} (\mathbf{W} - \mathbf{I}) + (\tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)})) \right\|_{\mathbb{F}}^2 \\ &= \left\| \gamma (\mathbf{G}^{(t)} - \mathbf{V}^{(t)}) (\mathbf{W} - \mathbf{I}) + \gamma (\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top) (\mathbf{W} - \mathbf{I}) + (\tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)})) \right\|_{\mathbb{F}}^2 \\ &\stackrel{(i)}{\leq} 3\gamma^2 \mathcal{C} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + 3\gamma^2 \mathcal{C} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 + 3 \left\| \tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)}) \right\|_{\mathbb{F}}^2 \\ &\stackrel{(ii)}{\leq} 3\gamma^2 \mathcal{C} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + 3\gamma^2 \mathcal{C} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 + 3 \left\| \nabla F(\mathbf{X}^{(t+1)}) - \nabla F(\mathbf{X}^{(t)}) \right\|_{\mathbb{F}}^2 + \frac{6n\sigma^2}{b} \\ &\stackrel{(iii)}{\leq} 3\gamma^2 \mathcal{C} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + 3\gamma^2 \mathcal{C} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 + 3L^2 \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_{\mathbb{F}}^2 + \frac{6n\sigma^2}{b} \\ &\stackrel{(iv)}{\leq} 3\gamma^2 \mathcal{C} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + (3\gamma^2 \mathcal{C} + 9L^2 \eta^2) \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 \\ &\quad + 9\gamma^2 \mathcal{C} L^2 \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\mathbb{F}}^2 + 9\gamma^2 \mathcal{C} L^2 \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 + 9L^2 \eta^2 n \|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{6n\sigma^2}{b}, \end{aligned}$$

where (i) comes from Young's inequality (cf. (C.3)) and basic facts of matrix norm (cf. (C.9)), (ii) comes from the bounded variance assumption (Assumption 7), (iii) comes from the smoothness assumption (Assumption 5), and (iv) follows from (C.12). Combining the above inequality with (C.13), we have

$$\begin{aligned} \left\| \mathbf{V}^{(t+1)} - \mathbf{G}^{(t+1)} \right\|_{\mathbb{F}}^2 &\leq \left(1 - \frac{\rho}{2}\right) \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + \frac{2}{\rho} \left\| \mathbf{V}^{(t+1)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 \\ &\leq \left(1 - \frac{\rho}{2} + \frac{6\gamma^2 \mathcal{C}}{\rho}\right) \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_{\mathbb{F}}^2 + \frac{6\gamma^2 \mathcal{C} + 18L^2 \eta^2}{\rho} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 \\ &\quad + \frac{18\gamma^2 \mathcal{C} L^2}{\rho} \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_{\mathbb{F}}^2 + \frac{18\gamma^2 \mathcal{C} L^2}{\rho} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 + \frac{18L^2 \eta^2 n}{\rho} \|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{12n\sigma^2}{b\rho}. \end{aligned}$$

Plugging in the definitions of $\Omega_i^{(t)}$, we obtain (C.8b).

Bounding $\Omega_3^{(t)}$ in (C.8c) To bound the consensus error $\left\| \mathbf{X}^{(t+1)} - \bar{\mathbf{x}}^{(t+1)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2$, by the update rule of BEER (cf. Line 3), we have

$$\begin{aligned} &\left\| \mathbf{X}^{(t+1)} - \bar{\mathbf{x}}^{(t+1)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 \\ &= \left\| \mathbf{X}^{(t)} + \gamma \mathbf{H}^{(t)} (\mathbf{W} - \mathbf{I}) - \eta \mathbf{V}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top + \eta \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 \\ &\stackrel{(i)}{=} \left\| \mathbf{X}^{(t)} \tilde{\mathbf{W}} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top + \gamma (\mathbf{H}^{(t)} - \mathbf{X}^{(t)}) (\mathbf{W} - \mathbf{I}) - \eta \mathbf{V}^{(t)} + \eta \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_{\mathbb{F}}^2 \end{aligned}$$

$$\begin{aligned}
&\stackrel{\text{(ii)}}{\leq} (1+\beta)(1-\gamma(1-\alpha)) \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \left(1 + \frac{1}{\beta}\right) \left(2\gamma^2 \left\| (\mathbf{H}^{(t)} - \mathbf{X}^{(t)})(\mathbf{W} - \mathbf{I}) \right\|_F^2 + 2\eta^2 \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2\right) \\
&\stackrel{\text{(iii)}}{\leq} \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \left(1 + \frac{2}{\gamma(1-\alpha)}\right) \left(2\gamma^2 \left\| (\mathbf{H}^{(t)} - \mathbf{X}^{(t)})(\mathbf{W} - \mathbf{I}) \right\|_F^2 + 2\eta^2 \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2\right) \\
&\stackrel{\text{(iv)}}{\leq} \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \left(1 + \frac{2}{\gamma(1-\alpha)}\right) \left(2\gamma^2 C \left\| \mathbf{H}^{(t)} - \mathbf{X}^{(t)} \right\|_F^2 + 2\eta^2 \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2\right) \\
&\leq \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\gamma C}{(1-\alpha)} \left\| \mathbf{H}^{(t)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{6\eta^2}{\gamma(1-\alpha)} \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2,
\end{aligned}$$

where (i) follows from the definition $\widetilde{\mathbf{W}} = \mathbf{I} + \gamma(\mathbf{W} - \mathbf{I})$, (ii) follows from applying Young's inequality twice and Lemma 7, i.e.

$$\left\| \mathbf{X}^{(t)} \widetilde{\mathbf{W}} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F \leq (1 - \gamma(1 - \alpha)) \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F,$$

(iii) follows by choosing $\beta = \gamma(1 - \alpha)/2$, and (iv) uses the definition of C (cf. (C.9)). Plugging in the definitions of $\Omega_i^{(t)}$, we obtain (C.8c).

Bounding $\Omega_4^{(t)}$ in (C.8d) First, note that

$$\begin{aligned}
\left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top \right\|_F^2 &= \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top + \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top \right\|_F^2 \\
&= \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 - n \left\| \bar{\mathbf{v}}^{(t+1)} - \bar{\mathbf{v}}^{(t)} \right\|_2^2 \\
&\leq \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2.
\end{aligned}$$

Thus by the update rule of BEER (cf. Line 6), we have

$$\begin{aligned}
&\left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top \right\|_F^2 \\
&\leq \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 \\
&= \left\| \mathbf{V}^{(t)} + \gamma \mathbf{G}^{(t+1)} (\mathbf{W} - \mathbf{I}) + \tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)}) - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 \\
&= \left\| (\mathbf{V}^{(t)} \widetilde{\mathbf{W}} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top) + \gamma (\mathbf{G}^{(t+1)} - \mathbf{V}^{(t)}) (\mathbf{W} - \mathbf{I}) + (\tilde{\nabla}_b F(\mathbf{X}^{(t+1)}) - \tilde{\nabla}_b F(\mathbf{X}^{(t)})) \right\|_F^2 \\
&\stackrel{\text{(i)}}{\leq} \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \left(1 + \frac{2}{\gamma(1-\alpha)}\right) \left(2\gamma^2 C \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_F^2 + 2L^2 \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{4n\sigma^2}{b}\right) \\
&\stackrel{\text{(ii)}}{\leq} \left(1 - \frac{\gamma(1-\alpha)}{2}\right) \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\gamma C}{(1-\alpha)} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_F^2 + \frac{6L^2}{\gamma(1-\alpha)} \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{12n\sigma^2}{b\gamma(1-\alpha)} \\
&\leq \left(1 - \frac{\gamma(1-\alpha)}{2} + \frac{18L^2\eta^2}{\gamma(1-\alpha)}\right) \left\| \mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\gamma C}{(1-\alpha)} \left\| \mathbf{G}^{(t)} - \mathbf{V}^{(t)} \right\|_F^2 \\
&\quad + \frac{18\gamma CL^2}{(1-\alpha)} \left\| \mathbf{X}^{(t)} - \mathbf{H}^{(t)} \right\|_F^2 + \frac{18\gamma CL^2}{(1-\alpha)} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{18n\eta^2 L^2}{\gamma(1-\alpha)} \left\| \bar{\mathbf{v}}^{(t)} \right\|_2^2 + \frac{12n\sigma^2}{b\gamma(1-\alpha)},
\end{aligned}$$

where (i) and (ii) are obtained similarly as the derivation of (C.8c), and the last line follows from (C.12).

Thus, we can get (C.8d) by plugging in the definitions of $\Omega_i^{(t)}$ and conclude the proof. \square

C.3 Proof of Theorem 7

This proof makes use of Lemma 8 and Lemma 9 to construct a Lyapunov function and then demonstrates its descending property using a linear system argument.

Step 1: establishing a descent property of the function value First, we have the following inequality captures the “descent” of the function value.

$$\begin{aligned}
f(\bar{\mathbf{x}}^{(t+1)}) &\stackrel{(i)}{\leq} f(\bar{\mathbf{x}}^{(t)}) - \eta \langle \bar{\mathbf{v}}^{(t)}, \nabla f(\bar{\mathbf{x}}^{(t)}) \rangle + \frac{\eta^2 L}{2} \|\bar{\mathbf{v}}^{(t)}\|_2^2 \\
&= f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 - \frac{\eta}{2} \|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)}) - \bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{\eta^2 L}{2} \|\bar{\mathbf{v}}^{(t)}\|_2^2 \\
&= f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)}) - \bar{\mathbf{v}}^{(t)}\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \|\bar{\mathbf{v}}^{(t)}\|_2^2 \\
&\stackrel{(ii)}{\leq} f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta}{2n^2} \|\nabla F(\bar{\mathbf{x}}^{(t)})\mathbf{1}_n - \tilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \|\bar{\mathbf{v}}^{(t)}\|_2^2 \\
&= f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta}{2n^2} \|\nabla F(\bar{\mathbf{x}}^{(t)})\mathbf{1}_n - \nabla F(\mathbf{X}^{(t)})\mathbf{1}_n\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \|\bar{\mathbf{v}}^{(t)}\|_2^2 \\
&\quad + \frac{\eta}{2n^2} \|\nabla F(\mathbf{X}^{(t)})\mathbf{1}_n - \tilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n\|_2^2 \\
&\quad + \frac{\eta}{n^2} \langle \nabla F(\mathbf{X}^{(t)})\mathbf{1}_n - \tilde{\nabla}_b F(\mathbf{X}^{(t)})\mathbf{1}_n, \nabla f(\bar{\mathbf{x}}^{(t)})\mathbf{1}_n - \nabla F(\mathbf{X}^{(t)})\mathbf{1}_n \rangle,
\end{aligned}$$

where (i) comes from the L -smooth assumption (Assumption 5), (ii) comes from Lemma 8. Take expectation on both sides, and using the bounded variance assumption (Assumption 7) and independence of stochastic samples, we get

$$\begin{aligned}
\mathbb{E}f(\bar{\mathbf{x}}^{(t+1)}) &\leq \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta}{2n^2} \mathbb{E}\|\nabla F(\bar{\mathbf{x}}^{(t)})\mathbf{1}_n - \nabla F(\mathbf{X}^{(t)})\mathbf{1}_n\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \mathbb{E}\|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{\eta\sigma^2}{2bn} \\
&\stackrel{(i)}{\leq} \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta}{2n} \mathbb{E}\|\nabla F(\mathbf{X}^{(t)}) - \nabla F(\bar{\mathbf{x}}^{(t)})\|_F^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \mathbb{E}\|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{\eta\sigma^2}{2bn} \\
&\stackrel{(ii)}{\leq} \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta L^2}{2n} \mathbb{E}\|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)}\mathbf{1}_n^\top\|_F^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \mathbb{E}\|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{\eta\sigma^2}{2bn},
\end{aligned}$$

where (i) comes from Young’s inequality, and (ii) comes from the L -smooth assumption (Assumption 5) again. Finally, by substituting definitions of $\Omega_3^{(t)}$ and $\Omega_5^{(t)}$, we reach

$$\mathbb{E}f(\bar{\mathbf{x}}^{(t+1)}) \leq \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta L^2}{2n} \Omega_3^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2}\right) \Omega_5^{(t)} + \frac{\eta\sigma^2}{2bn}. \quad (\text{C.14})$$

Step 2: constructing the Lyapunov function By representing

$$\Omega^{(t)} = [\Omega_1^{(t)} \quad \Omega_2^{(t)} \quad \Omega_3^{(t)} \quad \Omega_4^{(t)}]^\top, \quad (\text{C.15})$$

Lemma 9 can be written more compactly as

$$\mathbf{\Omega}^{(t+1)} \leq \underbrace{\begin{bmatrix} 1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho} & 0 & \frac{6\gamma^2 C}{\rho} & \frac{6\eta^2}{\rho} \\ \frac{18\gamma^2 CL^2}{\rho} & 1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho} & \frac{18\gamma^2 CL^2}{\rho} & \frac{6\gamma^2 C + 18L^2\eta^2}{\rho} \\ \frac{6\gamma C}{(1-\alpha)} & 0 & 1 - \frac{\gamma(1-\alpha)}{2} & \frac{6\eta^2}{\gamma(1-\alpha)} \\ \frac{18\gamma CL^2}{(1-\alpha)} & \frac{6\gamma C}{(1-\alpha)} & \frac{18\gamma CL^2}{(1-\alpha)} & 1 - \frac{\gamma(1-\alpha)}{2} + \frac{18L^2\eta^2}{\gamma(1-\alpha)} \end{bmatrix}}_{=:A} \mathbf{\Omega}^{(t)} + \underbrace{\begin{bmatrix} \frac{6n\eta^2}{\rho} \\ \frac{18L^2\eta^2 n}{\rho} \\ 0 \\ \frac{18n\eta^2 L^2}{\gamma(1-\alpha)} \end{bmatrix}}_{=:b_1} \mathbf{\Omega}_5^{(t)} + \underbrace{\begin{bmatrix} 0 \\ \frac{12n}{\rho} \\ 0 \\ \frac{12n}{\gamma(1-\alpha)} \end{bmatrix}}_{=:b_2} \frac{\sigma^2}{b}. \quad (\text{C.16})$$

Define the Lyapunov function

$$\begin{aligned} \Phi^{(t)} &= \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - f^* + \frac{c_1 L}{n} \cdot \Omega_1^{(t)} + \frac{c_2(1-\alpha)^2}{nL} \cdot \Omega_2^{(t)} + \frac{c_3 L}{n} \cdot \Omega_3^{(t)} + \frac{c_4(1-\alpha)^2}{nL} \Omega_4^{(t)} \\ &= \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - f^* + \mathbf{s}^\top \mathbf{\Omega}^{(t)}, \end{aligned} \quad (\text{C.17})$$

where

$$\mathbf{s} = \begin{bmatrix} c_1 L & c_2(1-\alpha)^2 & c_3 L & c_4(1-\alpha)^2 \\ n & nL & n & nL \end{bmatrix}$$

for some constants c_1, c_2, c_3, c_4 that will be specified later.

