

# GT-STORM: Taming Sample, Communication, and Memory Complexities in Decentralized Non-Convex Learning

Xin Zhang<sup>1</sup>, Jia Liu<sup>2</sup>, Zhengyuan Zhu<sup>1</sup>, and Elizabeth S. Bentley<sup>3</sup>

<sup>1</sup>Department of Statistics, Iowa State University

<sup>2</sup>Department of Electrical and Computer Engineering, The Ohio State University

<sup>3</sup>Information Directorate, Air Force Research Laboratory

## ABSTRACT

Decentralized nonconvex optimization has received increasing attention in recent years in machine learning due to its advantages in system robustness, data privacy, and implementation simplicity. However, three fundamental challenges in designing decentralized optimization algorithms are how to reduce their sample, communication, and memory complexities. In this paper, we propose a gradient-tracking-based stochastic recursive momentum (GT-STORM) algorithm for efficiently solving nonconvex optimization problems. We show that to reach an  $\epsilon^2$ -stationary solution, the total number of sample evaluations of our algorithm is  $\tilde{O}(m^{1/2}\epsilon^{-3})$  and the number of communication rounds is  $\tilde{O}(m^{-1/2}\epsilon^{-3})$ , which improve the  $O(\epsilon^{-4})$  costs of sample evaluations and communications for the existing decentralized stochastic gradient algorithms. We conduct extensive experiments with a variety of learning models, including non-convex logistical regression and convolutional neural networks, to verify our theoretical findings. Collectively, our results contribute to the state of the art of theories and algorithms for decentralized network optimization.

## CCS CONCEPTS

• **Computing methodologies** → **Distributed algorithms**; *Machine learning*; • **Networks** → *Network performance analysis*.

## KEYWORDS

Network Consensus Optimization, Stochastic Variance Reduction, Gradient Tracking

### ACM Reference Format:

Xin Zhang<sup>1</sup>, Jia Liu<sup>2</sup>, Zhengyuan Zhu<sup>1</sup>, and Elizabeth S. Bentley<sup>3</sup>. 2020. GT-STORM: Taming Sample, Communication, and Memory Complexities in Decentralized Non-Convex Learning. In *MobiHoc '20: ACM Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*, June 30 – July 03, 2020, Shanghai, China. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

In recent years, machine learning has witnessed enormous success in many areas, including image processing, natural language processing, online recommender systems, just to name a few. From

a mathematical perspective, training machine learning models amounts to solving an optimization problem. However, with the rapidly increasing dataset sizes and the high dimensionality and the non-convex hardness of the training problem (e.g., due to the use of deep neural networks), training large-scale machine learning models by a single centralized machine has become inefficient and unscalable. To address the efficiency and scalability challenges, an effective approach is to leverage *decentralized* computational resources in a computing network, which could follow a parameter server (PS)-worker architecture [1–3] or fully decentralized peer-to-peer network structure [4, 5]. Also, thanks to the robustness to single-point-of-failure, data privacy, and implementation simplicity, decentralized learning over computing networks has attracted increasing interest recently, and has been applied in various science and engineering areas (including dictionary learning [6], multi-agent systems [7, 8], multi-task learning [9, 10], information retrieval [11], energy allocation [12], etc.).

In the fast growing literature of decentralized learning over networks, a classical approach is the so-called network consensus optimization, which traces its roots to the seminal work by Tsitsiklis in 1984 [13]. Recently, network consensus optimization has gained a lot of renewed interest owing to the elegant decentralized subgradient descent method (DSGD) proposed by Nedic and Ozdaglar [4], which has been applied in decentralized learning due to its simple algorithmic structure and good convergence performance. In network-consensus-based decentralized learning, a set of geographically distributed computing nodes collaborate to train a common learning model. Each node holds a local dataset that may be too large to be sent to a centralized location due to network communication limits, or cannot be shared due to privacy/security risks. A distinctive feature of network-consensus-based decentralized learning is that there is a lack of a dedicated PS. As a result, each node has to exchange information with its local neighbors to reach a consensus on a global optimal learning model.

Despite its growing significance in practice, the design of high-performance network-consensus-based decentralized learning faces three fundamental *conflicting* complexities, namely *sample, communication, and memory complexities*. First, due to the high dimensionality of most deep learning models, it is impossible to leverage beyond first-order (stochastic) gradient information to compute the update direction in each iteration. The variability of a stochastic gradient is strongly influenced by the number of training samples in its mini-batch. However, the more training samples in a mini-batch, the higher computational cost of the stochastic gradient. Second, by using fewer training samples in each iteration to trade for a lower computational cost, the resulting stochastic gradient

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

*MobiHoc '20, June 30 – July 03, 2020, Shanghai, China*

© 2020 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-XXXX-X/18/06.

<https://doi.org/10.1145/1122445.1122456>

unavoidably has a larger variance, which further leads to more iterations (hence communication rounds) to reach a certain training accuracy (i.e., slower convergence). The low communication efficiency is particularly problematic in many wireless edge networks, where the communication links could be low-speed and highly unreliable. Lastly, in many mobile edge-computing environments, the mobile devices could be severely limited by hardware resources (e.g., CPU/GPU, memory) and they cannot afford reserving a large memory space to run a very sophisticated decentralized learning algorithm that has too many intermediate variables.

Due to the above fundamental trade-off between sample, communication, and computing resource costs, the notions of sample, communication, and memory complexities (to be formally defined in Section 2) become three of the most important measures in assessing the performances of decentralized learning algorithms. However, in the literature, most existing works have achieved low complexities in some of these measures, but not all (see Section 2 for in-depth discussions). The limitations of these existing works motivate to ask the following question: *Could we design a decentralized learning algorithm that strikes a good balance between sample complexity and communication complexity?* In this paper, we answer the above question positively by proposing a new GT-STORM algorithm (gradient-tracking-based stochastic recursive momentum) that achieves low sample, communication, and memory complexities. Our main results and contributions are summarized as follows:

- Unlike existing approaches, our proposed GT-STORM algorithm adopts a new estimator, which is updated with a consensus mixing of the neighboring estimators of the last iteration, which helps improve the global gradient estimation. Our method achieves the nice features of previous works [14–17] while avoiding their pitfalls. To some extent, our GT-STORM algorithm can be viewed as an indirect way of integrating the stochastic gradient method, variance reduction method, and gradient tracking method.
- We provide a detailed convergence analysis and complexity analysis. Under some mild assumptions and parameter conditions, our algorithm enjoys an  $\tilde{O}(T^{-2/3})$  convergence rate. Note that this rate is much faster than the rate of  $O(T^{1/2})$  for the classic decentralized stochastic algorithms, e.g., DSGD [18], PSGD [5] and GNSD [17]. Also, we show that to reach an  $\epsilon^2$ -stationary solution, the total number of sample evaluations of our algorithm is  $\tilde{O}(m^{1/2}\epsilon^{-3})$  and the communication round is  $\tilde{O}(m^{-1/2}\epsilon^{-3})$ .
- We conduct extensive experiments to examine the performance of our algorithm, including both a non-convex logistic regression model on the LibSVM datasets and convolutional neural network models on MNIST and CIFAR-10 datasets. Our experiments show that our algorithm outperforms two state-of-the-art decentralized learning algorithms [5, 17]. These experiments corroborate our theoretical results.

The rest of the paper is organized as follows. In Section 2, we first provide the preliminaries of network consensus optimization and discuss related works with a focus on sample, communication, and memory complexities. In Section 3, we present our proposed GT-STORM algorithm, as well as its communication, sample, and memory complexity analysis. We provide numerical results in Section 4 to verify the theoretical results of our GT-STORM algorithm. Lastly in Section 5, we provide concluding remarks.

## 2 PRELIMINARIES AND RELATED WORK

To facilitate our technical discussions, in Section 2.1, we first provide an overview on network consensus optimization and formally define the notions of sample, communication, and memory complexities of decentralized optimization algorithms for network consensus optimization. Then, in Section 2.2, we first review centralized stochastic first-order optimization algorithms for solving non-convex learning problems from a historical perspective and with a focus on sample, communication, and memory complexities. Here, we introduce several acceleration techniques that motivate our GT-STORM algorithmic design. Lastly, we review the recent developments of optimization algorithms for decentralized learning and compare them with our work.

### 2.1 Network Consensus Optimization

As mentioned in Section 1, in decentralized learning, there are a set of geographically distributed computing nodes forming a network. In this paper, we represent such a network by an undirected connected network  $\mathcal{G} = (\mathcal{N}, \mathcal{L})$ , where  $\mathcal{N}$  and  $\mathcal{L}$  are the sets of nodes and edges, respectively, with  $|\mathcal{N}| = m$ . Each node can communicate with their neighbors via the edges in  $\mathcal{L}$ . The goal of decentralized learning is to use the nodes to *distributively* and *collaboratively* solve a network-wide optimization problem as follows:

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^p} \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}), \quad (1)$$

where each local objective function  $f_i(\mathbf{x}) \triangleq \mathbb{E}_{\zeta \sim \mathcal{D}_i} f_i(\mathbf{x}; \zeta)$  is only observable to node  $i$  and not necessarily convex. Here,  $\mathcal{D}_i$  represents the distribution of the dataset at node  $i$ , and  $f_i(\mathbf{x}; \zeta)$  represents a loss function that evaluates the discrepancy between the learning model's output and the ground truth of a training sample  $\zeta$ . To solve Problem (1) in a decentralized fashion, a common approach is to rewrite Problem (1) in the following equivalent form:

$$\begin{aligned} & \text{Minimize} && \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}_i) \\ & \text{subject to} && \mathbf{x}_i = \mathbf{x}_j, \quad \forall (i, j) \in \mathcal{L}, \end{aligned} \quad (2)$$

where  $\mathbf{x} \triangleq [\mathbf{x}_1^\top, \dots, \mathbf{x}_m^\top]^\top$  and  $\mathbf{x}_i$  is an introduced local copy at node  $i$ . In Problem (2), the constraints ensure that the local copies at all nodes are equal to each other, hence the term “consensus.” Thus, Problems (1) and (2) share the same solutions. The main goal of network consensus optimization is to design an algorithm to attain an  $\epsilon^2$ -stationary point  $\mathbf{x}$  defined as follows:

$$\underbrace{\left\| \frac{1}{m} \sum_{i=1}^m \nabla f_i(\bar{\mathbf{x}}) \right\|^2}_{\text{Global gradient magnitude}} + \underbrace{\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2}_{\text{Consensus error}} \leq \epsilon^2, \quad (3)$$

where  $\bar{\mathbf{x}} \triangleq \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$  denotes the global average across all nodes. Different from the traditional  $\epsilon^2$ -stationary point in centralized optimization problems, the metric in Eq. (3) has two terms: the first term is the gradient magnitude for the (non-convex) global objective and the second term is the average consensus error of all local copies. To date, many decentralized algorithms have been developed to compute the  $\epsilon^2$ -stationary point (see Section 2.2). However, most

of these algorithms suffer limitations in sample, communication, and memory complexities. In what follows, we formally state the definitions of sample, communication, and memory complexities used in the literature (see, e.g., [19]):

**DEFINITION 1 (SAMPLE COMPLEXITY).** *The sample complexity is defined as the total number of the incremental first-order oracle (IFO) calls required across all the nodes to find an  $\epsilon^2$ -stationary point defined in Eq. (3), where one IFO call evaluates a pair of  $(f_i(\mathbf{x}; \zeta), \nabla f_i(\mathbf{x}; \zeta))$  on a sample  $\zeta \sim \mathcal{D}_i$  and parameter  $\mathbf{x} \in \mathbb{R}^p$  at node  $i$ .*

**DEFINITION 2 (COMMUNICATION COMPLEXITY).** *The communication complexity is defined as the total rounds of communications required to find an  $\epsilon^2$ -stationary point defined in Eq. (3), where each node can send and receive a  $p$ -dimensional vector with its neighboring nodes in one communication round.*

**DEFINITION 3 (MEMORY COMPLEXITY).** *The memory complexity is defined as total dimensionality of all intermediate variables in the algorithm run by a node to find an  $\epsilon^2$ -stationary point in Eq. (3).*

To make sense of these three complexity metrics into perspective, consider the standard centralized gradient descent (GD) method as an example. Note that the GD algorithm has an  $O(1/T)$  convergence rate for non-convex optimization, which suggests  $O(\epsilon^{-2})$  communication complexity. Also, it takes a full gradient evaluation in each iteration, i.e.,  $O(n)$  per-iteration sample complexity, where  $n$  is the total number of samples. This implies  $O(n\epsilon^{-2})$  sample complexity to converge to an  $\epsilon^2$ -stationary point. Hence, the sample complexity of GD is high if the dataset size  $n$  is large.

