ELIMINATING STEADY-STATE OSCILLATIONS IN DISTRIBUTED OPTIMIZATION AND LEARNING VIA ADAPTIVE STEPSIZE

Anonymous authors

000

001

002

004

006

008

009

010 011 012

013

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032 033 034

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Distributed stochastic optimization and learning is gaining increasing traction due to its ability to enable large-scale data processing and model training across multiple agents without the need for centralized coordination. However, existing distributed stochastic optimization and learning approaches, such as distributed SGD and their variants, generally face a dilemma in stepsize selection: a small stepsize leads to low convergence speed, whereas a large stepsize often incurs pronounced steady-state oscillations, which prevents the algorithm from achieving stable convergence accuracy. In this paper, we propose an adaptive stepsize approach for distributed stochastic optimization and learning that can eliminate steady-state oscillations and ensure fast convergence. Such guarantees are unattained by existing adaptive stepsize approaches, even in centralized optimization and learning. We prove that our proposed algorithm achieves linear convergence with respect to the iteration number, and that the convergence error decays sublinearly with the batch size of sampled data points. In the specific case in terms of deterministic distributed optimization with exact gradients accessible to agents, we prove that our proposed algorithm linearly converges to an exact optimal solution. Moreover, we quantify that the computational complexity of the proposed algorithm is on the order of $\mathcal{O}(\log(\epsilon^{-1}))$, which matches the existing results on adaptive stepsize approaches for centralized optimization/learning. Experimental results on machine learning benchmarks confirm the effectiveness of our proposed approach.

1 Introduction

With the advance of the era of big data, distributed stochastic optimization and learning methods have attracted increasing attention due to their unique ability to leverage the computational power of multiple devices to accelerate training (Nedic & Ozdaglar, 2009; Yang & Johansson, 2010; Shamir & Srebro, 2014; Lian et al., 2017; Nedić & Liu, 2018; Yang et al., 2019; Kim et al., 2024; Hu et al., 2024). Unlike centralized optimization and learning methods (Wang & Elia, 2011; Andrychowicz et al., 2016; Ruder, 2016) that typically rely on a central server to aggregate local model parameters or data from all participating agents, distributed methods allow each agent to collaboratively learn a global model using only its own local dataset and information exchanged with neighboring agents, without the assistance of any centralized server or aggregator (Scaman et al., 2018; Liu et al., 2020; Yang et al., 2022).

However, existing distributed stochastic optimization/learning approaches often face a dilemma in stepsize selection (Jacobs, 1988; Schaul et al., 2013; Wei et al., 2020; Zhuang et al., 2020; Li et al., 2024a; Huang et al., 2024b; Crawshaw et al., 2025). Specifically, an excessively small stepsize may lead to an overly low convergence speed (Srivastava & Nedic, 2011; Lin et al., 2023; Sharifnassab et al., 2024), whereas an excessively large stepsize often causes pronounced steady-state oscillations or overshoot, which prevents the algorithm from achieving stable convergence accuracy (Andriushchenko et al., 2023; Huang et al., 2024a). Recently, several adaptive or automatic stepsize approaches have been proposed for centralized optimization and learning (Fletcher, 2005; Kingma, 2014; Rolinek & Martius, 2018; Li & Orabona, 2019; Malitsky & Mishchenko, 2019;

Kavis et al., 2022; Jiang & Stich, 2023; Malitsky & Mishchenko, 2024). However, these approaches generally rely on a centralized server to coordinate computation that are impractical in a fully distributed setting where no centralized server/aggregator exists to determine a common stepsize across all agents (Nedić et al., 2018). Although some works have attempted to extend adaptive stepsize approaches to distributed optimization and learning (Nazari et al., 2022; Carnevale et al., 2022; Xie et al., 2022; Ramezani-Kebrya et al., 2023; Chen & Wang, 2024; Kuruzov et al., 2024; Saravanos et al., 2024), most of them still either require a centralized server to collect local model parameters/stepsizes from all agents (Ramezani-Kebrya et al., 2023; Chen & Wang, 2024; Kuruzov et al., 2024), or are limited to scenarios where agents must have access to accurate gradients of the objective functions (Carnevale et al., 2022; Xie et al., 2022; Saravanos et al., 2024) for stepsize adjustment. The only exception is the work in Nazari et al. (2022), which achieves adaptive stepsize adjustments in distributed online learning by normalizing the gradient using an accumulated sum of historical gradient values. However, this approach leads to a rapidly decaying stepsize, which in turn results in slow convergence in the later stages of the algorithm (see our experimental results in Fig. 5 in Appendix C.3 for details). To the best of our knowledge, no existing adaptive stepsize approaches can ensure fast and stable convergence in fully distributed stochastic optimization/learning.

Our contributions are summarized as follows:

- 1. We propose an adaptive stepsize algorithm for fully distributed stochastic optimization and learning. This is in stark contrast to existing adaptive stepsize approaches, which either rely on a centralized server to coordinate a common stepsize across all agents (in, e.g., Ramezani-Kebrya et al. (2023); Kim et al. (2024); Chen & Wang (2024); Kuruzov et al. (2024)), or require that agents have access to accurate gradients of the objective functions (Carnevale et al., 2022; Xie et al., 2022; Saravanos et al., 2024)—which, however, are often hard to obtain in real-world applications where the randomness in sampled data results in only noisy gradients being accessible to agents. To the best of our knowledge, this is the first adaptive (non-monotone decreasing) stepsize approach for fully distributed stochastic optimization/learning, without the need for accurate gradients or the assistance of any centralized servers.
- 2. Our adaptive stepsize algorithm can eliminate steady-state oscillations and ensure stable convergence accuracy in the later stages of the algorithm. This is unattained by most existing adaptive stepsize approaches even in centralized optimization and learning (Fletcher, 2005; Li & Orabona, 2019; Kavis et al., 2022; Jiang & Stich, 2023). The key enabler is our novel design of the stepsize update rule, which allows each agent to dynamically adjust its individual stepsizes based on locally estimated curvature of the global objective function. This provides each agent with large stepsizes in the early stages to accelerate convergence, and extremely small stepsizes near the global optimum to ensure stable convergence accuracy (see our experimental results in Figs. 1(d)-1(f) and Figs. 2(d)-2(f) for details). Furthermore, since stable convergence accuracy is achieved in the later stages of our algorithm, we can also provide a clear stopping criterion¹ for each agent in distributed optimization and learning, which is rarely addressed in the state-of-the-art literature.
- 3. In addition to eliminating steady-state oscillations, we also establish the convergence rate and computational complexity of our algorithm for both stochastic and deterministic distributed optimization and learning, which is different from existing adaptive stepsize results in, e.g., McMahan & Streeter (2014); Yang et al. (2019); Crawshaw et al. (2025) that focus solely on deterministic cases where accurate gradients of objective functions are accessible to agents. For distributed stochastic optimization/learning, we prove that our algorithm achieves linear convergence with respect to the number of algorithm iterations, and that the convergence error decays sublinearly with the batch size of sampled data points. For the deterministic case, we prove that our algorithm linearly converges to an exact optimal solution.
- 4. We systematically quantify that the computational complexity of our algorithm is on the order of $\mathcal{O}(\log(\epsilon^{-1}))$ for both stochastic and deterministic cases, which matches the existing results on adaptive stepsize approaches for centralized optimization and learning in, e.g., (Kavis et al., 2022; Yang & Ma, 2023).
- 5. We conduct experimental evaluations using several machine learning benchmark datasets, including the "MNIST" dataset, the "CIFAR-10" dataset, and the "CIFAR-100" dataset. The

¹We use the "stopping criterion" to denote the condition that determines when each agent in a distributed stochastic optimization and learning algorithm terminates its iterations (Vlachos, 2008; Ding et al., 2025).

results confirm the effectiveness of our algorithm in terms of both test accuracy and steadystate convergence performance.

2 RELATED WORK

Distributed stochastic optimization and learning. Distributed stochastic optimization methods have been widely employed in modern machine learning (Yang, 2013; Xin et al., 2020; Nedic, 2020; Guo et al., 2020; Pu et al., 2020; Allen-Zhu et al., 2020; Khaled & Jin, 2023; Song et al., 2025). However, most existing methods require all agents to share a common stepsize that is either fixed (Pu & Nedić, 2021; Koloskova et al., 2021; Nguyen et al., 2023; Song et al., 2024) or diminishing (Jakovetic et al., 2018; Dieuleveut & Patel, 2019; Li et al., 2024b; Lee et al., 2025). The fixed stepsize causes pronounced overshoot or oscillations near the global optimal solution (Pu & Nedić, 2021; Koloskova et al., 2021; Nguyen et al., 2023), whereas diminishing stepsizes often lead to an overly low convergence speed, both of which prevent the algorithm from achieving stable convergence accuracy (as shown in our experimental results in Fig. 1 and Fig. 2). Given these limitations, designing an adaptive stepsize approach that allows each participating agent to adaptively adjust its individual stepsizes is a promising direction for improving convergence speed and ensuring stable learning performance in distributed stochastic optimization and learning.

Adaptive stepsize approaches. Several adaptive stepsize approaches have been proposed for centralized optimization and learning (Fletcher, 2005; Kingma, 2014; Rolinek & Martius, 2018; Li & Orabona, 2019; Malitsky & Mishchenko, 2019; Kavis et al., 2022; Jiang & Stich, 2023; Malitsky & Mishchenko, 2024). However, these methods typically consider a single agent setting where learning is performed with only one adaptive stepsize adjustment. This makes them inapplicable to fully distributed stochastic optimization and learning, where multiple agents cooperatively perform learning and each agent has its own stepsize updates. Moreover, the existing adaptive stepsize approaches often lead to steady-state oscillations, which prevent stable convergence accuracy in the later stages of the algorithm and hinder the determination of a clear stopping criterion (as shown in our experimental results in Fig. 2). Although some efforts have attempted to extend adaptive stepsize approaches to distributed optimization and learning (Nazari et al., 2022; Carnevale et al., 2022; Xie et al., 2022; Ramezani-Kebrya et al., 2023; Chen & Wang, 2024; Kuruzov et al., 2024; Saravanos et al., 2024), most of them still rely on a centralized server to collect local model parameters/stepsizes from all agents to coordinate a stepsize (Ramezani-Kebrya et al., 2023; Chen & Wang, 2024; Kuruzov et al., 2024), or are limited to scenarios where accurate gradients of the objective functions must be accessible to agents (Carnevale et al., 2022; Xie et al., 2022; Saravanos et al., 2024), both of which are impractical in a fully distributed and stochastic setting. The only exception is the recent work in Nazari et al. (2022), which achieves stepsize adjustments in distributed stochastic optimization and learning. However, its approach parallels adaptive gradient methods (e.g., ADAM in Kingma (2014)), which makes the stepsizes decay rapidly in practical neural-network training, thereby leading to a low convergence speed in the later stages of the algorithm (as shown in our experimental results in Fig. 5 in Appendix C.3). To the best of our knowledge, no adaptive stepsize approaches have been reported for distributed stochastic optimization and learning that can ensure both fast convergence and stable steady-state performance.

Notations: We use \mathbb{R}^n to denote the n-dimensional real Euclidean space and $\mathbb{N}(\mathbb{N}^+)$ to denote the set of nonnegative (positive) integers. We write $\mathbf{0}_n$ and $\mathbf{1}_n$ for n-dimensional column vectors of all zeros and all ones, respectively; in both cases we suppress the dimension when clear from the context. We use $\langle x,y\rangle$ to denote the inner product of two vectors and $\|\cdot\|$ to denote the Euclidean norm of a vector. We write $\mathbb{E}[x]$ for the expected value of a random variable x. We use $[a]_+ = \max\{0,a\}$ to refer to the maximum of 0 and a for any real number a and the convention $\frac{a}{0} = +\infty$ for any a > 0. We denote the set of m agents as [m] and add an overbar to a letter to represent the average of m agents, e.g., $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$.

3 PROBLEM FORMULATION

We consider m agents that cooperatively learn a common optimal model parameter x^* to the following stochastic optimization problem (Sundhar Ram et al., 2010; Lian et al., 2017; Chen & Wang,

2024):

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x), \quad f_i(x) = \mathbb{E}_{\xi_i \sim \mathcal{P}_i}[l(x, \xi_i)]. \tag{1}$$

Here, the local objective function $f_i(x) : \mathbb{R}^n \to \mathbb{R}$ represents the mathematical expectation of agent i's loss function $l(x, \xi_i)$, where ξ_i denotes the agent i's data sample drawn from distribution \mathcal{P}_i .

In real-world applications, since the data distribution \mathcal{P}_i is typically unknown to each agent i, it can only have access to a noisy estimate on the gradient of $f_i(x)$ (Pu & Nedić, 2021; Nazari et al., 2022; Kim et al., 2024). In other words, at each iteration t, each agent i independently and identically samples $|\mathcal{B}|$ data points (also called a batch size of $|\mathcal{B}|$) from its local distribution \mathcal{P}_i and computes a noisy gradient estimate $g_i^t(x) = \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \nabla l(x, \xi_{ij}^t)$, where ξ_{ij}^t is the jth sampled data collected by agent i at iteration t. Based on the gradient estimate $g_i^t(x)$ and communication with its neighbors, each agent i performs distributed training. We make the following standard assumption about $f_i(x)$ and $g_i^t(x)$:

Assumption 1. For any agent $i \in [m]$, its local objective function $f_i(x)$ is μ -strongly convex and L-smooth. The gradient estimate $g_i^t(x)$ is unbiased with bounded variance σ^2 , i.e., $\mathbb{E}[g_i^t(x)] = \nabla f_i(x)$ and $\mathbb{E}[\|g_i^t(x) - \nabla f_i(x)\|^2] \leq \sigma^2$ hold for any $x \in \mathbb{R}^n$ and $t \geq 0$.

In Assumption 1, the strong convexity of $f_i(x)$ is used to ensure linear convergence, which is commonly used in the existing literature (Ivkin et al., 2019; Hou et al., 2021; Akhavan et al., 2021; Wang et al., 2023; Yang & Ma, 2023; He et al., 2024; Er et al., 2024).

We describe the local interaction among agents using a weight matrix $W = \{w_{ij}\} \in \mathbb{R}^{m \times m}$, where $w_{ij} > 0$ if agent i and agent j can directly communicate with each other, and $w_{ij} = 0$ otherwise. The neighboring set of agent i is defined as $\mathcal{N}_i = \{j \in [m] | w_{ij} > 0\}$, which includes itself. Furthermore, we make the following assumption on matrix W:

Assumption 2. The matrix $W \in \mathbb{R}^{m \times m}$ is symmetric and satisfies $\mathbf{1}_m^\top W = \mathbf{1}_m^\top$, $W \mathbf{1}_m = \mathbf{1}_m$, and $\rho \triangleq \|W - \frac{\mathbf{1}_m \mathbf{1}_m^\top}{m}\| < 1$.

Existing distributed optimization and learning approaches typically require the stepsize to be either fixed (Pu & Nedić, 2021; Koloskova et al., 2021; Nguyen et al., 2023; Song et al., 2024) or diminishing (Jakovetic et al., 2018; Dieuleveut & Patel, 2019; Li et al., 2024b; Lee et al., 2025). However, the use of a fixed stepsize often suffers from error/bias terms proportional to the stepsize (Yuan et al., 2016), which can cause pronounced overshoot or persistent oscillations near the global optimum, thereby compromising convergence stability in the later stages of the algorithm (as shown in our experimental results in Fig. 1). Although employing a diminishing stepsize can asymptotically eliminate such errors and ensure stable steady-state convergence, it often results in an undesirably low convergence speed, which is problematic for applications requiring fast convergence (Nedic & Ozdaglar, 2009; Jakovetic et al., 2018; Dieuleveut & Patel, 2019; Lee et al., 2025). Given these limitations, we aim to develop an adaptive stepsize approach for distributed stochastic optimization and learning, enabling each agent to adaptively adjust its stepsize during algorithm iterations to achieve both fast convergence and stable steady-state performance.

4 ALGORITHM DESIGN

In this section, we propose an adaptive stepsize approach for distributed stochastic optimization and learning that ensures both fast convergence and stable steady-state performance. The proposed approach is summarized in Algorithm 1, which is implemented in a fully distributed manner.

In Algorithm 1, Lines 3-7 execute a consensus-based gradient descent step for agent i. Lines 8, 11, and 14 update $y_{i,1}^{t+1}$ to track $\frac{1}{m}\sum_{i=1}^m g_i^t(x_i^{t+1})$, which serves to approximate the global gradient $\frac{1}{m}\sum_{i=1}^m \nabla f_i(x_i^{t+1})$. Lines 9, 12, and 14 update an auxiliary variable $y_{i,2}^t$ to track $\frac{1}{m}\sum_{i=1}^m g_i^t(x_i^t)$, which serves to approximate the global gradient $\frac{1}{m}\sum_{i=1}^m \nabla f_i(x_i^t)$. With this understanding, we let each agent i locally estimate the curvature of the global objective function in Line 15. Based on this estimate, each agent i's adaptive stepsize update rule is given in Line 15 and Line 16.