By (C.16) from Lemma 9 and the descent property (C.14), we have

$$\begin{aligned} \Phi^{(t+1)} &= \mathbb{E}f(\bar{\mathbf{x}}^{(t+1)}) - f^* + \mathbf{s}^\top \mathbf{\Omega}^{(t+1)} \\ &\leq \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - f^* - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta L^2}{2n} \Omega_3^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} \right) \Omega_5^{(t)} + \frac{\eta \sigma^2}{2bn} \\ &\quad + \mathbf{s}^\top \left(A \mathbf{\Omega}^{(t)} + \Omega_5^{(t)} \mathbf{b}_1 + \frac{\sigma^2}{b} \mathbf{b}_2 \right) \\ &\leq \Phi^{(t)} - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} \right) \Omega_5^{(t)} + \frac{\eta \sigma^2}{2bn} + (\mathbf{s}^\top A - \mathbf{s}^\top + \mathbf{q}^\top) \mathbf{\Omega}^{(t)} + \mathbf{s}^\top (\Omega_5^{(t)} \mathbf{b}_1 + \mathbf{b}_2 \frac{\sigma^2}{b}) \\ &= \Phi^{(t)} - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + (\mathbf{s}^\top A - \mathbf{s}^\top + \mathbf{q}^\top) \mathbf{\Omega}^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} - \mathbf{s}^\top \mathbf{b}_1 \right) \Omega_5^{(t)} + \left(\frac{\eta}{2n} + \mathbf{s}^\top \mathbf{b}_2 \right) \frac{\sigma^2}{b}, \end{aligned} \quad (\text{C.18})$$

where $\mathbf{q} = [0 \quad 0 \quad \frac{\eta L^2}{2n} \quad 0]^\top$. For a moment we assume that there exist some constants $c_1, c_2, c_3, c_4 > 0$ such that

$$\mathbf{s}^\top (A - I) + \mathbf{q}^\top \leq \mathbf{0}, \quad (\text{C.19a})$$

$$\frac{\eta}{2} - \frac{\eta^2 L}{2} - \mathbf{s}^\top \mathbf{b}_1 \geq 0, \quad (\text{C.19b})$$

leading to

$$\Phi^{(t+1)} \leq \Phi^{(t)} - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \left(\frac{\eta}{2n} + \mathbf{s}^\top \mathbf{b}_2 \right) \frac{\sigma^2}{b} \leq \Phi^{(t)} - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{36c_4 \sigma^2}{c_\gamma L b \rho}.$$

The proof is thus completed by recursing the above relation over $t = 0, \dots, T-1$.

Step 3: verifying (C.19) It boils down to verify (C.19) is feasible, and it is equivalent to verify there exist parameters $c_1, c_2, c_3, c_4, \gamma, \eta > 0$ satisfying the following matrix inequality:

$$\begin{bmatrix} \mathbf{I} - \mathbf{A}^\top \\ -\mathbf{b}_1 \end{bmatrix} \text{diag} \left[\frac{L}{n}, \frac{(1-\alpha)^2}{nL}, \frac{L}{n}, \frac{(1-\alpha)^2}{nL} \right] \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} \geq \begin{bmatrix} \mathbf{q} \\ \frac{\eta^2 L}{2} - \frac{\eta}{2} \end{bmatrix}.$$

Note that by choosing $\gamma = c_\gamma(1-\alpha)\rho, \eta = c_\eta\gamma(1-\alpha)^2/L$, and setting $c_\gamma \leq \frac{1}{6\sqrt{C}}$ and $c_\eta \leq \frac{1}{9}$, we get

$$1 - \frac{\rho}{2} + \frac{6\gamma^2 C}{\rho} \leq 1 - \frac{\rho}{4}, \quad 1 - \frac{\gamma(1-\alpha)}{2} + \frac{18L^2\eta^2}{\gamma(1-\alpha)} \leq 1 - \frac{\gamma(1-\alpha)}{4}. \quad (\text{C.20})$$

Now, it suffices to show that there exist $c_1, c_2, c_3, c_4, c_\gamma, c_\eta > 0$ such that the following inequalities are satisfied:

$$\begin{bmatrix} \frac{\rho L}{4n} & -\frac{18c_\gamma^2\rho(1-\alpha)^4 L}{n} & -\frac{6c_\gamma\rho L}{n} & -\frac{18Cc_\gamma\rho(1-\alpha)^2 L}{n} \\ 0 & \frac{\rho(1-\alpha)^2}{4nL} & 0 & -\frac{6Cc_\gamma\rho(1-\alpha)^2}{nL} \\ -\frac{6c_\gamma(1-\alpha)\gamma CL}{n} & -\frac{18c_\gamma(1-\alpha)\gamma L}{n} & \frac{\gamma(1-\alpha)L}{2n} & -\frac{18C\gamma(1-\alpha)L}{n} \\ -\frac{6c_\eta^2 c_\gamma \gamma(1-\alpha)^3}{nL} & -\frac{6C\gamma(1-\alpha)^3(1+3c_\eta^2(1-\alpha)^4)}{nL} & -\frac{6c_\eta^2 \gamma(1-\alpha)^3}{nL} & \frac{\gamma(1-\alpha)^3}{4nL} \\ -12c_\eta c_\gamma(1-\alpha)^3 \frac{\eta}{2} & -36c_\eta c_\gamma(1-\alpha)^5 \frac{\eta}{2} & 0 & -36c_\eta(1-\alpha) \frac{\eta}{2} \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} \geq \begin{bmatrix} 0 \\ 0 \\ \frac{c_\eta \gamma(1-\alpha)L}{2n} \\ 0 \\ (c_\eta c_\gamma \rho(1-\alpha)^3 - 1) \frac{\eta}{2} \end{bmatrix}.$$

Given $\rho \leq 1, (1-\alpha) \leq 1$, this can be further reduced to show the existence of $c_1, c_2, c_3, c_4, c_\gamma, c_\eta > 0$ such that

$$\begin{bmatrix} 1 & -72Cc_\gamma^2 & -24Cc_\gamma & -72Cc_\gamma \\ 0 & 1 & 0 & -24Cc_\gamma \\ -12Cc_\gamma & -35Cc_\gamma & 1 & -36C \\ -24c_\eta^2 c_\gamma & -24c_\gamma(1+3c_\eta^2) & -24c_\eta^2 & 1 \\ -12c_\eta c_\gamma & -36c_\eta c_\gamma & 0 & -36c_\eta \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} \geq \begin{bmatrix} 0 \\ 0 \\ c_\eta \\ 0 \\ -1 + c_\eta c_\gamma \end{bmatrix}.$$

This can be easily verified by noting that as long as c_η and c_γ are set sufficiently small, it is straightforward to find feasible c_1, c_2, c_3, c_4 .

C.4 Proof of Theorem 8

The proof strategy of Theorem 8 is similar to that of Theorem 7. However, we need to create a slightly different linear system to improve convergence under the new assumption.

Denote $\kappa := L/\mu$. Taking the same Lyapunov function $\Phi^{(t)}$ in (C.17), by the same argument of Appendix C.3 up to (C.18), we have

$$\Phi^{(t+1)} \leq \mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - f^* - \frac{\eta}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta L^2}{2n} \Omega_3^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} \right) \Omega_5^{(t)} + \frac{\eta \sigma^2}{2bn} + \mathbf{s}^\top \left(\mathbf{A} \Omega^{(t)} + \Omega_5^{(t)} \mathbf{b}_1 + \frac{\sigma^2}{b} \mathbf{b}_2 \right)$$

$$\begin{aligned}
&\leq (1 - \eta\mu)(\mathbb{E}f(\bar{\mathbf{x}}^{(t)}) - f^*) + (\mathbf{s}^\top \mathbf{A} + \mathbf{q}^\top) \boldsymbol{\Omega}^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} - \mathbf{s}^\top \mathbf{b}_1 \right) \Omega_5^{(t)} + \left(\frac{\eta}{2n} + \mathbf{s}^\top \mathbf{b}_2 \right) \frac{\sigma^2}{b} \\
&= (1 - \eta\mu) \Phi^{(t)} + \left(\mathbf{s}^\top \mathbf{A} - (1 - \eta\mu) \mathbf{s}^\top + \mathbf{q}^\top \right) \boldsymbol{\Omega}^{(t)} - \left(\frac{\eta}{2} - \frac{\eta^2 L}{2} - \mathbf{s}^\top \mathbf{b}_1 \right) \Omega_5^{(t)} + \left(\frac{\eta}{2n} + \mathbf{s}^\top \mathbf{b}_2 \right) \frac{\sigma^2}{b},
\end{aligned}$$

where $\mathbf{q} = [0 \quad 0 \quad \frac{\eta L^2}{2n} \quad 0]^\top$, and the second inequality follows from the PL condition (Assumption 6). If we can establish that there exist there exist some constants c_1, c_2, c_3, c_4 such that

$$\mathbf{s}^\top (\mathbf{A} - (1 - \eta\mu) \mathbf{I}) + \mathbf{q}^\top \leq 0, \quad (\text{C.21a})$$

$$\frac{\eta}{2} - \frac{\eta^2 L}{2} - \mathbf{s}^\top \mathbf{b}_1 \geq 0, \quad (\text{C.21b})$$

we arrive at

$$\Phi^{(t+1)} \leq (1 - \eta\mu) \Phi^{(t)} + \left(\frac{\eta}{2n} + \mathbf{s}^\top \mathbf{b}_2 \right) \frac{\sigma^2}{b} \leq (1 - \eta\mu) \Phi^{(t)} + \frac{36c_4 \sigma^2}{c_\gamma L b \rho}.$$

Recurring the above relation then complete the proof.

It then boils down to establish (C.21). By similar arguments as Appendix C.3, in view of (C.20) and $\rho \leq 1, (1 - \alpha) \leq 1$, it is sufficient to show there exist constants $c_1, c_2, c_3, c_4, c_\gamma, c_\eta > 0$ such that

$$\begin{bmatrix}
1 - \frac{4c_\eta c_\gamma}{\kappa} & -72C c_\gamma^2 & -24C c_\gamma & -72C c_\gamma \\
0 & 1 - \frac{4c_\eta c_\gamma}{\kappa} & 0 & -24C c_\gamma \\
-12C c_\gamma & -35C c_\gamma & 1 - \frac{2c_\eta}{\kappa} & -36C \\
-24c_\eta^2 c_\gamma & -24c_\gamma(1 + 3c_\eta^2) & -24c_\eta^2 & 1 - \frac{4c_\eta}{\kappa} \\
-12c_\eta c_\gamma & -36c_\eta c_\gamma & 0 & -36c_\eta
\end{bmatrix}
\begin{bmatrix}
c_1 \\
c_2 \\
c_3 \\
c_4
\end{bmatrix}
\geq
\begin{bmatrix}
0 \\
0 \\
c_\eta \\
0 \\
-1 + c_\eta c_\gamma
\end{bmatrix}.$$

This can be easily verified by noting that as long as c_η and c_γ are set sufficiently small, it is straightforward to find feasible c_1, c_2, c_3, c_4 .

Appendix D

Appendix for Chapter 5

D.1 Proof of Theorem 9

This section proves Theorem 9 in the following steps: 1) define privacy loss and moment generating function, 2) define mechanisms and sub-mechanisms, 3) bound overall moment generating function and show the choice of perturbation variance satisfies all conditions.