In contrast, consider the classical stochastic gradient descent (SGD) algorithm that is widely used in machine learning. The basic idea of SGD is to lower the gradient evaluation cost by using only a mini-batch of samples in each iteration. However, due to the sample randomness in mini-batches, the convergence rate of SGD for non-convex optimization is reduced to  $O(1/\sqrt{T})$  [20–22]. Thus, to reach an  $\epsilon^2$ -stationary point  $\mathbf{x}$  with  $\|\nabla f(\mathbf{x})\|^2 \leq \epsilon^2$ , SGD has  $O(\epsilon^{-4})$  sample complexity, which could be either higher or lower than the  $O(n\epsilon^{-2})$  sample complexity of the GD method, depending on the relationship between  $n$  and  $\epsilon$ . Also, for  $p$ -dimensional problems, both GD and SGD have memory complexity  $p$ , since they only need a  $p$ -dimensional vector to store (stochastic) gradients.

## 2.2 Related Work

### 1) Centralized First-Order Methods with Low Complexities:

Now, we review several state-of-the-art low-complexity centralized stochastic first-order methods that are related to our GT-STORM algorithm. To reduce the overall sample and communication complexities of the standard GD and SGD algorithms, a natural approach is variance reduction. Earlier works following this approach include SVRG [23, 24], SAGA [25] and SCSG [26]. These algorithms has an overall sample complexity of  $O(n + n^{2/3}\epsilon^{-2})$ . A more recent variance reduction method is the stochastic path-integrated differential estimator (SPIDER) [27], which is based on the SARAH gradient estimator developed by Nguyen *et al.* [28]. SPIDER further lowers the sample complexity to  $O(n + \sqrt{n}\epsilon^{-2})$ , which attains the  $\Omega(\sqrt{n}\epsilon^{-2})$  theoretical lower bound for finding an  $\epsilon^2$ -stationary point for  $n = O(\epsilon^{-4})$ . More recently, to improve the small step-size  $O(\epsilon L^{-1})$  in SPIDER, a variant called SpiderBoost was proposed

in [29], which allows a larger constant step-size  $O(L^{-1})$  while keeping the same  $O(n + \sqrt{n}\epsilon^{-2})$  sample complexity. It should be noted, however, that the significantly improved sample complexity of SPIDER/SpiderBoost is due to a restrictive assumption that a universal Lipschitz smoothness constant exists for all local objectives  $f(\cdot; \zeta_i) \forall i$ . This means that the objectives are “similar” and there are no “outliers” in the training samples. Meanwhile, to obtain the optimal communication complexity, SpiderBoost require a (nearly) full gradient every  $\sqrt{n}$  iterations and a mini-batch of stochastic gradient evaluation with batch size  $\sqrt{n}$  in each iteration.

To overcome the above limitations, a hybrid stochastic gradient descent (Hybrid-SGD) method is recently proposed in [14], where a convex combination of the SARAH estimator [28] and an unbiased stochastic gradient is used as the gradient estimator. The Hybrid-SGD method relaxes the universal Lipschitz constant assumption in SpiderBoost to an average Lipschitz smoothness assumption. Moreover, it only requires two samples to evaluate the gradient per iteration. As a result, Hybrid-SGD has a  $O(\epsilon^{-3})$  sample complexity that is *independent* of dataset size. Although Hybrid-SGD is for centralized optimization, the interesting ideas therein motivate our GT-STORM approach for *decentralized learning* following a similar token. Interestingly, we show that in decentralized settings, our GT-STORM method can further improve the gradient evaluation to only *one sample* per iteration, while not degrading the communication complexity order. Lastly, we remark that all algorithms above have memory complexity at least  $2p$  for  $p$ -dimensional problems. In contrast, GT-STORM enjoys a  $p$  memory complexity.

**2) Decentralized Optimization Algorithms** In the literature, many decentralized learning optimization algorithms have been proposed to solve Problem (1), e.g., first-order methods [4, 16, 30, 31], prime-dual methods [32, 33], Newton-type methods [34, 35] (see in [36, 37] for comprehensive surveys). In this paper, we consider decentralized first-order methods for the non-convex network consensus optimization in (2). In the literature, the convergence rate of the well-known decentralized gradient descent (DGD) algorithm [4] was studied in [38], which showed that DGD with a constant step-size converges with an  $O(1/T)$  rate to a step-size-dependent error ball around a stationary point. Later, a gradient tracking (GT) method was proposed in [16] to find an  $\epsilon^2$ -stationary point with an  $O(1/T)$  convergence rate under constant step-sizes. However, these methods require a full gradient evaluation per iteration, which yields  $O(n\epsilon^{-2})$  sample complexity. To reduce the per-iteration sample complexity, stochastic gradients are adopted in the decentralized optimization, e.g., DSGD [18], PSGD [5], GNSD [17]. Due to the randomness in stochastic gradients, the convergence rate is reduced to  $O(1/\sqrt{T})$ . Thus, the sample and communication complexities of these stochastic methods are  $O(\epsilon^{-4})$  and  $O(m^{-1}\epsilon^{-4})$ , two orders of magnitude higher than their deterministic counterparts. To overcome the limitations in stochastic methods, a natural idea is to use variance reduction techniques similar to those for centralized optimization to reduce the sample and communication complexities for the non-convex network consensus optimization. So far, existing works on the decentralized stochastic variance reduction methods include DSA [39], diffusion-AVRG [40] and GT-SAGA [41] etc., all of which focus on convex problems. To our knowledge, the decentralized gradient estimation and tracking (D-GET) algorithm

in [19] is the only work for non-convex optimization. D-GET integrates the decentralized gradient tracking [17] and the SpiderBoost gradient estimator [29] to obtain  $O(mn + m\sqrt{n}\epsilon^{-2})$  dataset-size-dependent sample complexity and  $O(\epsilon^{-2})$  communication complexity. Recall that the sample and communication complexities of GT-STORM are  $O(m^{1/2}\epsilon^{-3})$  and  $O(m^{-1/2}\epsilon^{-3})$ , respectively. Thus, if dataset size  $n = \Omega(\epsilon^{-2})$ , D-GET has a higher sample complexity than GT-STORM. As an example, when  $\epsilon = 10^{-2}$ ,  $n$  is on the order of  $10^4$ , which is common in modern machine learning datasets. Also, the memory complexity of D-GET is  $2p$  as opposed to the  $p$  memory complexity of GT-STORM. This implies a huge saving with GT-STORM if  $p$  is large, e.g.,  $p \approx 10^6$  in many deep learning models.

### 3 A GRADIENT-TRACKING STOCHASTIC RECURSIVE MOMENTUM ALGORITHM

In this section, we introduce our gradient-tracking-based stochastic recursive momentum (GT-STORM) algorithm for solving Problem (2) in Section 3.1. Then, we will state the main theoretical results and their proofs in Sections 3.2 and 3.3, respectively.

#### 3.1 The GT-STORM Algorithm

In the literature, a standard starting point to solve Problem (2) is to reformulate the problem as [4]:

$$\begin{aligned} & \text{Minimize} \quad \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}_i) \\ & \text{subject to} \quad (\mathbf{W} \otimes \mathbf{I}_p) \mathbf{x} = \mathbf{x}, \end{aligned} \quad (4)$$

where  $\mathbf{I}_p$  denotes the  $p$ -dimensional identity matrix, the operator  $\otimes$  denotes the Kronecker product, and  $\mathbf{W} \in \mathbb{R}^{m \times m}$  is often referred to as the consensus matrix. We let  $[\mathbf{W}]_{ij}$  represent the element in the  $i$ -th row and the  $j$ -th column in  $\mathbf{W}$ . For Problems (4) and (2) to be equivalent,  $\mathbf{W}$  should satisfy the following properties:

- (a) *Doubly Stochastic*:  $\sum_{i=1}^m [\mathbf{W}]_{ij} = \sum_{j=1}^m [\mathbf{W}]_{ij} = 1$ .
- (b) *Symmetric*:  $[\mathbf{W}]_{ij} = [\mathbf{W}]_{ji}, \forall i, j \in \mathcal{N}$ .
- (c) *Network-Defined Sparsity Pattern*:  $[\mathbf{W}]_{ij} > 0$  if  $(i, j) \in \mathcal{L}$ ; otherwise  $[\mathbf{W}]_{ij} = 0, \forall i, j \in \mathcal{N}$ .

The above properties imply that the eigenvalues of  $\mathbf{W}$  are real and can be sorted as  $-1 < \lambda_m \leq \dots \leq \lambda_2 < \lambda_1 = 1$ . We define the second-largest eigenvalue in magnitude of  $\mathbf{W}$  as  $\lambda \triangleq \max\{|\lambda_2|, |\lambda_m|\}$  for the further notation convenience. It can be seen later that  $\lambda$  plays an important role in the step-size selection and the algorithm's convergence rate.

As mentioned in Section 2.1, our GT-STORM algorithm is inspired by the GT method [16, 42] for reducing consensus error and the recursive variance reduction (VR) methods [27, 29] developed for centralized optimization. Specifically, in the centralized GT method, an estimator  $\mathbf{y}$  is introduced to track the global gradient:

$$\mathbf{y}_t = \mathbf{W}\mathbf{y}_{t-1} + \mathbf{g}_t - \mathbf{g}_{t-1}, \quad (5)$$

where  $\mathbf{g}_t$  is the gradient estimation in the  $t$ th iteration. Meanwhile, to reduce the stochastic error, a gradient estimator  $\mathbf{v}$  in VR methods is updated recursively based on a *double-loop* structure as follows:

$$\mathbf{v}_t = \mathbf{v}_{t-1} + \nabla f(\mathbf{x}_t; \zeta_t) - \nabla f(\mathbf{x}_{t-1}; \zeta_t), \quad \text{if } \text{mod}(t, q) \neq 0, \quad (6)$$

where  $\nabla f(\mathbf{x}; \zeta)$  is the stochastic gradient dependent on parameter  $\mathbf{x}$  and a data sample  $\zeta$ , and  $q$  is the number of the inner loop iterations. On the other hand, if  $\text{mod}(t, q) = 0$ ,  $\mathbf{v}_t$  takes a full gradient. Note that these two estimators have a *similar* structure: Both are *recursively* updating the previous estimation based on the difference of the gradient estimations between two consecutive iterations (i.e., momentum). This motivates us to consider the following question: *Could we somehow “integrate” these two methods to develop a new decentralized gradient estimator to track the global gradient and reduce the stochastic error at the same time?* Unfortunately, the GT and VR estimators can not be combined straightforwardly. The major challenge lies in the structural difference in the outer loop iteration (i.e.,  $\text{mod}(t, q) = 0$ ), where the VR estimator requires full gradient and does not follow the recursive updating structure.

Surprisingly, in this paper, we show that there exists an “indirect” way to achieve the salient features of both GT and VR. Our approach is to abandon the double-loop structure of VR and pursue a *single-loop* structure. Yet, this single-loop structure should still be able to reduce the variance and consistently track the global gradient. Specifically, we introduce a parameter  $\beta_t \in [0, 1]$  in the recursive update and integrate it with a consensus step as follows:

$$\mathbf{v}_{i,t} = \beta_t \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{v}_{j,t-1} + \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t}) - \beta_t \nabla f_i(\mathbf{x}_{i,t-1}; \zeta_{i,t}), \quad (7)$$

where  $\mathbf{x}_{i,t}$ ,  $\mathbf{v}_{i,t}$  and  $\zeta_{i,t}$  are the parameter, gradient estimator, and random sample in the  $t$ th iteration at node  $i$ , respectively. Note that the estimator reduces to the classical stochastic gradient estimator when  $\beta_t = 0$ . On the other hand, if we set  $\beta_t = 1$ , the estimator becomes the (stochastic) gradient tracking estimator based on a single sample (implying low sample complexity). Then, the key to the success of our GT-STORM design lies in meticulously choosing parameter  $\beta_t$  to mimic the gradient estimator technique in centralized optimization [14, 15]. Lastly, the local parameters can be updated by the conventional decentralized stochastic gradient descent step:

$$\mathbf{x}_{i,t+1} = \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{x}_{j,t} - \eta_t \mathbf{v}_{i,t}, \quad (8)$$

where  $\eta_t$  is the step-size in iteration  $t$ . To summarize, we state our algorithm in Algorithm 1 as follows.

**Algorithm 1:** Gradient-Tracking-based Stochastic Recursive Momentum Algorithm (GT-STORM).

#### Initialization:

1. Choose  $T > 0$  and let  $t = 1$ . Set  $\mathbf{x}_{i,0} = \mathbf{x}^0$  at node  $i$ . Calculate  $\mathbf{v}_{i,0} = \nabla f_i(\mathbf{x}_{i,0}; \zeta_{i,0})$  at node  $i$ .