The key enabler for us to ensure stable steady-state convergence is our meticulously designed stepsize update rule. More specifically, our stepsize update rule enables each agent to locally estimate

Algorithm 1 Adaptive stepsize design for distributed stochastic optimization and learning (from agent i's perspective)

```
1: Input: x_i^0 \in \mathbb{R}^n, y_{i,1}^0 = g_i^{-1}(x_i^0) = g_i^0(x_i^0), y_{i,2}^{-1} = g_i^{-1}(x_i^{-1}) = \mathbf{0}_n, \eta_i^0 > 0, \beta \in (0, 1.36), r \in (0, 1), M \in \mathbb{N}^+, \text{ and } T \in \mathbb{N}^+.

2: for t = 0, 1, \dots, T do

3: x_i^{t+1}(0) = x_i^t - \eta_i^t y_{i,1}^t

4: for q = 0, 1, \dots, M - 1 do

5: x_i^{t+1}(q+1) = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{t+1}(q)

6: end for

7: x_i^{t+1} = x_i^{t+1}(M)

8: y_{i,1}^{t+1}(0) = y_{i,1}^t + g_i^t(x_i^{t+1}) - g_i^{t-1}(x_i^t)

9: y_{i,2}^t(0) = y_{i,2}^{t-1} + g_i^t(x_i^t) - g_i^{t-1}(x_i^{t-1})

10: for q = 0, 1, \dots, M - 1 do

11: y_{i,1}^{t+1}(q+1) = \sum_{j \in \mathcal{N}_i} w_{ij} y_{j,1}^{t+1}(q)

12: y_{i,2}^t(q+1) = \sum_{j \in \mathcal{N}_i} w_{ij} y_{j,2}^t(q)

13: end for

14: y_{i,1}^{t+1} = y_{i,1}^{t+1}(M) and y_{i,2}^t = y_{i,2}^t(M)

15: L_i^{t+1} = \frac{\|y_{i,1}^{t+1} - y_{i,2}^t\|}{\|x_i^{t+1} - x_i^t\|}

16: \eta_i^{t+1} = \min \left\{ \beta \eta_i^t, \frac{7\sqrt{r}}{10} \frac{\eta_i^t}{\sqrt{[(\eta_i^t L_i^{t+1})^2 - 1]_+}} \right\}

17: end for
```

the curvature of the global objective function. In this way, each agent's stepsize can be adapted to large values in the early stages of the algorithm, and to extremely small values near the global optimum (as shown in our experimental results in Figs. 1(d)-1(f) and Figs. 2(d)-2(f)). Therefore, our design avoids the slow convergence caused by small diminishing stepsizes used in, e.g., Jakovetic et al. (2018); Dieuleveut & Patel (2019); Li et al. (2024b); Lee et al. (2025) and eliminate the oscillations arising from fixed stepsizes in e.g., Pu & Nedić (2021); Koloskova et al. (2021); Nguyen et al. (2023); Song et al. (2024).

It is worth noting that our algorithm is fundamentally different from existing adaptive stepsize methods in e.g., Malitsky & Mishchenko (2019; 2024); Kim et al. (2024); Chen & Wang (2024), which explicitly require a centralized server to coordinate stepsize adjustment, which is infeasible in fully distributed settings in the absence of a centralized server. Furthermore, our design is also different from existing adaptive stepsize approaches for deterministic distributed optimization/learning in Carnevale et al. (2022); Xie et al. (2022); Saravanos et al. (2024), which typically require access to exact gradients of objective functions—such exact gradients are often unattainable in real-world applications where only noisy gradient estimates are available to each agent (Lian et al., 2017).

In Algorithm 1, we provide *optional* inner-consensus-loop iterations for $x_i^t, y_{i,1}^t$, and $y_{i,2}^t$. This design is intended to accelerate consensus among agents and improve the accuracy of global gradient tracking, thereby guaranteeing linear convergence (see Theorem 1 for details). In practical machine learning applications, the number of inner-consensus-loop iterations M can be chosen as any positive integer. For example, we can simply select M=1 (in which case Algorithm 1 reduces to a single-loop algorithm) to minimize the computational and communication costs of our algorithm. In fact, our experimental results in Fig. 3(c) show that the test accuracy of our algorithm remains comparable even with M=1.

5 Convergence results

In this section, we prove that Algorithm 1 can ensure linear convergence with respect to the number of iterations T, and the convergence error decreases sublinearly with the batch size of sampled data. The results are summarized in Theorem 1, whose proof can be found in Appendix B.2.

Theorem 1. Under Assumptions 1 and 2, for any $T \ge 0$ and batch size $|\mathcal{B}| > 0$, if the number of inner-consensus-loop iterations M satisfies $M \ge M_0$ with detailed forms of M_0 given in equa-

tion 70 in Appendix B.2, the iterates x_i^t generated by Algorithm 1 satisfy

$$\mathbb{E}\left[\left\|x_i^T - x^*\right\|^2\right] \le \mathcal{O}\left(\gamma^T\right) + \mathcal{O}\left(\frac{\sigma^2}{|\mathcal{B}|}\right),\tag{2}$$

where the convergence rate γ is given by $\gamma = \max\Big\{1 - \frac{\mu^2}{4L}, \, \frac{91}{92}\Big\}.$

Theorem 1 proves that Algorithm 1 linearly converges to an optimal solution to problem 1 with the optimization error decreasing as the batch size of sampled data $|\mathcal{B}|$ increases. It is worth noting that the bound $\mathcal{O}\left(\frac{\sigma^2}{|\mathcal{B}|}\right)$ in Theorem 1, caused by finite batch size of sampled data, inherently exists in all stochastic optimization approaches with finite samples (Yuan et al., 2022; Sharma et al., 2023). Although variance reduction techniques (Reddi et al., 2016; Fang et al., 2018) and diminishing stepsize methods (Nedic & Ozdaglar, 2009) can be used to mitigate the influence of this term in distributed stochastic optimization and learning, their successful implementation heavily relies on the assumption of a fixed upper bound on the stepsizes, which is hard to satisfy when each agent's stepsize is dynamic and adaptive over iterations.

In Theorem 1, we consider a stochastic scenario in which each agent can only access to noisy gradient estimates (which are computed based on data sampled from an unknown data distribution \mathcal{P}_i). Next, we consider a deterministic scenario in which each agent can access to accurate gradients. The convergence result of Algorithm 1 in the deterministic scenario is summarized in the following Theorem 2, whose proof is given in Appendix B.3.

Theorem 2. Under Assumptions 1 and 2, for any $T \ge 0$, if the number of inner-consensus-loop iterations M satisfies $M \ge M_0$ with detailed forms of M_0 given in equation 70 of Appendix B.2, the iterates x_i^t generated by Algorithm 1 with deterministic gradients satisfy

$$\mathbb{E}\left[\left\|x_i^T - x^*\right\|^2\right] \le \mathcal{O}\left(\gamma^T\right),\tag{3}$$

where the convergence rate γ is given by $\gamma = \max \left\{ 1 - \frac{\mu^2}{4L}, \frac{91}{92} \right\}$.

Theorem 2 proves that when we consider distributed optimization and learning in a deterministic scenario, Algorithm 1 converges to an exact solution to problem in equation 1 with a linear convergence rate, which matches existing convergence results on adaptive stepsizes for centralized optimization and learning (Li & Orabona, 2019; Malitsky & Mishchenko, 2019; Kavis et al., 2022; Malitsky & Mishchenko, 2024). Moreover, this is also stronger than the convergence results achieved by existing distributed optimization methods with diminishing stepsizes (Jakovetic et al., 2018; Dieuleveut & Patel, 2019; Li et al., 2024b; Lee et al., 2025), which guarantee only sublinear convergence rates.

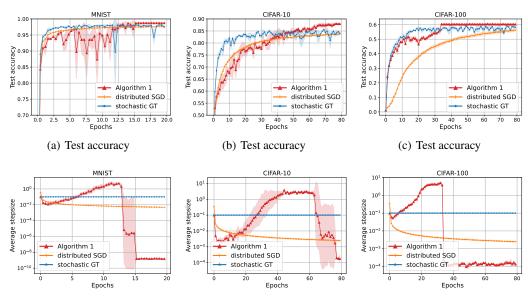
Furthermore, to give a more intuitive description of the computational complexity, we define an ϵ -solution to problem in equation 1 as follows.

Definition 1 (Lian et al. (2017)). For some integer T > 0, if $\mathbb{E}[\|x_i^T - x^*\|^2] \le \epsilon$ holds, then we say that the sequence $\{x_i^t\}$ can reach an ϵ -solution to the problem in equation 1.

Building on Theorem 1 and Theorem 2, we have the following corollary.

Corollary 1. Under Assumptions 1 and 2, for any $\epsilon > 0$, Algorithm 1 with noisy gradient estimates requires at most $\mathcal{O}((2|\mathcal{B}|+3M+3)\log(\epsilon^{-1}))$ gradient evaluation to obtain an ϵ -solution, and Algorithm 1 with accurate gradients requires at most $\mathcal{O}((2M+3)\log(\epsilon^{-1}))$ gradient evaluation to obtain an ϵ -solution.

In Corollary 1, the low bound on the number of inner-consensus-loop iterations M in Algorithm 1 is a fixed constant, which is different from the existing distributed optimization results in, e.g., Berahas et al. (2019); Li et al. (2020) which have the inner-loop iteration number increasing with the outer-loop iterations, and hence have a higher computational complexity of the order of $\mathcal{O}((\log(\epsilon^{-1}))^2)$. Moreover, the computational complexity of our Algorithm 1 matches the adaptive stepsize results on centralized learning in, e.g., Malitsky & Mishchenko (2019; 2024) and the convergence results on distributed optimization in, e.g., Chen & Wang (2024); Kuruzov et al. (2024). This is also less than the convergence results in, e.g., Jakovetic et al. (2018); Dieuleveut & Patel (2019); Li et al. (2024b) with diminishing stepsizes which have a computation complexity of the order of $\mathcal{O}(\epsilon^{-1})$.



(d) Average stepsize across agents (e) Average stepsize across agents (f) Average stepsize across agents

Figure 1: Test-accuracy and average-stepsize (across five agents) evolutions of Algorithm 1, distributed SGD in Jakovetic et al. (2018), and stochastic GT in Pu & Nedić (2021) on the "MNIST", "CIFAR-10", and "CIFAR-100" datasets, respectively. The 95% confidence intervals were computed from three independent runs with random seeds 42, 1010, and 2024.

6 EXPERIMENTS

In this section, we evaluate the performance of our proposed Algorithm 1 on image classification tasks using representative benchmark datasets, including the "MNIST" dataset (Deng, 2012), the "CIFAR-10" dataset (Krizhevsky et al., 2010), and "CIFAR-100" dataset (DeVries & Taylor, 2017). All these tasks involve nonsmooth and nonconvex objective functions, which are intended to show the effectiveness of our algorithm beyond the settings of strong convexity or smoothness. Due to the space limitations, we leave the experimental results on logistic regression with strongly convex and smooth loss functions to Appendix C.3. In all experiments, we considered five agents connected in a ring, where each agent communicates only with its two immediate neighbors. For the coupling matrix W, we set $w_{ii} = 0.4$ for all agent i, $w_{ij} = 0.3$ if agents i and j are neighbors, and $w_{ij} = 0$ otherwise. For each experiment, we considered heterogeneous data distribution, with each agent i randomly sampling 40% data points from the class i and sampling 60% data points from each remaining class. We evaluated the performance of our proposed algorithm through the following three cases: 1) we compared Algorithm 1 with existing distributed stochastic optimization/learning approaches, including distributed SGD in Jakovetic et al. (2018) with diminishing stepsize and the stochastic gradient-tracking (called stochastic GT) in Pu & Nedić (2021) with fixed stepsize; 2) we compared Algorithm 1 with existing adaptive stepsize approaches for centralized learning, including the well-known ADAM in Kingma (2014) and the adaptive SGD in Malitsky & Mishchenko (2024); and 3) to evaluate the effect of the coefficients β and r in the stepsize update rule (i.e., Line 16 in Algorithm 1) and the number of inner-consensus-loop iterations M in Algorithm 1 on convergence accuracy, we test the convergence performance of Algorithm 1 under different β , r, and M, respectively. The detailed experimental settings are given in Appendix C.1 and Appendix C.2, and additional experimental results on comparison of Algorithm 1 and distributed ADAM in Nazari et al. (2022) are provided in Appendix C.3. The code for all experiments is available online².

Comparison with existing distributed stochastic optimization approaches. We trained convolutional neural networks (CNNs) with two, four, and five layers on the "MNIST", "CIFAR-10", and "CIFAR-100" datasets, respectively. We conducted training for 20 epochs on the "MNIST" dataset

²https://anonymous.4open.science/r/DASGD-71D1/README.md

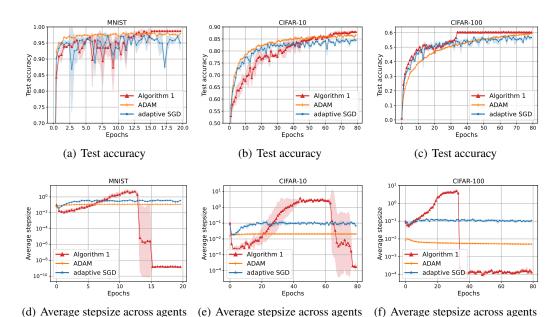


Figure 2: Test-accuracy and average-stepsize (across five agents) evolutions of Algorithm 1, ADAM

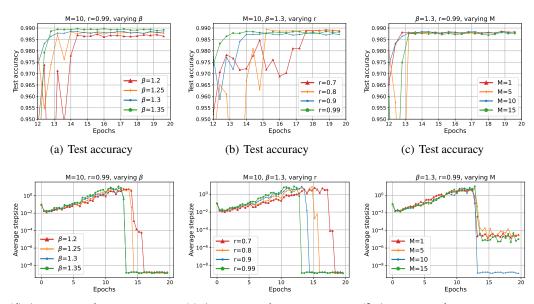
in Kingma (2014), and adaptive SGD in Malitsky & Mishchenko (2024) on the "MNIST", "CIFAR-10", and "CIFAR-100" datasets, respectively. The 95% confidence intervals were computed from three independent runs with random seeds 42, 1010, and 2024.

and 80 epochs on the "CIFAR-10" and "CIFAR-100" datasets, using a batch size of 128. The step-size for distributed SGD was set as $\eta_i = \frac{0.1}{(t+1)^{0.5}}$ and for stochastic GT was set as $\eta_i = 0.1$. Both of them represent the best-performing stepsizes we could find in our comparison. In fact, during our tuning process, we obverse that setting $\eta = 0.01$ for stochastic GT results in overly slow convergence, whereas setting $\eta = 1$ leads to divergent behaviors. For Algorithm 1, we set the coefficients β and r in stepsize update rule as $\beta = 1.3$ and r = 0.99, and the number of inner-loop iterations as M = 10. (The test accuracies of Algorithm 1 under different β , r, and M are provided in Figs. 3(a), 3(b), and 3(c), respectively.)

Fig. 1(a) to Fig. 1(c) show that our proposed Algorithm 1 achieves the highest test accuracy and a more stable steady-state convergence compared with distributed SGD in Jakovetic et al. (2018) and stochastic GT in Pu & Nedić (2021). The early-stage oscillations in test accuracy of Algorithm 1 are mainly attributable to the adaptive process of stepsize adjustments. Compared with distributed SGD with diminishing stepsizes, stochastic GT with a fixed stepsize achieves faster convergence, however, it suffers from larger steady-state oscillations. In contrast, our proposed algorithm eliminates steady-state oscillations, and hence, ensures fast convergence. This is achieved because our proposed adaptive stepsize rule allows each agent to take large stepsizes in the early stages of Algorithm 1 and extremely small stepsizes near the global optimum in the later stages, as shown in Fig. 1(d) to Fig. 1(f). These results further imply a clear stopping criterion for each agent in the implementation of our Algorithm 1. Specifically, we can preset a constant $\tau > 0$ (e.g., $\tau = 10^{-9}$ in the "MNIST" experiment) for all agents, and once an agent i's stepsize η_i^t falls below τ , it can terminate training, which does not compromise the global learning accuracy.

Comparison with existing adaptive stepsize approaches. Since adaptive stepsize approaches are rarely reported in a fully distributed setting without a centralized server/aggregator, we compared the convergence performance of Algorithm 1 with that of existing adaptive stepsize approaches for centralized learning, including ADAM in Kingma (2014) and the adaptive SGD in Malitsky & Mishchenko (2019; 2024). This comparison is challenging because centralized methods can perform training directly on aggregated data, while our approach in Algorithm 1 operates in a fully distributed manner where each agent can only perform local computations and neighboring communication.

Fig. 2(a) to Fig. 2(c) show that Algorithm 1 has a higher test accuracy than both ADAM and adaptive SGD, even without the assistance of any centralized server/aggregator. This finding is noteworthy, as



(d) Average stepsize across agents (e) Average stepsize across agents (f) Average stepsize across agents

Figure 3: Test-accuracy and average-stepsize (across five agents) evolutions of Algorithm 1 under different parameters β , r, and M, respectively.

it empirically demonstrate that our fully distributed learning approach with heterogeneous adaptive stepsizes among agents can accelerate learning compared with centralized methods with a single adaptive stepsize. Furthermore, Fig. 2(d) to Fig. 2(f) once again confirm that our adaptive stepsize approach provides agents with large stepsizes in the early stages and small stepsizes in the convergence stages, thereby facilitating better performance than existing centralized counterparts.