Moment generating function Let o and aux denote an outcome and an auxiliary input, respectively. Then, we can define the privacy loss of an outcome o on neighboring dataset \mathbf{Z} and \mathbf{Z}_i as

$$c(o; \mathcal{M}, \text{aux}, \mathbf{Z}, \mathbf{Z}_i) = \log \frac{\mathbb{P}(\mathcal{M}(\text{aux}, \mathbf{Z}) = o)}{\mathbb{P}(\mathcal{M}(\text{aux}, \mathbf{Z}_i) = o)},$$

and its log moment generating functions as

$$\alpha_i^{\mathcal{M}}(\lambda; \text{aux}, \mathbf{Z}, \mathbf{Z}_i) = \log \mathbb{E}_{o \sim \mathcal{M}(\text{aux}, \mathbf{Z})} [\exp(\lambda c(o; \mathcal{M}, \text{aux}, \mathbf{Z}, \mathbf{Z}_i))].$$

Take maximum over conditions, the unconditioned log moment generating function is

$$\hat{\alpha}_i^{\mathcal{M}}(\lambda) = \max_{\text{aux}, \mathbf{Z}, \mathbf{Z}_i} \alpha_i^{\mathcal{M}}(\lambda; \text{aux}, \mathbf{Z}, \mathbf{Z}_i).$$

Sub-mechanisms Definition 7 defines the LDP mechanism, but it is not enough to model decentralized algorithms. To model the perturbation operation happens on agent i at time t , we define a sub-mechanism as $\mathcal{M}_i^{(t)} : \mathcal{D} \rightarrow \mathcal{R}$, where $i \in [n], t \in [T]$, which can be understood as the perturbation added on agent i at time t . In addition, we define another mechanism $\mathcal{C} : \mathcal{R} \rightarrow \mathcal{R}$ to model the compression operator and $\mathcal{C} \circ \mathcal{M}_i^{(t)}$ to represent the full update at an agent, and use \mathcal{M} to represent the full algorithm.

Proof of LDP The overall log moment generating function for agent i can be bounded using [LZLC22, Lemma 2] as

$$\hat{\alpha}_i^{\mathcal{M}}(\lambda) \leq \sum_{t=1}^T \hat{\alpha}_i^{\mathcal{C} \circ \mathcal{M}_i^{(t)}}(\lambda) \leq \sum_{t=1}^T \hat{\alpha}_i^{\mathcal{M}_i^{(t)}}(\lambda).$$

Let $q = \frac{b}{m}$ denote the probability each data sample is chose. For agent i and $\lambda > 0$, assume $q \leq \frac{\tau}{16\sigma_p}$ and $\lambda \leq \frac{\sigma_p^2}{\tau^2} \log \frac{\tau}{q\sigma_p}$. We can apply [ACG⁺16, Lemma 3] to bound each $\hat{\alpha}_i^{\mathcal{M}_i^{(t)}}(\lambda)$ as

$$\hat{\alpha}_i^{\mathcal{M}_i^{(t)}}(\lambda) \leq \frac{q^2 \lambda (\lambda + 1) \tau^2}{(1 - q) \sigma_p^2} + O\left(\frac{q^3 \lambda^3 \tau^3}{\sigma_p^3}\right) = O\left(\frac{q^2 \lambda^2 \tau^2}{\sigma_p^2}\right).$$

To conclude the proof, we can verify there exists some λ that satisfies the following inequalities when choosing $\sigma_p = \frac{\tau q \sqrt{T \log(1/\delta)}}{\epsilon}$ and $q = \frac{b}{m}$,

$$\begin{aligned} \left(\frac{Tq\tau\lambda}{\sigma_p}\right)^2 &\leq \frac{\lambda\epsilon}{2}, \\ \exp(-\lambda\epsilon/2) &\leq \delta, \\ \lambda &\leq \frac{\sigma_p^2}{\tau^2} \log \frac{\tau}{q\sigma_p}. \end{aligned}$$

D.2 Proof of Theorem 10

This section proves Theorem 10 in the following 4 subsections: Appendix D.2.1 derives the descent inequality, Appendices D.2.2 and D.2.3 create two linear systems to bound the sum of consensus errors in the descent inequality, and finally Appendix D.2.4 specifies hyper parameters to obtain convergence rate.

To reuse this section's results in Appendix D.3, we assume Assumption 10 in deriving descent lemma and linear systems, and lift this assumption when computing convergence rate in Appendix D.2.4 using $\sigma_g \leq 2\tau$.

D.2.1 Function value descent

Using Taylor expansion, and taking expectation conditioned on time t ,

$$\begin{aligned} \mathbb{E}_t[f(\bar{\mathbf{x}}^{(t+1)}) - f(\bar{\mathbf{x}}^{(t)})] &\leq \mathbb{E}_t\langle \nabla f(\bar{\mathbf{x}}^{(t)}), -\eta \bar{\mathbf{v}}^{(t+1)} \rangle + \frac{L}{2} \mathbb{E}_t \|\eta \bar{\mathbf{v}}^{(t+1)}\|_2^2 \\ &= -\eta \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \mathbb{E}_t[\bar{\mathbf{v}}^{(t+1)}] \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\ &= -\eta \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \mathbb{E}_t[\bar{\mathbf{g}}_\tau^{(t+1)} + \bar{\mathbf{e}}^{(t+1)}] \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2, \end{aligned}$$

where the last equality is due to $\bar{\mathbf{v}}^{(t)} = \bar{\mathbf{g}}_p^{(t)}$ that can be proved by induction.

Because $\mathbb{E}_t[\mathbf{e}_i^{(t)}] = \mathbf{0}_d$ and stochastic gradients are unbiased,

$$\mathbb{E}_t[f(\bar{\mathbf{x}}^{(t+1)}) - f(\bar{\mathbf{x}}^{(t)})]$$

$$\begin{aligned}
&= -\eta \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \nabla F(\mathbf{X}^{(t)}) \left(\frac{1}{n} \mathbf{1}_n \right) \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&= \frac{\eta}{2} \left(\|\nabla f(\bar{\mathbf{x}}^{(t)}) - \nabla F(\mathbf{X}^{(t)}) \left(\frac{1}{n} \mathbf{1}_n \right)\|_2^2 - \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 - \|\nabla F(\mathbf{X}^{(t)}) \left(\frac{1}{n} \mathbf{1}_n \right)\|_2^2 \right) + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\leq -\frac{\eta}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta L^2}{2n} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 - \frac{\eta}{2} \|\nabla F(\mathbf{X}^{(t)}) \left(\frac{1}{n} \mathbf{1}_n \right)\|_2^2, \tag{D.1}
\end{aligned}$$

where the last inequality is due to Assumption 2.

Let $\Delta = \mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^*$. Take full expectation and average (D.1) over $t = 1, \dots, T$, the expected utility can be bounded by

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 &\leq \frac{2\Delta}{\eta T} + \frac{1}{T} \cdot \frac{L^2}{n} \sum_{t=1}^T \mathbb{E} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2 \\
&\quad + \frac{1}{T} \cdot \eta L \sum_{t=1}^T \mathbb{E} \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) \left(\frac{1}{n} \mathbf{1}_n \right)\|_2^2. \tag{D.2}
\end{aligned}$$

D.2.2 Sum of variable consensus errors

This subsection creates a linear system to bound $\sum_{t=1}^T \mathbb{E} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2$ by $\sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2$ and $\sum_{t=1}^T \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2$. To simplify notations, let $\widehat{\mathbf{W}} = \mathbf{I}_n + \gamma(\mathbf{W} - \mathbf{I}_n)$, and denote the mixing rate of $\widehat{\mathbf{W}}$ by $\hat{\alpha} = \|\widehat{\mathbf{W}} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top\right)\|_{\text{op}}$. Lemma 10 analyzes the mixing rate of the regularized mixing matrix.

Lemma 10 (Mixing rate of regularized mixing matrix). *Assuming $0 < \gamma \leq 1$. The mixing rate of $\widehat{\mathbf{W}}$ can be bounded as*

$$\hat{\alpha} \leq 1 + \gamma(\alpha - 1). \tag{D.3}$$

Proof of Lemma 10. Let $\lambda_1 = 1 > \lambda_2 \geq \dots \geq \lambda_n > -1$ denote the eigenvalues of \mathbf{W} . Corresponding eigenvalues of $\widehat{\mathbf{W}}$ are $1 + \gamma(\lambda_i - 1)$, $i = 1, \dots, n$.

The mixing rate of $\widehat{\mathbf{W}}$ is

$$\begin{aligned}
\hat{\alpha} &= \max \{ |1 + \gamma(\lambda_2 - 1)|, |1 + \gamma(\lambda_n - 1)| \} \\
&\leq \max \{ |1 - \gamma| + \gamma|\lambda_2|, |1 - \gamma| + \gamma|\lambda_n| \} \\
&= 1 + \gamma(\alpha - 1).
\end{aligned}$$

□

Variable consensus error

Take expectation conditioned on time t , and use Young's inequality, the variable consensus error can be bounded as

$$\mathbb{E}_t \|\mathbf{X}^{(t+1)} - \bar{\mathbf{x}}^{(t+1)} \mathbf{1}_n^\top\|_F^2$$

$$\begin{aligned}
&= \mathbb{E}_t \left\| \left(\mathbf{X}^{(t)} + \gamma \mathbf{Q}_x^{(t+1)} (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t+1)} \right) \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&= \mathbb{E}_t \left\| \left(\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right) \left(\widehat{\mathbf{W}} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) + \gamma \left(\mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t+1)} \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\leq \frac{2}{1 + \hat{\alpha}^2} \left\| \left(\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right) \left(\widehat{\mathbf{W}} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\quad + \frac{2}{1 - \hat{\alpha}^2} \mathbb{E}_t \left\| \gamma \left(\mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t+1)} \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\leq \frac{2}{1 + \hat{\alpha}^2} \left\| \left(\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right) \left(\widehat{\mathbf{W}} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\quad + \frac{4}{1 - \hat{\alpha}^2} \mathbb{E}_t \left\| \gamma \left(\mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) \right\|_F^2 + \frac{4}{1 - \hat{\alpha}^2} \mathbb{E}_t \left\| \eta \mathbf{V}^{(t+1)} \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\stackrel{(i)}{\leq} \frac{2\hat{\alpha}^2}{1 + \hat{\alpha}^2} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{16\gamma^2}{1 - \hat{\alpha}^2} \mathbb{E}_t \left\| \mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{4\eta^2}{1 - \hat{\alpha}^2} \mathbb{E}_t \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top \right\|_F^2 \\
&\stackrel{(ii)}{\leq} \hat{\alpha} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{16(1 - \rho)\gamma^2}{1 - \hat{\alpha}} \left\| \mathbf{Q}_x^{(t)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{4\eta^2}{1 - \hat{\alpha}} \mathbb{E}_t \left\| \mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top \right\|_F^2, \tag{D.4}
\end{aligned}$$

where (i) is obtained by $\|\mathbf{W} - \mathbf{I}_n\|_{\text{op}} \leq 2$, (ii) uses $2\hat{\alpha} \leq 1 + \hat{\alpha}^2$, $1 - \hat{\alpha} \leq 1 - \hat{\alpha}^2$ and Definition 5.

Variable quantization error

Assume γ satisfies the following inequality (which will be verified in Appendix D.2.4)

$$\gamma^2 \leq \frac{\rho^2}{96(1 - \rho)}. \tag{D.5}$$

Take expectation conditioned on time t , the variable quantization error can be decomposed and bounded as

$$\begin{aligned}
&\mathbb{E}_t \left\| \mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t+1)} \right\|_F^2 \\
&= \mathbb{E}_t \left\| \mathbf{Q}_x^{(t)} + \mathcal{C}(\mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)}) - \mathbf{X}^{(t+1)} \right\|_F^2 \\
&= \mathbb{E}_t \left\| \mathcal{C}(\mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)}) - (\mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)}) - (\mathbf{X}^{(t+1)} - \mathbf{X}^{(t)}) \right\|_F^2 \\
&\stackrel{(i)}{\leq} \frac{2}{1 + (1 - \rho)} \mathbb{E}_t \left\| \mathcal{C}(\mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)}) - (\mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)}) \right\|_F^2 + \frac{2}{1 - (1 - \rho)} \mathbb{E}_t \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 \\
&\stackrel{(ii)}{\leq} \frac{2(1 - \rho)}{1 + (1 - \rho)} \left\| \mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)} \right\|_F^2 + \frac{2}{\rho} \mathbb{E}_t \left\| \mathbf{X}^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 \\
&= \frac{2(1 - \rho)}{1 + (1 - \rho)} \left\| \mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)} \right\|_F^2 \\
&\quad + \frac{2}{\rho} \mathbb{E}_t \left\| \gamma \left(\mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) + \gamma \left(\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right) (\mathbf{W} - \mathbf{I}_n) - \eta \mathbf{V}^{(t+1)} \right\|_F^2 \\
&\stackrel{(iii)}{\leq} \left(1 - \frac{\rho}{2} \right) \left\| \mathbf{X}^{(t)} - \mathbf{Q}_x^{(t)} \right\|_F^2 + \frac{24\gamma^2}{\rho} \mathbb{E}_t \left\| \mathbf{Q}_x^{(t+1)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{24\gamma^2}{\rho} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\eta^2}{\rho} \mathbb{E}_t \left\| \mathbf{V}^{(t+1)} \right\|_F^2 \\
&\leq \left(1 - \frac{\rho}{2} + \frac{24(1 - \rho)\gamma^2}{\rho} \right) \left\| \mathbf{Q}_x^{(t)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{24\gamma^2}{\rho} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\eta^2}{\rho} \mathbb{E}_t \left\| \mathbf{V}^{(t+1)} \right\|_F^2 \\
&\stackrel{(iv)}{\leq} \left(1 - \frac{\rho}{4} \right) \left\| \mathbf{Q}_x^{(t)} - \mathbf{X}^{(t)} \right\|_F^2 + \frac{24\gamma^2}{\rho} \left\| \mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top \right\|_F^2 + \frac{6\eta^2}{\rho} \mathbb{E}_t \left\| \mathbf{V}^{(t+1)} \right\|_F^2, \tag{D.6}
\end{aligned}$$

where (i) is obtained by applying Young's inequality, (ii) uses Definition 5, (iii) uses the fact $\|\mathbf{W} - \mathbf{I}_n\|_{\text{op}} \leq 2$, and (iv) uses (D.5).