#### Main Loop:

2. In the  $t$ -th iteration, each node sends  $\mathbf{x}_{i,t-1}$  and local gradient estimator  $\mathbf{v}_{i,t-1}$  to its neighbors. Meanwhile, upon the reception of all neighbors' information, each node performs the following:
  - a) Update local parameter:  $\mathbf{x}_{i,t} = \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{x}_{j,t-1} - \eta_{t-1} \mathbf{v}_{i,t-1}$ .
  - b) Update local gradient estimator:  $\mathbf{v}_{i,t} = \beta_t \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{v}_{j,t-1} + \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t}) - \beta_t \nabla f_i(\mathbf{x}_{i,t-1}; \zeta_{i,t})$ .
3. Stop if  $t > T$ ; otherwise, let  $t \leftarrow t + 1$  and go to Step 2.

Two remarks for Algorithm 1 are in order. First, thanks to the single-loop structure, GT-STORM is easier to implement compared to the low-sample-complexity D-GET [19] method, which has in a

double-loop structure. Second, GT-STORM only requires  $p$  memory space due to the use of only one intermediate vector  $\mathbf{v}$  at each node. In contrast, the memory complexity of D-GET is  $2p$  (cf.  $\mathbf{y}$  and  $\mathbf{v}$  in [19]). This 50% saving is huge particularly for deep learning models, where the number of parameters could be in the range of millions.

### 3.2 Main Theoretical Results

In this section, we will establish the complexity properties of the proposed GT-STORM algorithm. For better readability, we state the main theorem and its corollary in this section and provide the intermediate lemmas to Section 3.3. We start with the following assumptions on the global and local objectives:

**ASSUMPTION 1.** The objective function  $f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x})$  with  $f_i(\mathbf{x}) = \mathbb{E}_{\zeta \sim \mathcal{D}_i} f_i(\mathbf{x}; \zeta)$  satisfies the following assumptions:

- (a) Boundedness from below: There exists a finite lower bound  $f^* = \inf_{\mathbf{x}} f(\mathbf{x}) > -\infty$ ;
- (b)  $L$ -average smoothness:  $f_i(\cdot; \zeta_i)$  is  $L$ -average smooth on  $\mathbb{R}^p$ , i.e., there exists a positive constant  $L$ , such that  $\mathbb{E}_{\zeta \sim \mathcal{D}_i} [\|\nabla f_i(\mathbf{x}; \zeta) - \nabla f_i(\mathbf{y}; \zeta)\|^2] \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p, i \in [m]$ ;
- (c) Bounded variance: There exists a constant  $\sigma \geq 0$  such that  $\mathbb{E}_{\zeta \sim \mathcal{D}_i} [\|\nabla f_i(\mathbf{x}; \zeta) - \nabla f_i(\mathbf{x})\|^2] \leq \sigma^2, \forall \mathbf{x} \in \mathbb{R}^p, i \in [m]$ ;
- (d) Bounded gradient: There exists a constant  $G \geq 0$  such that  $\mathbb{E}_{\zeta \sim \mathcal{D}_i} [\|\nabla f_i(\mathbf{x}; \zeta)\|^2] \leq G^2, \forall \mathbf{x} \in \mathbb{R}^p, i \in [m]$ .

In the above assumptions, (a) and (c) are standard in the stochastic non-convex optimization literature; (b) is an expected Lipschitz smoothness condition over the data distribution, which implies the conventional global Lipschitz smoothness [20] by the Jensen's inequality. Note that (b) is weaker than the individual Lipschitz smoothness in [19, 27, 29]: if there exists an outlier data sample, then the individual objective function might have a very large smoothness parameter while the average smoothness can still be small; (d) is equivalent to the Lipschitz continuity assumption, which is also commonly used for non-convex stochastic algorithms [43–45] and is essential for analyzing the decentralized gradient descent method [18, 30, 38].<sup>1</sup>

For convenience, in the subsequent analysis, we define  $\tilde{\mathbf{W}} = \mathbf{W} \otimes \mathbf{I}_m$ ,  $\mathbf{g}_{i,t} = \nabla f_i(\mathbf{x}_{i,t})$ ,  $\mathbf{u}_{i,t} = \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t})$ ,  $\mathbf{w}_{i,t} = \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t}) - \nabla f_i(\mathbf{x}_{i,t-1}; \zeta_{i,t})$  and  $\mathbf{a}_t = [\mathbf{a}_{1,t}^\top, \dots, \mathbf{a}_{m,t}^\top]^\top$  and  $\tilde{\mathbf{a}}_t = \frac{1}{m} \sum_{i=1}^m \mathbf{a}_{i,t}$ , for  $\mathbf{a} \in \{\mathbf{x}, \mathbf{u}, \mathbf{w}, \mathbf{v}, \mathbf{g}\}$ . Then, the algorithm can be compactly rewritten in the following matrix-vector form:

$$\mathbf{x}_t = \tilde{\mathbf{W}} \mathbf{x}_{t-1} - \eta_{t-1} \mathbf{v}_{t-1}, \quad (9)$$

$$\mathbf{v}_t = \beta_t \tilde{\mathbf{W}} \mathbf{v}_{t-1} + \beta_t \mathbf{w}_t + (1 - \beta_t) \mathbf{u}_t. \quad (10)$$

Furthermore, since  $\mathbf{1}^\top \mathbf{W} = \mathbf{1}^\top$ , we have  $\bar{\mathbf{x}}_t = \bar{\mathbf{x}}_{t-1} - \eta_{t-1} \bar{\mathbf{v}}_{t-1}$ ,  $\bar{\mathbf{v}}_t = \beta_t \bar{\mathbf{v}}_{t-1} + \beta_t \bar{\mathbf{w}}_t + (1 - \beta_t) \bar{\mathbf{u}}_t$ . We first state the convergence result for Algorithm 1 as follows:

**THEOREM 1.** Under Assumption 1 and with the positive constants  $c_0$  and  $c_1$  satisfying  $1 - (1 + c_1)\lambda^2 - \frac{1}{c_0} > 0$ , if we set  $\eta_t = \tau/(\omega + t)^{1/3}$  and  $\beta_{t+1} = 1 - \rho \eta_t^2$ , with  $\tau > 0$ ,  $\omega \geq \max\{2, \tau^3 / \min\{k_1^3, k_2^3, k_3^3\}\}$  and

$\rho = 2/(3\tau^3) + 32L^2$ , then we have the following result for Algorithm 1:

$$\begin{aligned} & \min_{t \in [T]} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\ & \leq \frac{2(f(\bar{\mathbf{x}}_0) - f(\bar{\mathbf{x}}^*))}{\tau(T+1)^{2/3}} + \frac{2c_0 \mathbb{E}[\|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2]}{m\tau(T+1)^{2/3}} \\ & \quad + \frac{(\omega - 1)\sigma^2}{16mL^2\tau^2(T+1)^{2/3}} + \frac{\rho^2\sigma^2 \ln(\omega + T - 1)}{8mL^2(T+1)^{2/3}} \\ & \quad + \frac{12(1 + \frac{1}{c_1})c_0\tau^{1/3}G^2\rho^2}{(\omega - 1)^{1/3}(T+1)^{2/3}} + O\left(\frac{c_3\omega}{\tau T^{5/3}}\right), \end{aligned} \quad (11)$$

where  $c_3 = \max\{1, \omega/(m\tau^2), \tau^{4/3}/\omega^{1/3}, \tau \ln(\omega + T)/m\}$ , and the constants  $k_1$ ,  $k_2$  and  $k_3$  are:

$$k_1 = 1/\left(2L + 32(1 + \frac{1}{c_1})c_0L^2\right), \quad (12)$$

$$k_2 = \left(1 - (1 + c_1)\lambda^2\right)/\left(1 + \frac{1}{c_1} + \frac{1}{c_0}\right), \quad (13)$$

$$k_3 = \sqrt{\left(1 - (1 + c_1)\lambda^2 - \frac{1}{c_0}\right)/\left(\frac{2}{3\tau^3} + \frac{2L^2 + 1}{2c_0}\right)}. \quad (14)$$

In Theorem 1,  $c_0$  and  $c_1$  are two constants depending on the network topology, which in turn will affect the step-size and convergence: with a sparse network, i.e.,  $\lambda$  is close to but not exactly one (recall that  $\lambda = \max\{|\lambda_2|, |\lambda_m|\}$ ). In order for  $1 - (1 + c_1)\lambda^2 - \frac{1}{c_0} > 0$  to hold,  $c_0$  needs to be large and  $c_1$  needs to be close to zero, which leads to small  $k_1$ ,  $k_2$  and  $k_3$ . Note that the step-size  $\eta_t$  is of the order  $O(t^{-1/3})$ , which is larger than the  $O(t^{-1/2})$  order for the classical decentralized SGD algorithms. With this larger step-size, the convergence rate is  $O(t^{-2/3})$  and faster than the rate  $O(t^{-1/2})$  for the decentralized SGD algorithms. Based on Theorem 1, we have the sample and communication complexity results for Algorithm 1:

**COROLLARY 2.** Under the conditions in Theorem 1, if  $\tau = O(m^{1/3})$  and  $\omega = O(m^{4/3})$ , then to achieve an  $\epsilon^2$ -stationary solution, the total communication rounds are on the order of  $\tilde{O}(m^{-1/2}\epsilon^{-3})$  and the total samples evaluated across the network is on the order of  $\tilde{O}(m^{1/2}\epsilon^{-3})$ .

### 3.3 Proofs of the Theoretical Results

Due to space limitation, we provide a proof sketch for Theorem 1 here and relegate the details to the appendices. First, we bound the error of gradient estimator  $\mathbb{E}[\|\mathbf{v}_t - \mathbf{g}_t\|^2]$  as follows:

**LEMMA 1 (ERROR OF GRADIENT ESTIMATOR).** Under Assumption 1 and with  $\mathbf{v}_t$  defined in (10), it holds that  $\mathbb{E}[\|\tilde{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2] \leq \beta_t^2 \mathbb{E}[\|\tilde{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2] + \frac{2\beta_t^2 L^2}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2] + \frac{2(1 - \beta_t)\sigma^2}{m}$ .

It can be seen that the upper bound depends on the error in the previous step with a factor  $\beta_t^2$ . This will be helpful when we construct a potential function. Then, according to the algorithm updates (9)–(10), we show the following descent inequality:

**LEMMA 2 (DESCENT LEMMA).** Under Assumption 1, Algorithm 1 satisfies:  $\mathbb{E}[f(\bar{\mathbf{x}}_{t+1})] - \mathbb{E}[f(\bar{\mathbf{x}}_t)] \leq -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \times \mathbb{E}[\|\tilde{\mathbf{v}}_t\|^2] + \eta_t \mathbb{E}[\|\tilde{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2] + \frac{L^2\eta_t}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2]$ .

We remark that the right-hand-side (RHS) of the above inequality contains the consensus error of local parameters  $\sum_{t=0}^T \mathbb{E}[\|\mathbf{x}_t -$

<sup>1</sup>Note that under the assumption (b), as long as the parameter  $\mathbf{x}$  is bounded, (d) is satisfied.

$1 \otimes \bar{\mathbf{x}}_t\|^2$ ), which makes the analysis more difficult than that of the centralized optimization. Next, we prove the contraction of iterations in the following lemma, which is useful in analyzing the decentralized gradient tracking algorithms.