The effects of β , r, and M on convergence accuracy. We evaluate the test accuracies of Algorithm 1 under different coefficients β and r in the stepsize update rule (i.e., Line 16 in Algorithm 1) and the number of inner-loop iterations M in Algorithm 1, respectively. We ran this experiment on the "MNIST" dataset over 20 epochs, with a batch size of 128 and a random seed as 42.

Fig. 3(a), Fig. 3(b), Fig. 3(d), and Fig. 3(e) imply that larger β and r lead to faster convergence and earlier stopping in Algorithm 1. This result is intuitively consistent, as large β and r contribute to larger stepsizes before convergence stages (as shown in Fig. 3(d) and Fig. 3(e)), which in turn leads to a higher convergence speed. Furthermore, Fig. 3(c) and Fig. 3(f) show that the number of inner-consensus-loop iterations M has a negligible effect on convergence accuracy and the stopping criterion. Hence, in practical machine learning tasks, we can set M=1 (so that Algorithm 1 reduces to a single-loop algorithm) to minimize computational and communication costs of Algorithm 1.

7 CONCLUSION

In this paper, we have proposed an adaptive stepsize approach for distributed stochastic optimization and learning without the assistance of any centralized server/aggregator or the need for accurate gradients. This is nontrivial, because existing adaptive stepsize approaches either rely on a centralized server to coordinate stepsizes among agents, or are limited to deterministic scenarios where agents have access to accurate gradients of the objective functions. Moreover, our approach can eliminate steady-state oscillations, and hence, ensures fast convergence. This stands in stark contrast to most existing adaptive stepsize approaches that often incur steady-state oscillations near the global optimal solution, and thereby preventing the algorithm from achieving stable convergence accuracy. In addition, we have systematically characterized the convergence rates of our algorithm for both stochastic and deterministic distributed optimization, and quantified the computational complexities for gradient evaluations on both cases. Experimental results on image classifications using three benchmark datasets confirm the advantages of the proposed approach over existing counterparts.

Ethics statement. All authors declare no conflicts of interest and no ethical issues in this work.

Reproducibility statement. All authors confirm the reproducibility of both the theoretical and experimental results. The code for all experiments is available online at https://anonymous.4open.science/r/DASGD-71D1/README.md. Detailed descriptions of the experimental settings and implementation details are provided in the main text and Appendix. Theoretical assumptions are clearly stated, and complete proofs of all results are included in the Appendix.

REFERENCES

- Arya Akhavan, Massimiliano Pontil, and Alexandre Tsybakov. Distributed zero-order optimization under adversarial noise. *Advances in Neural Information Processing Systems*, 34:10209–10220, 2021.
- Zeyuan Allen-Zhu, Faeze Ebrahimianghazani, Jerry Li, and Dan Alistarh. Byzantine-resilient nonconvex stochastic gradient descent. In *International Conference on Learning Representations*,
- Maksym Andriushchenko, Aditya Vardhan Varre, Loucas Pillaud-Vivien, and Nicolas Flammarion. SGD with large step sizes learns sparse features. In *International Conference on Machine Learning*, pp. 903–925. PMLR, 2023.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. *Advances in Neural Information Processing Systems*, 29, 2016.
- Albert S. Berahas, Raghu Bollapragada, Nitish Shirish Keskar, and Ermin Wei. Balancing communication and computation in distributed optimization. *IEEE Transactions on Automatic Control*, 64(8):3141–3155, 2019.
- Guido Carnevale, Francesco Farina, Ivano Notarnicola, and Giuseppe Notarstefano. GTAdam: Gradient tracking with adaptive momentum for distributed online optimization. *IEEE Transactions on Control of Network Systems*, 10(3):1436–1448, 2022.
- Ziqin Chen and Yongqiang Wang. Locally differentially private distributed online learning with guaranteed optimality. *IEEE Transactions on Automatic Control*, 2024.
- Michael Crawshaw, Blake Woodworth, and Mingrui Liu. Constant stepsize local GD for logistic regression: Acceleration by instability. In *International Conference on Machine Learning*, 2025.
- Li 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.
- Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- Aymeric Dieuleveut and Kumar Kshitij Patel. Communication trade-offs for local-SGD with large step size. *Advances in Neural Information Processing Systems*, 32, 2019.
- Xuhui Ding, Yuchen Xu, Gaoyang Li, Kai Yang, Jinhong Yuan, and Jianping An. Design and performance evaluation for BILCM-ID system with improved stopping criterion. *IEEE Transactions on Vehicular Technology*, 74(4):6779–6784, 2025.
- Guner Dilsad Er, Sebastian Trimpe, and Michael Muehlebach. Distributed event-based learning via ADMM. In *International Conference on Machine Learning*, 2024.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in Neural Information Processing Systems*, 31, 2018.
- Roger Fletcher. On the barzilai-borwein method. In *Optimization and Control with Applications*, pp. 235–256. Springer, 2005.

- Zhishuai Guo, Mingrui Liu, Zhuoning Yuan, Li Shen, Wei Liu, and Tianbao Yang. Communication efficient distributed stochastic auc maximization with deep neural networks. In *International Conference on Machine Learning*, pp. 3864–3874. PMLR, 2020.
 - Yutong He, Jie Hu, Xinmeng Huang, Songtao Lu, Bin Wang, and Kun Yuan. Distributed bilevel optimization with communication compression. In *International Conference on Machine Learning*, 2024.
 - Charlie Hou, Kiran K Thekumparampil, Giulia Fanti, and Sewoong Oh. FeDChain: Chained algorithms for near-optimal communication cost in federated learning. In *International Conference on Learning Representations*, 2021.
 - Jie Hu, Vishwaraj Doshi, et al. Accelerating distributed stochastic optimization via self-repellent random walks. In *International Conference on Learning Representations*, 2024.
 - Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Booster: Tackling harmful fine-tuning for large language models via attenuating harmful perturbation. In *International Conference on Learning Representations*, 2024a.
 - Yan Huang, Xiang Li, Yipeng Shen, Niao He, and Jinming Xu. Achieving near-optimal convergence for distributed minimax optimization with adaptive stepsizes. *Advances in Neural Information Processing Systems*, 37:19740–19782, 2024b.
 - Nikita Ivkin, Daniel Rothchild, Enayat Ullah, Ion Stoica, Raman Arora, et al. Communication-efficient distributed SGD with sketching. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Robert A Jacobs. Increased rates of convergence through learning rate adaptation. *Neural Networks*, 1(4):295–307, 1988.
 - Dusan Jakovetic, Dragana Bajovic, Anit Kumar Sahu, and Soummya Kar. Convergence rates for distributed stochastic optimization over random networks. In *IEEE Conference on Decision and Control*, pp. 4238–4245. IEEE, 2018.
 - Xiaowen Jiang and Sebastian U Stich. Adaptive SGD with polyak stepsize and line-search: Robust convergence and variance reduction. *Advances in Neural Information Processing Systems*, 36: 26396–26424, 2023.
 - Ali Kavis, Kfir Yehuda Levy, and Volkan Cevher. High probability bounds for a class of nonconvex algorithms with AdaGrad stepsize. In *International Conference on Learning Representations*, 2022.
 - Ahmed Khaled and Chi Jin. Faster federated optimization under second-order similarity. In *International Conference on Learning Representations*, 2023.
 - Junhyung Lyle Kim, Mohammad Taha Toghani, César A Uribe, and Anastasios Kyrillidis. Adaptive federated learning with auto-tuned clients. In *International Conference on Learning Representa*tions, 2024.
 - Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - Anastasiia Koloskova, Tao Lin, and Sebastian U Stich. An improved analysis of gradient tracking for decentralized machine learning. *Advances in Neural Information Processing Systems*, 34: 11422–11435, 2021.
 - Alex Krizhevsky, Geoff Hinton, et al. Convolutional deep belief networks on CIFAR-10. *Unpublished Manuscript*, 40(7):1–9, 2010.
 - Ilya Kuruzov, Gesualdo Scutari, and Alexander Gasnikov. Achieving linear convergence with parameter-free algorithms in decentralized optimization. *Advances in Neural Information Processing Systems*, 37:96011–96044, 2024.

- Su Hyeong Lee, Manzil Zaheer, and Tian Li. Efficient distributed optimization under heavy-tailed noise. In *International Conference on Learning Representations*, 2025.
 - Hanmin Li, Avetik Karagulyan, and Peter Richtarik. Det-CGD: Compressed gradient descent with matrix stepsizes for non-convex optimization. In *International Conference on Learning Repre*sentations, 2024a.
 - Huan Li, Cong Fang, Wotao Yin, and Zhouchen Lin. Decentralized accelerated gradient methods with increasing penalty parameters. *IEEE Transactions on Signal Processing*, 68:4855–4870, 2020.
 - Junyi Li, Feihu Huang, and Heng Huang. FedDA: Faster adaptive gradient methods for federated constrained optimization. In *International Conference on Learning Representations*, 2024b.
 - Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *International Conference on Artificial Intelligence and Statistics*, pp. 983–992. PMLR, 2019.
 - Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? A case study for decentralized parallel stochastic gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.
 - Jihao Andreas Lin, Shreyas Padhy, Javier Antorán, Austin Tripp, Alexander Terenin, Csaba Szepesvári, José Miguel Hernández-Lobato, and David Janz. Stochastic gradient descent for gaussian processes done right. In *International Conference on Learning Representations*, 2023.
 - Xiaorui Liu, Yao Li, Rongrong Wang, Jiliang Tang, and Ming Yan. Linear convergent decentralized optimization with compression. In *International Conference on Learning Representations*, 2020.
 - Yura Malitsky and Konstantin Mishchenko. Adaptive gradient descent without descent. *arXiv* preprint arXiv:1910.09529, 2019.
 - Yura Malitsky and Konstantin Mishchenko. Adaptive proximal gradient method for convex optimization. *Advances in Neural Information Processing Systems*, 37:100670–100697, 2024.
 - Brendan McMahan and Matthew Streeter. Delay-tolerant algorithms for asynchronous distributed online learning. *Advances in Neural Information Processing Systems*, 27, 2014.
 - Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. Dadam: A consensus-based distributed adaptive gradient method for online optimization. *IEEE Transactions on Signal Processing*, 70:6065–6079, 2022.
 - Angelia Nedic. Distributed gradient methods for convex machine learning problems in networks: Distributed optimization. *IEEE Signal Processing Magazine*, 37(3):92–101, 2020.
 - Angelia Nedić and Ji Liu. Distributed optimization for control. *Annual Review of Control, Robotics, and Autonomous Systems*, 1(1):77–103, 2018.
 - Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
 - Angelia Nedić, Alex Olshevsky, Wei Shi, and César A Uribe. Geometrically convergent distributed optimization with uncoordinated step-sizes. In *American Control Conference*, pp. 3950–3955. IEEE, 2017.
 - Angelia Nedić, Alex Olshevsky, and Michael G Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976, 2018.
- Edward Duc Hien Nguyen, Sulaiman A Alghunaim, Kun Yuan, and César A Uribe. On the performance of gradient tracking with local updates. In *IEEE Conference on Decision and Control*, pp. 4309–4313. IEEE, 2023.
 - Shi Pu and Angelia Nedić. Distributed stochastic gradient tracking methods. *Mathematical Programming*, 187(1):409–457, 2021.

- Shi Pu, Alex Olshevsky, and Ioannis Ch Paschalidis. Asymptotic network independence in distributed stochastic optimization for machine learning: Examining distributed and centralized stochastic gradient descent. *IEEE Signal Processing Magazine*, 37(3):114–122, 2020.
 - Ali Ramezani-Kebrya, Kimon Antonakopoulos, Igor Krawczuk, Justin Deschenaux, and Volkan Cevher. Distributed extra-gradient with optimal complexity and communication guarantees. In *International Conference on Learning Representations*, 2023.
 - Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In *International Conference on Machine Learning*, pp. 314–323. PMLR, 2016.
 - Michal Rolinek and Georg Martius. L4: Practical loss-based stepsize adaptation for deep learning. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
 - Augustinos D Saravanos, Hunter Kuperman, Alex Oshin, Arshiya Taj Abdul, Vincent Pacelli, and Evangelos A Theodorou. Deep distributed optimization for large-scale quadratic programming. In *International Conference on Learning Representations*, 2024.
 - Kevin Scaman, Francis Bach, Sébastien Bubeck, Laurent Massoulié, and Yin Tat Lee. Optimal algorithms for non-smooth distributed optimization in networks. *Advances in Neural Information Processing Systems*, 31, 2018.
 - Tom Schaul, Sixin Zhang, and Yann LeCun. No more pesky learning rates. In *International Conference on Machine Learning*, 2013.
 - Ohad Shamir and Nathan Srebro. Distributed stochastic optimization and learning. In *Annual Allerton Conference on Communication, Control, and Computing*, pp. 850–857. IEEE, 2014.
 - Arsalan Sharifnassab, Saber Salehkaleybar, and Richard Sutton. Metaoptimize: A framework for optimizing step sizes and other meta-parameters. In *International Conference on Machine Learning*, 2024.
 - Rohan Sharma, Kaiyi Ji, and Changyou Chen. AUC-CL: A batchsize-robust framework for self-supervised contrastive representation learning. In *International Conference on Learning Representations*, 2023.
 - Yilong Song, Peijin Li, Bin Gao, and Kun Yuan. Distributed retraction-free and communication-efficient optimization on the stiefel manifold. In *International Conference on Machine Learning*, 2025.
 - Zhuoqing Song, Lei Shi, Shi Pu, and Ming Yan. Optimal gradient tracking for decentralized optimization. *Mathematical Programming*, 207(1):1–53, 2024.
 - Kunal Srivastava and Angelia Nedic. Distributed asynchronous constrained stochastic optimization. *IEEE Journal of Selected Topics in Signal Processing*, 5(4):772–790, 2011.
 - S Sundhar Ram, Angelia Nedić, and Venugopal V Veeravalli. Distributed stochastic subgradient projection algorithms for convex optimization. *Journal of Optimization Theory and Applications*, 147(3):516–545, 2010.
 - Kemal Tutuncu, Ilkay Cinar, Ramazan Kursun, and Murat Koklu. Edible and poisonous mushrooms classification by machine learning algorithms. In *Mediterranean Conference on Embedded Computing*, pp. 1–4. IEEE, 2022.
 - Andreas Vlachos. A stopping criterion for active learning. *Computer Speech and Language*, 22(3): 295–312, 2008. ISSN 0885-2308.
 - Jing Wang and Nicola Elia. A control perspective for centralized and distributed convex optimization. In *IEEE Conference on Decision and Control and European Control*, pp. 3800–3805. IEEE, 2011.

- Xiaoyu Wang, Mikael Johansson, and Tong Zhang. Generalized polyak step size for first order optimization with momentum. In *International Conference on Machine Learning*, pp. 35836–35863. PMLR, 2023.
 - Kaixuan Wei, Angelica Aviles-Rivero, Jingwei Liang, Ying Fu, Carola-Bibiane Schönlieb, and Hua Huang. Tuning-free plug-and-play proximal algorithm for inverse imaging problems. In *International Conference on Machine Learning*, pp. 10158–10169. PMLR, 2020.
 - Siyu Xie, Masoud H Nazari, George Yin, et al. Adaptive step size selection in distributed optimization with observation noise and unknown stochastic target variation. *Automatica*, 135:109940, 2022.
 - Ran Xin, Soummya Kar, and Usman A Khan. Decentralized stochastic optimization and machine learning: As unified variance-reduction framework for robust performance and fast convergence. *IEEE Signal Processing Magazine*, 37(3):102–113, 2020.
 - Bo Yang and Mikael Johansson. Distributed optimization and games: A tutorial overview. *Networked Control Systems*, pp. 109–148, 2010.
 - Shuoguang Yang, Xuezhou Zhang, and Mengdi Wang. Decentralized gossip-based stochastic bilevel optimization over communication networks. *Advances in Neural Information Processing Systems*, 35:238–252, 2022.
 - Tao Yang, Xinlei Yi, Junfeng Wu, Ye Yuan, Di Wu, Ziyang Meng, Yiguang Hong, Hong Wang, Zongli Lin, and Karl H Johansson. A survey of distributed optimization. *Annual Reviews in Control*, 47:278–305, 2019.
 - Tianbao Yang. Trading computation for communication: Ddistributed stochastic dual coordinate ascent. *Advances in Neural Information Processing Systems*, 26, 2013.
 - Zhuang Yang and Li Ma. Adaptive step size rules for stochastic optimization in large-scale learning. *Statistics and Computing*, 33(2):45, 2023.
 - Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 26(3):1835–1854, 2016.
 - Zhuoning Yuan, Yuexin Wu, Zi-Hao Qiu, Xianzhi Du, Lijun Zhang, Denny Zhou, and Tianbao Yang. Provable stochastic optimization for global contrastive learning: Small batch does not harm performance. In *International Conference on Machine Learning*, pp. 25760–25782. PMLR, 2022.
 - Juntang Zhuang, Tommy Tang, Yifan Ding, Sekhar C Tatikonda, Nicha Dvornek, Xenophon Papademetris, and James Duncan. Adabelief optimizer: Aadapting stepsizes by the belief in observed gradients. *Advances in Neural Information Processing Systems*, 33:18795–18806, 2020.