Linear system

Let $\mathbf{e}_1^{(t)} = \begin{bmatrix} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_{\text{F}}^2 \\ \|\mathbf{Q}_x^{(t)} - \mathbf{X}^{(t)}\|_{\text{F}}^2 \end{bmatrix}$, we can take full expectation and rewrite (D.4) and (D.6) in matrix form as

$$\begin{aligned} \mathbb{E}[\mathbf{e}_1^{(t+1)}] &\leq \begin{bmatrix} \hat{\alpha} & \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \\ \frac{24\gamma^2}{\rho} & 1 - \frac{\rho}{4} \end{bmatrix} \mathbb{E}[\mathbf{e}_1^{(t)}] + \begin{bmatrix} \frac{4\eta^2}{1-\hat{\alpha}} \mathbb{E}\|\mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top\|_{\text{F}}^2 \\ \frac{6\eta^2}{\rho} \mathbb{E}\|\mathbf{V}^{(t+1)}\|_{\text{F}}^2 \end{bmatrix} \\ &:= \mathbf{G}_1 \mathbb{E}[\mathbf{e}_1^{(t)}] + \mathbf{b}_1^{(t)}. \end{aligned} \quad (\text{D.7})$$

We can compute $(\mathbf{I}_n - \mathbf{G}_1)^{-1}$ and verify all its entries are positive:

$$\begin{aligned} (\mathbf{I}_n - \mathbf{G}_1)^{-1} &= \frac{1}{(1-\hat{\alpha}) \cdot \frac{\rho}{4} - \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \cdot \frac{24\gamma^2}{\rho}} \begin{bmatrix} \frac{\rho}{4} & \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \\ \frac{24\gamma^2}{\rho} & 1 - \hat{\alpha} \end{bmatrix} \\ &\leq \frac{1}{\frac{1}{8}(1-\hat{\alpha})\rho} \begin{bmatrix} \frac{\rho}{4} & \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \\ \frac{24\gamma^2}{\rho} & 1 - \hat{\alpha} \end{bmatrix}, \end{aligned} \quad (\text{D.8})$$

where we assume the following inequality to prove (D.8), which will be validated in Appendix D.2.4:

$$(1-\hat{\alpha}) \cdot \frac{\rho}{4} - \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \cdot \frac{24\gamma^2}{\rho} \geq \frac{1}{8}(1-\hat{\alpha})\rho. \quad (\text{D.9})$$

Sum expected error vectors $\mathbb{E}[\mathbf{e}_1^{(t)}]$ over $t = 1, \dots, T$,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[\mathbf{e}_1^{(t)}] &\leq \sum_{t=1}^T (\mathbf{G}_1 \mathbb{E}[\mathbf{e}_1^{(t-1)}] + \mathbf{b}_1^{(t-1)}) \\ &\leq \mathbf{G}_1 \sum_{t=1}^T \mathbb{E}[\mathbf{e}_1^{(t)}] + \mathbf{G}_1 \mathbb{E}[\mathbf{e}_1^{(0)}] + \sum_{t=1}^T \mathbf{b}_1^{(t-1)}. \end{aligned}$$

Reorganize terms, multiply $(\mathbf{I}_n - \mathbf{G}_1)^{-1}$ on both sides and use $\mathbf{e}_1^{(0)} = \mathbf{0}_2$, the sum of error vectors can be bounded as

$$\sum_{t=1}^T \mathbb{E}[\mathbf{e}_1^{(t)}] \leq (\mathbf{I}_n - \mathbf{G}_1)^{-1} \sum_{t=0}^{T-1} \mathbf{b}_1^{(t)}. \quad (\text{D.10})$$

The sum of consensus error can be computed as

$$\begin{aligned} &\sum_{t=1}^T \mathbb{E}\|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_{\text{F}}^2 \\ &\leq \begin{bmatrix} 1 & 0 \end{bmatrix} (\mathbf{I}_n - \mathbf{G})^{-1} \sum_{t=0}^{T-1} \mathbf{b}_1^{(t)} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\frac{1}{8}(1-\hat{\alpha})\rho} \left[\frac{\rho}{4} \quad \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \right] \left[\begin{array}{c} \frac{4\eta^2}{1-\hat{\alpha}} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 \\ \frac{6\eta^2}{\rho} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)}\|_F^2 \end{array} \right] \\
&= \frac{\eta^2}{\frac{1}{8}(1-\hat{\alpha})\rho} \left(\frac{\rho}{1-\hat{\alpha}} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \cdot \frac{6}{\rho} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)}\|_F^2 \right) \\
&\stackrel{(i)}{=} \frac{8\eta^2}{(1-\hat{\alpha})\rho} \left(\frac{\rho}{1-\hat{\alpha}} + \frac{96(1-\rho)\gamma^2}{(1-\hat{\alpha})\rho} \right) \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{768(1-\rho)\gamma^2\eta^2}{(1-\hat{\alpha})^2\rho^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_F^2 \\
&\stackrel{(ii)}{\leq} \frac{16\eta^2}{(1-\hat{\alpha})^2} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{768(1-\rho)\gamma^2\eta^2}{(1-\hat{\alpha})^2\rho^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_F^2, \tag{D.11}
\end{aligned}$$

where we use the equality $\|\mathbf{V}^{(t)}\|_F^2 = \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + n\|\bar{\mathbf{v}}^{(t)}\|_2^2$ for (i) and use (D.5) for (ii).

D.2.3 Sum of gradient consensus errors

This section creates a linear system to bound the sum of gradient consensus error $\sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2$ by $\sum_{t=1}^T \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2$ and constant terms.

Gradient consensus error

Take expectation conditioned on time t and reorganize terms, the gradient consensus error can be expanded as

$$\begin{aligned}
&\mathbb{E}_t \|\mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top\|_F^2 \\
&= \mathbb{E}_t \left\| \left(\mathbf{V}^{(t)} + \gamma \mathbf{Q}_v^{(t+1)} (\mathbf{W} - \mathbf{I}_n) + \mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)} \right) \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&= \mathbb{E}_t \left\| \left(\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top \right) \left(\widehat{\mathbf{W}} - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) + \gamma \left(\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) + \left(\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)} \right) \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2.
\end{aligned}$$

Then, take full expectation, use the update formula and Young's inequality similarly to (D.4),

$$\begin{aligned}
&\mathbb{E} \|\mathbf{V}^{(t+1)} - \bar{\mathbf{v}}^{(t+1)} \mathbf{1}_n^\top\|_F^2 \\
&\leq \frac{2\hat{\alpha}^2}{1+\hat{\alpha}^2} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 \\
&\quad + \frac{2}{1-\hat{\alpha}^2} \mathbb{E} \left\| \gamma \left(\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) + \left(\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)} \right) \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\leq \hat{\alpha} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{4}{1-\hat{\alpha}} \mathbb{E} \left\| \gamma \left(\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)} \right) (\mathbf{W} - \mathbf{I}_n) \right\|_F^2 \\
&\quad + \frac{4}{1-\hat{\alpha}} \mathbb{E} \left\| \left(\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)} \right) \left(\mathbf{I}_n - \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \right) \right) \right\|_F^2 \\
&\stackrel{(i)}{\leq} \hat{\alpha} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{16\gamma^2}{1-\hat{\alpha}} \mathbb{E} \|\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)}\|_F^2 + \frac{4}{1-\hat{\alpha}} \mathbb{E} \|\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)}\|_F^2 \\
&\stackrel{(ii)}{\leq} \hat{\alpha} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \mathbb{E} \|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_F^2 + \frac{16n(\tau^2 + \sigma_p^2 d)}{1-\hat{\alpha}}, \tag{D.12}
\end{aligned}$$

where (i) is proved using the facts $\|\mathbf{W} - \mathbf{I}_n\|_{\text{op}} \leq 2$ and $\|\mathbf{I}_n - (\frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top)\|_{\text{op}} \leq 1$, (ii) is due to Definition 5 and

$$\begin{aligned} \mathbb{E}\|\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)}\|_{\text{F}}^2 &\leq 2\mathbb{E}\|\mathbf{G}_p^{(t+1)}\|_{\text{F}}^2 + 2\mathbb{E}\|\mathbf{G}_p^{(t)}\|_{\text{F}}^2 \\ &= 2(\mathbb{E}\|\mathbf{G}_\tau^{(t+1)}\|_{\text{F}}^2 + n\sigma_p^2d) + 2(\mathbb{E}\|\mathbf{G}_\tau^{(t)}\|_{\text{F}}^2 + n\sigma_p^2d) \\ &\leq 4n(\tau^2 + \sigma_p^2d). \end{aligned} \quad (\text{D.13})$$

Gradient quantization error

$$\begin{aligned} \mathbb{E}_t\|\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t+1)}\|_{\text{F}}^2 &= \mathbb{E}_t\|(\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)}) - (\mathbf{V}^{(t+1)} - \mathbf{V}^{(t)})\|_{\text{F}}^2 \\ &\leq \frac{2}{1 + (1-\rho)}\mathbb{E}_t\|\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 + \frac{2}{1 - (1-\rho)}\mathbb{E}_t\|\mathbf{V}^{(t+1)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 \\ &\leq \frac{2(1-\rho)}{2-\rho}\|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 + \frac{2}{\rho}\mathbb{E}_t\|\gamma\mathbf{Q}_v^{(t+1)}(\mathbf{W} - \mathbf{I}_n) + \mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)}\|_{\text{F}}^2 \\ &\leq \frac{2(1-\rho)}{2-\rho}\|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 + \frac{6\gamma^2}{\rho}\mathbb{E}_t\|(\mathbf{Q}_v^{(t+1)} - \mathbf{V}^{(t)})(\mathbf{W} - \mathbf{I}_n)\|_{\text{F}}^2 \\ &\quad + \frac{6\gamma^2}{\rho}\|\mathbf{V}^{(t)}(\mathbf{W} - \mathbf{I}_n)\|_{\text{F}}^2 + \frac{6}{\rho}\mathbb{E}_t\|\mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)}\|_{\text{F}}^2 \\ &\stackrel{\text{(i)}}{\leq} \left(1 - \frac{\rho}{2} + \frac{24\gamma^2(1-\rho)}{\rho}\right)\|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 + \frac{24\gamma^2}{\rho}\|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)}\mathbf{1}_n^\top\|_{\text{F}}^2 + \frac{24n(\tau^2 + \sigma_p^2d)}{\rho} \\ &\stackrel{\text{(ii)}}{\leq} \left(1 - \frac{\rho}{4}\right)\|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 + \frac{24\gamma^2}{\rho}\|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)}\mathbf{1}_n^\top\|_{\text{F}}^2 + \frac{24n(\tau^2 + \sigma_p^2d)}{\rho}, \end{aligned} \quad (\text{D.14})$$

where we use (D.13) and the fact $\frac{2(1-\rho)}{2-\rho} = 1 - \frac{\rho}{2-\rho} \geq 1 - \frac{\rho}{2}$ when $\rho \geq 0$ to reach (i) and use (D.5) to reach (ii).

Linear system

Let $\mathbf{e}_2^{(t)} = \begin{bmatrix} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)}\mathbf{1}_n^\top\|_{\text{F}}^2 \\ \|\mathbf{Q}_v^{(t)} - \mathbf{V}^{(t)}\|_{\text{F}}^2 \end{bmatrix}$. We can write (D.12) and (D.14) in matrix form as

$$\begin{aligned} \mathbb{E}[\mathbf{e}_2^{(t+1)}] &\leq \begin{bmatrix} \hat{\alpha} & \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \\ \frac{24\gamma^2}{\rho} & 1 - \frac{\rho}{4} \end{bmatrix} \mathbb{E}[\mathbf{e}_2^{(t)}] + \begin{bmatrix} \frac{16n(\tau^2 + \sigma_p^2d)}{1-\hat{\alpha}} \\ \frac{24n(\tau^2 + \sigma_p^2d)}{\rho} \end{bmatrix} \\ &:= \mathbf{G}_2\mathbb{E}[\mathbf{e}_2^{(t)}] + \mathbf{b}_2^{(t)}. \end{aligned}$$

Because $\mathbf{G}_2 = \mathbf{G}_1$, we can use the same argument as in Appendix D.2.2, and use (D.8) to prove

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}\|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)}\mathbf{1}_n^\top\|_{\text{F}}^2 &\leq \begin{bmatrix} 1 & 0 \end{bmatrix} (\mathbf{I}_n - \mathbf{G}_2)^{-1} \left(\mathbb{E}[\mathbf{e}_2^{(0)}] + \sum_{t=0}^{T-1} \mathbf{b}_2^{(t)} \right) \\ &\leq \frac{1}{\frac{1}{8}(1-\hat{\alpha})\rho} \begin{bmatrix} \frac{\rho}{4} & \frac{16(1-\rho)\gamma^2}{1-\hat{\alpha}} \end{bmatrix} \begin{bmatrix} \frac{16Tn(\tau^2 + \sigma_p^2d)}{1-\hat{\alpha}} \\ \frac{24Tn(\tau^2 + \sigma_p^2d)}{\rho} \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \frac{Tn(\tau^2 + \sigma_p^2 d)}{\frac{1}{8}(1 - \hat{\alpha})\rho} \cdot \left(\frac{\rho}{4} \cdot \frac{16}{1 - \hat{\alpha}} + \frac{16(1 - \rho)\gamma^2}{1 - \hat{\alpha}} \cdot \frac{24}{\rho} \right) \\
&\leq \frac{Tn(\tau^2 + \sigma_p^2 d)}{\frac{1}{8}(1 - \hat{\alpha})\rho} \cdot \left(\frac{4\rho}{1 - \hat{\alpha}} + \frac{4\rho}{1 - \hat{\alpha}} \right) \\
&= \frac{64}{(1 - \hat{\alpha})^2} Tn(\tau^2 + \sigma_p^2 d), \tag{D.15}
\end{aligned}$$

where we use (D.5) to prove the last inequality.