**LEMMA 3 (ITERATES CONTRACTION).** *The following contraction properties of the iterates produced by Algorithm 1 hold:*

$$\|\mathbf{x}_{t+1} - 1 \otimes \bar{\mathbf{x}}_{t+1}\|^2 \leq (1 + c_1)\lambda^2 \|\mathbf{x}_t - 1 \otimes \bar{\mathbf{x}}_t\|^2 + (1 + \frac{1}{c_1})\eta_t^2 \|\mathbf{v}_t - 1 \otimes \bar{\mathbf{v}}_t\|^2, \quad (15)$$

$$\|\mathbf{v}_{t+1} - 1 \otimes \bar{\mathbf{v}}_{t+1}\|^2 \leq (1 + c_1)\beta_{t+1}^2 \lambda^2 \|\mathbf{v}_t - 1 \otimes \bar{\mathbf{v}}_t\|^2 + 2(1 + \frac{1}{c_1})(\beta_{t+1}^2 \|\mathbf{w}_{t+1}\|^2 + (1 - \beta_{t+1})^2 \|\mathbf{u}_{t+1}\|^2), \quad (16)$$

where  $c_1$  is a positive constant. Additionally, we have

$$\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \leq 8\|\mathbf{x}_t - 1 \otimes \bar{\mathbf{x}}_t\|^2 + 4\eta_t^2 \|\mathbf{v}_t - 1 \otimes \bar{\mathbf{v}}_t\|^2 + 4\eta_t^2 m \|\bar{\mathbf{v}}_t\|^2. \quad (17)$$

Finally, we define a potential function in (18), based on which we prove the convergence bound:

**LEMMA 4. (Convergence of Potential Function)** *Define the following potential function:*

$$H_t = \mathbb{E}[f(\bar{\mathbf{x}}_t) + \frac{1}{32L^2\eta_{t-1}} \|\bar{\mathbf{g}}_t - \bar{\mathbf{v}}_t\|^2 + \frac{c_0}{m\eta_{t-1}} \|\mathbf{x}_t - 1 \otimes \bar{\mathbf{x}}_t\|^2 + \frac{c_0}{m} \|\mathbf{v}_t - 1 \otimes \bar{\mathbf{v}}_t\|^2], \quad (18)$$

where  $c_0$  is a positive constant. Under Assumption 1, if we set  $\eta_t = \tau/(\omega + t)^{1/3}$  and  $\beta_{t+1} = 1 - \rho\eta_t^2$ , where  $\tau, \omega \geq 2, \rho = 2/(3\tau^3) + 32L^2$  are three constants, then it holds that:

$$H_{t+1} - H_t \leq -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{\rho^2 \sigma^2 \eta_t^3}{16mL^2} + 2(1 + \frac{1}{c_1})c_0 G^2 \rho^2 \eta_t^4 - \frac{c_0 C_1}{m\eta_t} [\|\mathbf{x}_t - 1 \otimes \bar{\mathbf{x}}_t\|^2] - \frac{c_0 C_2}{m} \mathbb{E}[\|\mathbf{v}_t - 1 \otimes \bar{\mathbf{v}}_t\|^2] - \frac{C_3 \eta_t}{4} \mathbb{E}[\|\bar{\mathbf{v}}_t\|^2], \quad (19)$$

where  $C_1, C_2$ , and  $C_3$  are following constants:  $C_1 = 1 - (1 + c_1)\lambda^2 - \frac{1}{2c_0} - 16(1 + \frac{1}{c_1})L^2\eta_t - (\frac{2}{3\tau^3} + \frac{L^2}{c_0})\eta_t^2$ ,  $C_2 = 1 - (1 + c_1)\lambda^2 - (1 + \frac{1}{c_1})\eta_t - \frac{\eta_t}{4c_0} - 8(1 + \frac{1}{c_1})L^2\eta_t^2$ ,  $C_3 = 1 - 2L\eta_t - 32(1 + \frac{1}{c_1})c_0 L^2 \eta_t$ .

Finally, by properly selecting the parameters, constants  $C_1, C_2$  and  $C_3$  can be made non-negative, which leads to Theorem 1.

## 4 EXPERIMENTAL RESULTS

In this section, we conduct experiments using several non-convex machine learning problems to evaluate the performance of our method. In particular, we compare our algorithm with the following state-of-art *single-loop* algorithms:

- DSGD [4, 18, 30]: Each node performs:  $\mathbf{x}_{i,t+1} = \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{x}_{j,t} - \eta \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t})$ , where the stochastic gradient  $\nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t})$  corresponds to random sample  $\zeta_{i,t}$ . Then, each node exchanges the local parameter  $\mathbf{x}_{i,t}$  with its neighbors.
- GNSD [17]: Each node keeps two variables  $\mathbf{x}_{i,t}$  and  $\mathbf{y}_{i,t}$ . The local parameter  $\mathbf{x}_{i,t}$  is updated as  $\mathbf{x}_{i,t+1} = \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{x}_{j,t} - \eta \mathbf{y}_{i,t}$  and the tracked gradient  $\mathbf{y}_{i,t}$  is updated as  $\mathbf{y}_{i,t+1} = \sum_{j \in \mathcal{N}_i} [\mathbf{W}]_{ij} \mathbf{y}_{j,t} + \nabla f_i(\mathbf{x}_{i,t+1}; \zeta_{i,t+1}) - \nabla f_i(\mathbf{x}_{i,t}; \zeta_{i,t})$ .

Here, we compare with the above two classes of stochastic algorithms because they all employ a single-loop structure and do not require full gradient evaluations. We note that it is hard to have a direct and fair comparison with D-GET [19] numerically, since D-GET uses full gradients and has a double-loop structure.

**Network Model:** The communication graph  $\mathcal{G}$  is generated by the Erdős-Rényi graph with different edge connectivity probability  $p_c$  and number of nodes  $m$ . We set  $m = 10$  and the edge connectivity probability as  $p_c = 0.5$ . The consensus matrix is chosen as  $\mathbf{W} = \mathbf{I} - \frac{2}{3\lambda_{\max}(\mathbf{L})}\mathbf{L}$ , where  $\mathbf{L}$  is the Laplacian matrix of  $\mathcal{G}$ , and  $\lambda_{\max}(\mathbf{L})$  denotes the largest eigenvalue of  $\mathbf{L}$ .

**1) Non-convex logistic regression:** In our first experiment, we consider the binary logistic regression problem with a non-convex regularizer [14, 29, 46]:

$$\min_{\mathbf{x} \in \mathbb{R}^d} -\frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n [y_{ij} \log(\frac{1}{1 + e^{-\mathbf{x}^\top \zeta_{ij}}}) + (1 - y_{ij}) \log(\frac{e^{-\mathbf{x}^\top \zeta_{ij}}}{1 + e^{-\mathbf{x}^\top \zeta_{ij}}})] + \alpha \sum_{i=1}^d \frac{\mathbf{x}_i^2}{1 + \mathbf{x}_i^2}, \quad (20)$$

where the label  $y_{ij} \in \{0, 1\}$ , the feature  $\zeta_{ij} \in \mathbb{R}^d$  and  $\alpha = 0.1$ .

**1-a) Datasets:** We consider three commonly used binary classification datasets from LibSVM: *a9a*, *rcv1.binary* and *ijcnn1*. The *a9a* dataset has 32561 samples, 123 features, the *rcv1.binary* dataset has 20242 samples, 47236 features, and the *ijcnn1* dataset has 49990 samples, 22 features. We evenly divide the dataset into  $m$  sub-datasets corresponding to the  $m$  nodes.

**1-b) Parameters:** For all algorithms, we set the batch size as one and the initial step-size  $\eta_0$  is tuned by searching over the grid  $\{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1.0\}$ . For DSGD and GNSD, the step-size is set to  $\eta_t = \eta_0 / \sqrt{1 + 0.1t}$ , which is on the order of  $O(t^{-1/2})$  following the state-of-the-art theoretical result [17]. For GT-STORM, the step-size is set as  $\eta_t = \eta_0 / \sqrt[3]{1 + 0.1t}$ , which is on the order of  $O(t^{-1/3})$  as specified in our theoretical result. In addition, we choose the parameter  $\rho$  for GT-STORM as  $1/\eta_0^2$ , so that  $\beta_1 = 0$  in the first step.

**1-c) Results:** We first compare the convergence rates of the algorithms. We adopt the consensus loss defined in the left-hand-side (LHS) of (3) as the criterion. After tuning, the best initial step-sizes are 0.1, 0.5 and 0.2 for *a9a*, *ijcnn1* and *rcv1.binary*, respectively. The results are shown in Figs. 1–3. It can be seen that our algorithm has a better performance: for *a9a* and *rcv1.binary* datasets, all algorithms reach almost the same accuracy but our algorithm has a faster speed; for *ijcnn1* dataset, our algorithm outperforms other methods both in the speed and accuracy.

Next, we examine the effect of the parameter  $\rho$  on our algorithm. We focus on the *a9a* dataset and fix the initial step-size as  $\eta_0 = 0.1$ . We choose  $\rho$  from  $\{10^{-1}, 10^0, 10^1, 10^2\}$ . Note that  $\rho = 10^2$  is corresponding to the case  $\rho = 1/\eta_0^2$ . The results are shown in Fig. 4. It can be seen that the case  $\rho = 10^1$  has the best performance, which is followed by the case  $\rho = 10^2$ . Also, as  $\rho$  decreases, the convergence speed becomes slower (see the cases  $\rho = 10^{-1}$  and  $10^0$ ).

In addition, we examine the effect of the network topology. We first fix the number of workers as  $m = 10$  and change the edge connectivity probability  $p_c$  from 0.35 to 0.9. Note that with a smaller

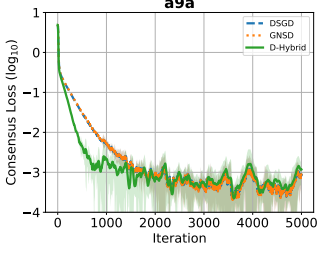


Figure 1: Non-convex logistic regression on LibSVM: a9a.

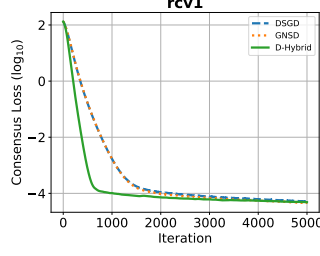


Figure 2: Non-convex logistic regression on LibSVM: rcv1.

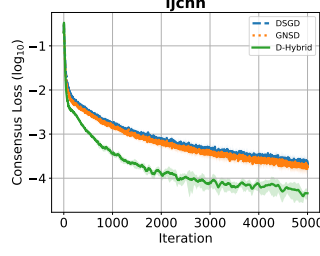


Figure 3: Non-convex logistic regression on LibSVM: ijcnn.

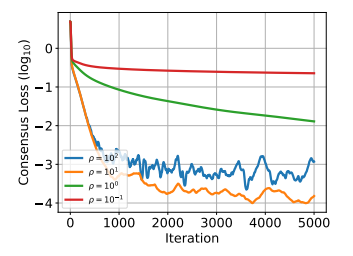
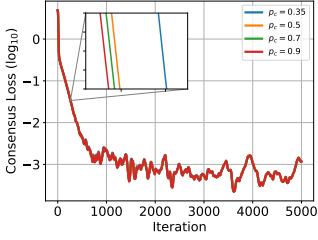
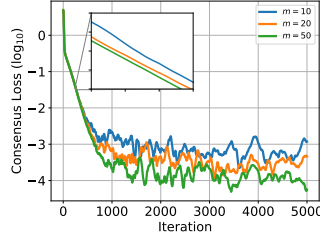
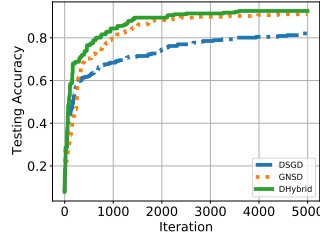
Figure 4: Non-convex logistic regression: The effect of  $\rho$ .Figure 5: Non-convex logistic regression: The effect of  $p_c$ .Figure 6: Non-convex logistic regression: The effect of  $m$ .

Figure 7: CNN experimental results on MNIST dataset.

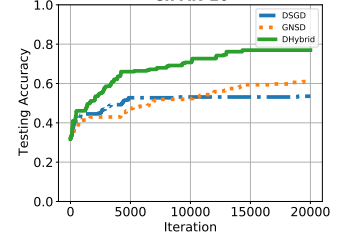


Figure 8: CNN experimental results on CIFAR-10 dataset.

$p_c$ , the network becomes sparser. We set  $\eta_0 = 0.1$  and  $\rho = 10^2$ . The results are shown in Fig. 5. Under different  $p_c$ -values, our algorithm has a similar performance in terms of convergence speed and accuracy. But with a larger  $p_c$ -values i.e., a denser network, the convergence speed slightly increases (see the zoom-in view in Fig. 6). Then, we fix the the edge connectivity probability  $p_c = 0.5$  but change the number of workers  $m$  from 10 to 50. We show the results in Fig. 6. It can be seen that with more workers, the algorithm converges faster and reaches a better accuracy.

**2) Convolutional neural networks** We use all three algorithms to train a convolutional neural network (CNN) model for image classification on MNIST and CIFAR-10 datasets. We adopt the same network topology as in the previous experiment. We use a non-identically distributed data partition strategy: the  $i$ th machine can access the data with the  $i$ th label. We fix the initial step-size as  $\eta_0 = 0.01$  for all three algorithms and the remaining settings are the same as in the previous experiment.

**2-a) Learning Models:** For MNIST, the adopted CNN model has two convolutional layers (first of size  $1 \times 16 \times 5$  and then of size  $16 \times 32 \times 5$ ), each of which is followed by a max-pooling layer with size  $2 \times 2$ , and then a fully connected layer. The ReLU activation is used for the two convolutional layers and the “softmax” activation is applied at the output layer. The batch size is 64 for the CNN training on MNIST. For CIFAR-10, we apply the CNN model with two convolutional layers (first of size  $3 \times 6 \times 5$  and then of size  $6 \times 16 \times 5$ ). Each of the convolutional layers is followed by a max-pooling layer of size  $2 \times 2$ , and then three fully connected layers. The ReLU activation is used for the two convolutional layers and the first two fully connected layers, and the “softmax” activation is applied at the output layer. The batch size is chosen as 128 for the CNN training on CIFAR-10.