APPENDIX

- Section A: Notations
- Section B: Results of Algorithm 1
 - B.1 Technical lemmas
 - B.2 Proof of Theorem 1
 - B.3 Proof of Theorem 2
 - B.4 Proof of Corollary 1
- Section C: Experimental setups and additional experimental results
 - C.1 Benchmark datasets
 - C.2 Experimental setups
 - C.3 Additional experimental results

A NOTATIONS

For the sake of notational simplicity, we introduce some additional notations. We use \mathbb{R}^+ to denote the set of positive real numbers and use $\mathbf{X}^t \triangleq \operatorname{col}(x_1^t, x_2^t, \cdots, x_m^t) \in \mathbb{R}^{mn}$ to denote the stacked model parameters of all agents. We also use \otimes to denote the Kronecker product. We use $\overline{x}^t(q)$ to denote the average of all agents' model parameters at the qth inner iteration of the tth outer iteration. We define $\mathcal{F}_t \triangleq \{\xi_{i,s}|i=1,\cdots,m \text{ and } s=0,\cdots,t\}$, where $\xi_{i,t}$ represents the data point sampled by agent i at the tth iteration. For further notational simplicity, we define $\overline{x}^t = \frac{1}{m} \sum_{i=1}^m x_i^t, \overline{\eta^t y_1^t} = \frac{1}{m} \sum_{i=1}^m \eta_i^t y_{1,i}^t, \eta_{\max}^t = \max_{i \in [m]} \eta_i^t, \eta_{\max} = \max_{t \in \mathbb{N}} \eta_{\max}^t, \overline{\eta}^t = \frac{1}{m} \sum_{i=1}^m \eta_i^t, \overline{y}_1^t = \frac{1}{m} \sum_{i=1}^m y_{1,i}^t, \overline{y}_2^t = \frac{1}{m} \sum_{i=1}^m y_{2,i}^t, \hat{x}_i^t = x_i^t - \overline{x}^t, \hat{y}_{1,i}^t = y_{1,i}^t - \overline{y}_1^t, \text{ and } \hat{y}_{2,i}^t = y_{2,i}^t - \overline{y}_2^t.$

B RESULTS OF ALGORITHM 1

B.1 TECHNICAL LEMMAS

We introduce the following three lemmas to characterize the consensus errors of Algorithm 1.

Lemma 1. Under Assumptions 1 and 2, the following inequality holds for Algorithm 1:

$$\mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^t\|^2] \le \frac{45}{46} \mathbb{E}[\|\overline{x}^t - \overline{x}^{t-1}\|^2] - \mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^t\|^2] + \frac{125}{31\beta} \mathbb{E}[\overline{\eta}^t (f(\overline{x}^{t-1}) - f(\overline{x}^t))] + 2\delta_3^t,$$
(4)

where the constant δ_3^t is given by

$$\delta_{3}^{t} = \frac{125}{124m} \left(\frac{L^{2}\eta_{\text{max}}}{a_{6}} + \frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1+a_{7})}{50a_{7}} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t}\|^{2}]$$

$$+ \frac{125}{124m} \left(\frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1+a_{7})}{50a_{7}} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t-1}\|^{2}]$$

$$+ \frac{\eta_{\text{max}}^{2}}{m} \left(\frac{187}{62} + \frac{125}{62a_{2}} + \frac{1}{a_{1}} \right) \sum_{i=1}^{m} \left(\mathbb{E}[\|\hat{y}_{1,i}^{t}\|^{2}] + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^{2}] \right)$$

$$+ \frac{125\eta_{\text{max}}^{2}}{124m} \left(\frac{1}{a_{3}} + 1 + \frac{1}{a_{8}} + \frac{2}{a_{2}} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^{2}]$$

$$+ \left(\left(1 - \frac{1}{a_{3}} \right) \left(2 + a_{9} + \frac{1}{a_{9}} \right) + \frac{1}{a_{6}m} \right) \frac{\sigma^{2}}{|\mathcal{B}|},$$

$$(5)$$

with $b_1 \triangleq \frac{2\beta^2}{m}(1-a_3)(a_4-1)(1-\frac{1}{a_5}) + 2\beta^2(1-a_3)(1-\frac{1}{a_4}), a_1 \in (0, \frac{1}{124}), a_2 \in \mathbb{R}^+, a_3 \in \mathbb{R}^+, a_4, a_5, a_6, a_7 \in (0, 1), a_8 = \frac{1-124a_1}{250}.$

Proof. According to Line 3 in Algorithm 1, we have

$$\overline{x}^{t+1} = \overline{x}^{t+1}(0) = \overline{x}^t - \overline{\eta^t y_1^t},\tag{6}$$

with $\overline{x}^t = \frac{1}{m} \sum_{i=1}^m x_i^t$ and $\overline{\eta^t y_1^t} = \frac{1}{m} \sum_{i=1}^m \eta_i^t y_{1,i}^t$. Since $\overline{\eta}^t \overline{y}_1^t = \frac{1}{m^2} \sum_{i=1}^m \eta_i^t \sum_{j=1}^m y_{1,j}^t$ holds, we obtain the following inequality:

$$\mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^{t}\|^{2} | \mathcal{F}_{t}] = \|\overline{\eta^{t}} y_{1}^{t}\|^{2} = \|\overline{\eta^{t}} y_{1}^{t} - \overline{\eta}^{t} \overline{y}_{1}^{t} + \overline{\eta}^{t} \overline{y}_{1}^{t}\|^{2} \\
\leq \left(1 + \frac{1}{a_{1}}\right) \|\overline{\eta^{t}} y_{1}^{t} - \overline{\eta}^{t} \overline{y}_{1}^{t}\|^{2} + (1 + a_{1}) \|\overline{\eta}^{t} \overline{y}_{1}^{t}\|^{2} \\
= \left(1 + \frac{1}{a_{1}}\right) \|\frac{1}{m} \sum_{i=1}^{m} (\eta_{i}^{t} y_{1,i}^{t} - \eta_{i}^{t} \overline{y}_{1}^{t}) \|^{2} + (1 + a_{1}) \|\overline{\eta}^{t} \overline{y}_{1}^{t}\|^{2} \\
\leq \left(1 + \frac{1}{a_{1}}\right) \frac{1}{m} \sum_{i=1}^{m} \|\eta_{i}^{t} y_{1,i}^{t} - \eta_{i}^{t} \overline{y}_{1}^{t} \|^{2} + (1 + a_{1}) \|\overline{\eta}^{t} \overline{y}_{1}^{t} \|^{2} \\
= \left(1 + \frac{1}{a_{1}}\right) \frac{1}{m} \sum_{i=1}^{m} (\eta_{i}^{t})^{2} \|y_{1,i}^{t} - \overline{y}_{1}^{t} \|^{2} + (1 + a_{1}) \|\overline{\eta}^{t} \overline{y}_{1}^{t} \|^{2} \\
\leq \left(1 + \frac{1}{a_{1}}\right) \frac{1}{m} \sum_{i=1}^{m} \eta_{\max}^{2} \|\hat{y}_{1,i}^{t}\|^{2} + (1 + a_{1}) \|\overline{\eta}^{t} \overline{y}_{1}^{t}\|^{2}, \\$$

where $\mathcal{F}_t=\{\xi_{i,s}|i=1,\dots,N;\ s=0,\dots,t\}$ with $\xi_{i,t}$ denoting the data point sampled by agent i at iteration t. Here, we have used the inequality $\|a+b\|^2 \leq (1+\frac{1}{\alpha})\|a\|^2 + (1+\alpha)\|b\|^2$ for any $\alpha>0$ and $a,\ b\in\mathbb{R}^n$ in the first inequality and the inequality $\|\frac{1}{m}\sum_{i=1}^m a_i\|^2 \leq \frac{1}{m}\sum_{i=1}^m \|a_i\|^2$ for any $a_i\in\mathbb{R},\ i=1,\cdots,m$ in the second inequality. By choosing $a_1\in(0,\frac{1}{124})$ and applying the relation $\|a\|^2=\|a-b\|^2-\|b\|^2+2\langle a,b\rangle$ to equation 7, the term $\|\overline{\eta}^t\overline{y}_1^t\|$ can be bounded by

$$\|\overline{\eta}^t \overline{y}_1^t\|^2 = \|\overline{\eta}^t (\overline{y}_1^t - \overline{y}_2^{t-1})\|^2 - \|\overline{\eta}^t \overline{y}_2^{t-1}\|^2 + 2\langle \overline{\eta}^t \overline{y}_1^t, \overline{\eta}^t \overline{y}_2^{t-1}\rangle. \tag{8}$$

The first term on the right-hand side of equation 8 satisfies

$$\|\overline{\eta}^{t}(\overline{y}_{1}^{t} - \overline{y}_{2}^{t-1})\|^{2} = \left\|\frac{1}{m}\sum_{i=1}^{m}\eta_{i}^{t}(\overline{y}_{1}^{t} - \overline{y}_{2}^{t-1})\right\|^{2}$$

$$\leq \frac{1}{m}\sum_{i=1}^{m}\|\eta_{i}^{t}(\overline{y}_{1}^{t} - \overline{y}_{2}^{t-1})\|^{2} = \frac{1}{m}\sum_{i=1}^{m}\|\eta_{i}^{t}(y_{1,i}^{t} - y_{2,i}^{t-1}) - \eta_{i}^{t}(\hat{y}_{1,i}^{t} - \hat{y}_{2,i}^{t-1})\|^{2}$$

$$\leq \frac{1}{m}\sum_{i=1}^{m}(1+a_{2})\|\eta_{i}^{t}(y_{1,i}^{t} - y_{2,i}^{t-1})\|^{2} + \frac{1}{m}\sum_{i=1}^{m}\left(1 + \frac{1}{a_{2}}\right)\|\eta_{i}^{t}(\hat{y}_{1,i}^{t} - \hat{y}_{2,i}^{t-1})\|^{2},$$

$$(9)$$

for any $a_2 \in \mathbb{R}^+$. According to $L_i^t = \frac{\|y_{1,i}^t - y_{2,i}^{t-1}\|}{\|x_i^t - x_i^{t-1}\|}$ and $\|a + b\| \le 2\|a\|^2 + 2\|b\|^2$, equation 9 can be rewritten as follows:

$$\mathbb{E}[\|\overline{\eta}^{t}(\overline{y}_{1}^{t} - \overline{y}_{2}^{t-1})\|^{2}] \\
\leq \mathbb{E}\left[\frac{1}{m}\sum_{i=1}^{m}(1+a_{2})(\eta_{i}^{t}L_{i}^{t})^{2}\|x_{i}^{t} - x_{i}^{t-1}\|^{2} + \frac{2}{m}\sum_{i=1}^{m}\eta_{\max}^{2}(1+\frac{1}{a_{2}})(\|\hat{y}_{1,i}^{t}\|^{2} + \|\hat{y}_{2,i}^{t-1}\|^{2})\right] \\
\leq \mathbb{E}\left[\frac{1}{m}\sum_{i=1}^{m}(1+a_{2})(1+a_{7})(\eta_{i}^{t}L_{i}^{t})^{2}\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \mathbb{E}\left[\frac{2}{m}\sum_{i=1}^{m}\eta_{\max}^{2}\left(1+\frac{1}{a_{2}}\right)(\|\hat{y}_{1,i}^{t}\|^{2} + \|\hat{y}_{2,i}^{t-1}\|^{2})\right] \\
+ \frac{2\eta_{\max}^{2}L^{2}}{m}\sum_{i=1}^{m}(1+a_{2})\left(1+\frac{1}{a_{7}}\right)\mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2} + \|\hat{x}_{i}^{t-1}\|^{2}\right], \tag{10}$$

for any $a_7 \in \mathbb{R}^+$. Here, we have used the L-smoothness of each f_i in the second inequality.

 We proceed to estimate a lower bound on the second term on the right-hand side of equation 8.

$$\|\overline{\eta}^{t}\overline{y}_{2}^{t-1}\|^{2} = (\overline{\eta}^{t})^{2}\|\overline{y}_{2}^{t-1}\|^{2} = (\overline{\eta}^{t})^{2}\|\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1} + \overline{y}_{1}^{t-1}\|^{2}$$

$$\geq (1 - a_{3})(\overline{\eta}^{t})^{2}\|\overline{y}_{1}^{t-1}\|^{2} + \left(1 - \frac{1}{a_{3}}\right)(\overline{\eta}^{t})^{2}\|\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}\|^{2},$$
(11)

for any $a_3 \in (0,1)$, where in the derivation we have used the inequality $||a+b||^2 \ge (1-\frac{1}{\alpha})||a||^2 + (1-\alpha)||b||^2$, for any $a \in (0,1)$. Since the relationship $1-\frac{1}{a_3} < 0$ holds, the second term on the right-hand side of equation 11 satisfies

$$\left(1 - \frac{1}{a_3}\right) (\overline{\eta}^t)^2 \|\overline{y}_2^{t-1} - \overline{y}_1^{t-1}\|^2
\geq \left(1 - \frac{1}{a_3}\right) (\eta_{\max})^2 \|\frac{1}{m} \sum_{i=1}^m (g_i^{t-1}(x_i^{t-1}) - g_i^{t-2}(x_i^{t-1}))\|^2
\geq \left(1 - \frac{1}{a_3}\right) (\eta_{\max})^2 \frac{1 + a_9}{m} \sum_{i=1}^m \|g_i^{t-1}(x_i^{t-1}) - \nabla f_i(x_i^{t-1})\|^2
+ \left(1 - \frac{1}{a_3}\right) (\eta_{\max})^2 \frac{a_9 + 1}{a_9 m} \sum_{i=1}^m \|\nabla f_i(x_i^{t-1}) - g_i^{t-2}(x_i^{t-1})\|^2.$$
(12)

The first term on the right-hand side of equation 12 satisfies

$$\mathbb{E}[\|g_{i}^{t-1}(x_{i}^{t-1}) - \nabla f_{i}(x_{i}^{t-1})\|^{2}] \\
= \mathbb{E}\left[\left\|\frac{1}{m}\sum_{j=1}^{m}(\nabla l_{i}(x_{i}^{t-1},\xi_{i,j}^{t-1}) - \nabla f_{i}(x_{i}^{t-1}))\right\|^{2}\right] \\
= \frac{1}{m^{2}}\mathbb{E}\left[\sum_{j=1}^{m}\left\|\nabla l_{i}(x_{i}^{t-1},\xi_{i,j}^{t-1}) - \nabla f_{i}(x_{i}^{t-1})\right\|^{2}\right] \\
= \frac{1}{m^{2}}\sum_{i=1}^{k}\frac{\sigma^{2}}{|\mathcal{B}|} = \frac{\sigma^{2}}{m|\mathcal{B}|},$$
(13)

where in the derivation we have used Assumption 1.

The second term on the right-hand side of equation 12 satisfies

$$\mathbb{E}[\|\nabla f_i(x_i^{t-1}) - g_i^{t-2}(x_i^{t-1})\|^2] \le \frac{\sigma^2}{|\mathcal{B}|}.$$
 (14)

By using the inequality $(\sum_{i=1}^n a_i)^2 \ge \sum_{i=1}^n a_i^2$ for any nonnegative constants a_1, \ldots, a_n , the first term on the right-hand side of equation 11 satisfies

$$\mathbb{E}\left[(\overline{\eta}^{t})^{2} \| \overline{y}_{1}^{t-1} \|^{2}\right] = \mathbb{E}\left[\left(\overline{\eta}^{t} \| \overline{y}_{1}^{t-1} \|\right)^{2}\right] = \mathbb{E}\left[\left(\frac{1}{m} \sum_{i=1}^{m} \eta_{i}^{t} \| \overline{y}_{1}^{t-1} \|\right)^{2}\right] \\
= \mathbb{E}\left[\left(\frac{1}{m} \sum_{i=1}^{m} \eta_{i}^{t} \| y_{1,i}^{t-1} - \hat{y}_{1,i}^{t-1} \|\right)^{2}\right] \\
\geq \frac{1}{m^{2}} \sum_{i=1}^{m} \mathbb{E}\left[(\eta_{i}^{t})^{2} \| y_{1,i}^{t-1} - \hat{y}_{1,i}^{t-1} \|^{2}\right] \\
\geq \frac{1 - a_{10}}{m^{2}} \sum_{i=1}^{m} \mathbb{E}\left[(\eta_{i}^{t})^{2} \| y_{1,i}^{t-1} \|^{2}\right] + \frac{a_{10} - 1}{a_{10} m^{2}} \sum_{i=1}^{m} \mathbb{E}\left[(\eta_{i}^{t})^{2} \| \hat{y}_{1,i}^{t-1} \|^{2}\right], \tag{15}$$

for $a_{10} \in (0,1)$.