With (D.15), we can bound (D.11) by

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2 \\
&\leq \frac{16\eta^2}{(1 - \hat{\alpha})^2} \sum_{t=1}^T \mathbb{E} \|\mathbf{V}^{(t)} - \bar{\mathbf{v}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{768(1 - \rho)\gamma^2\eta^2}{(1 - \hat{\alpha})^2\rho^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_F^2 \\
&\leq \frac{16\eta^2}{(1 - \hat{\alpha})^2} \cdot \frac{64}{(1 - \hat{\alpha})^2} Tn(\tau^2 + \sigma_p^2 d) + \frac{768(1 - \rho)\gamma^2\eta^2}{(1 - \hat{\alpha})^2\rho^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_F^2 \\
&\leq \frac{1024\eta^2}{(1 - \hat{\alpha})^4} Tn(\tau^2 + \sigma_p^2 d) + \frac{8\eta^2}{(1 - \hat{\alpha})^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_F^2, \tag{D.16}
\end{aligned}$$

where we use (D.5) again to prove the last inequality.

D.2.4 Convergence rate

Note bounded gradient assumption can imply Assumption 10 for some $\sigma_g \leq 2\tau$, we can bound the expected norm of average gradient estimate as

$$\begin{aligned}
\mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2 &= \mathbb{E} \|\bar{\mathbf{g}}_p^{(t)}\|_2^2 \\
&= \mathbb{E} \|\bar{\mathbf{g}}_\tau^{(t)}\|_2^2 + \frac{\sigma_p^2 d}{n} \\
&\leq \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n)\|_2^2 + \frac{\sigma_g^2}{b} + \frac{\sigma_p^2 d}{n} \\
&\leq \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n)\|_2^2 + \frac{4\tau^2}{b} + \frac{\sigma_p^2 d}{n}. \tag{D.17}
\end{aligned}$$

We assume

$$\eta L \leq \frac{1}{8} (1 - \hat{\alpha})^{\frac{4}{3}}. \tag{D.18}$$

Using (D.16) (D.17), expected utility (D.2) can be bounded by

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 &\leq \frac{2\Delta}{\eta T} + \frac{1}{T} \cdot \frac{L^2}{n} \sum_{t=1}^T \mathbb{E} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2 \\
&\quad + \frac{1}{T} \sum_{t=1}^T \eta L \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2 - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n)\|_2^2
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{2\Delta}{\eta T} + \frac{1}{T} \cdot \frac{L^2}{n} \left(\frac{1024\eta^2}{(1-\hat{\alpha})^4} Tn(\tau^2 + \sigma_p^2 d) + \frac{8\eta^2}{(1-\hat{\alpha})^2} \sum_{t=1}^T n \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2 \right) \\
&\quad + \frac{1}{T} \sum_{t=1}^T \eta L \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2 - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n)\|_2^2 \\
&\stackrel{(i)}{\leq} \frac{2\Delta}{\eta T} + \frac{1024\eta^2 L^2}{(1-\hat{\alpha})^4} (\tau^2 + \sigma_p^2 d) + \frac{2\eta L}{(1-\hat{\alpha})^{\frac{4}{3}} T} \sum_{t=1}^T \mathbb{E} \|\bar{\mathbf{v}}^{(t)}\|_2^2 - \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n)\|_2^2 \\
&\leq \frac{2\Delta}{\eta T} + \frac{1024\eta^2 L^2}{(1-\hat{\alpha})^4} (\tau^2 + \sigma_p^2 d) + \frac{2\eta L}{(1-\hat{\alpha})^{\frac{4}{3}}} \left(\frac{4\tau^2}{b} + \frac{\sigma_p^2 d}{n} \right) \\
&\stackrel{(ii)}{=} \frac{2\Delta}{\eta T} + \frac{1024\eta^2 L^2 \tau^2}{(1-\hat{\alpha})^4} (1 + T\phi_m^2) + \frac{8\eta L \tau^2}{(1-\hat{\alpha})^{\frac{4}{3}}} (1 + T\phi_m^2) \\
&\stackrel{(iii)}{=} \frac{2\Delta}{\eta T} + \frac{2048\eta^2 L^2 \tau^2}{(1-\hat{\alpha})^4} + \frac{16\eta L \tau^2}{(1-\hat{\alpha})^{\frac{4}{3}}} \tag{D.19}
\end{aligned}$$

where we use (D.18) for (i), substitute $b = 1$ and $\sigma_p^2 d = \left(\frac{\tau \sqrt{T \log(1/\delta)}}{m\epsilon} \right)^2 = T\tau^2 \phi_m^2$ for (ii), and substitute $T = \phi_m^{-2}$ for (iii).

We set the step size as

$$\eta = \frac{\gamma^{\frac{4}{3}} (1-\alpha)^{\frac{4}{3}}}{32} \cdot \frac{\phi_m}{L},$$

(D.19) can be further bounded as

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 &\leq \frac{64L\Delta\phi_m}{\gamma^{\frac{4}{3}} (1-\alpha)^{\frac{4}{3}}} + \frac{2\tau^2 \phi_m^2}{(1-\hat{\alpha})^{\frac{4}{3}}} + \frac{\tau^2 \phi_m}{2} \\
&\leq \frac{64L\Delta\phi_m}{\gamma^{\frac{4}{3}} (1-\alpha)^{\frac{4}{3}}} + \frac{3\tau^2 \phi_m}{(1-\hat{\alpha})^{\frac{4}{3}}} \\
&\leq \frac{67\phi_m}{\gamma^{\frac{4}{3}} (1-\alpha)^{\frac{4}{3}}} \max\{\tau^2, L\Delta\},
\end{aligned}$$

where we use the assumption in (5.2) that $\phi_m < 1$ to prove the second inequality, and use Lemma 10 to reach the last inequality.

Lastly, set the hyper parameter γ as

$$\gamma = \frac{1}{100} (1-\alpha)\rho.$$

We can now verify conditions (D.5), (D.9) and the condition on η are all met to conclude the proof:

$$\gamma^2 \leq \frac{\rho^2}{10000} \quad \Rightarrow \quad \tag{D.5}$$

$$\gamma^4 = \gamma^2 \cdot \frac{(1-\alpha)^2 \rho^2}{10000} \leq \frac{(1-\hat{\alpha})^2 \rho^2}{10000} \quad \Rightarrow \quad \tag{D.9}$$

$$\eta L \leq \frac{(1-\hat{\alpha})^{\frac{4}{3}}}{32} \quad \Rightarrow \quad \tag{D.18}$$

D.3 Proof of Theorem 11

This section proves Theorem 11 in 2 subsections. Appendix D.3.1 derives the descent inequality using results from Appendices D.2.2 and D.2.3. Appendix D.3.2 first assumes all expected gradient norm $\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2$ are greater than a threshold ν (i.e. $\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \geq \nu$ for all $t = 1, \dots, T$), then specifies parameters and proves the average of expected gradient norm is smaller than that threshold $\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \leq \nu$, which contradicts the assumption hence proves the algorithm reaches $\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \leq \nu$ within T steps.

D.3.1 Function value descent

Let $\delta_i^{(t)} = \frac{\tau}{\tau + \|\mathbf{g}_i^{(t)}\|_2}$ and $\delta^{(t)} = \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2}$. Similar to Appendix D.2.1, use Taylor expansion and take expectation conditioned on t , we can expand the function value descent as

$$\begin{aligned}
& \mathbb{E}_t [f(\bar{\mathbf{x}}^{(t+1)}) - f(\bar{\mathbf{x}}^{(t)})] \\
& \leq \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), -\eta \bar{\mathbf{v}}^{(t+1)} \rangle + \frac{L}{2} \mathbb{E}_t \|\eta \bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
& = -\eta \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \bar{\mathbf{g}}_p^{(t+1)} \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
& = -\eta \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \bar{\mathbf{g}}_\tau^{(t+1)} \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
& = -\eta \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) \rangle + \eta \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)} \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
& = -\eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \eta \mathbb{E}_t \langle \nabla f(\bar{\mathbf{x}}^{(t)}), \text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)} \rangle + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
& \leq -\eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \eta \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \mathbb{E}_t \|\text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)}\|_2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2. \tag{D.20}
\end{aligned}$$

The $\mathbb{E}_t \|\text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)}\|_2$ term in (D.20) is the error introduced by gradient clipping, which can be analyzed by splitting it to 4 terms as following

$$\begin{aligned}
& \mathbb{E}_t \|\text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)}\|_2 \\
& = \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \frac{\tau}{\tau + \|\mathbf{g}_i^{(t)}\|_2} \mathbf{g}_i^{(t)} - \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \nabla f(\bar{\mathbf{x}}^{(t)}) \right\|_2 \\
& = \mathbb{E}_t \left\| \frac{1}{n} \sum_{i=1}^n \left(\frac{\tau}{\tau + \|\mathbf{g}_i^{(t)}\|_2} \mathbf{g}_i^{(t)} - \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \mathbf{g}_i^{(t)} \right) \right. \\
& \quad + \frac{1}{n} \sum_{i=1}^n \left(\frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \mathbf{g}_i^{(t)} - \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \nabla f_i(\mathbf{x}_i^{(t)}) \right) \\
& \quad + \frac{1}{n} \sum_{i=1}^n \left(\frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \nabla f_i(\mathbf{x}_i^{(t)}) - \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \nabla f_i(\mathbf{x}_i^{(t)}) \right) \\
& \quad \left. + \frac{1}{n} \sum_{i=1}^n \left(\frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \nabla f_i(\mathbf{x}_i^{(t)}) - \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \nabla f(\bar{\mathbf{x}}^{(t)}) \right) \right\|_2
\end{aligned}$$

$$\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \left(\frac{\tau}{\tau + \|\mathbf{g}_i^{(t)}\|_2} - \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \right) \mathbf{g}_i^{(t)} \right\|_2 \quad (\text{D.21})$$

$$+ \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} (\mathbf{g}_i^{(t)} - \nabla f_i(\mathbf{x}_i^{(t)})) \right\|_2 \quad (\text{D.22})$$

$$+ \frac{1}{n} \sum_{i=1}^n \left\| \left(\frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} - \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \right) \nabla f_i(\mathbf{x}_i^{(t)}) \right\|_2 \quad (\text{D.23})$$

$$+ \left\| \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \left(\frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^{(t)}) - \nabla f(\bar{\mathbf{x}}^{(t)}) \right) \right\|_2. \quad (\text{D.24})$$

Next, we bound each term separately using triangle inequality, Assumptions 9 and 10.