**2-b) Results:** Fig. 7 illustrates the testing accuracy of different algorithms versus iterations on MNIST and CIFAR-10 datasets. It can be seen from Fig. 7 that on the MNIST dataset, GNSD and GT-STORM have similar performance, but our GT-STORM maintains a faster speed and a better prediction accuracy. Compared with DSGD, our GT-STORM can gain about 10% more accuracy. On the CIFAR-10 dataset (see Fig. 8), the performances of DSGD and GNSD deteriorate, while GT-STORM can achieve a better accuracy. Specifically, the accuracy of GT-STORM is around 15% higher than that of GNSD and 25% higher than that of DSGD.

## 5 CONCLUSION

In this paper, we proposed a gradient-tracking-based stochastic recursive momentum (GT-STORM) algorithm for decentralized non-convex optimization, which enjoys low sample, communication, and memory complexities. Our algorithm fuses the gradient tracking estimator and the variance reduction estimator and has a simple single-loop structure. Thus, it is more practical compared to existing works (e.g. GT-SAGA/SVRG and D-GET) in the literature. We have also conducted extensive numerical studies to verify the performance of our method, including non-convex logistic regression and neural networks. The numerical results show that our method outperforms the state-of-the-art methods when training on the large datasets. Our results in this work contribute to the increasingly important field of decentralized network training.

## REFERENCES

- [1] B. Recht, C. Re, S. Wright, and F. Niu, “Hogwild: A lock-free approach to parallelizing stochastic gradient descent,” in *Advances in neural information processing systems*, 2011, pp. 693–701.
- [2] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, “Parallelized stochastic gradient descent,” in *Advances in neural information processing systems*, 2010, pp. 2595–2603.



- [3] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang *et al.*, "Large scale distributed deep networks," in *Advances in neural information processing systems*, 2012, pp. 1223–1231.
- [4] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Transactions on Automatic Control*, vol. 54, no. 1, p. 48, 2009.
- [5] 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*, 2017, pp. 5330–5340.
- [6] J. Chen, Z. J. Towfic, and A. H. Sayed, "Dictionary learning over distributed models," *IEEE Transactions on Signal Processing*, vol. 63, no. 4, pp. 1001–1016, 2014.
- [7] Y. Cao, W. Yu, W. Ren, and G. Chen, "An overview of recent progress in the study of distributed multi-agent coordination," *IEEE Transactions on Industrial Informatics*, vol. 9, no. 1, pp. 427–438, 2012.
- [8] K. Zhou and S. I. Roumeliotis, "Multirobot active target tracking with combinations of relative observations," *IEEE Transactions on Robotics*, vol. 27, no. 4, pp. 678–695, 2011.
- [9] W. Wang, J. Wang, M. Kolar, and N. Srebro, "Distributed stochastic multi-task learning with graph regularization," *arXiv preprint arXiv:1802.03830*, 2018.
- [10] X. Zhang, J. Liu, and Z. Zhu, "Distributed linear model clustering over networks: A tree-based fused-lasso admm approach," *arXiv preprint arXiv:1905.11549*, 2019.
- [11] K. Ali and W. Van Stam, "Tivo: Making show recommendations using a distributed collaborative filtering architecture," in *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, 2004, pp. 394–401.
- [12] Z. Jiang, K. Mukherjee, and S. Sarkar, "On consensus-disagreement tradeoff in distributed optimization," in *2018 Annual American Control Conference (ACC)*. IEEE, 2018, pp. 571–576.
- [13] J. N. Tsitsiklis, "Problems in decentralized decision making and computation." Massachusetts Inst of Tech Cambridge Lab for Information and Decision Systems, Tech. Rep., 1984.
- [14] Q. Tran-Dinh, N. H. Pham, D. T. Phan, and L. M. Nguyen, "Hybrid stochastic gradient descent algorithms for stochastic nonconvex optimization," *arXiv preprint arXiv:1905.05920*, 2019.
- [15] A. Cutkosky and F. Orabona, "Momentum-based variance reduction in non-convex sgd," in *Advances in Neural Information Processing Systems*, 2019, pp. 15 210–15 219.
- [16] P. Di Lorenzo and G. Scutari, "Next: In-network nonconvex optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.
- [17] S. Lu, X. Zhang, H. Sun, and M. Hong, "GNSD: a gradient-tracking based non-convex stochastic algorithm for decentralized optimization," in *2019 IEEE Data Science Workshop, DSW 2019*. Institute of Electrical and Electronics Engineers Inc., 2019, pp. 315–321.
- [18] Z. Jiang, A. Balu, C. Hegde, and S. Sarkar, "Collaborative deep learning in fixed topology networks," in *Advances in Neural Information Processing Systems*, 2017, pp. 5904–5914.
- [19] H. Sun, S. Lu, and M. Hong, "Improving the sample and communication complexity for decentralized non-convex optimization: A joint gradient estimation and tracking approach," *ICML 2020*, 2019.
- [20] S. Ghadimi and G. Lan, "Stochastic first-and zeroth-order methods for nonconvex stochastic programming," *SIAM Journal on Optimization*, vol. 23, no. 4, pp. 2341–2368, 2013.
- [21] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.
- [22] P. Zhou, X. Yuan, and J. Feng, "New insight into hybrid stochastic gradient descent: Beyond with-replacement sampling and convexity," in *Advances in Neural Information Processing Systems*, 2018, pp. 1234–1243.
- [23] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," in *Advances in neural information processing systems*, 2013, pp. 315–323.
- [24] S. J. Reddi, A. Hefny, S. Sra, B. Póczos, and A. Smola, "Stochastic variance reduction for nonconvex optimization," in *International conference on machine learning*, 2016, pp. 314–323.
- [25] 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*, 2014, pp. 1646–1654.
- [26] 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*, 2017, pp. 2348–2358.
- [27] 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*, 2018, pp. 689–699.
- [28] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takác, "Sarah: A novel method for machine learning problems using stochastic recursive gradient," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2613–2621.
- [29] 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*, 2019, pp. 2403–2413.
- [30] K. Yuan, Q. Ling, and W. Yin, "On the convergence of decentralized gradient descent," *SIAM Journal on Optimization*, vol. 26, no. 3, pp. 1835–1854, 2016.
- [31] W. Shi, Q. Ling, G. Wu, and W. Yin, "Extra: An exact first-order algorithm for decentralized consensus optimization," *SIAM Journal on Optimization*, vol. 25, no. 2, pp. 944–966, 2015.
- [32] H. Sun and M. Hong, "Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms," *IEEE Transactions on Signal processing*, vol. 67, no. 22, pp. 5912–5928, 2019.
- [33] J. F. Mota, J. M. Xavier, P. M. Aguiar, and M. Püschel, "D-admm: A communication-efficient distributed algorithm for separable optimization," *IEEE Transactions on Signal Processing*, vol. 61, no. 10, pp. 2718–2723, 2013.
- [34] A. Mokhtari, W. Shi, Q. Ling, and A. Ribeiro, "A decentralized second-order method with exact linear convergence rate for consensus optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 4, pp. 507–522, 2016.
- [35] M. Eisen, A. Mokhtari, and A. Ribeiro, "Decentralized quasi-newton methods," *IEEE Transactions on Signal Processing*, vol. 65, no. 10, pp. 2613–2628, 2017.
- [36] A. Nedić, A. Olshevsky, and M. G. Rabbat, "Network topology and communication-computation tradeoffs in decentralized optimization," *Proceedings of the IEEE*, vol. 106, no. 5, pp. 953–976, 2018.
- [37] T.-H. Chang, M. Hong, H.-T. Wai, X. Zhang, and S. Lu, "Distributed learning in the non-convex world: From batch to streaming data, and beyond," *arXiv preprint arXiv:2001.04786*, 2020.
- [38] J. Zeng and W. Yin, "On nonconvex decentralized gradient descent," *IEEE Transactions on signal processing*, vol. 66, no. 11, pp. 2834–2848, 2018.
- [39] A. Mokhtari and A. Ribeiro, "DSA: decentralized double stochastic averaging gradient algorithm," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 2165–2199, 2016.
- [40] K. Yuan, B. Ying, J. Liu, and A. H. Sayed, "Variance-reduced stochastic learning by networked agents under random reshuffling," *IEEE Transactions on Signal Processing*, vol. 67, no. 2, pp. 351–366, 2018.
- [41] R. Xin, U. A. Khan, and S. Kar, "Variance-reduced decentralized stochastic optimization with gradient tracking," *arXiv preprint arXiv:1909.11774*, 2019.
- [42] A. Nedić, A. Olshevsky, and W. Shi, "A geometrically convergent method for distributed optimization over time-varying graphs," in *2016 IEEE 55th Conference on Decision and Control (CDC)*. IEEE, 2016, pp. 1023–1029.
- [43] Y. Zhou, Y. Liang, and H. Zhang, "Generalization error bounds with probabilistic guarantee for sgd in nonconvex optimization," *arXiv preprint arXiv:1802.06903*, 2018.
- [44] S. P. Karimireddy, Q. Rebjock, S. U. Stich, and M. Jaggi, "Error feedback fixes SignSGD and other gradient compression schemes," *arXiv preprint arXiv:1901.09847*, 2019.
- [45] A. Koloskova, S. U. Stich, and M. Jaggi, "Decentralized stochastic optimization and gossip algorithms with compressed communication," *arXiv preprint arXiv:1902.00340*, 2019.
- [46] Z. Wang, Y. Zhou, Y. Liang, and G. Lan, "Cubic regularization with momentum for nonconvex optimization," *arXiv preprint arXiv:1810.03763*, 2018.
- [47] X. Zhang, J. Liu, Z. Zhu, and E. S. Bentley, (2020) GT-STORM: taming sample, communication, and memory complexities in decentralized non-convex learning. [Online]. Available: [https://kevinliu-osu-ece.github.io/publications/GT-STORM\\_TR.pdf](https://kevinliu-osu-ece.github.io/publications/GT-STORM_TR.pdf)

## A ADDITIONAL EXPERIMENT DETAILS

In our simulation, the communication graph  $\mathcal{G}$  is generated by the Erdős-Rényi graph with different edge connectivity probability  $p_c$  and nodes number  $m$ . We set  $m = 10$  and  $p_c = 0.5$ . The generated graph is shown in Figure 9.

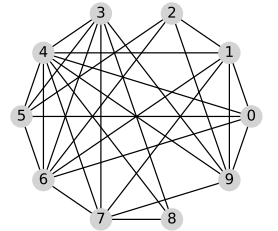


Figure 9: The network topology generated by the Erdős-Rényi graph.

### A.1 Nonconvex Logistic Regression

In Section 4.1, we set the step-size as  $\eta_t = \eta_0 / \sqrt{1 + 0.1 \times t}$  for DSGD and GNSD, while  $\eta_t =$



$\eta_0/\sqrt[3]{1+0.1 \times t}$  for GT-STORM. It can be noted that the step-size adopted for GT-STORM is diminishing slower than those for DSGD and GNSD, though the choices are following the theoretical results. Thus, here we apply the step-size as  $\eta_t = \eta_0/\sqrt[3]{1+0.1 \times t}$  for all the three algorithms. We tune the initial step-size  $\eta_0$  by searching the grid  $\{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1.0\}$ . After tuning, the best initial step-sizes are 0.1, 0.5 and 0.2 for *a9a*, *ijcnn1* and *rcv1.binary*, respectively. We show the results in Figure 10. It can be seen that with a larger step-size, though the convergence is faster for DSGD and GNSD at the beginning, the accuracy is unsatisfactory (e.g. *a9a* and *rcv1.binary*). Also, with the same step-size, our algorithm performs much better than the other two.

## A.2 Convolutional Neural Networks

Here we show the testing loss and accuracy for the CNN models on the MNIST and CIFAR-10 datasets in Figure 11-12. In all experiment results, our algorithm has a better performance: a higher accuracy and a smaller loss. The final testing accuracy results are summarized in Table 1.

**Table 1: The results of testing accuracy with the CNN models trained by different algorithms.**

Dataset		DSGD	GNSD	GT-STORM
MNIST	I.D.	0.9102	0.9102	0.9375
	N.D.	0.8203	0.9102	0.9257
CIFAR-10	I.D.	0.6093	0.6016	0.7734
	N.D.	0.5352	0.6133	0.7695

## B PROOF OF MAIN RESULTS

Due to space limitation, we provide key proof steps of the key lemmas and theorems in this appendix. We refer readers to [47] for the complete proofs.