We estimate a lower bound on the first term on the right-hand side of equation 15 as follows:

$$\sum_{i=1}^{m} \mathbb{E}\left[(\eta_{i}^{t})^{2} \| y_{1,i}^{t-1} \|^{2} \right] = \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| x_{i}^{t}(0) - x_{i}^{t-1} \|^{2} \right] \\
\geq (1 - a_{4}) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| \overline{x}^{t} - \overline{x}^{t-1} \|^{2} \right] \\
+ \left(1 - \frac{1}{a_{4}} \right) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| \hat{x}_{i}^{t}(0) - \hat{x}_{i}^{t-1} \|^{2} \right] \\
\geq (1 - a_{4})(1 - a_{5}) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| \overline{x}^{t} - \overline{x}^{t-1} \|^{2} \right] \\
+ (1 - a_{4})a_{5} \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| \overline{x}^{t} - \overline{x}^{t-1} \|^{2} \right] \\
+ (1 - \frac{1}{a_{4}}) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}} \right)^{2} \| \hat{x}_{i}^{t}(0) - \hat{x}_{i}^{t-1} \|^{2} \right],$$
(16)

for any $a_4 \in (1, +\infty)$ and $a_5 \in (1, +\infty)$.

We estimate a lower bound on the third term on the right-hand side of equation 16 as follows:

$$\sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\hat{x}_{i}^{t}(0) - \hat{x}_{i}^{t-1}\|^{2}\right] \\
\geq (1 - a_{11}) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\hat{x}_{i}^{t}(0)\|^{2}\right] + \left(1 - \frac{1}{a_{11}}\right) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\hat{x}_{i}^{t-1}\|^{2}\right] \\
\geq (1 - a_{11}) \left(\frac{\eta_{\max}}{\eta_{\min}}\right)^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2}\right] + \left(1 - \frac{1}{a_{11}}\right) \sum_{i=1}^{m} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\hat{x}_{i}^{t-1}\|^{2}\right], \tag{17}$$

for any $a_{11} \in (0,1)$, where we have used the inequality $\sum_{i=1}^m \|\hat{x}_i^t(0)\|^2 \ge \rho^{2M} \sum_{i=1}^m \|\hat{x}_i^t(0)\|^2 \ge \sum_{i=1}^m \|\hat{x}_i^t\|^2$ in the last inequality. This inequality will be proved in the subsequent Lemma 3.

Finally, using inequalities equation 11 -equation 17, we arrive at

$$\mathbb{E}[\|\overline{\eta}^{t}\overline{y}_{2}^{t-1}\|^{2}]$$

$$\geq (1 - a_{3})(1 - a_{10})(1 - a_{4})(1 - a_{5})\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right]$$

$$+ (1 - a_{3})(1 - a_{10})(1 - a_{4})a_{5}\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right]$$

$$+ (1 - a_{3})(1 - a_{10})(1 - \frac{1}{a_{4}})(1 - a_{11})\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|\hat{x}_{i}^{t}\|^{2}\right]$$

$$+ (1 - a_{3})(1 - a_{10})\left(1 - \frac{1}{a_{4}}\right)\left(1 - \frac{1}{a_{11}}\right)\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|\hat{x}_{i}^{t-1}\|^{2}\right]$$

$$+ (1 - a_{3})\left(1 - \frac{1}{a_{10}}\right)\frac{1}{m^{2}}\sum_{i=1}^{m}(\eta_{i}^{t})^{2}\mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right]$$

$$+ \left(1 - \frac{1}{a_{3}}\right)\left(2 + a_{9} + \frac{1}{a_{9}}\right)\frac{\eta_{\max}^{2}\sigma^{2}}{|\mathcal{B}|}.$$
(18)

The inequality in equation 18 implies

$$-\mathbb{E}[\|\overline{\eta}^{t}\overline{y}_{2}^{t-1}\|^{2}] \leq -(1-a_{3})(1-a_{10})(1-a_{4})(1-a_{5})\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|\overline{x}^{t}-\overline{x}^{t-1}\|^{2}\right] + (1-a_{3})(1-a_{10})(a_{4}-1)a_{5}\frac{\beta^{2}}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\|\overline{x}^{t}-\overline{x}^{t-1}\|^{2}\right] + \frac{b_{1}}{m}\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2}\right] + \frac{b_{2}}{m}\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + \frac{b_{3}\eta_{\max}^{2}}{m}\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + \frac{b_{4}\eta_{\max}^{2}\sigma^{2}}{|\mathcal{B}|},$$

$$(1-a_{2})(1-a_{10})(1-a_{10})(1-a_{10})\beta^{2} + \frac{b_{1}\eta_{\min}^{2}\sigma^{2}}{m}(1-a_{20})(1-a_{10})(1-a_{10})\beta^{2}$$

$$\begin{array}{lll} \text{with} & b_1 & = & \frac{(1-a_3)(1-a_{10})\left(1-\frac{1}{a_4}\right)(1-a_{11})\beta^2}{m}, & b_2 & = & \frac{(1-a_3)(1-a_{10})\left(1-\frac{1}{a_4}\right)\left(\frac{1}{a_{11}}-1\right)\beta^2}{m}, & b_3 & = & \frac{(1-a_3)\left(\frac{1}{a_{10}}-1\right)}{m}, & \text{and} & b_4 & = & \left(\frac{1}{a_3}-1\right)\left(2+a_9+\frac{1}{a_9}\right). \end{array}$$

Next, we estimate an upper bound for the last term on the right-hand side of equation 8 as follows:

$$2\mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \overline{\eta}^{t} \overline{y}_{2}^{t-1} \rangle\right] = \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1} + \overline{y}_{1}^{t-1}) \rangle\right]$$

$$= \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} y_{1,i}^{t-1} \rangle\right] - \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} \hat{y}_{1,i}^{t-1} \rangle\right] + \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}) \rangle\right]$$

$$= \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{x}_{i}^{t}(0)) \rangle\right] - \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}) \rangle\right]$$

$$+ \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}) \rangle\right]$$

$$\leq \frac{2}{m} \sum_{i=1}^{m} \mathbb{E}\left[\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \eta_{i}^{t} (\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}) \rangle\right]$$

$$+ \mathbb{E}\left[\frac{\eta_{\max}^{2}}{a_{8}} \|\overline{y}_{2}^{t-1} - \overline{y}_{1}^{t-1}\|^{2}\right], \qquad (20)$$

with $a_8 = \frac{1-124a_1}{250} > 0$, where we have used the inequality $2\langle a,b\rangle \leq \frac{1}{\alpha}\|a\|^2 + \alpha\|b\|^2$ in the last inequality.

Next, we need to transform the first term on the right-hand side of equation 20.

$$2\mathbb{E}\left[\left\langle \overline{\eta}^{t}\overline{y}_{1}^{t}, \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(x_{i}^{t-1} - x_{i}^{t})\right\rangle\right]$$

$$= 2\mathbb{E}\left[\left\langle \overline{\eta}^{t}\overline{y}_{1}^{t}, \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] + 2\mathbb{E}\left[\left\langle \overline{\eta}^{t}\overline{y}_{1}^{t}, \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\hat{x}_{i}^{t-1} - \hat{x}_{i}^{t})\right\rangle\right]$$

$$\leq \frac{2}{m}\sum_{j=1}^{m}\mathbb{E}\left[\left\langle \eta_{j}^{t}\overline{y}_{1}^{t}, \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] + a_{8}\mathbb{E}\left[\left\|\overline{\eta}^{t}\overline{y}_{1}^{t}\right\|^{2}\right] + \frac{2\beta^{2}}{a_{8}}\mathbb{E}\left[\left\|\hat{x}_{i}^{t-1}\right\|^{2}\right] + \frac{2\beta^{2}}{a_{8}}\mathbb{E}\left[\left\|\hat{x}_{i}^{t}\right\|^{2}\right].$$

$$(21)$$

Then, we estimate an upper bound on the first term on the right-hand side of equation 21 as follows:

$$\mathbb{E}\left[\left\langle \eta_{j}^{t}\overline{y}_{1}^{t}, \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] = \mathbb{E}\left[\eta_{j}^{t}\left\langle \overline{y}_{1}^{t} - \nabla f(\overline{x}^{t}), \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] \\
+ \mathbb{E}\left[\eta_{j}^{t}\left\langle \nabla f(\overline{x}^{t}), \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] \\
\leq \mathbb{E}\left[\eta_{j}^{t}\left\langle \overline{y}_{1}^{t} - \nabla f(\overline{x}^{t}), \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t})\right\rangle\right] + \frac{1}{\beta}\mathbb{E}\left[\eta_{j}^{t}(f(\overline{x}^{t-1}) - f(\overline{x}^{t}))\right], \tag{22}$$

with $\nabla f(\overline{x}^t) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(\overline{x}^t)$, where we have used the convexity of the function f(x) and the relationship $\frac{\eta_i^t}{\eta_i^{t-1}} \leq \frac{1}{\beta}$ for any given t in the last inequality.

The first term on the right-hand side of equation 22 can be bounded by

$$\mathbb{E}\left[\left\langle \overline{y}_{1}^{t} - \nabla f(\overline{x}^{t}), \frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}(\overline{x}^{t-1} - \overline{x}^{t}) \right\rangle\right] \\
\leq \frac{1}{2a_{6}} \mathbb{E}\left[\|\overline{y}_{1}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right] + \frac{a_{6}}{2} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right] \\
= \frac{1}{2a_{6}} \mathbb{E}\left[\left\|\frac{1}{m} \sum_{k=1}^{m} (\nabla f_{k}(x_{k}^{t}) - \nabla f_{k}(\overline{x}^{t}))\right\|^{2}\right] + \frac{a_{6}}{2} \mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2} \|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right] + \frac{\sigma^{2}}{2|\mathcal{B}|ma_{6}} \\
\leq \frac{1}{2ma_{6}} \sum_{k=1}^{m} \mathbb{E}\left[\|\nabla f_{k}(x_{k}^{t}) - \nabla f_{k}(\overline{x}^{t})\|^{2}\right] + \frac{a_{6}\beta^{2}}{2} \mathbb{E}\left[\|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right] + \frac{\sigma^{2}}{2|\mathcal{B}|ma_{6}} \\
\leq \frac{L^{2}}{2ma_{6}} \sum_{k=1}^{m} \mathbb{E}\left[\|\hat{x}_{k}^{t}\|^{2}\right] + \frac{a_{6}\beta^{2}}{2} \mathbb{E}\left[\|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right] + \frac{\sigma^{2}}{2|\mathcal{B}|ma_{6}}, \tag{23}$$

for any $a_6 \in \mathbb{R}^+$, where L is from the L-smoothness of f_i given in Assumption 1. Here, we have used the inequality $\langle a,b \rangle \leq \frac{1}{2\alpha}\|a\|^2 + \frac{\alpha}{2}\|b\|^2$ for any $\alpha>0$ and $a,\ b\in\mathbb{R}^n$ in the first inequality, and the L-smoothness of $f_i(x)$ in the last inequality.

Combining equation 20 and equation 23, we obtain the following inequality:

$$2\mathbb{E}\left[\left\langle \overline{\eta}^{t} \overline{y}_{1}^{t}, \overline{\eta}^{t} \overline{y}_{1}^{t-1} \right\rangle\right] \leq 2a_{8} \mathbb{E}\left[\left\| \overline{\eta}^{t} \overline{y}_{1}^{t} \right\|^{2}\right] + \frac{\eta_{\max}^{2}}{a_{8} m} \sum_{i=1}^{m} \mathbb{E}\left[\left\| \hat{y}_{1,i}^{t-1} \right\|^{2}\right] + \frac{a_{6} \beta^{2}}{m} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t} \| \overline{x}^{t-1} - \overline{x}^{t} \|^{2}\right] + \frac{2\beta^{2}}{a_{8} m} \sum_{i=1}^{m} \mathbb{E}\left[\left\| \hat{x}_{i}^{t-1} \right\|^{2}\right] + \frac{1}{m} \sum_{i=1}^{m} \left(\frac{L^{2} \eta_{\max}}{a_{6}} + \frac{2\beta^{2}}{a_{8}}\right) \mathbb{E}\left[\left\| \hat{x}_{i}^{t} \right\|^{2}\right] + \frac{2}{m\beta} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t} (f(\overline{x}^{t-1}) - f(\overline{x}^{t}))\right].$$

$$(24)$$

Substituting equation 10, equation 19, and equation 24 into equation 8, we obtain

$$(1 - 2a_{8})\mathbb{E}\left[\|\overline{\eta}^{t}\overline{y}_{1}^{t}\|^{2}\right]$$

$$\leq \frac{1}{m}\sum_{i=1}^{m}(1 + a_{2})\mathbb{E}\left[(\eta_{i}^{t}L_{i}^{t})^{2}\|x_{i}^{t} - x_{i}^{t-1}\|^{2}\right] + \frac{a_{6}\beta^{2}}{m}\sum_{i=1}^{m}\mathbb{E}\left[\eta_{i}^{t}\|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right]$$

$$+ \frac{2\beta}{m}\sum_{i=1}^{m}\mathbb{E}\left[\eta_{i}^{t}(f(\overline{x}^{t-1}) - f(\overline{x}^{t}))\right]$$

$$- (1 - a_{3})(1 - a_{4})(1 - a_{5})\frac{1}{m^{2}}\sum_{i=1}^{m}\mathbb{E}\left[\left(\frac{\eta_{i}^{t}}{\eta_{i}^{t-1}}\right)^{2}\|x_{i}^{t} - x_{i}^{t-1}\|^{2}\right] + \delta_{1}^{t},$$

$$(25)$$

where the constant δ_1^t is given by

$$\delta_{1}^{t} = \frac{1}{m} \sum_{i=1}^{m} \left(\frac{L^{2} \eta_{\text{max}}}{a_{6}} + \frac{2\beta^{2}}{a_{8}} + b_{1} \right) \mathbb{E} \left[\|\hat{x}_{i}^{t}\|^{2} \right] + \frac{1}{m} \sum_{i=1}^{m} \left(\frac{2\beta^{2}}{a_{8}} + b_{1} \right) \mathbb{E} \left[\|\hat{x}_{i}^{t-1}\|^{2} \right]$$

$$+ \frac{1}{m} \sum_{i=1}^{m} \eta_{\text{max}}^{2} \left(2 + \frac{2}{a_{2}} \right) \mathbb{E} \left[\|\hat{y}_{1,i}^{t}\|^{2} \right] + \frac{1}{m} \sum_{i=1}^{m} \left(\frac{1}{a_{3}} + 1 + \frac{1}{a_{8}} + \frac{2}{a_{2}} \right) \mathbb{E} \left[\eta_{\text{max}}^{2} \|\hat{y}_{1,i}^{t-1}\|^{2} \right]$$

$$+ \frac{\sigma^{2}}{2ma_{6}}.$$

$$(26)$$

By choosing some appropriate a_2 , a_3 , a_4 , and a_5 such that the following inequality holds:

$$(1+a_2)(\eta_i^t L_i^t)^2 - (1-a_3)(1-a_4)(1-a_5)\frac{(\eta_i^t)^2}{m(\eta_i^{t-1})^2} \le \frac{49}{100},$$

we have the following relationship:

$$\eta_i^t \le \frac{7\sqrt{r}}{10} \frac{\eta_i^{t-1}}{\sqrt{[(\eta_i^{t-1}L_i^t)^2 - 1]_+}}. (27)$$

Based on the above analysis, we can rewrite equation 25 as follows:

$$(1 - 2a_8)\mathbb{E}\left[\|\overline{\eta}^t \overline{y}_1^t\|^2\right] \le \frac{49}{100m} \sum_{i=1}^m \mathbb{E}\left[\|x_i^t - x_i^{t-1}\|^2\right] + \frac{a_6 \beta^2}{m} \sum_{i=1}^m \mathbb{E}\left[\eta_i^t \|\overline{x}^{t-1} - \overline{x}^t\|^2\right] + \frac{2\beta}{m} \sum_{i=1}^m \mathbb{E}\left[\eta_i^t (f(\overline{x}^{t-1}) - f(\overline{x}^t))\right] + \delta_1^t.$$
(28)

By using the inequality $\|a+b\|^2 \le (1+\frac{1}{\alpha})\|a\|^2 + (1+\alpha)\|b\|^2$ for any $\alpha>0$ and $a,\ b\in\mathbb{R}^n$, and the equation $x_i^t=\overline{x}^t+\hat{x}_i^t$, we can rewrite equation 28 as follows:

$$(1 - 2a_8)\mathbb{E}\left[\|\overline{\eta}^t \overline{y}_1^t\|^2\right] \le \left(\frac{49(1 + a_7)}{100} + \frac{a_6 \beta^2}{m} \sum_{i=1}^m \eta_i^t\right) \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] + \frac{2}{m\beta} \sum_{i=1}^m \mathbb{E}\left[\eta_i^t (f(\overline{x}^{t-1}) - f(\overline{x}^t))\right] + \delta_1^t + \frac{49(1 + a_7)}{100ma_7} \sum_{i=1}^m \mathbb{E}\left[\|\hat{x}_i^t - \hat{x}_i^{t-1}\|^2\right] \le \left(\frac{49(1 + a_7)}{100} + a_6 \beta^2 \eta_{\text{max}}\right) \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] + \frac{2}{m\beta} \sum_{i=1}^m \mathbb{E}\left[\eta_i^t (f(\overline{x}^{t-1}) - f(\overline{x}^t))\right] + \delta_2^t,$$
(29)

for any $a_7 \in \mathbb{R}^+$, where the constant δ_2^t is given by $\delta_2^t = \delta_1^t + \frac{49(1+a_7)}{50ma_7} \sum_{i=1}^m \mathbb{E}\left[\|\hat{x}_i^t\|^2 + \|\hat{x}_i^{t-1}\|^2\right]$.