Bound the first term (D.21) as

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \left(\frac{\tau}{\tau + \|\mathbf{g}_i^{(t)}\|_2} - \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \right) \mathbf{g}_i^{(t)} \right\|_2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \frac{\tau(\|\mathbf{g}_i^{(t)}\|_2 - \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2)}{(\tau + \|\mathbf{g}_i^{(t)}\|_2)(\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2)} \mathbf{g}_i^{(t)} \right\|_2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left(\left| \|\mathbf{g}_i^{(t)}\|_2 - \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2 \right| \cdot \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \cdot \frac{\|\mathbf{g}_i^{(t)}\|_2}{\tau + \|\mathbf{g}_i^{(t)}\|_2} \right) \\ &\leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left| \|\mathbf{g}_i^{(t)}\|_2 - \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2 \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n \sqrt{\mathbb{E}_t (\|\mathbf{g}_i^{(t)}\|_2 - \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2)^2} \\ &= \frac{1}{n} \sum_{i=1}^n \sqrt{\mathbb{E}_t (\|\mathbf{g}_i^{(t)}\|_2^2 + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2^2 - 2\|\mathbf{g}_i^{(t)}\|_2 \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2)} \\ &\leq \frac{1}{n} \sum_{i=1}^n \sqrt{\mathbb{E}_t (\|\mathbf{g}_i^{(t)}\|_2^2 + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2^2 - 2\langle \mathbf{g}_i^{(t)}, \nabla f_i(\mathbf{x}_i^{(t)}) \rangle)} \\ &= \frac{1}{n} \sum_{i=1}^n \sqrt{\mathbb{E}_t \|\mathbf{g}_i^{(t)} - \nabla f_i(\mathbf{x}_i^{(t)})\|_2^2} \\ &\leq \frac{\sigma_g}{\sqrt{b}}. \end{aligned} \quad (\text{D.25})$$

Bound the second term (D.22) as

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_t \left\| \frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} (\mathbf{g}_i^{(t)} - \nabla f_i(\mathbf{x}_i^{(t)})) \right\|_2 \leq \frac{1}{n} \sum_{i=1}^n \frac{\tau \sigma_g / \sqrt{b}}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} \leq \frac{\sigma_g}{\sqrt{b}}. \quad (\text{D.26})$$

Bound the third term (D.23) as

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \left\| \left(\frac{\tau}{\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2} - \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \right) \nabla f_i(\mathbf{x}_i^{(t)}) \right\|_2 \\ &= \frac{1}{n} \sum_{i=1}^n \left\| \frac{\tau(\|\nabla f_i(\mathbf{x}_i^{(t)})\|_2 - \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2)}{(\tau + \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2)(\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2)} \nabla f_i(\mathbf{x}_i^{(t)}) \right\|_2 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{n} \sum_{i=1}^n \frac{\tau \left| \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2 - \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \right|}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \\
&\leq \frac{1}{n} \sum_{i=1}^n \delta^{(t)} \left| \|\nabla f_i(\mathbf{x}_i^{(t)})\|_2 - \|\nabla f_i(\bar{\mathbf{x}}^{(t)})\|_2 \right| + \frac{1}{n} \sum_{i=1}^n \delta^{(t)} \left| \|\nabla f_i(\bar{\mathbf{x}}^{(t)})\|_2 - \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \right| \\
&\leq \frac{1}{n} \sum_{i=1}^n \delta^{(t)} L \|\mathbf{x}_i^{(t)} - \bar{\mathbf{x}}^{(t)}\|_2 + \frac{1}{n} \sum_{i=1}^n \delta^{(t)} \cdot \frac{1}{12} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\
&\leq \frac{\delta^{(t)} L}{\sqrt{n}} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F + \frac{1}{12} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2,
\end{aligned} \tag{D.27}$$

where we use $\delta^{(t)} \leq 1$ to reach the last inequality.

Bound (D.24) as

$$\begin{aligned}
&\left\| \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \left(\frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^{(t)}) - \nabla f(\bar{\mathbf{x}}^{(t)}) \right) \right\|_2 \\
&= \frac{\tau}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \|\nabla F(\mathbf{X}^{(t)}) (\frac{1}{n} \mathbf{1}_n) - \nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\
&\leq \frac{\delta^{(t)} L}{\sqrt{n}} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F.
\end{aligned} \tag{D.28}$$

Using (D.25), (D.26), (D.27) and (D.28), the function value descent inequality (D.20) becomes

$$\begin{aligned}
\mathbb{E}_t [f(\bar{\mathbf{x}}^{(t+1)}) - f(\bar{\mathbf{x}}^{(t)})] &\leq -\eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\quad + \eta \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \mathbb{E}_t \|\text{Clip}_\tau(\nabla f(\bar{\mathbf{x}}^{(t)})) - \bar{\mathbf{g}}_\tau^{(t+1)}\|_2 \\
&\leq -\eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\quad + \eta \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \left(\frac{2\sigma_g}{\sqrt{b}} + \frac{1}{12} \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 + \frac{2\delta^{(t)} L}{\sqrt{n}} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F \right) \\
&= -\frac{11}{12} \eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\quad + \eta \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \left(\frac{2\sigma_g}{\sqrt{b}} + \frac{2\delta^{(t)} L}{\sqrt{n}} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F \right) \\
&\leq -\frac{5}{12} \eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{\eta^2 L}{2} \mathbb{E}_t \|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\quad + \frac{2\eta\sigma_g}{\sqrt{b}} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 + \frac{2\delta^{(t)} \eta L^2}{n} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2,
\end{aligned} \tag{D.29}$$

where the last inequality is due to

$$\begin{aligned}
&\eta \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \cdot \frac{2\delta^{(t)} L}{\sqrt{n}} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F \\
&\leq \eta \delta^{(t)} \cdot 2 \cdot \sqrt{\frac{1}{2} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2} \cdot \sqrt{\frac{2L^2}{n} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2} \\
&\leq \frac{1}{2} \eta \delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2 + \frac{2\delta^{(t)} \eta L^2}{n} \|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2.
\end{aligned}$$

D.3.2 Convergence rate

Different from (D.17), with the use of gradient clipping operator, we can only bound the expected norm of average gradient estimate as

$$\begin{aligned}
\mathbb{E}\|\bar{\mathbf{g}}^{(t)}\|_2^2 &= \mathbb{E}\|\bar{\mathbf{g}}_p^{(t)}\|_2^2 \\
&= \mathbb{E}\|\bar{\mathbf{g}}_\tau^{(t)} + \bar{\mathbf{e}}^{(t)}\|_2^2 \\
&= \mathbb{E}\|\bar{\mathbf{g}}_\tau^{(t)}\|_2^2 + \mathbb{E}\|\bar{\mathbf{e}}^{(t)}\|_2^2 \\
&\leq \tau^2 + \frac{\sigma_p^2 d}{n}.
\end{aligned} \tag{D.30}$$

Let $\Delta = \mathbb{E}[f(\bar{\mathbf{x}}^{(0)})] - f^*$. The techniques used is similar to that used in Appendix D.2, so that we can reuse results from Appendices D.2.2 and D.2.3, namely (D.16), in the following proof. Take full expectation and use (D.30), sum (D.29) over $t = 1, \dots, T$,

$$\begin{aligned}
-\Delta &\leq -\frac{5\eta}{12} \sum_{t=1}^T \mathbb{E}(\delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2) + \frac{2\eta\sigma}{\sqrt{b}} \sum_{t=1}^T \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\
&\quad + \frac{2\eta L^2}{n} \sum_{t=1}^T \mathbb{E}\|\mathbf{X}^{(t)} - \bar{\mathbf{x}}^{(t)} \mathbf{1}_n^\top\|_F^2 + \frac{\eta^2 L}{2} \sum_{t=1}^T \mathbb{E}\|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&\leq -\frac{5\eta}{12} \sum_{t=1}^T \mathbb{E}(\delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2) + \frac{2\eta\sigma}{\sqrt{b}} \sum_{t=1}^T \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\
&\quad + \frac{2\eta L^2}{n} \left(\frac{8\eta^2}{(1-\hat{\alpha})^2} \sum_{t=1}^T n \mathbb{E}\|\bar{\mathbf{v}}^{(t)}\|_2^2 + \frac{1024\eta^2}{(1-\hat{\alpha})^4} T n (\tau^2 + \sigma_p^2 d) \right) + \frac{\eta^2 L}{2} \sum_{t=1}^T \mathbb{E}\|\bar{\mathbf{v}}^{(t+1)}\|_2^2 \\
&= -\frac{5\eta}{12} \sum_{t=1}^T \mathbb{E}(\delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2) + \frac{2\eta\sigma}{\sqrt{b}} \sum_{t=1}^T \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\
&\quad + \frac{16\eta^3 L^2}{(1-\hat{\alpha})^2} \left(\tau^2 + \frac{\sigma_p^2 d}{n} \right) + \frac{2048\eta^3 L^2}{(1-\hat{\alpha})^4} T (\tau^2 + \sigma_p^2 d) + \frac{\eta^2 L T}{2} \left(\tau^2 + \frac{\sigma_p^2 d}{n} \right).
\end{aligned} \tag{D.31}$$

Compared to the descent inequality in Appendix D.2.4, (D.31) has a unique terms: $\mathbb{E}(\delta^{(t)} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2)$ which is expected norm of clipped gradients. To be able to analyze this terms, we need to use convexity and monotonicity from Lemma 11.

Lemma 11. Let $g(x) = \frac{x}{c+x}$ and $h(x) = xg(x) = \frac{x^2}{c+x}$. When $x \geq 0$, $g(x)$ and $h(x)$ increase monotonically, while $g(x)$ is concave and $h(x)$ is convex.

Proof of Lemma 11. It is sufficient to prove Lemma 11 by evaluating the first-order and second-order derivatives of $g(x)$ and $h(x)$.

Because $g'(x) = \frac{(c+x)-x}{(c+x)^2} = \frac{c}{(c+x)^2} > 0$ and $h'(x) = g(x) + xg'(x) \geq 0$, $g(x)$ and $h(x)$ increase monotonically.

$g(x)$ is concave because $g''(x) = \frac{2c(c+x)}{(c+x)^4} = -\frac{2c}{(c+x)^3} < 0$.

$h(x)$ is convex because $h''(x) = 2g'(x) + xg''(x) = \frac{2c}{(c+x)^2} - \frac{2cx}{(c+x)^3} = \frac{2c^2}{(c+x)^3} > 0$. \square

Next, we substitute $\tau = \nu$ (cf. Theorem 11), and assume the following inequality

$$\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \geq \nu. \quad (\text{D.32})$$

By Lemma 11, the expectation of clipped gradients can be bounded as

$$\begin{aligned} \mathbb{E}\left(\delta^{(t)}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2\right) &= \mathbb{E}\left(\frac{\tau\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2^2}{\tau + \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2}\right) \\ &\geq \frac{\tau(\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2)^2}{\tau + \mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2} \\ &\geq \frac{\tau\nu}{\tau + \nu}\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\ &= \frac{\nu}{2}\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2. \end{aligned} \quad (\text{D.33})$$

Using (D.30), (D.33) and Assumptions 9 and 10, we can further bound the RHS of (D.31) as

$$\begin{aligned} -\Delta &\leq -\frac{5\eta\nu}{24}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 + \frac{2\eta\sigma_g}{\sqrt{b}}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \\ &\quad + \frac{16\eta^3L^2}{(1-\hat{\alpha})^2}\left(\tau^2 + \frac{\sigma_p^2d}{n}\right) + \frac{2048\eta^3L^2}{(1-\hat{\alpha})^4}T(\tau^2 + \sigma_p^2d) + \frac{\eta^2LT}{2}\left(\tau^2 + \frac{\sigma_p^2d}{n}\right) \\ &\leq -\frac{\eta\nu}{8}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 + \frac{2048T\eta^3L^2}{(1-\hat{\alpha})^4}(\tau^2 + \sigma_p^2d) + \frac{3T\eta^2L}{1-\hat{\alpha}}\left(\tau^2 + \frac{\sigma_p^2d}{n}\right) \\ &= -\frac{\eta\nu}{8}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 + \frac{2048T\eta^3L^2\tau^2}{(1-\hat{\alpha})^4}(1 + Tb^2\phi_m^2) + \frac{3T\eta^2L\tau^2}{1-\hat{\alpha}}(1 + Tb^2\phi_m^2), \end{aligned} \quad (\text{D.34})$$

where we use (D.18) and $b = \left(\frac{24\sigma_g}{\nu}\right)^2$ to prove the last inequality.

Reorganize terms, (D.34) can be further bounded as

$$\begin{aligned} \frac{1}{T}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 &\leq \frac{8\Delta}{\eta\nu T} + \frac{16384\eta^2L^2\nu}{(1-\hat{\alpha})^4}(1 + Tb^2\phi_m^2) + \frac{24\eta L\nu}{1-\hat{\alpha}}(1 + Tb^2\phi_m^2) \\ &\stackrel{(i)}{\leq} \frac{8\Delta b\phi_m}{\eta\nu} + \frac{32768\eta^2L^2\nu}{(1-\hat{\alpha})^4} + \frac{48\eta L\nu}{1-\hat{\alpha}} \\ &\stackrel{(ii)}{\leq} \frac{8\Delta b\phi_m}{\eta\nu} + \frac{4096\eta L\nu}{(1-\hat{\alpha})^{\frac{8}{3}}} + \frac{48\eta L\nu}{1-\hat{\alpha}}, \\ &\leq \frac{8\Delta b\phi_m}{\eta\nu} + \frac{4144\eta L\nu}{(1-\hat{\alpha})^{\frac{8}{3}}}, \end{aligned} \quad (\text{D.35})$$

where we substitute $T = (b\phi_m)^{-1}$ to prove (i), use (D.18) for (ii), and use $(1-\hat{\alpha})^{-1} \leq (1-\hat{\alpha})^{-8/3}$ to prove the last inequality.

Set $\eta = \frac{\gamma^{\frac{4}{3}}(1-\alpha)^{\frac{4}{3}}}{24} \cdot \sqrt{\frac{\Delta}{L}} \cdot \frac{\sqrt{b\phi_m}}{\tau}$ and $\gamma = \frac{1}{100}(1-\alpha)\rho$, and use $b = \left(\frac{24\sigma_g}{\nu}\right)^2$, (D.35) can be further bounded

as

$$\frac{1}{T}\sum_{t=1}^T\mathbb{E}\|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 \leq \frac{192\sqrt{L\Delta b\phi_m}}{(1-\hat{\alpha})^{\frac{4}{3}}} + \frac{173\sqrt{L\Delta b\phi_m}}{(1-\hat{\alpha})^{\frac{4}{3}}}$$

$$\begin{aligned}
&= \frac{365\sqrt{L\Delta b\phi_m}}{(1-\hat{\alpha})^{\frac{4}{3}}} \\
&= \frac{8760\sqrt{L\Delta\phi_m}}{(1-\hat{\alpha})^{\frac{4}{3}}} \cdot \frac{\sigma_g}{\nu}.
\end{aligned} \tag{D.36}$$

Choosing $\nu = \sqrt{\frac{8760\sigma_g\sqrt{L\Delta\phi_m}}{(1-\hat{\alpha})^{\frac{4}{3}}}}$, (D.36) simplifies to $\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 < \nu$, which further implies $\exists t \in [T]$ such that $\mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(t)})\|_2 < \nu$. However, this contradicts the assumption (D.32), which leads to the convergence results in the theorem.