### B.1 Proof for Lemma 1

PROOF. Recall that  $\bar{\mathbf{v}}_t = \beta_t(\bar{\mathbf{v}}_{t-1} + \bar{\mathbf{w}}_t) + (1-\beta_t)\bar{\mathbf{u}}_t$ , then

$$\begin{aligned} \|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2 &= \|\beta_t(\bar{\mathbf{v}}_{t-1} + \bar{\mathbf{w}}_t) + (1-\beta_t)\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t\|^2 \\ &= \|\beta_t(\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}) + \beta_t(\bar{\mathbf{w}}_t - \bar{\mathbf{g}}_t + \bar{\mathbf{g}}_{t-1}) + (1-\beta_t)(\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t)\|^2 \\ &= \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \|\beta_t(\bar{\mathbf{w}}_t - \bar{\mathbf{g}}_t + \bar{\mathbf{g}}_{t-1}) + (1-\beta_t)(\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t)\|^2 \\ &\quad + 2\langle \bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}, \beta_t(\bar{\mathbf{w}}_t - \bar{\mathbf{g}}_t + \bar{\mathbf{g}}_{t-1}) + (1-\beta_t)(\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t) \rangle \end{aligned} \quad (21)$$

Note that  $\mathbb{E}_{\zeta_t}[\bar{\mathbf{w}}_t] = \bar{\mathbf{g}}_t - \bar{\mathbf{g}}_{t-1}$  and  $\mathbb{E}_{\zeta_t}[\bar{\mathbf{u}}_t] = \bar{\mathbf{g}}_t$ . Taking expectation with respect to  $\zeta_t$ , we have:

$$\begin{aligned} \mathbb{E}_{\zeta_t}[\|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2] &\stackrel{(a)}{\leq} \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \mathbb{E}_{\zeta_t}[\|\beta_t(\bar{\mathbf{w}}_t - \bar{\mathbf{g}}_t + \bar{\mathbf{g}}_{t-1}) + (1-\beta_t)(\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t)\|^2] \\ &\stackrel{(b)}{\leq} \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \mathbb{E}_{\zeta_t}[2\beta_t^2\|\bar{\mathbf{w}}_t - \bar{\mathbf{g}}_t + \bar{\mathbf{g}}_{t-1}\|^2 + 2(1-\beta_t)^2\|\bar{\mathbf{u}}_t - \bar{\mathbf{g}}_t\|^2] \\ &\stackrel{(c)}{\leq} \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + 2\beta_t^2\mathbb{E}_{\zeta_t}[\|\bar{\mathbf{w}}_t\|^2] + \frac{2(1-\beta_t)^2\sigma^2}{m} \\ &\stackrel{(d)}{\leq} \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \frac{2\beta_t^2}{m} \sum_{i=1}^m \mathbb{E}_{\zeta_t}[\|\mathbf{w}_{i,t}\|^2] + \frac{2(1-\beta_t)^2\sigma^2}{m} \end{aligned}$$

$$\begin{aligned} &\stackrel{(e)}{\leq} \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \frac{2\beta_t^2L^2}{m} \sum_{i=1}^m \mathbb{E}_{\zeta_t}[\|\mathbf{x}_{i,t} - \mathbf{x}_{i,t-1}\|^2] + \frac{2(1-\beta_t)^2\sigma^2}{m} \\ &= \beta_t^2\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2 + \frac{2\beta_t^2L^2}{m} \mathbb{E}_{\zeta_t}[\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2] + \frac{2(1-\beta_t)^2\sigma^2}{m}, \end{aligned} \quad (22)$$

where (a) is because the cross term has the expectation as zero; (b) is by  $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$ ; (c) is by  $\mathbb{E}[\|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2] \leq \mathbb{E}[\|\mathbf{X}\|^2]$  and Assumption 1 (c); (d) is because of the Jensen's inequality; (e) is by the  $L$ -average smoothness. Thus, taking the full expectation, it holds that

$$\begin{aligned} \mathbb{E}[\|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2] &\leq \beta_t^2\mathbb{E}[\|\bar{\mathbf{v}}_{t-1} - \bar{\mathbf{g}}_{t-1}\|^2] \\ &\quad + \frac{2\beta_t^2L^2}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2] + \frac{2(1-\beta_t)^2\sigma^2}{m}. \end{aligned} \quad (23)$$

□

### B.2 Proof for Lemma 2

PROOF. From the  $L$ -smoothness of  $f$ , we have:

$$\begin{aligned} f(\bar{\mathbf{x}}_{t+1}) &\leq f(\bar{\mathbf{x}}_t) - \eta_t \langle \nabla f(\bar{\mathbf{x}}_t), \bar{\mathbf{v}}_t \rangle + \frac{L\eta_t^2}{2} \|\bar{\mathbf{v}}_t\|^2 \\ &= f(\bar{\mathbf{x}}_t) - \frac{\eta_t}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \|\bar{\mathbf{v}}_t\|^2 + \frac{\eta_t}{2} \|\bar{\mathbf{v}}_t - \nabla f(\bar{\mathbf{x}}_t)\|^2 \\ &\leq f(\bar{\mathbf{x}}_t) - \frac{\eta_t}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \|\bar{\mathbf{v}}_t\|^2 \\ &\quad + \eta_t \|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2 + \eta_t \|\bar{\mathbf{g}}_t - \nabla f(\bar{\mathbf{x}}_t)\|^2 \\ &\stackrel{(a)}{\leq} f(\bar{\mathbf{x}}_t) - \frac{\eta_t}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \|\bar{\mathbf{v}}_t\|^2 \\ &\quad + \eta_t \|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2 + \frac{\eta_t}{m} \sum_{i=1}^m \|\mathbf{g}_{i,t} - \nabla f_i(\bar{\mathbf{x}}_t)\|^2 \\ &\stackrel{(b)}{\leq} f(\bar{\mathbf{x}}_t) - \frac{\eta_t}{2} \|\nabla f(\bar{\mathbf{x}}_t)\|^2 - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \|\bar{\mathbf{v}}_t\|^2 \\ &\quad + \eta_t \|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2 + \frac{L^2\eta_t}{m} \|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2, \end{aligned} \quad (24)$$

where  $\bar{\mathbf{g}}_t = \frac{1}{m} \sum_{i=1}^m \mathbf{g}_{i,t}$  and  $\mathbf{g}_{i,t} = \nabla f_i(\mathbf{x}_{i,t})$ , (a) is because of the Jensen's inequality and (b) is by the  $L$ -average smoothness. Take the full expectation on the above inequality:

$$\begin{aligned} \mathbb{E}[f(\bar{\mathbf{x}}_{t+1})] - \mathbb{E}[f(\bar{\mathbf{x}}_t)] &\leq \frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] - \left(\frac{\eta_t}{2} - \frac{L\eta_t^2}{2}\right) \mathbb{E}[\|\bar{\mathbf{v}}_t\|^2] \\ &\quad + \eta_t \mathbb{E}[\|\bar{\mathbf{v}}_t - \bar{\mathbf{g}}_t\|^2] + \frac{L^2\eta_t}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2]. \end{aligned} \quad (25)$$

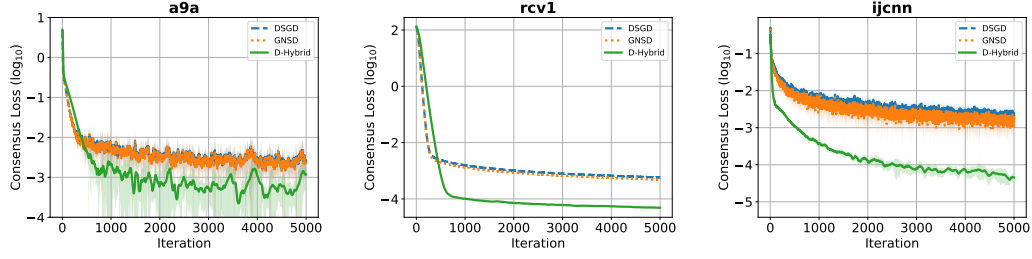
□

### B.3 Proof for Lemma 3

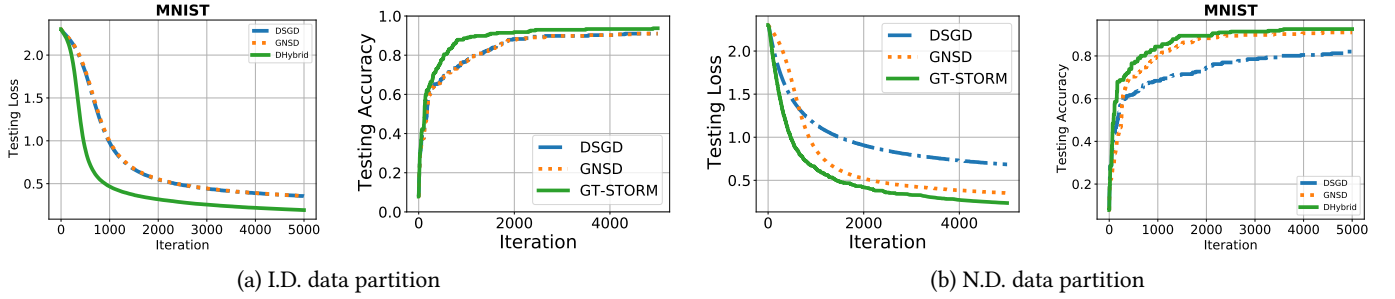
PROOF. First for the iterate  $\mathbf{x}_t$ , we have the following contraction:

$$\|\tilde{\mathbf{W}}\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2 = \|\tilde{\mathbf{W}}(\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t)\|^2 \leq \lambda^2 \|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2, \quad (26)$$

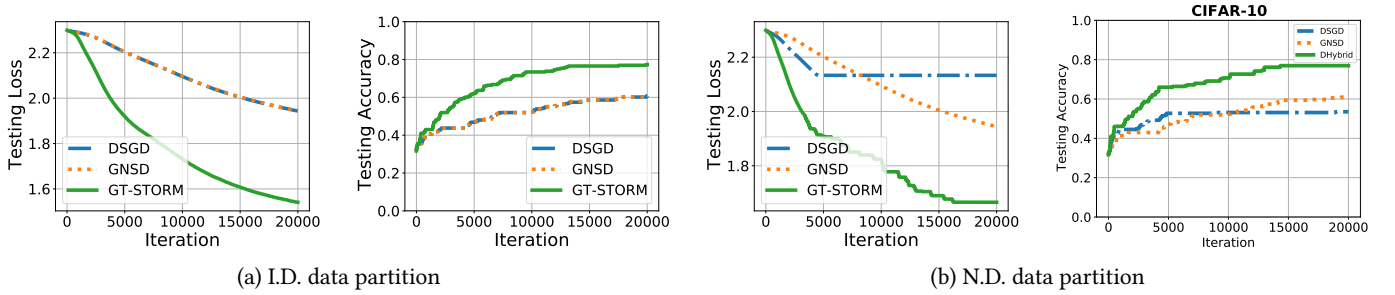
This is because  $\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t$  is orthogonal to  $\mathbf{1}$ , which is the eigenvector corresponding to the largest eigenvalue of  $\tilde{\mathbf{W}}$ , and  $\lambda =$



**Figure 10: Experimental results for nonconvex logistic regression on LibSVM datasets: the three algorithms adopt a diminishing step-size  $\eta_t = \eta_0 / \sqrt[3]{1 + 0.1 \times t}$ ; the curves are averaged over 10 random trials and the shaded regions represent the standard deviation.**



**Figure 11: Experimental results for CNN on MNIST dataset.**



**Figure 12: Experimental results for CNN on CIFAR-10 dataset.**

$\max\{|\lambda_2|, |\lambda_m|\}$ . Recall that  $\bar{\mathbf{x}}_{t+1} = \bar{\mathbf{x}}_t - \eta_t \bar{\mathbf{v}}_t$ , hence,

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{x}}_{t+1}\|^2 &= \|\tilde{\mathbf{W}}\mathbf{x}_t - \eta_t \mathbf{v}_t - \mathbf{1} \otimes (\bar{\mathbf{x}}_t - \eta_t \bar{\mathbf{v}}_t)\|^2 \\ &\leq (1 + c_1) \|\tilde{\mathbf{W}}\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2 + (1 + \frac{1}{c_1}) \eta_t^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2 \\ &\leq (1 + c_1) \lambda^2 \|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2 + (1 + \frac{1}{c_1}) \eta_t^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2. \end{aligned} \quad (27)$$

Similarly to (27), we have:

$$\begin{aligned} &\|\mathbf{v}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{v}}_{t+1}\|^2 \\ &= \|\beta_{t+1}(\tilde{\mathbf{W}}\mathbf{v}_t + \mathbf{w}_{t+1}) + (1 - \beta_{t+1})\mathbf{u}_{t+1} - \mathbf{1} \otimes (\beta_{t+1}(\bar{\mathbf{v}}_t + \bar{\mathbf{w}}_{t+1}) + (1 - \beta_{t+1})\bar{\mathbf{u}}_{t+1})\|^2 \\ &\leq (1 + c_1) \beta_{t+1}^2 \lambda^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2 + (1 + \frac{1}{c_1}) \|\beta_{t+1}(\mathbf{w}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{w}}_{t+1}) + (1 - \beta_{t+1})(\mathbf{u}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{u}}_{t+1})\|^2 \end{aligned}$$

$$\begin{aligned} &\leq (1 + c_1) \beta_{t+1}^2 \lambda^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2 + 2(1 + \frac{1}{c_1}) (\beta_{t+1}^2 \|\mathbf{w}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{w}}_{t+1}\|^2 + (1 - \beta_{t+1})^2 \|\mathbf{u}_{t+1} - \mathbf{1} \otimes \bar{\mathbf{u}}_{t+1}\|^2) \\ &\stackrel{(a)}{\leq} (1 + c_1) \beta_{t+1}^2 \lambda^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2 + 2(1 + \frac{1}{c_1}) (\beta_{t+1}^2 \|\mathbf{w}_{t+1}\|^2 + (1 - \beta_{t+1})^2 \|\mathbf{u}_{t+1}\|^2) \end{aligned} \quad (28)$$

where (a) is due to  $\|\mathbf{I} - \frac{1}{m} \mathbf{1} \mathbf{1}^\top\| \leq 1$ . Lastly, according to the updating equation (9) in main paper, it holds

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 &= \|\tilde{\mathbf{W}}\mathbf{x}_t - \eta_t \mathbf{v}_t - \mathbf{x}_t\|^2 \\ &= \|(\tilde{\mathbf{W}} - \mathbf{I})\mathbf{x}_t - \eta_t \mathbf{v}_t\|^2 \leq 2\|(\tilde{\mathbf{W}} - \mathbf{I})\mathbf{x}_t\|^2 + 2\eta_t^2 \|\mathbf{v}_t\|^2 \\ &= 2\|(\tilde{\mathbf{W}} - \mathbf{I})(\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t)\|^2 + 2\eta_t^2 \|\mathbf{v}_t\|^2 \\ &\stackrel{(a)}{\leq} 8\|(\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t)\|^2 + 4\eta_t^2 \|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2 + 4\eta_t^2 m \|\bar{\mathbf{v}}_t\|^2, \end{aligned} \quad (29)$$

where (a) is because  $\|\tilde{\mathbf{W}} - \mathbf{I}\| \leq 2$ .