By substituting equation 29 into equation 7, we obtain

$$\mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^{t}\|^{2}\right] \\
\leq c_{1}\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \frac{2(1+a_{1})}{m\beta(1-2a_{8})} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t}(f(\overline{x}^{t-1}) - f(\overline{x}^{t}))\right] + \delta_{3}^{t}, \tag{30}$$

with $\delta_3^t = \frac{1+a_1}{1-2a_8} \delta_2^t + (1+\frac{1}{a_1}) \frac{1}{m} \sum_{i=1}^m \eta_{\max}^2 \|\hat{y}_{1,i}^t\|^2$ and $c_1 = \frac{1+a_1}{1-2a_8} (\frac{49(1+a_7)}{100} + a_6 \beta^2 \eta_{\max})$. To ensure that $c_1 < \frac{1}{2}$, we select sufficiently small a_6 , a_7 , and a_8 such that the following relations hold:

$$\frac{49(1+a_7)}{100} + a_6 \beta^2 \eta_{\text{max}} = \frac{495}{1000},$$

$$a_8 = \frac{1-124a_1}{250} > 0,$$
(31)

where a_8 exists due to $a_1 < \frac{1}{124}$ given in the lemma statement. Then, we have

$$c_1 \le \frac{495}{1000} \times \frac{1+a_1}{1-2a_8} < \frac{1}{2},\tag{32}$$

where we have used the relationship $\frac{495(1+a_1)}{500} < 1 - 2a_8$.

Substituting the second equation in equation 31 and equation 32 into equation 30, we obtain

$$\mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^{t}\|^{2}\right] \\
\leq \frac{495}{496 \times 2} \mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \frac{125}{62m\beta} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t}(f(\overline{x}^{t-1}) - f(\overline{x}^{t}))\right] + \delta_{3}^{t}. \tag{33}$$

Multiplying both sides of equation 33 by C = 2 leads to

$$\mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] \le \frac{45}{46} \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] - \mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] + \frac{125}{31\beta} \mathbb{E}\left[\overline{\eta}^t (f(\overline{x}^{t-1}) - f(\overline{x}^t))\right] + 2\delta_3^t, \tag{34}$$

which implies Lemma 1.

 Lemma 2. Under Assumptions 1 and 2, the following inequality holds for Algorithm 1:

$$\mathbb{E}[\|\overline{x}^{t+1} - x^*\|^2] + \mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^t\|^2] + \left(2 + \frac{125\beta}{31}\right) \overline{\eta}^t \mathbb{E}[f(\overline{x}^t) - f(x^*)] \\
\leq \left(1 - \frac{\mu\eta_{\min}}{2} + a_9\mu\eta_{\max}\right) \mathbb{E}[\|\overline{x}^t - x^*\|^2] + \frac{45}{46} \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] \\
+ \gamma \left(2 + \frac{125\beta}{31}\right) \mathbb{E}\left[\overline{\eta}^{t-1}(f(\overline{x}^{t-1}) - f(x^*))\right] + \delta_5^t, \tag{35}$$

for any $\gamma \in (0,1)$, where the constant δ_5^t is given by

$$\begin{split} \delta_{5}^{t} &= \frac{125}{62m} \left(\frac{L^{2} \eta_{\text{max}}}{a_{6}} + \frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1 + a_{7})}{50a_{7}} + \frac{124L^{2} \eta_{\text{max}}}{125a_{9}\mu} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t}\|^{2}] \\ &+ \frac{125}{62m} \left(\frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1 + a_{7})}{50a_{7}} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t-1}\|^{2}] \\ &+ \frac{\eta_{\text{max}}^{2}}{m} \left(\frac{187}{31} + \frac{125}{31a_{2}} + \frac{2}{a_{1}} + \frac{4}{\mu\eta_{\text{min}}} \right) \sum_{i=1}^{m} \left(\mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^{2} + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^{2}] \right) \\ &+ \frac{125\eta_{\text{max}}^{2}}{62m} \left(\frac{1}{a_{3}} + 1 + \frac{1}{a_{8}} + \frac{2}{a_{2}} \right) \sum_{i=1}^{m} \mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^{2}] + \frac{2\eta_{\text{max}}L^{2}\sigma^{2}}{a_{9}\mu} \\ &+ 2\left(\left(1 - \frac{1}{a_{3}} \right) \left(2 + \frac{2\eta_{\text{max}}}{\eta_{\text{min}}} + \frac{\eta_{\text{min}}}{2\eta_{\text{max}}} \right) + \frac{1}{a_{6}m} \right) \frac{\sigma^{2}}{|\mathcal{B}|}. \end{split}$$

Proof. According to the dynamics of x_i^t in Algorithm 1, we have

$$\mathbb{E}\left[\|\overline{x}^{t+1} - x^*\|^2\right] = \mathbb{E}\left[\|\overline{x}^t - \overline{\eta^t y^t} - x^*\|^2\right]$$

$$= \mathbb{E}\left[\|\overline{x}^t - x^*\|^2\right] + \mathbb{E}\left[\|\overline{\eta^t y^t}\|^2\right] - 2\mathbb{E}\left[\langle \overline{x}^t - x^*, \overline{\eta^t y^t}\rangle\right]$$

$$= \mathbb{E}\left[\|\overline{x}^t - x^*\|^2\right] + \mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] - 2\mathbb{E}\left[\langle \overline{x}^t - x^*, \overline{\eta^t y^t}\rangle\right],$$
(37)

with $\overline{\eta^t y^t} := \frac{1}{N} \sum_{i=1}^N \eta_i^t y_{1,i}^t$.

The third term on the right-hand side of equation 37 satisfies:

$$2\mathbb{E}\left[\langle \overline{x}^{t} - x^{*}, \overline{\eta^{t}} \overline{y^{t}} \rangle\right] = -2\langle \overline{x}^{t} - x^{*}, \frac{1}{m} \sum_{i=1}^{m} \eta_{i}^{t} y_{1,i}^{t} \rangle$$

$$\leq -2\mathbb{E}\left[\overline{\eta}^{t} (f(\overline{x}^{t}) - f(x^{*})) - \mu \eta_{\min} \|\overline{x}^{t} - x^{*}\|^{2}\right]$$

$$-2\mathbb{E}\left[\left\langle \overline{x}^{t} - x^{*}, \frac{1}{m} \sum_{i=1}^{m} \eta_{i}^{t} (y_{1,i}^{t} - \nabla f(\overline{x}^{t})) \right\rangle\right],$$
(38)

where in the derivation we have used the μ -strong convexity of f(x) and $\eta_{\min} = \min_{i \in [m], t \in \mathbb{N}^+} \eta_i^t$. Furthermore, since the function f is L-smoothness, the minimum of the stepsizes exists.

By using the Cauchy-Schwarz inequality, the third term on the right-hand side of inequality equation 38 satisfies

$$\begin{split} &-2\mathbb{E}\left[\left\langle \overline{x}^t - x^*, \frac{1}{m}\sum_{i=1}^m \eta_i^t(y_{1,i}^t - \nabla f(\overline{x}^t))\right\rangle\right] \\ &= -\frac{2}{m}\sum_{i=1}^m \mathbb{E}\left[\left\langle \sqrt{\eta_i^t}(\overline{x}^t - x^*), \sqrt{\eta_i^t}(y_{1,i}^t - \nabla f(\overline{x}^t))\right\rangle\right] \\ &\leq \frac{2}{m}\sqrt{\left(\sum_{i=1}^m \mathbb{E}\left[\left\|\sqrt{\eta_i^t}(\overline{x}^t - x^*)\right\|^2\right]\right)\left(\sum_{i=1}^m \mathbb{E}\left[\left\|\sqrt{\eta_i^t}(y_{1,i}^t - \nabla f(\overline{x}^t)\right\|^2\right)\right]}. \end{split}$$

By applying the inequality $2\langle a,b\rangle \leq \alpha \|a\|^2 + \frac{1}{\alpha} \|b\|^2$ for any $\alpha > 0$ and $a, b \in \mathbb{R}^n$ to equation 39, we obtain

$$-2\mathbb{E}\left[\left\langle \overline{x}^{t} - x^{*}, \frac{1}{m} \sum_{i=1}^{m} \eta_{i}^{t}(y_{1,i}^{t} - \nabla f(\overline{x}^{t})) \right\rangle\right]$$

$$\leq \frac{a_{9}\mu}{m} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t} \| \overline{x}^{t} - x^{*} \|^{2}\right] + \frac{1}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\eta_{i}^{t} \| y_{1,i}^{t} - \nabla f(\overline{x}^{t}) \|^{2}\right]$$

$$\leq a_{9}\mu \eta_{\max} \mathbb{E}\left[\|\overline{x}^{t} - x^{*}\|^{2}\right] + \frac{\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|y_{1,i}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right].$$

The second term on the right-hand side of equation 39 satisfies

$$\frac{\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|y_{1,i}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right] = \frac{\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|y_{1,i}^{t} - \overline{y}_{1}^{t} + \overline{y}_{1}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right] \\
\leq \frac{2\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|y_{1,i}^{t} - \overline{y}_{1}^{t}\|^{2}\right] + \frac{2\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|\overline{y}_{1}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right].$$
(39)

By using the Lipschitz continuity of ∇f from Assulption 1, we have

$$\begin{split} &\sum_{i=1}^m \mathbb{E}\left[\|\overline{y}_1^t - \nabla f(\overline{x}^t)\|^2\right] \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}\left[\|g_i^t(x_i^t) - \nabla f_i(\overline{x}^t)\|^2\right] \\ &= \frac{\sigma^2}{B} + \frac{1}{m} \sum_{i=1}^m \mathbb{E}\left[\|\nabla f_i^t(x_i^t) - \nabla f_i(\overline{x}^t)\|^2\right] \leq \frac{\sigma^2}{|\mathcal{B}|} + \frac{L^2}{m} \sum_{i=1}^m \mathbb{E}\left[\|\hat{x}_i^t\|^2\right]. \end{split}$$

Substituting equation 40 into equation 39 leads to

$$\frac{\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|y_{1,i}^{t} - \nabla f(\overline{x}^{t})\|^{2}\right] \\
\leq \frac{2\eta_{\max}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t}\|^{2}\right] + \frac{2\eta_{\max}L^{2}}{ma_{9}\mu} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2}\right] + \frac{2\eta_{\max}L^{2}\sigma^{2}}{|\mathcal{B}|a_{9}\mu}.$$
(40)

By substituting equation 38 to equation 40 into equation 37, we obtain

$$\mathbb{E}\left[\|\overline{x}^{t+1} - x^*\|^2\right] \le (1 - \mu \eta_{\min} + a_9 \mu \eta_{\max}) \mathbb{E}\left[\|\overline{x}^t - x^*\|^2\right] + \mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] - 2\mathbb{E}\left[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))\right] + \delta_4^t,$$

$$(41)$$

with
$$\delta_4^t = \frac{2\eta_{\max}}{ma_9\mu} \sum_{i=1}^m \mathbb{E}\left[\|\hat{y}_{1,i}^t\|^2\right] + \frac{2\eta_{\max}L^2}{ma_9\mu} \sum_{i=1}^m \mathbb{E}\left[\|\hat{x}_i^t\|^2\right] + \frac{2\eta_{\max}L^2\sigma^2}{|\mathcal{B}|a_9\mu}$$

By adding both sides of equation 4 in Lemma 1 and equation 41, we obtain

$$\mathbb{E}\left[\|\overline{x}^{t+1} - x^*\|^2\right] + \mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] + \left(2 + \frac{125\beta}{31}\right) \mathbb{E}\left[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))\right] \\
\leq (1 - \mu\eta_{\min} + a_9\mu\eta_{\max}) \mathbb{E}\left[\|\overline{x}^t - x^*\|^2\right] + \frac{45}{46} \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] \\
+ \frac{125\beta}{31} \mathbb{E}\left[\overline{\eta}^t (f(\overline{x}^{t-1}) - f(x^*))\right] + \delta_5^t, \tag{42}$$

with $\delta_5^t = \delta_4^t + 2\delta_3^t$.

By setting $\beta \in (1, 1.36)$ and using Line 16 in Algorithm 1, we have $\frac{125\beta}{31}\overline{\eta}^t \leq \frac{125\beta^2}{31}\overline{\eta}^{t-1} = \gamma_1\left(2 + \frac{125\beta}{31}\right)\overline{\eta}^{t-1}$ for some $\gamma_1 \in \left(0, \frac{91}{92}\right)$, which implies the following inequality:

$$\frac{125\beta}{31}\overline{\eta}^t \leq \gamma_1 \left(2 + \frac{125\beta}{31}\right)\overline{\eta}^{t-1} \quad \text{and} \quad \frac{125\beta^2}{31} \leq \gamma_1 \left(2 + \frac{125\beta}{31}\right). \tag{43}$$

By letting $a_9 = \frac{\eta_{\min}}{2\eta_{\max}}$, we have

$$\mathbb{E}\left[\|\overline{x}^{t+1} - x^*\|^2\right] + \mathbb{E}\left[\|\overline{x}^{t+1} - \overline{x}^t\|^2\right] + \left(2 + \frac{125\beta}{31}\right) \mathbb{E}\left[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))\right] \\
\leq \left(1 - \frac{\mu\eta_{\min}}{2}\right) \mathbb{E}\left[\|\overline{x}^t - x^*\|^2\right] + \frac{45}{46} \mathbb{E}\left[\|\overline{x}^t - \overline{x}^{t-1}\|^2\right] \\
+ \mathbb{E}\left[\gamma_1 \left(2 + \frac{125\beta}{31}\right) \overline{\eta}^{t-1} (f(\overline{x}^{t-1}) - f(x^*))\right] + \delta_5^t, \tag{44}$$

(47)

which proves Lemma 2.

Lemma 3. Under Assumptions 1 and 2, the following inequality holds for Algorithm 1:

$$\sum_{i=1}^{m} \mathbb{E}\left[\hat{x}_{i}^{t}\|^{2}\right] \leq \rho^{2M} \left(12\eta_{\max}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + (24\eta_{\max}^{2} L^{2} + 3) \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + 48m\eta_{\max}^{2} L^{2}\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \frac{12\eta_{\max}^{2}\sigma^{2}}{|\mathcal{B}|}\right), \tag{45}$$

$$\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t}\|^{2}\right] \leq \rho^{2M} \left(18L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2}\right] + 18L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + 3\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 18mL^{2}\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + 3\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 18mL^{2}\mathbb{E}\left[\|\hat{y}_{2,i}^{t-1}\|^{2}\right] + 3\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{2,i}^{t-1}\|^{2}\right] + 18L^{2}\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + 3\sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{2,i}^{t-1}\|^{2}\right] + 18mL^{2}\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \frac{36m\sigma^{2}}{|\mathcal{B}|}\right), \tag{47}$$

where $\rho < 1$ is from Assumption 2 and M is the number of inner-consensus-loop iterations from Algorithm 1.

Proof. According to Line 5 in Algorithm 1, we have

$$\mathbf{X}^{t}(q+1) = (W \otimes I_{n})\mathbf{X}^{t}(q), \ q = 0, 1, \dots, M-1,$$
 (48)

where $W \in \mathbb{R}^{m \times m}$ is the adjacency matrix given in Assumption 2. Since the relationship $\bar{x}^t(q) =$ $\frac{1}{m}\sum_{i=1}^{m}x_{i}^{t}(q)$ holds, we have

$$\overline{x}^{t}(q+1) = \frac{1}{m} \sum_{i=1}^{m} x_{i}^{t}(q+1) = \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} w_{ij} x_{j}^{t}(q) = \frac{1}{m} \sum_{j=1}^{m} x_{j}^{t}(q) = \overline{x}^{t}(q), \tag{49}$$

where we have used Assumption 2 in the derivation.