Lastly, we can verify conditions (D.5), (D.9) and (D.18) are all met to conclude the proof:

$$\gamma^2 \leq \frac{\rho^2}{10000} \Rightarrow \tag{D.5}$$

$$\gamma^4 = \gamma^2 \cdot \frac{(1-\alpha)^2 \rho^2}{10000} \leq \frac{(1-\hat{\alpha})^2 \rho^2}{10000} \Rightarrow \tag{D.9}$$

$$\begin{aligned}
\eta L &= \frac{\gamma^{\frac{4}{3}}(1-\alpha)^{\frac{4}{3}}}{24} \cdot \frac{\sqrt{bL\Delta\phi_m}}{\tau} \\
&= \gamma^{\frac{4}{3}}(1-\alpha)^{\frac{4}{3}} \cdot \frac{\sigma_g\sqrt{L\Delta\phi_m}}{\nu^2} \\
&= \gamma^{\frac{4}{3}}(1-\alpha)^{\frac{4}{3}} \cdot \frac{(1-\hat{\alpha})^{\frac{4}{3}}}{8760} \Rightarrow \tag{D.18}
\end{aligned}$$

Bibliography

- [ABCP13] M. E. Andrés, N. E. Bordenabe, K. Chatzikokolakis, and C. Palamidessi. Geo-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security, CCS '13*, pages 901–914, New York, NY, USA, November 2013. Association for Computing Machinery.
- [ACCÖ21] S. Asoodeh, W.-N. Chen, F. P. Calmon, and A. Özgür. Differentially private federated learning: An information-theoretic perspective. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pages 344–349, July 2021.
- [ACG⁺16] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 308–318, October 2016.
- [AGL⁺17] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. QSGD: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [AHJ⁺18] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli. The convergence of sparsified gradient methods. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [ALBR19] M. Assran, N. Loizou, N. Ballas, and M. Rabbat. Stochastic gradient push for distributed deep learning. In *International Conference on Machine Learning*, pages 344–353. PMLR, 2019.
- [AS14] M. Arioli and J. Scott. Chebyshev acceleration of iterative refinement. *Numerical Algorithms*, 66(3):591–608, 2014.
- [AZ17] Z. Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017*, page 1200–1205, New York, NY, USA, 2017. Association for Computing Machinery.

- [AZH16] Z. Allen-Zhu and E. Hazan. Variance reduction for faster non-convex optimization. In *International Conference on Machine Learning*, pages 699–707. PMLR, 2016.
- [BBKW19] A. S. Berahas, R. Bollapragada, N. S. Keskar, and E. Wei. Balancing communication and computation in distributed optimization. *IEEE Transactions on Automatic Control*, 64(8):3141–3155, 2019.
- [BBP13] Y. Bengio, N. Boulanger-Lewandowski, and R. Pascanu. Advances in optimizing recurrent networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8624–8628, May 2013.
- [BBW21] A. S. Berahas, R. Bollapragada, and E. Wei. On the convergence of nested decentralized gradient methods with multiple consensus and gradient steps. *IEEE Transactions on Signal Processing*, 69:4192–4203, 2021.
- [BJ13] P. Bianchi and J. Jakubowicz. Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization. *IEEE Transactions on Automatic Control*, 58(2):391–405, February 2013.
- [BPC⁺11] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends[®] in Machine learning*, 3(1):1–122, 2011.
- [BT89] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and distributed computation: numerical methods*, volume 23. Prentice hall Englewood Cliffs, NJ, 1989.
- [BT09] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [CABP13] K. Chatzikokolakis, M. E. Andrés, N. E. Bordenabe, and C. Palamidessi. Broadening the scope of differential privacy using metrics. In E. De Cristofaro and M. Wright, editors, *Privacy Enhancing Technologies*, Lecture Notes in Computer Science, pages 82–102, Berlin, Heidelberg, 2013. Springer.
- [CEBM22] E. Cyffers, M. Even, A. Bellet, and L. Massoulié. Muffliato: Peer-to-peer privacy amplification for decentralized optimization and averaging. In *Advances in Neural Information Processing Systems*, 2022.

- [CL11] C.-C. Chang and C.-J. Lin. LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):1–27, May 2011.
- [CSU⁺19] A. Cheu, A. Smith, J. Ullman, D. Zeber, and M. Zhilyaev. Distributed differential privacy via shuffling. In Y. Ishai and V. Rijmen, editors, *Advances in Cryptology – EUROCRYPT 2019*, Lecture Notes in Computer Science, pages 375–403, Cham, 2019. Springer International Publishing.
- [CWH20] X. Chen, S. Z. Wu, and M. Hong. Understanding gradient clipping in private SGD: A geometric perspective. In *Advances in Neural Information Processing Systems*, volume 33, pages 13773–13782. Curran Associates, Inc., 2020.
- [CZC⁺20] S. Cen, H. Zhang, Y. Chi, W. Chen, and T.-Y. Liu. Convergence of distributed stochastic variance reduced methods without sampling extra data. *IEEE Transactions on Signal Processing*, 68:3976–3989, 2020.
- [DBL14] A. Defazio, F. Bach, and S. Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pages 1646–1654, 2014.
- [Den12] L. Deng. The MNIST database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29(6):141–142, 2012.
- [DJW13] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438, October 2013.
- [DKX⁺22] R. Das, S. Kale, Z. Xu, T. Zhang, and S. Sanghavi. Beyond uniform Lipschitz condition in differentially private optimization. *arXiv preprint arXiv:2206.10713*, 2022.
- [DLC⁺22] J. Du, S. Li, X. Chen, S. Chen, and M. Hong. Dynamic differential-privacy preserving sgd. *arXiv preprint arXiv:2111.00173*, 2022.
- [DLS16] P. Di Lorenzo and G. Scutari. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.
- [DMNS06] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In S. Halevi and T. Rabin, editors, *Theory of Cryptography*, Lecture Notes in Computer Science, pages 265–284, Berlin, Heidelberg, 2006. Springer.

- [DR14] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, August 2014.
- [DSZOR15] C. M. De Sa, C. Zhang, K. Olukotun, and C. Ré. Taming the wild: A unified analysis of Hogwild-style algorithms. *Advances in Neural Information Processing Systems*, 28, 2015.
- [Dwo08] C. Dwork. Differential privacy: A survey of results. In M. Agrawal, D. Du, Z. Duan, and A. Li, editors, *Theory and Applications of Models of Computation*, Lecture Notes in Computer Science, pages 1–19, Berlin, Heidelberg, 2008. Springer.
- [FGW21] J. Fan, Y. Guo, and K. Wang. Communication-efficient accurate statistical estimation. *Journal of the American Statistical Association*, 0(0):1–11, August 2021.
- [FKT20] V. Feldman, T. Koren, and K. Talwar. Private stochastic convex optimization: Optimal rates in linear time. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing, STOC 2020*, pages 439–449, New York, NY, USA, June 2020. Association for Computing Machinery.
- [FLLZ18] C. Fang, C. J. Li, Z. Lin, and T. Zhang. SPIDER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems*, pages 687–697, 2018.
- [FSG⁺21] I. Fatkhullin, I. Sokolov, E. Gorbunov, Z. Li, and P. Richtárik. EF21 with bells & whistles: Practical algorithmic extensions of modern error feedback. *arXiv preprint arXiv:2110.03294*, 2021.
- [GBLR21] E. Gorbunov, K. P. Burlachenko, Z. Li, and P. Richtarik. MARINA: Faster non-convex distributed learning with compression. In *Proceedings of the 38th International Conference on Machine Learning*, pages 3788–3798. PMLR, July 2021.
- [GF20] L. V. Gambuzza and M. Frasca. Distributed control of multiconsensus. *IEEE Transactions on Automatic Control*, 66(5):2032–2044, 2020.
- [HAD⁺21] A. Hashemi, A. Acharya, R. Das, H. Vikalo, S. Sanghavi, and I. S. Dhillon. On the benefits of multiple gossip steps in communication-constrained decentralized federated learning. *IEEE Transactions on Parallel and Distributed Systems*, 2021.
- [HDJ⁺20] X. Huang, Y. Ding, Z. L. Jiang, S. Qi, X. Wang, and Q. Liao. DP-FL: A novel differentially private federated learning framework for the unbalanced data. *World Wide Web*, 23(4):2529–2545, July 2020.

- [HSZ⁺22] Y. Huang, Y. Sun, Z. Zhu, C. Yan, and J. Xu. Tackling data heterogeneity: A new unified framework for decentralized sgd with sample-induced topology. In *Proceedings of the 39th International Conference on Machine Learning*, pages 9310–9345. PMLR, June 2022.
- [INS⁺19] R. Iyengar, J. P. Near, D. Song, O. Thakkar, A. Thakurta, and L. Wang. Towards practical differentially private convex optimization. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 299–316, May 2019.
- [IW22] C. Iakovidou and E. Wei. S-NEAR-DGD: A flexible distributed stochastic gradient method for inexact communication. *IEEE Transactions on Automatic Control*, 2022.
- [JZ13] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pages 315–323, 2013.
- [KDG03] D. Kempe, A. Dobra, and J. Gehrke. Gossip-based computation of aggregate information. In *44th Annual IEEE Symposium on Foundations of Computer Science, 2003. Proceedings.*, pages 482–491, October 2003.
- [KFI17] S. Kanai, Y. Fujiwara, and S. Iwamura. Preventing gradient explosions in gated recurrent units. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [KFJ18] S. Khirirat, H. R. Feyzmahdavian, and M. Johansson. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- [KKJ⁺21] D. Kovalev, A. Koloskova, M. Jaggi, P. Richtarik, and S. Stich. A linearly convergent algorithm for decentralized optimization: Sending less bits for free! In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pages 4087–4095. PMLR, 13–15 Apr 2021.
- [KLL16] J. Kim, J. K. Lee, and K. M. Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1646–1654, 2016.
- [KMR15] J. Konečný, B. McMahan, and D. Ramage. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*, 2015.
- [KMY⁺16] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.

- [KRS19] S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *International Conference on Machine Learning*, pages 3252–3261. PMLR, 2019.
- [KS19] A. Koloskova, S. Stich, and M. Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3478–3487. PMLR, 09–15 Jun 2019.
- [Lam01] L. Lamberg. Confidentiality and privacy of electronic medical records. *JAMA*, 285(24):3075–3076, June 2001.
- [LBZR21] Z. Li, H. Bao, X. Zhang, and P. Richtárik. Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International Conference on Machine Learning*, pages 6286–6295. PMLR, 2021.
- [LCCC20] B. Li, S. Cen, Y. Chen, and Y. Chi. Communication-efficient distributed optimization in networks with gradient tracking and variance reduction. *Journal of Machine Learning Research*, 21(180):1–51, 2020.
- [LDS21] Y. Lu and C. De Sa. Optimal complexity in decentralized training. In *International Conference on Machine Learning*, pages 7111–7123. PMLR, 2021.
- [LFYL20] H. Li, C. Fang, W. Yin, and Z. Lin. Decentralized accelerated gradient methods with increasing penalty parameters. *IEEE Transactions on Signal Processing*, 68:4855–4870, 2020.
- [Li19] Z. Li. SSRGD: Simple stochastic recursive gradient descent for escaping saddle points. In *Advances in Neural Information Processing Systems*, pages 1523–1533, 2019.
- [LJB⁺95] Y. LeCun, L. D. Jackel, L. Bottou, C. Cortes, J. S. Denker, H. Drucker, I. Guyon, U. A. Muller, E. Sackinger, and P. Simard. Learning algorithms for classification: A comparison on handwritten digit recognition. *Neural networks: the statistical mechanics perspective*, 261(276):2, 1995.
- [LJC17] L. Lei, C. Ju, J. Chen, and M. I. Jordan. Non-convex finite-sum optimization via SCSG methods. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [LKQR20] Z. Li, D. Kovalev, X. Qian, and P. Richtarik. Acceleration for compressed gradient descent in distributed and federated optimization. In *Proceedings of the 37th International Conference on Machine Learning*, pages 5895–5904. PMLR, November 2020.

- [LL18] Z. Li and J. Li. A simple proximal stochastic gradient method for nonsmooth nonconvex optimization. In *Advances in Neural Information Processing Systems*, pages 5569–5579, 2018.
- [LLC22] B. Li, Z. Li, and Y. Chi. DESTRESS: Computation-optimal and communication-efficient decentralized nonconvex finite-sum optimization. *SIAM Journal on Mathematics of Data Science*, 4(3):1031–1051, September 2022.
- [LLHP22] Y. Liao, Z. Li, K. Huang, and S. Pu. A compressed gradient tracking method for decentralized optimization with linear convergence. *IEEE Transactions on Automatic Control*, 67(10):5622–5629, 2022.
- [LLP23] Y. Liao, Z. Li, and S. Pu. A linearly convergent robust compressed Push-Pull method for decentralized optimization. *arXiv preprint arXiv:2303.07091*, 2023.
- [LR20] Z. Li and P. Richtárik. A unified analysis of stochastic gradient methods for nonconvex federated optimization. *arXiv preprint arXiv:2006.07013*, 2020.
- [LR21a] Z. Li and P. Richtarik. CANITA: Faster rates for distributed convex optimization with communication compression. In *Advances in Neural Information Processing Systems*, volume 34, pages 13770–13781. Curran Associates, Inc., 2021.
- [LR21b] Z. Li and P. Richtárik. ZeroSARAH: Efficient nonconvex finite-sum optimization with zero full gradient computation. *arXiv preprint arXiv:2103.01447*, 2021.
- [LSY19] Z. Li, W. Shi, and M. Yan. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates. *IEEE Transactions on Signal Processing*, 67(17):4494–4506, 2019.
- [LY22] L. Luo and H. Ye. An optimal stochastic algorithm for decentralized nonconvex finite-sum optimization. *arXiv preprint arXiv:2210.13931*, 2022.
- [LZLC22] Z. Li, H. Zhao, B. Li, and Y. Chi. SoteriaFL: A unified framework for private federated learning with communication compression. In *Advances in Neural Information Processing Systems*, volume 35, pages 4285–4300. Curran Associates, Inc., 2022.
- [LZZ⁺17] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340, 2017.