□

Recall the result from Lemma 2:

#### B.4 Proof for Lemma 4

PROOF. First, with  $\eta_t = \tau/(\omega + t)^{1/3}$ , we have:

$$\begin{aligned} \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} &= \frac{1}{\tau} ((\omega + t)^{\frac{1}{3}} - (\omega + t - 1)^{\frac{1}{3}}) \\ &\stackrel{(a)}{\leq} \frac{1}{3\tau} \cdot \frac{1}{(\omega + t - 1)^{\frac{2}{3}}} = \frac{1}{3\tau^3} \eta_{t-1}^2 \\ &\stackrel{(b)}{\leq} \frac{1}{3\tau} \cdot \frac{1}{(\frac{\omega}{2} + t)^{\frac{2}{3}}} \leq \frac{1}{3\tau} \cdot \frac{2^{\frac{2}{3}}}{(\omega + t)^{\frac{2}{3}}} \\ &\leq \frac{2^{\frac{2}{3}}}{3\tau^3} \cdot \frac{\tau^2}{(\omega + t)^{\frac{2}{3}}} \leq \frac{2}{3\tau^3} \eta_t, \end{aligned} \quad (30)$$

where (a) is by  $(x + y)^{1/3} - x^{1/3} \leq yx^{-2/3}/3$  and (b) is by  $\omega \geq 2$ .

Then, we give the following three contractions:

i) for  $\mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2]$ , we have:

$$\begin{aligned} &\frac{1}{\eta_t} \mathbb{E}[\|\tilde{\mathbf{v}}_{t+1} - \tilde{\mathbf{g}}_{t+1}\|^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] \\ &\stackrel{(a)}{\leq} \left( \frac{\beta_{t+1}^2}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] + \frac{2\beta_{t+1}^2 L^2}{m\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + \frac{2(1 - \beta_{t+1})^2 \sigma^2}{m\eta_t} \\ &\stackrel{(b)}{\leq} \left( \frac{1 - \rho\eta_t^2}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] + \frac{2\beta_{t+1}^2 L^2}{m\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + \frac{2\rho^2 \sigma^2 \eta_t^3}{m} \\ &\stackrel{(c)}{\leq} \left( \frac{2}{3\tau^3} - \rho \right) \eta_t \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] + \frac{2\beta_{t+1}^2 L^2}{m\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + \frac{2\rho^2 \sigma^2 \eta_t^3}{m} \\ &\stackrel{(d)}{=} -32\eta_t L^2 \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] + \frac{2\beta_{t+1}^2 L^2}{m\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + \frac{2\rho^2 \sigma^2 \eta_t^3}{m} \end{aligned} \quad (31)$$

where (a) is from Lemma 1, (b) is by  $\beta_{t+1} = 1 - \rho\eta_t^2 < 1$ , (c) is by (30) and (d) is by the setting  $\rho = 2/(3\tau^3) + 32L^2$ .

ii) for  $\mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2]$ , we have:

$$\begin{aligned} &\frac{1}{\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{1} \otimes \tilde{\mathbf{x}}_{t+1}\|^2] - \frac{1}{\eta_{t-1}} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] \\ &\stackrel{(a)}{\leq} \left( \frac{(1+c_1)\lambda^2}{\eta_t} - \frac{1}{\eta_{t-1}} \right) [\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] + (1 + \frac{1}{c_1}) \eta_t \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\stackrel{(b)}{\leq} \left( \frac{(1+c_1)\lambda^2 - 1}{\eta_t} + \frac{2}{3\tau^3} \right) \eta_t [\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] + (1 + \frac{1}{c_1}) \eta_t \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \end{aligned} \quad (32)$$

where (a) is by Lemma 3 and (b) is by (30).

iii) for  $\mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2]$ , we have:

$$\begin{aligned} &\mathbb{E}[\|\mathbf{v}_{t+1} - \mathbf{1} \otimes \tilde{\mathbf{v}}_{t+1}\|^2] - \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\stackrel{(a)}{\leq} ((1+c_1)\beta_{t+1}^2 \lambda^2 - 1) \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\quad + 2(1 + \frac{1}{c_1}) (\beta_{t+1}^2 \mathbb{E}[\|\mathbf{w}_{t+1}\|^2] + (1 - \beta_{t+1})^2 \mathbb{E}[\|\mathbf{u}_{t+1}\|^2]) \\ &\stackrel{(b)}{\leq} ((1+c_1)\beta_{t+1}^2 \lambda^2 - 1) \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\quad + 2(1 + \frac{1}{c_1}) (\beta_{t+1}^2 L^2 \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + mG^2 \rho^2 \eta_t^4) \end{aligned} \quad (33)$$

where (a) is by Lemma 3 and (b) is by  $\beta_{t+1} = 1 - \rho\eta_t^2$  and Assumption 1.

$$\begin{aligned} \mathbb{E}[f(\tilde{\mathbf{x}}_{t+1})] - \mathbb{E}[f(\tilde{\mathbf{x}}_t)] &\leq -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\tilde{\mathbf{x}}_t)\|^2] - \left( \frac{\eta_t}{2} - \frac{L\eta_t^2}{2} \right) \mathbb{E}[\|\tilde{\mathbf{v}}_t\|^2] \\ &\quad + \eta_t \mathbb{E}[\|\tilde{\mathbf{v}}_t - \tilde{\mathbf{g}}_t\|^2] + \frac{L^2 \eta_t}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2]. \end{aligned} \quad (34)$$

Then, with the result in i), we have:

$$\begin{aligned} &\mathbb{E}[f(\tilde{\mathbf{x}}_{t+1}) + \frac{1}{32L^2\eta_t} \|\tilde{\mathbf{g}}_{t+1} - \tilde{\mathbf{v}}_{t+1}\|^2] - \mathbb{E}[f(\tilde{\mathbf{x}}_t) + \frac{1}{32L^2\eta_{t-1}} \|\tilde{\mathbf{g}}_t - \tilde{\mathbf{v}}_t\|^2] \\ &\leq -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\tilde{\mathbf{x}}_t)\|^2] - \left( \frac{\eta_t}{2} - \frac{L\eta_t^2}{2} \right) \mathbb{E}[\|\tilde{\mathbf{v}}_t\|^2] + \frac{\beta_{t+1}^2}{16m\eta_t} \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] \\ &\quad + \frac{L^2 \eta_t}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] + \frac{\rho^2 \sigma^2 \eta_t^3}{16mL^2} \end{aligned} \quad (35)$$

Next, with the results in ii) and iii), we have:

$$\begin{aligned} &\mathbb{E}[\frac{1}{\eta_t} \|\mathbf{x}_{t+1} - \mathbf{1} \otimes \tilde{\mathbf{x}}_{t+1}\|^2 + \|\mathbf{v}_{t+1} - \mathbf{1} \otimes \tilde{\mathbf{v}}_{t+1}\|^2] \\ &\quad - \mathbb{E}[\frac{1}{\eta_{t-1}} \|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2 + \|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\leq -\left( \frac{1 - (1+c_1)\lambda^2}{\eta_t} - \frac{2\eta_t}{3\tau^3} \right) \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] \\ &\quad - \left( 1 - (1+c_1)\beta_{t+1}^2 \lambda^2 - (1 + \frac{1}{c_1}) \eta_t \right) \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\quad + 2(1 + \frac{1}{c_1}) \beta_{t+1}^2 L^2 \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] + 2(1 + \frac{1}{c_1}) mG^2 \rho^2 \eta_t^4. \end{aligned} \quad (36)$$

Thus, for the defined potential function:

$$\begin{aligned} H_t &= \mathbb{E}[f(\tilde{\mathbf{x}}_t) + \frac{1}{32L^2\eta_{t-1}} \|\tilde{\mathbf{g}}_t - \tilde{\mathbf{v}}_t\|^2] \\ &\quad + \frac{c_0}{m\eta_{t-1}} \|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2 + \frac{c_0}{m} \|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2, \end{aligned} \quad (37)$$

its differential can be calculated as

$$\begin{aligned} H_{t+1} - H_t &\leq -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\tilde{\mathbf{x}}_t)\|^2] - \left( \frac{\eta_t}{2} - \frac{L\eta_t^2}{2} \right) \mathbb{E}[\|\tilde{\mathbf{v}}_t\|^2] \\ &\quad + \left( \frac{\beta_{t+1}^2}{16m\eta_t} + 2(1 + \frac{1}{c_1}) \frac{c_0 \beta_{t+1}^2 L^2}{m} \right) \mathbb{E}[\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2] \\ &\quad - \left( 1 - (1+c_1)\lambda^2 - \frac{2\eta_t^2}{3\tau^3} - \frac{L^2 \eta_t^2}{c_0} \right) \times \frac{c_0}{m\eta_t} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \tilde{\mathbf{x}}_t\|^2] \\ &\quad - \left( 1 - (1+c_1)\beta_{t+1}^2 \lambda^2 - (1 + \frac{1}{c_1}) \eta_t \right) \times \frac{c_0}{m} \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \tilde{\mathbf{v}}_t\|^2] \\ &\quad + \frac{\rho^2 \sigma^2 \eta_t^3}{16mL^2} + 2(1 + \frac{1}{c_1}) c_0 G^2 \rho^2 \eta_t^4 \\ &\stackrel{(a)}{\leq} -\frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\tilde{\mathbf{x}}_t)\|^2] + \frac{\rho^2 \sigma^2 \eta_t^3}{16mL^2} + 2(1 + \frac{1}{c_1}) c_0 G^2 \rho^2 \eta_t^4 \end{aligned}$$

$$\begin{aligned}
& -\left(1-(1+c_1)\lambda^2 - \frac{1}{2c_0} - 16\left(1+\frac{1}{c_1}\right)L^2\eta_t - \frac{2\eta_t^2}{3\tau^3}\right. \\
& \left. - \frac{L^2\eta_t^2}{c_0}\right) \times \frac{c_0}{m\eta_t} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
& -\left(1-(1+c_1)\lambda^2 - \left(1+\frac{1}{c_1}\right)\eta_t - \frac{\eta_t}{4c_0}\right. \\
& \left. - 8\left(1+\frac{1}{c_1}\right)L^2\eta_t^2\right) \times \frac{c_0}{m} \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2] \\
& -\left(1-2L\eta_t - 32\left(1+\frac{1}{c_1}\right)c_0L^2\eta_t\right) \times \frac{\eta_t}{4} \mathbb{E}[\|\bar{\mathbf{v}}_t\|^2]. \tag{38}
\end{aligned}$$

where (a) follows from plugging the result for  $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2$  from Lemma 3 and  $\beta_{t+1} < 1$ .  $\square$

## B.5 Proof for Theorem 1

PROOF. From Lemma 4, we have:

$$\begin{aligned}
& \sum_{t=0}^T \frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] \leq H_0 - H_{T+1} + \sum_{t=0}^T \frac{\rho^2 \sigma^2 \eta_t^3}{16mL^2} \\
& + \sum_{t=0}^T 2\left(1+\frac{1}{c_1}\right)c_0G^2\rho^2\eta_t^4 - \sum_{t=0}^T \frac{c_0C_1}{m\eta_t} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
& - \sum_{t=0}^T \frac{c_0C_2}{m} \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2] - \sum_{t=0}^T \frac{C_3\eta_t}{4} \mathbb{E}[\|\bar{\mathbf{v}}_t\|^2]. \tag{39}
\end{aligned}$$