By using the definition $\bar{\mathbf{X}}^t = \operatorname{col}(\bar{x}^t, \dots, \bar{x}^t) \in \mathbb{R}^{mn}$ and equation 49, we have

$$\bar{\mathbf{X}}^t = (W \otimes I_n)\bar{\mathbf{X}}^t. \tag{50}$$

By defining $\Delta^t(q) \triangleq \mathbf{X}^t(q) - \bar{\mathbf{X}}^t$ and subtracting equation 50 from equation 48, we obtain

$$\Delta^{t}(q+1) = \mathbf{X}^{t}(q+1) - \mathbf{X}_{1}^{t} = (W \otimes I_{n})\Delta^{t}(q)$$

$$\Delta^{t}(q+1) = (W \otimes I_{n})\Delta^{t}(q).$$

Since W is a doubly stochastic matrix, there must exist an orthogonal matrix $\Phi \in \mathbb{R}^{m \times m}$ such that W satisfies the following transformation:

$$\Phi^{\top} W \Phi = \operatorname{diag}\{1, \lambda_2, \dots, \lambda_m\},\tag{51}$$

with $|\lambda_i| < 1$, i = 2, ..., m. The first column of Φ is given by $\frac{1}{\sqrt{m}} \mathbf{1}_n$, which corresponds to the eigenvalue 1 of W. By further considering the following transformation:

$$\Delta_1^t(q) = (\Phi^\top \otimes I_n) \Delta^t(q), \tag{52}$$

with $\Delta_1^t(q) = [\sigma_1^t(q); \sigma_2^t(q); \dots; \sigma_m^t(q)] \in \mathbb{R}^{mn}$, we have

$$\sigma_i^t(q) = \sum_{j=1}^m \Phi_{ij}^\top (x_j^t(q) - \overline{x}^t), \tag{53}$$

where Φ_{j}^{\top} denotes the element in the *i*th row and *j*th column of the matrix Φ^{\top} . By using $\sigma_1^t(q) = \frac{1}{\sqrt{m}} \sum_{j=1}^m (x_j^t(q) - \overline{x}^t) = \mathbf{0}$, equation 51 can be rewritten as follows:

$$\Delta_1^t(q+1) = (\operatorname{diag}\{1, \lambda_2, \dots, \lambda_m\} \otimes I_n) \Delta_1^t(q). \tag{54}$$

Since the relationship $\sigma_1^t(q) = 0$ holds, equation 54 implies

$$\sigma_i^t(q+1) = \lambda_i \sigma_i^t(q) \le \rho \sigma_i^t(q) \le \rho^{q+1} \sigma_i^t(0), \tag{55}$$

with $\rho = \max\{|\lambda_2|, \cdots, |\lambda_m|\} < 1$. According to equation 55, we have

$$\|\Delta^t(M)\|^2 \le \rho^{2M} \|\Delta^t(0)\|^2,\tag{56}$$

which further implies

$$\sum_{i=1}^{m} \|x_i^t - \overline{x}^t\|^2 \le \rho^{2M} \sum_{i=1}^{m} \|x_i^t(0) - \overline{x}^t\|^2.$$
 (57)

By using an argument similar to the derivation of equation 57, we obtain

$$\sum_{i=1}^{m} \|y_{1,i}^{t} - \overline{y}_{1}^{t}\|^{2} \leq \rho^{2M} \sum_{i=1}^{m} \|y_{1,i}^{t}(0) - \overline{y}_{1}^{t}\|^{2},$$

$$\sum_{i=1}^{m} \|y_{2,i}^{t} - \overline{y}_{2}^{t}\|^{2} \leq \rho^{2M} \sum_{i=1}^{m} \|y_{2,i}^{t}(0) - \overline{y}_{2}^{t}\|^{2}.$$
(58)

Using equation 57, we have

$$\sum_{i=1}^{m} \|x_{i}^{t} - \overline{x}^{t}\|^{2} \leq \rho^{2M} \sum_{i=1}^{m} \|x_{i}^{t}(0) - \overline{x}^{t}\|^{2}$$

$$= \rho^{2M} \sum_{i=1}^{m} \|x_{i}^{t}(0) - x_{i}^{t-1} + x_{i}^{t-1} - \overline{x}^{t-1} + \overline{x}^{t-1} - \overline{x}^{t}\|^{2}$$

$$\leq 3\rho^{2M} \left(\sum_{i=1}^{m} \|x_{i}^{t}(0) - x_{i}^{t-1}\|^{2} + \sum_{i=1}^{m} \|x_{i}^{t-1} - \overline{x}^{t-1}\|^{2} + \sum_{i=1}^{m} \|\overline{x}^{t-1} - \overline{x}^{t}\|^{2} \right)$$

$$= 3\rho^{2M} \left(\sum_{i=1}^{m} \|x_{i}^{t}(0) - x_{i}^{t-1}\|^{2} + \sum_{i=1}^{m} \|\hat{x}_{i}^{t-1}\|^{2} + m\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2} \right), \tag{59}$$

where we have used the relationship $\|a+b+c\|^2 \le 3\|a\|^2 + 3\|b\|^2 + 3\|c\|^2$ in the second inequality.

We estimate an upper bound on the first term on the right-hand side of equation 59 as follows:

$$\sum_{i=1}^{m} \|x_{i}^{t}(0) - x_{i}^{t-1}\|^{2} = \sum_{i=1}^{m} \|\eta_{i}^{t-1}y_{1,i}^{t-1}\|^{2} \le 2 \sum_{i=1}^{m} \|\eta_{i}^{t-1}\hat{y}_{1,i}^{t-1}\|^{2} + 2 \sum_{i=1}^{m} \|\eta_{i}^{t-1}\overline{y}_{1}^{t-1}\|^{2}
= 2\eta_{\max}^{2} \sum_{i=1}^{m} \|\hat{y}_{1,i}^{t-1}\|^{2} + 2m\eta_{\max}^{2} \left\| \frac{1}{m} \sum_{i=1}^{m} (g_{i}^{t-1}(x_{i}^{t-1}) - \nabla f(x^{*})) \right\|^{2}.$$
(60)

By using the following inequality and equation 60

$$\mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^{m}(g_{i}^{t-1}(x_{i}^{t-1}) - \nabla f(x^{*}))\right\|^{2}\right] \leq \frac{\sigma^{2}}{|\mathcal{B}|m} + \frac{1}{m}\mathbb{E}\left[\left\|\sum_{i=1}^{m}(\nabla f_{i}(x_{i}^{t-1}) - \nabla f_{i}(x^{*}))\right\|^{2}\right]$$
$$\leq \frac{\sigma^{2}}{|\mathcal{B}|m} + \frac{L^{2}}{m}\sum_{i=1}^{m}\mathbb{E}\left[\left\|x_{i}^{t-1} - x^{*}\right\|^{2}\right],$$

we obtain the following relationship:

$$\sum_{i=1}^{m} \mathbb{E}\left[\|x_{i}^{t}(0) - x_{i}^{t-1}\|^{2}\right] \\
1354 \qquad \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 2\eta_{\max}^{2} L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|x_{i}^{t-1} - x^{*}\|^{2}\right] + 2\eta_{\max}^{2} \sigma^{2} \\
1357 \qquad \leq 2\eta_{\max}^{2} \sum_{i=1}^{m} \|\hat{y}_{1,i}^{t-1}\|^{2} + 4\eta_{\max}^{2} L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] \\
+ 4m\eta_{\max}^{2} L^{2} \mathbb{E}\left[\|\overline{x}^{t-1} - x^{*}\|^{2}\right] + 2\eta_{\max}^{2} \sigma^{2} \\
1361 \qquad + 4m\eta_{\max}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 4\eta_{\max}^{2} L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + 8m\eta_{\max}^{2} L^{2} (\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] \\
1362 \qquad \leq 2\eta_{\max}^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 4\eta_{\max}^{2} L^{2} \sum_{i=1}^{m} \mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right] + 8m\eta_{\max}^{2} L^{2} (\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] \\
+ \mathbb{E}\left[\|\overline{x}^{t} - x^{*}\|^{2}\right]) + \frac{2\eta_{\max}^{2} \sigma^{2}}{|\mathcal{B}|}.$$
(61)

The third term on the right-hand side of equation 59 satisfies

$$m\|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}$$

$$= m\|\overline{\eta^{t-1}y^{t-1}}\|^{2} = m\|\frac{1}{m}\sum_{i=1}^{m}\eta_{i}^{t-1}y_{i}^{t-1}\|^{2}$$

$$\leq \sum_{i=1}^{m}\|\eta_{i}^{t-1}y_{i}^{t-1}\|^{2} \leq \eta_{\max}^{2}\sum_{i=1}^{m}\|y_{i}^{t-1}\|^{2} \leq 2\eta_{\max}^{2}\sum_{i=1}^{m}\|\hat{y}_{1,i}^{t-1}\|^{2} + 2m\eta_{\max}^{2}\|\overline{y}_{1}^{t-1}\|^{2}$$

$$= 2\eta_{\max}^{2}\left(\sum_{i=1}^{m}\|\hat{y}_{1,i}^{t-1}\|^{2} + m\left\|\frac{1}{m}\sum_{i=1}^{m}(g_{i}^{t-1}(x_{i}^{t-1}) - f(x^{*}))\right\|^{2}\right),$$
(62)

with $\overline{\eta^{t-1}y^{t-1}}=\frac{1}{m}\sum_{i=1}^m\eta_i^{t-1}y_{1,i}^{t-1}$. Substituting equation 60 into equation 63 leads to

$$m\mathbb{E}\left[\|\overline{x}^{t-1} - \overline{x}^{t}\|^{2}\right]$$

$$\leq 2\eta_{\max}^{2}\left(\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + L^{2}\sum_{i=1}^{m}\mathbb{E}\left[\|x_{i}^{t-1} - x^{*}\|^{2}\right]\right) + 2\eta_{\max}^{2}\sigma^{2}$$

$$\leq 2\eta_{\max}^{2}\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{y}_{1,i}^{t-1}\|^{2}\right] + 4\eta_{\max}^{2}L^{2}\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right]$$

$$+8m\eta_{\max}^{2}L^{2}\left(\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \mathbb{E}\left[\|\overline{x}^{t} - x^{*}\|^{2}\right]\right) + \frac{2\eta_{\max}^{2}\sigma^{2}}{|\mathcal{B}|}.$$

$$(63)$$

By substituting equation 61 and equation 63 into equation 59, we arrive at equation 45.

By using equation 58, we have

$$\sum_{i=1}^{m} \|y_{1,i}^{t} - \overline{y}_{1}^{t}\|^{2}
\leq \rho^{2M} \sum_{i=1}^{m} \|y_{1,i}^{t}(0) - \overline{y}_{1}^{t}\|^{2}
= \rho^{2M} \sum_{i=1}^{m} \|y_{1,i}^{t}(0) - y_{1,i}^{t-1} + y_{1,i}^{t-1} - \overline{y}_{1}^{t-1} + \overline{y}_{1}^{t-1} - \overline{y}_{1}^{t}\|^{2}
\leq 3\rho^{2M} \left(\sum_{i=1}^{m} \|y_{1,i}^{t}(0) - y_{1,i}^{t-1}\|^{2} + \sum_{i=1}^{m} \|\hat{y}_{1,i}^{t-1}\|^{2} + m\|\overline{y}_{1}^{t} - \overline{y}_{1}^{t-1}\|^{2} \right).$$
(64)

The first term on the right-hand side of equation 64 satisfies

$$\sum_{i=1}^{m} \mathbb{E} \left[\| y_{1,i}^{t}(0) - y_{1,i}^{t-1} \|^{2} \right] = \sum_{i=1}^{m} \mathbb{E} \left[\| g_{i}^{t-1}(x_{i}^{t}) - g_{i}^{t-2}(x_{i}^{t-1}) \|^{2} \right]$$

$$\leq \sum_{i=1}^{m} \mathbb{E} \left[\| \nabla f_{i}(x_{i}^{t}) - \nabla f_{i}(x_{i}^{t-1}) \|^{2} \right] + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

$$\leq L^{2} \sum_{i=1}^{m} \mathbb{E} \left[\| x_{i}^{t} - x_{i}^{t-1} \|^{2} \right] + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

$$\leq L^{2} \sum_{i=1}^{m} \mathbb{E} \left[\| x_{i}^{t} - \overline{x}^{t-1} \|^{2} \right] + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

$$= L^{2} \sum_{i=1}^{m} \mathbb{E} \left[\| x_{i}^{t} - \overline{x}^{t} + \overline{x}^{t} - \overline{x}^{t-1} + \overline{x}^{t-1} - x_{i}^{t-1} \|^{2} \right] + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

$$\leq 3L^{2} \left(\sum_{i=1}^{m} \mathbb{E} \left[\| \hat{x}_{i}^{t} \|^{2} \right] + m\mathbb{E} \left[\| \overline{x}^{t} - \overline{x}^{t-1} \|^{2} \right] + \sum_{i=1}^{m} \mathbb{E} \left[\| \hat{x}_{i}^{t-1} \|^{2} \right] \right) + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

$$\leq 3L^{2} \left(\sum_{i=1}^{m} \mathbb{E} \left[\| \hat{x}_{i}^{t} \|^{2} \right] + m\mathbb{E} \left[\| \overline{x}^{t} - \overline{x}^{t-1} \|^{2} \right] + \sum_{i=1}^{m} \mathbb{E} \left[\| \hat{x}_{i}^{t-1} \|^{2} \right] \right) + \frac{2m\sigma^{2}}{|\mathcal{B}|}$$

The third term on the right-hand side of inequality equation 64 satisfies

$$\begin{split} m\mathbb{E}\left[\|\overline{y}_{1}^{t} - \overline{y}_{1}^{t-1}\|^{2}\right] &= m\mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^{m}(g_{i}^{t-1}(x_{i}^{t}) - g_{i}^{t-2}(x_{i}^{t-1}))\right\|^{2}\right] \\ &= m\mathbb{E}\left[\left\|\frac{1}{m}\sum_{i=1}^{m}(\nabla f_{i}(x_{i}^{t}) - \nabla f_{i}(x_{i}^{t-1}))\right\|^{2}\right] + \frac{2\sigma^{2}}{|\mathcal{B}|} \leq L^{2}\sum_{i=1}^{m}\mathbb{E}\left[\|x_{i}^{t} - x_{i}^{t-1}\|^{2}\right] + \frac{2\sigma^{2}}{|\mathcal{B}|} \\ &\leq 3L^{2}\left(\sum_{i=1}^{m}\mathbb{E}\left[\|\hat{x}_{i}^{t}\|^{2}\right] + m\mathbb{E}\left[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}\right] + \sum_{i=1}^{m}\mathbb{E}\left[\|\hat{x}_{i}^{t-1}\|^{2}\right)\right) + \frac{2\sigma^{2}}{|\mathcal{B}|}. \end{split}$$

By substituting equation 65 and equation 66 into equation 64, we arrive at equation 46.

The proof of equation 47 is similar to the derivation of equation 46, and thus is omitted here. \Box

B.2 Proof of Theorem 1

Proof of theorem 1: By setting $\alpha_1=1-\frac{\mu\eta_{\min}}{2},\ \alpha_2=\frac{45}{46},\ \alpha_3=\frac{125}{62m}(\frac{L^2\eta_{\max}}{a_6}+\frac{2\beta^2}{a_8}+b_1+\frac{49(1+a_7)}{50a_7}+\frac{124L^2\eta_{\max}}{125a_9\mu}),\ \alpha_4=\frac{125}{62m}(\frac{2\beta^2}{a_8}+b_1+\frac{49(1+a_7)}{50a_7}),\ \alpha_5=\frac{\eta_{\max}^2}{m}(\frac{187}{31}+\frac{125}{31a_2}+\frac{2}{a_1}+\frac{2}{a_9\mu\eta_{\max}}),\ \alpha_6=\frac{125\eta_{\max}^2}{62m}(\frac{1}{a_3}+1+\frac{1}{a_8}+\frac{2}{a_2}),\ \text{and}\ \alpha_7:=2\left(\left(1-\frac{1}{a_3}\right)\left(2+\frac{2\eta_{\max}}{\eta_{\min}}+\frac{\eta_{\min}}{2\eta_{\max}}\right)+\frac{1}{a_6m}\right)+\frac{2\eta_{\max}L^2}{a_9\mu},$ equation 35 can be rewritten as follows:

$$\mathbb{E}[\|\overline{x}^{t+1} - x^*\|^2] + \mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^t\|^2] + \left(2 + \frac{125\beta}{31}\right) \mathbb{E}[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))] \\
\leq \alpha_1 \mathbb{E}[\|\overline{x}^t - x^*\|^2] + \alpha_2 \mathbb{E}[\|\overline{x}^t - \overline{x}^{t-1}\|^2] + \gamma \left(2 + \frac{125\beta}{31}\right) \mathbb{E}[\overline{\eta}^{t-1} (f(\overline{x}^{t-1}) - f(x^*))] \\
+ \alpha_3 \sum_{i=1}^m \mathbb{E}[\|\hat{x}_i^t\|^2] + \alpha_4 \sum_{i=1}^m \mathbb{E}[\|\hat{x}_i^{t-1}\|^2] + \alpha_5 \sum_{i=1}^m \left(\mathbb{E}[\|\hat{y}_{1,i}^t\|^2] + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^2]\right) \\
+ \alpha_6 \sum_{i=1}^m \mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^2] + \frac{\alpha_7 \sigma^2}{|\mathcal{B}|}.$$
(66)

By using an argument similar to the derivation of equation 66, equation 45 and equation 46 can be rewritten as follows:

$$\sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t}\|^{2}] \leq \rho^{2M} \left(\alpha_{8} \sum_{i=1}^{m} \mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^{2}] + \alpha_{9} \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t-1}\|^{2}] \right) + \alpha_{10} \mathbb{E}[\|\overline{x}^{t} - x^{*}\|^{2}] + \alpha_{10} \mathbb{E}[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}] + \frac{\alpha_{8}\sigma^{2}}{|\mathcal{B}|} ,$$

$$\sum_{i=1}^{m} \left(\mathbb{E}[\|\hat{y}_{1,i}^{t}\|^{2}] + \mathbb{E}[\|\hat{y}_{2,i}^{t}\|^{2}] \right) \leq \rho^{2M} \left(\alpha_{11} \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t}\|^{2}] + \alpha_{11} \sum_{i=1}^{m} \mathbb{E}[\|\hat{x}_{i}^{t-1}\|^{2}] \right) + \alpha_{12} \sum_{i=1}^{m} \left(\mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^{2}] + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^{2}] \right) + \alpha_{13} \mathbb{E}[\|\overline{x}^{t} - \overline{x}^{t-1}\|^{2}] + \frac{\alpha_{14}\sigma^{2}}{|\mathcal{B}|},$$
(68)

with $\alpha_8=12\eta_{\max}^2$, $\alpha_9=24\eta_{\max}^2L^2+3$, $\alpha_{10}=48m\eta_{\max}^2L^2$, $\alpha_{11}=18L^2$, $\alpha_{12}=3$, $\alpha_{13}=18mL^2>0$, and $\alpha_{14}=72m$.