- [MGTR19] K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- [MMR⁺17] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pages 1273–1282, 2017.
- [MS23] T. Murata and T. Suzuki. DIFF2: Differential private optimization via gradient differences for nonconvex distributed learning. *arXiv preprint arXiv:2302.03884*, 2023.
- [NBD22] M. Noble, A. Bellet, and A. Dieuleveut. Differentially private federated learning on heterogeneous data. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, pages 10110–10145. PMLR, May 2022.
- [NLST17] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pages 2613–2621, 2017.
- [NO09] A. Nedić and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [NOP10] A. Nedić, A. Ozdaglar, and P. A. Parrilo. Constrained consensus and optimization in multi-agent networks. *IEEE Transactions on Automatic Control*, 55(4):922–938, 2010.
- [NOR18] A. Nedić, A. Olshevsky, and M. G. Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976, 2018.
- [NOS17] A. Nedić, A. Olshevsky, and W. Shi. Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633, 2017.
- [NvP⁺22] L. M. Nguyen, M. van Dijk, D. T. Phan, P. H. Nguyen, T.-W. Weng, and J. R. Kalagnanam. Finite-sum smooth optimization with SARAH. *Computational Optimization and Applications*, May 2022.
- [PLW20] T. Pan, J. Liu, and J. Wang. D-SPIDER-SFO: A decentralized optimization algorithm with faster convergence rate for nonconvex problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 1619–1626, 2020.

- [PMB13] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1310–1318. PMLR, May 2013.
- [Pol63] B. T. Polyak. Gradient methods for the minimisation of functionals. *USSR Computational Mathematics and Mathematical Physics*, 3(4):864–878, 1963.
- [QL18] G. Qu and N. Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 5(3):1245–1260, 2018.
- [RHS⁺16] S. J. Reddi, A. Hefny, S. Sra, B. Póczós, and A. Smola. Stochastic variance reduction for nonconvex optimization. In *International conference on machine learning*, pages 314–323, 2016.
- [RKR⁺16] S. J. Reddi, J. Konečný, P. Richtárik, B. Póczós, and A. Smola. AIDE: fast and communication efficient distributed optimization. *arXiv preprint arXiv:1608.06879*, 2016.
- [RLD⁺23] A. Reisizadeh, H. Li, S. Das, and A. Jadbabaie. Variance-reduced clipping for non-convex optimization. *arXiv preprint arXiv:2303.00883*, 2023.
- [RSF21] P. Richtárik, I. Sokolov, and I. Fatkhullin. EF21: A new, simpler, theoretically better, and practically faster error feedback. In *Advances in Neural Information Processing Systems*, volume 34, pages 4384–4396. Curran Associates, Inc., 2021.
- [RSPS16] S. J. Reddi, S. Sra, B. Póczós, and A. Smola. Fast incremental method for smooth nonconvex optimization. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 1971–1977. IEEE, 2016.
- [SBB⁺17] K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In *International Conference on Machine Learning*, pages 3027–3036, 2017.
- [SBB⁺18] K. Scaman, F. Bach, S. Bubeck, L. Massoulié, and Y. T. Lee. Optimal algorithms for non-smooth distributed optimization in networks. In *Advances in Neural Information Processing Systems*, pages 2740–2749, 2018.
- [SCJ18] S. U. Stich, J.-B. Cordonnier, and M. Jaggi. Sparsified SGD with memory. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

- [SDGD21] N. Singh, D. Data, J. George, and S. Diggavi. SQuARM-SGD: Communication-Efficient Momentum SGD for Decentralized Optimization. *IEEE Journal on Selected Areas in Information Theory*, 2(3):954–969, September 2021.
- [SFD⁺14] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *Interspeech 2014*, pages 1058–1062. ISCA, September 2014.
- [Sha07] D. Shah. Gossip algorithms. *Foundations and Trends[®] in Networking*, 3(1):1–125, 2007.
- [SK20] S. U. Stich and S. P. Karimireddy. The error-feedback framework: Better rates for sgd with delayed gradients and compressed updates. *J. Mach. Learn. Res.*, 21(1), jan 2020.
- [SLH20] H. Sun, S. Lu, and M. Hong. Improving the sample and communication complexity for decentralized non-convex optimization: Joint gradient estimation and tracking. In *International Conference on Machine Learning*, pages 9217–9228. PMLR, 2020.
- [SLWY15a] W. Shi, Q. Ling, G. Wu, and W. Yin. EXTRA: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [SLWY15b] W. Shi, Q. Ling, G. Wu, and W. Yin. A proximal gradient algorithm for decentralized composite optimization. *IEEE Transactions on Signal Processing*, 63(22):6013–6023, 2015.
- [SS19] G. Scutari and Y. Sun. Distributed nonconvex constrained optimization over time-varying digraphs. *Mathematical Programming*, 176(1-2):497–544, 2019.
- [SSZ14] O. Shamir, N. Srebro, and T. Zhang. Communication-efficient distributed optimization using an approximate Newton-type method. In *International conference on machine learning*, pages 1000–1008, 2014.
- [TGZ⁺18] H. Tang, S. Gan, C. Zhang, T. Zhang, and J. Liu. Communication compression for decentralized training. *Advances in Neural Information Processing Systems*, 31:7652–7662, 2018.
- [TLQ⁺19] H. Tang, X. Lian, S. Qiu, L. Yuan, C. Zhang, T. Zhang, and J. Liu. Deepsqueeze: Decentralization meets error-compensated compression. *arXiv preprint arXiv:1907.07346*, 2019.
- [TLY⁺18] H. Tang, X. Lian, M. Yan, C. Zhang, and J. Liu. D^2 : Decentralized training over decentralized data. In *International Conference on Machine Learning*, pages 4848–4856. PMLR, 2018.

- [TMHP20] H. Taheri, A. Mokhtari, H. Hassani, and R. Pedarsani. Quantized decentralized stochastic learning over directed graphs. In *Proceedings of the 37th International Conference on Machine Learning*, pages 9324–9333. PMLR, November 2020.
- [WJ21] J. Wang and G. Joshi. Cooperative SGD: A unified framework for the design and analysis of local-update SGD algorithms. *The Journal of Machine Learning Research*, 22(1):213:9709–213:9758, January 2021.
- [WJEG19] L. Wang, B. Jayaraman, D. Evans, and Q. Gu. Efficient privacy-preserving stochastic nonconvex optimization. *arXiv preprint arXiv:1910.13659*, 2019.
- [WJZ⁺19] Z. Wang, K. Ji, Y. Zhou, Y. Liang, and V. Tarokh. SpiderBoost and momentum: Faster variance reduction algorithms. In *Advances in Neural Information Processing Systems*, pages 2406–2416, 2019.
- [WXDX20] D. Wang, H. Xiao, S. Devadas, and J. Xu. On differentially private stochastic convex optimization with heavy-tailed data. In *Proceedings of the 37th International Conference on Machine Learning*, pages 10081–10091. PMLR, November 2020.
- [WYWH18] H.-T. Wai, Z. Yang, P. Z. Wang, and M. Hong. Multi-agent reinforcement learning via double averaging primal-dual optimization. In *Advances in Neural Information Processing Systems*, pages 9649–9660, 2018.
- [WYX17] D. Wang, M. Ye, and J. Xu. Differentially private empirical risk minimization revisited: Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017.
- [XB04] L. Xiao and S. Boyd. Fast linear iterations for distributed averaging. *Systems & Control Letters*, 53(1):65–78, September 2004.
- [XKK22a] R. Xin, U. A. Khan, and S. Kar. Fast decentralized nonconvex finite-sum optimization with recursive variance reduction. *SIAM Journal on Optimization*, 32(1):1–28, March 2022.
- [XKK22b] R. Xin, U. A. Khan, and S. Kar. A Fast Randomized Incremental Gradient Method for Decentralized Nonconvex Optimization. *IEEE Transactions on Automatic Control*, 67(10):5150–5165, October 2022.
- [XSKK19] R. Xin, A. K. Sahu, U. A. Khan, and S. Kar. Distributed stochastic optimization with gradient tracking over strongly-connected networks. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 8353–8358, 2019.

- [XXK17] C. Xi, R. Xin, and U. A. Khan. ADD-OPT: Accelerated distributed directed optimization. *IEEE Transactions on Automatic Control*, 63(5):1329–1339, 2017.
- [XYD19] H. Xiao, Y. Ye, and S. Devadas. Local differential privacy in decentralized optimization. *arXiv preprint arXiv:1902.06101*, 2019.
- [XZ14] L. Xiao and T. Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- [YCC⁺23] Y. Yan, J. Chen, P.-Y. Chen, X. Cui, S. Lu, and Y. Xu. Compressed decentralized proximal stochastic gradient method for nonconvex composite problems with heterogeneous data. *arXiv preprint arXiv:2302.14252*, 2023.
- [YGG17] Y. You, I. Gitman, and B. Ginsburg. Scaling SGD batch size to 32k for ImageNet training. Technical Report UCB/EECS-2017-156, EECS Department, University of California, Berkeley, Sep 2017.
- [YYZS18] K. Yuan, B. Ying, X. Zhao, and A. H. Sayed. Exact diffusion for distributed optimization and learning – part I: Algorithm development. *IEEE Transactions on Signal Processing*, 67(3):708–723, 2018.
- [YZLZ20] H. Ye, Z. Zhou, L. Luo, and T. Zhang. Decentralized accelerated proximal gradient descent. *Advances in Neural Information Processing Systems*, 33:18308–18317, 2020.
- [ZBLR21] H. Zhao, K. Burlachenko, Z. Li, and P. Richtárik. Faster rates for compressed federated learning with client-variance reduction. *arXiv preprint arXiv:2112.13097*, 2021.
- [ZCH⁺22] X. Zhang, X. Chen, M. Hong, S. Wu, and J. Yi. Understanding clipping for federated learning: Convergence and client-level differential privacy. In *Proceedings of the 39th International Conference on Machine Learning*, pages 26048–26067. PMLR, June 2022.
- [ZHSJ20] J. Zhang, T. He, S. Sra, and A. Jadbabaie. Why gradient clipping accelerates training: A theoretical justification for adaptivity. In *International Conference on Learning Representations*, March 2020.
- [ZJFW20] B. Zhang, J. Jin, C. Fang, and L. Wang. Improved analysis of clipping algorithms for non-convex optimization. In *Advances in Neural Information Processing Systems*, volume 33, pages 15511–15521. Curran Associates, Inc., 2020.

- [ZKL18] X. Zhang, M. M. Khalili, and M. Liu. Improving the privacy and accuracy of ADMM-based distributed algorithms. In *Proceedings of the 35th International Conference on Machine Learning*, pages 5796–5805. PMLR, July 2018.
- [ZKL20] X. Zhang, M. M. Khalili, and M. Liu. Recycled ADMM: Improving the privacy and accuracy of distributed algorithms. *IEEE Transactions on Information Forensics and Security*, 15:1723–1734, 2020.
- [ZKV⁺20] J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. J. Reddi, S. Kumar, and S. Sra. Why are adaptive methods good for attention models? *arXiv preprint arXiv:1912.03194*, October 2020.
- [ZLL⁺22] H. Zhao, B. Li, Z. Li, P. Richtárik, and Y. Chi. Beer: Fast $O(1/T)$ rate for decentralized nonconvex optimization with communication compression. *Advances in Neural Information Processing Systems*, 35, 2022.
- [ZM10] M. Zhu and S. Martínez. Discrete-time dynamic average consensus. *Automatica*, 46(2):322–329, 2010.
- [ZP23] Y. Zhou and S. Pu. Private and accurate decentralized optimization via encrypted and structured functional perturbation. *IEEE Control Systems Letters*, 7:1339–1344, 2023.
- [ZXC18] D. Zhou, P. Xu, and Q. Gu. Stochastic nested variance reduction for nonconvex optimization. *Advances in Neural Information Processing Systems*, 31:3921–3932, 2018.
- [ZY19] J. Zhang and K. You. Decentralized stochastic gradient tracking for non-convex empirical risk minimization. *arXiv preprint arXiv:1909.02712*, 2019.
- [ZZY⁺21] Y. Zhao, J. Zhao, M. Yang, T. Wang, N. Wang, L. Lyu, D. Niyato, and K.-Y. Lam. Local differential privacy-based federated learning for Internet of Things. *IEEE Internet of Things Journal*, 8(11):8836–8853, June 2021.