Note with  $\eta_t = \tau/(\omega + t)^{1/3}$  and  $\tau \geq 2$ , thus

$$\sum_{t=0}^T \eta_t^3 = \sum_{t=0}^T \frac{\tau}{\omega + t} \leq \int_{-1}^{T-1} \frac{\tau}{\omega + t} dt \leq \tau \ln(\omega + T - 1) \tag{40}$$

$$\sum_{t=0}^T \eta_t^4 = \sum_{t=0}^T \left(\frac{\tau}{\omega + t}\right)^{\frac{4}{3}} \leq \int_{-1}^{T-1} \left(\frac{\tau}{\omega + t}\right)^{\frac{4}{3}} dt \leq \frac{3\tau^{4/3}}{(\omega - 1)^{1/3}}. \tag{41}$$

Hence, since  $\eta_t$  is decreasing, we have:

$$\begin{aligned}
& \frac{\eta_T}{2} \sum_{t=0}^T \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_T - \mathbf{1} \otimes \bar{\mathbf{x}}_T\|^2] \\
& \leq \sum_{t=0}^T \frac{\eta_t}{2} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{\eta_t}{2m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
& \leq H_0 - H_{T+1} + \frac{\tau\rho^2\sigma^2\ln(\omega+T-1)}{16mL^2} \\
& + 6\left(1+\frac{1}{c_1}\right)c_0G^2\rho^2\frac{\tau^{4/3}}{(\omega-1)^{1/3}} - \sum_{t=0}^T \frac{2c_0C_1 - \eta_t^2}{2m\eta_t} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
& - \sum_{t=0}^T \frac{c_0C_2}{m} \mathbb{E}[\|\mathbf{v}_t - \mathbf{1} \otimes \bar{\mathbf{v}}_t\|^2] - \sum_{t=0}^T \frac{C_3\eta_t}{4} \mathbb{E}[\|\bar{\mathbf{v}}_t\|^2]. \tag{42}
\end{aligned}$$

Now, we show that by properly choosing  $\eta_t$ ,  $c_1$ , and  $c_0$ , the coefficients  $C_1 - \eta^2/2c_0$ ,  $C_2$  and  $C_3$  can be non-negative. Recall that:

$$C_1 = 1 - (1+c_1)\lambda^2 - \frac{1}{2c_0} - 16\left(1+\frac{1}{c_1}\right)L^2\eta_t - \left(\frac{2}{3\tau^3} + \frac{L^2}{c_0}\right)\eta_t^2, \tag{43}$$

$$C_2 = 1 - (1+c_1)\lambda^2 - \left(1+\frac{1}{c_1}\right)\eta_t - \frac{\eta_t}{4c_0} - 8\left(1+\frac{1}{c_1}\right)L^2\eta_t^2, \tag{44}$$

$$C_3 = 1 - 2L\eta_t - 32\left(1+\frac{1}{c_1}\right)c_0L^2\eta_t. \tag{45}$$

In order to have  $C_3 \geq 0$ , we have:

$$\eta_t \leq 1/\left(2L + 32\left(1+\frac{1}{c_1}\right)c_0L^2\right) := k_1. \tag{46}$$

With (46), it follows that:

$$C_2 \geq 1 - (1+c_1)\lambda^2 - \left(1+\frac{1}{c_1}\right)\eta_t - \frac{\eta_t}{2c_0}. \tag{47}$$

Thus,  $C_2 \geq 0$  if we set

$$\eta_t \leq \left(1 - (1+c_1)\lambda^2\right)/\left(1+\frac{1}{c_1} + \frac{1}{2c_0}\right) := k_2. \tag{48}$$

For  $C_1 - \eta^2/2c_0$ , it follows from (46) that:

$$C_1 - \frac{\eta^2}{2c_0} \geq 1 - (1+c_1)\lambda^2 - \frac{1}{c_0} - \left(\frac{2}{3\tau^3} + \frac{2L^2+1}{2c_0}\right)\eta_t^2. \tag{49}$$

By choosing

$$\eta_t \leq \sqrt{\left(1 - (1+c_1)\lambda^2 - \frac{1}{c_0}\right)/\left(\frac{2}{3\tau^3} + \frac{2L^2+1}{2c_0}\right)} := k_3, \tag{50}$$

$$\text{and } 0 < 1 - (1+c_1)\lambda^2 - \frac{3}{4c_0}, \tag{51}$$

we have  $C_1 - \eta^2/2c_0 \geq 0$ . To summarize, we need to set  $\eta_t \leq \min\{k_1, k_2, k_3\}$ . Since  $\eta_t$  is decreasing and  $\eta_0 = \tau/\omega^{1/3}$ , it implies that  $\omega \geq (\tau/\min\{k_1, k_2, k_3\})^3$ .

With the above parameter setting, we have:

$$\begin{aligned}
& \frac{\eta_T}{2} \sum_{t=0}^T \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_T - \mathbf{1} \otimes \bar{\mathbf{x}}_T\|^2] \leq H_0 - H_{T+1} \\
& + \frac{\tau\rho^2\sigma^2\ln(\omega+T-1)}{16mL^2} + 6\left(1+\frac{1}{c_1}\right)c_0G^2\rho^2\frac{\tau^{4/3}}{(\omega-1)^{1/3}}. \tag{52}
\end{aligned}$$

Multiplying both side of the above inequality by  $2/\eta_T(T+1)$ , we have:

$$\begin{aligned}
& \frac{1}{T+1} \sum_{t=0}^T \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_T - \mathbf{1} \otimes \bar{\mathbf{x}}_T\|^2] \\
& \leq \frac{2(H_0 - H_{T+1})}{\eta_T(T+1)} + \frac{\tau\rho^2\sigma^2\ln(\omega+T-1)}{8mL^2\eta_T(T+1)} + \frac{12\left(1+\frac{1}{c_1}\right)c_0\tau^{4/3}G^2\rho^2}{(\omega-1)^{1/3}\eta_T(T+1)}. \tag{53}
\end{aligned}$$

Note that

$$\begin{aligned}
H_0 &= \mathbb{E}[f(\bar{\mathbf{x}}_0) + \frac{1}{32L^2\eta_{-1}} \|\bar{\mathbf{g}}_0 - \bar{\mathbf{v}}_0\|^2 \\
&\quad + \frac{c_0}{m\eta_{-1}} \|\mathbf{x}_0 - \mathbf{1} \otimes \bar{\mathbf{x}}_0\|^2 + \frac{c_0}{m} \|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2] \\
&\stackrel{(a)}{=} \mathbb{E}[f(\bar{\mathbf{x}}_0) + \frac{1}{32L^2\eta_{-1}} \|\bar{\mathbf{g}}_0 - \bar{\mathbf{v}}_0\|^2 + \frac{c_0}{m} \|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2] \\
&\stackrel{(b)}{\leq} \mathbb{E}[f(\bar{\mathbf{x}}_0) + \frac{\sigma^2}{32mL^2\eta_{-1}} + \frac{c_0}{m} \|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2] \\
H_{T+1} &\geq \mathbb{E}[f(\bar{\mathbf{x}}_{T+1}) + \frac{c_0}{m\eta_T} \|\mathbf{x}_{T+1} - \mathbf{1} \otimes \bar{\mathbf{x}}_{T+1}\|^2 + \frac{c_0}{m} \|\mathbf{v}_{T+1} - \mathbf{1} \otimes \bar{\mathbf{v}}_{T+1}\|^2] \\
&\geq f(\bar{\mathbf{x}}^*),
\end{aligned}$$

where (a) is by  $\mathbf{x}_{i,0} = \mathbf{x}_0$  from line 1 in Algorithm 1 and (b) is by Assumption 1.

Hence, it follows that

$$\begin{aligned}
&\frac{1}{T+1} \sum_{t=0}^T \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
&\leq \frac{2(f(\bar{\mathbf{x}}_0) - f(\bar{\mathbf{x}}^*))}{\eta_T(T+1)} + \frac{2c_0\mathbb{E}[\|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2]}{m\eta_T(T+1)} \\
&\quad + \frac{(\omega-1)\sigma^2}{16mL^2\tau\eta_T(T+1)} + \frac{\tau\rho^2\sigma^2\ln(\omega+T-1)}{8mL^2\eta_T(T+1)} \\
&\quad + 12(1 + \frac{1}{c_1})c_0 \frac{\tau^{4/3}G^2\rho^2}{(\omega-1)^{1/3}\eta_T(T+1)}. \tag{54}
\end{aligned}$$

Since  $\eta_T = \tau/(\omega+T)^{1/3}$ , we have:

$$\begin{aligned}
&\min_{t \in [T]} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
&\leq \frac{2(f(\bar{\mathbf{x}}_0) - f(\bar{\mathbf{x}}^*))}{\tau(T+1)^{2/3}} + \frac{2c_0\mathbb{E}[\|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2]}{m\tau(T+1)^{2/3}} \\
&\quad + \frac{(\omega-1)\sigma^2}{16mL^2\tau^2(T+1)^{2/3}} + \frac{\rho^2\sigma^2\ln(\omega+T-1)}{8mL^2(T+1)^{2/3}} \\
&\quad + \frac{12(1 + \frac{1}{c_1})c_0\tau^{1/3}G^2\rho^2}{(\omega-1)^{1/3}(T+1)^{2/3}} + O\left(\frac{c_3\omega}{\tau T^{5/3}}\right). \tag{55}
\end{aligned}$$

where the  $O$ -notation is from  $(\omega+T)^{1/3} - (T+1)^{1/3} \leq (\omega-1)(T+1)^{-2/3}/3$  and  $c_3 = \max\{1, (\omega-1)/(m\tau^2), \tau^{4/3}/\omega^{1/3}, \tau\ln(\omega+T-1)/m\}$ .

□

## B.6 Proof for Corollary 2

PROOF. First, note that  $\omega \geq \max\{2, \tau^3/\min\{k_1^3, k_2^3, k_3^3\}\}$  holds with  $\tau = O(m^{1/3})$  and  $\omega = O(m^{4/3})$ . Plugging these parameters into Theorem 1 yields:

$$\begin{aligned}
&\min_{t \in [T]} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
&\leq O\left(\frac{2(f(\bar{\mathbf{x}}_0) - f(\bar{\mathbf{x}}^*))}{m^{1/3}(T+1)^{2/3}} + \frac{2c_0\mathbb{E}[\|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2]}{m^{4/3}(T+1)^{2/3}} + \frac{\sigma^2}{16L^2m^{1/3}(T+1)^{2/3}}\right. \\
&\quad \left. + \frac{\rho^2\sigma^2\ln(m^{4/3}+T)}{8mL^2(T+1)^{2/3}} + \frac{12(1 + \frac{1}{c_1})c_0G^2\rho^2}{m^{1/3}(T+1)^{2/3}} + \frac{c_3m}{T^{5/3}}\right). \tag{56}
\end{aligned}$$

With  $T \gg m^{4/3}$ , we have:

$$\begin{aligned}
&\min_{t \in [T]} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}_t)\|^2] + \frac{1}{m} \mathbb{E}[\|\mathbf{x}_t - \mathbf{1} \otimes \bar{\mathbf{x}}_t\|^2] \\
&\leq O\left(\frac{2(f(\bar{\mathbf{x}}_0) - f(\bar{\mathbf{x}}^*))}{m^{1/3}(T+1)^{2/3}} + \frac{2c_0\mathbb{E}[\|\mathbf{v}_0 - \mathbf{1} \otimes \bar{\mathbf{v}}_0\|^2]}{m^{4/3}(T+1)^{2/3}}\right. \\
&\quad \left. + \frac{\sigma^2}{16L^2m^{1/3}(T+1)^{2/3}} + \frac{\rho^2\sigma^2\ln T}{8mL^2(T+1)^{2/3}}\right. \\
&\quad \left. + \frac{12(1 + \frac{1}{c_1})c_0G^2\rho^2}{m^{1/3}(T+1)^{2/3}} + \frac{c_3}{m^{1/3}T^{2/3}}\right), \tag{57}
\end{aligned}$$

where  $c_3 = \max\{O(1), O(1/m^{1/3}), \ln(m^{4/3} + T)/m^{2/3}\}$ . The above result implies that the convergence rate is  $\tilde{O}(m^{-1/3}T^{-2/3})$ . Thus, to achieve an  $\epsilon^2$ -stationary solution, the total communication rounds needed are  $T = \tilde{O}(m^{-1/2}\epsilon^{-3})$ , and the total samples needed are  $mT = \tilde{O}(m^{1/2}\epsilon^{-3})$ . □