Multiplying inequalities equation 67 and equation 68 by K and then using equation 66 lead to

$$\mathbb{E}[\|\overline{x}^{t+1} - x^*\|^2] + \mathbb{E}[\|\overline{x}^{t+1} - \overline{x}^t\|^2] + \left(2 + \frac{125\beta}{31}\right) \mathbb{E}[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))] \\
+ (K - \alpha_3 - \rho^{2M}\alpha_{10}K) \sum_{i=1}^m \mathbb{E}[\|\hat{x}_i^t\|^2] + (K - \alpha_5) \sum_{i=1}^m \left(\mathbb{E}[\|\hat{y}_{1,i}^t\|^2] + \mathbb{E}[\|\hat{y}_{2,i}^t\|^2]\right) \\
\leq \left(\alpha_1 + \rho^{2M}\alpha_9K\right) \mathbb{E}[\|\overline{x}^t - x^*\|^2] + (\alpha_2 + \rho^{2M}K(\alpha_9 + \alpha_{12}))\mathbb{E}[\|\overline{x}^t - \overline{x}^{t-1}\|^2] \\
+ \gamma(2 + \frac{125\beta}{31})\mathbb{E}[\overline{\eta}^{t-1}(f(\overline{x}^{t-1}) - f(x^*))] + \left(\alpha_4 + \rho^{2M}K(\alpha_8 + \alpha_{10})\right) \sum_{i=1}^m \mathbb{E}[\|\hat{x}_i^{t-1}\|^2] \\
+ \left(\alpha_6 + \rho^{2M}K(\alpha_7 + \alpha_{11})\right) \sum_{i=1}^m \left(\mathbb{E}[\|\hat{y}_{1,i}^{t-1}\|^2] + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^2]\right) + (\alpha_7 + K\alpha_{14} + K\alpha_8) \frac{\sigma^2}{|\mathcal{B}|}. \tag{69}$$

By choosing sufficiently large K and M satisfying

$$K \ge \max\left\{\frac{2(92\alpha_4 + 91\alpha_3)}{91}, \frac{2(92\alpha_5 + 91\alpha_6)}{91}\right\},$$

$$M \ge \max\left\{\frac{\ln(1/2) - \ln(\alpha_8 + 2\alpha_{10})}{2\ln(\rho)}, \frac{\ln(1/2) - \ln(\alpha_7 + \alpha_{11})}{2\ln(\rho)}, \frac{\ln(1 - \alpha_1) - \ln(2) - \ln(\alpha_9 K)}{2\ln(\rho)}, \frac{\ln(1 - \alpha_2) - \ln(2) - \ln((\alpha_9 + \alpha_{12})K)}{2\ln(\rho)}\right\} \triangleq M_0,$$
(70)

the following inequalities always hold:

$$\alpha_{1} + \rho^{2M} \alpha_{9} K \leq \frac{1 + \alpha_{1}}{2} < 1,$$

$$\alpha_{2} + \rho^{2M} K(\alpha_{9} + \alpha_{12}) \leq \frac{1 + \alpha_{2}}{2} < 1,$$

$$\alpha_{4} + \rho^{2M} K(\alpha_{8} + \alpha_{10}) < \frac{91}{92} \left(K - \alpha_{3} - \rho^{2M} \alpha_{10} K \right),$$

$$\alpha_{6} + \rho^{2M} K(\alpha_{7} + \alpha_{11}) < \frac{91}{92} (K - \alpha_{5}).$$

$$(71)$$

According to the definition $\gamma = \max\left\{1 - \frac{\mu^2}{4L}, \frac{91}{92}\right\}$ in the theorem statement, we define an auxiliary function V(t+1) as follows:

$$V(t+1) = \mathbb{E}[\|\overline{x}^{t+1} - x^*\|^2 + \|\overline{x}^{t+1} - \overline{x}^t\|^2\} + \left(2 + \frac{125\beta}{31}\right) \mathbb{E}[\overline{\eta}^t (f(\overline{x}^t) - f(x^*))] + (K - \alpha_3 - \rho^{2M}\alpha_{10}K) \sum_{i=1}^m \mathbb{E}[\|\hat{x}_i^t\|^2\} + (K - \alpha_5) \left(\sum_{i=1}^m \mathbb{E}[\|\hat{y}_{1,i}^t\|^2] + \mathbb{E}[\|\hat{y}_{2,i}^{t-1}\|^2]\right).$$
(72)

Using equation 71, we obtain the following inequality:

1514
1515
$$V(t+1) \le \gamma V(t) + (\alpha_7 + K\alpha_{14} + K\alpha_8) \frac{\sigma^2}{|\mathcal{B}|},$$
1516

which is equivalent to

$$\left(V(t+1) - \frac{(\alpha_7 + K\alpha_{14} + K\alpha_8)\sigma^2}{(1-\gamma)|\mathcal{B}|}\right) \le \gamma \left(V(t) - \frac{(\alpha_7 + K\alpha_{14} + K\alpha_8)\sigma^2}{(1-\gamma)|\mathcal{B}|}\right). \tag{73}$$

Therefore, by using equation 73, we arrive at

$$V(t) \le \gamma^t V(0) + \frac{(\alpha_7 + K\alpha_{14} + K\alpha_8)\sigma^2}{(1 - \gamma)|\mathcal{B}|}.$$
 (74)

Moreover, since the relations $\overline{\eta}^{-1} = 0$, $\sum_{i=1}^{m} \|\hat{x}_{i}^{-1}\|^{2} = 0$ and $\sum_{i=1}^{m} \|\hat{y}_{1,i}^{-1}\|^{2} = 0$ hold, we have $V(0) = \|\overline{x}^{0} - x^{*}\|^{2} + \|\overline{x}^{0}\|^{2}$.

Furthermore, according to the definition of V(t) in equation 72, we arrive at

$$\mathbb{E}[\|x_{i}^{t} - x^{*}\|^{2}] = \mathbb{E}[\|x_{i}^{t} - \overline{x}^{t} + \overline{x}^{t} - x^{*}\|^{2}] \leq 2\mathbb{E}[\|\hat{x}_{i}^{t}\|^{2}] + 2\mathbb{E}[\|\overline{x}^{t} - x^{*}\|^{2}] \\
\leq \max \left\{ \frac{2}{K_{1} - \alpha_{3} - \rho^{2M_{2}}\alpha_{10}K_{1}}, 2\right] V(t) \\
\leq \max \left\{ \frac{2V(0)}{K_{1} - \alpha_{3} - \rho^{2M_{2}}\alpha_{10}K_{1}}, 2V(0)\right] \gamma^{t} \\
+ \left(\frac{(\alpha_{7} + K\alpha_{14} + K\alpha_{8})}{1 - \gamma} \max \left\{ \frac{2}{K_{1} - \alpha_{3} - \rho^{2M_{2}}\alpha_{10}K_{1}}, 2\right] \right) \frac{\sigma^{2}}{|\mathcal{B}|}, \tag{75}$$

which implies $\mathbb{E}\left[\|x_i^t - x^*\|^2\right] \leq \mathcal{O}\left(\gamma^t\right) + \mathcal{O}\left(\frac{\sigma^2}{|\mathcal{B}|}\right)$ and Theorem 1.

B.3 PROOF OF THEOREM 2

When accurate gradients are accessible to agents, Algorithm 1 reduces to the following algorithm.

Algorithm 2 Deterministic version of Algorithm 1 (from agent *i*'s perspective)

```
1550
                       1: Input: x_i^0 \in \mathbb{R}^n, y_i^0 = \nabla f_i(x_i^0), \eta_i^0 > 0, \beta \in (0, 1.36), r \in (0, 1), M \in \mathbb{N}^+, and T \in \mathbb{N}^+.
1551
                       2: for t = 0, 1, ..., T do
                                   x_i^{t+1}(0) = x_i^t - \eta_i^t y_i^t
\mathbf{for} \ q = 0, 1, \dots, M - 1 \ \mathbf{do}
x_i^{t+1}(q+1) = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{t+1}(q)
1552
1553
1554
                               end for x_i^{t+1} = x_i^{t+1}(M) y_i^{t+1}(0) = y_i^t + \nabla f_i(x_i^{t+1}) - \nabla f_i(x_i^t) for q = 0, 1, \dots, M-1 do y_i^{t+1}(q+1) = \sum_{j \in \mathcal{N}_i} w_{ij} y_j^{t+1}(q)
1555
1556
1557
1558
1559
1560
                   12: y_i^{t+1} = y_i^{t+1}(M)

13: L_i^{t+1} = \frac{\|y_i^{t+1} - y_i^t\|}{\|x_i^{t+1} - x_i^t\|}
1561
1562
                                   \eta_i^{t+1} = \min \left\{ \beta \eta_i^t, \frac{7\sqrt{r}}{10} \frac{\eta_i^t}{\sqrt{[(\eta_i^t L_i^{t+1})^2 - 1]_+}} \right\}
1564
1565
                    15: end for
```

Proof of Theorem 2: By using an argument similar to the derivation of equation 35, we obtain

$$\|\overline{x}^{t+1} - x^*\|^2 + \|\overline{x}^{t+1} - \overline{x}^t\|^2 + \left(2 + \frac{125\beta}{31}\right) \overline{\eta}^t f(\overline{x}^t) - f(x^*)$$

$$\leq (1 - \mu \eta_{\min} + a_9 \mu \eta_{\max}) \|\overline{x}^t - x^*\|^2 + \frac{45}{46} \|\overline{x}^t - \overline{x}^{t-1}\|^2$$

$$+ \gamma \left(2 + \frac{125\beta}{31}\right) \overline{\eta}^{t-1} (f(\overline{x}^{t-1}) - f(x^*)) + \delta_5^t, \tag{76}$$

for $\gamma \in (0,1)$, where the constant and δ_5^t is given by

$$\delta_{5}^{t} = \frac{125}{62m} \left(\frac{L^{2}\eta_{\text{max}}}{a_{6}} + \frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1+a_{7})}{50a_{7}} + \frac{124L^{2}\eta_{\text{max}}}{125a_{9}\mu} \right) \sum_{i=1}^{m} \|\hat{x}_{i}^{t}\|^{2}$$

$$+ \frac{125}{62m} \left(\frac{2\beta^{2}}{a_{8}} + b_{1} + \frac{49(1+a_{7})}{50a_{7}} \right) \sum_{i=1}^{m} \|\hat{x}_{i}^{t-1}\|^{2}$$

$$+ \frac{\eta_{\text{max}}^{2}}{m} \left(\frac{187}{31} + \frac{125}{31a_{2}} + \frac{2}{a_{1}} + \frac{2}{a_{9}\mu\eta_{\text{max}}} \right) \sum_{i=1}^{m} \|\hat{y}_{i}^{t}\|^{2}$$

$$+ \frac{125\eta_{\text{max}}^{2}}{62m} \left(\frac{1}{a_{3}} + 1 + \frac{1}{a_{8}} + \frac{2}{a_{2}} \right) \sum_{i=1}^{m} \|\hat{y}_{i}^{t-1}\|^{2}. \tag{77}$$

By using an argument similar to the derivations of equation 45 and equation 46, we have

$$\sum_{i=1}^{m} \|\hat{x}_{i}^{t}\|^{2} \leq \rho^{2M} \left(12\eta_{\max}^{2} \sum_{i=1}^{m} \|\hat{y}_{i}^{t-1}\|^{2} + (24\eta_{\max}^{2} L^{2} + 3) \sum_{i=1}^{m} \|\hat{x}_{i}^{t-1}\|^{2} + 48m\eta_{\max}^{2} L^{2} \|\overline{x}^{t} - x^{*}\|^{2} + 48m\eta_{\max}^{2} L^{2} \|\overline{x}^{t} - \overline{x}^{t-1}\|^{2} \right), \tag{78}$$

$$\sum_{i=1}^{m} \|\hat{y}_{1,i}^{t}\|^{2} \leq \rho^{2M} \left(18L^{2} \sum_{i=1}^{m} \|\hat{x}_{i}^{t}\|^{2} + 18L^{2} \sum_{i=1}^{m} \|\hat{x}_{i}^{t-1}\|^{2} + 3 \sum_{i=1}^{m} \|\hat{y}_{1,i}^{t-1}\|^{2} + 18mL^{2} \|\overline{x}^{t} - \overline{x}^{t-1}\|^{2} \right).$$
(79)

By using an argument similar to the derivation of equation 74 and constructing the following function:

$$V(t+1) = \|\overline{x}^{t+1} - x^*\|^2 + \|\overline{x}^{t+1} - \overline{x}^t\|^2 + \left(2 + \frac{125\beta}{31}\right) \overline{\eta}^t (f(\overline{x}^t) - f(x^*)) + (K - \alpha_3 - \rho^{2M} \alpha_{10} K) \sum_{i=1}^m \|\hat{x}_i^t\|^2 + (K - \alpha_5) \sum_{i=1}^m \|\hat{y}_i^t\|^2,$$
(80)

we obtain the following relationship:

$$V(t+1) \le \gamma V(t),\tag{81}$$

which implies $V(t) \leq \gamma^t V(0)$. Then, following an argument similar to the derivations of equation 75, we arrive at $\|x_i^t - x^*\|^2 \leq \mathcal{O}(\gamma^t)$, which proves Theorem 2.

B.4 PROOF OF COROLLARY 1

According to Theorem 1, the convergence rate of Algorithm 1 is $\mathcal{O}\left(\gamma^T\right) + \mathcal{O}\left(\frac{\delta^2}{|\mathcal{B}|}\right)$. Hence, to find an ϵ -optimal solution, the number of outer-loop iterations T needs to satisfy $T = \mathcal{O}(\log(\epsilon^{-1}))$. At each outer-loop iteration, Algorithm 1 requires $|\mathcal{B}|$ gradient evaluations at both $g_i^t(x_i^{t+1})$ and $g_i^t(x_i^t)$, resulting in a total of $2|\mathcal{B}|$ evaluations. Meanwhile, Lines 3, 8, and 9 in Algorithm 1 require M

gradient evaluations at $x_{i,1}^{t+1}(0)$, $y_{i,1}^{t+1}(0)$, and $y_{i,2}^t(0)$, Lines 5, 11, and 12 in Algorithm 1 require M gradient evaluations at $x_i^{t+1}(q)$, $y_{i,1}^{t+1}(q)$, and $y_{i,2}^t(q)$; and Lines 15 and 16 in Algorithm 1 each require one gradient evaluation at L_i^{t+1} and η_i^{t+1} , respectively. Based on the above discussion, we have that Algorithm 1 requires at most $2|\mathcal{B}| + 3M + 3$ gradient evaluations per outer-loop iteration t, leading to a computational complexity of $\mathcal{O}((2|\mathcal{B}|+3M+3)\log(\epsilon^{-1}))$ over T iterations. In the deterministic setting, Algorithm 1 reduces to Algorithm 2, which requires at most 2M+3 gradient evaluations per outer-loop iteration t, and thus has a computational complexity of $\mathcal{O}((2M+3)\log(\epsilon^{-1}))$ over T iterations.

C EXPERIMENTAL SETUPS AND ADDITIONAL EXPERIMENTAL RESULTS

C.1 BENCHMARK DATASETS

MNIST. The "MNIST" dataset is a benchmark dataset widely used in machine learning and computer vision (Deng, 2012). It typically consists of 70,000 grayscale images of handwritten digits (i.e., 0–9), with 60,000 used for training and 10,000 for testing. Each image has a size of 28×28 pixels, with the digit centered in the frame.

CIFAR-10. The "CIFAR-10" dataset consists of 60,000 color images of size 32×32 pixels in 10 classes, with 6,000 images per class (Krizhevsky et al., 2010). Among them, 50,000 images are used for training and 10,000 for testing. The dataset covers a diverse set of object categories, including airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks. Compared with the "MNIST" dataset, the "CIFAR-10" dataset poses a greater challenge due to its colored and natural images with larger intra-class variability.

CIFAR-100. The "CIFAR-100" dataset is a natural extension of the "CIFAR-10" dataset (DeVries & Taylor, 2017). It contains 60,000 color images of size 32×32 pixels, and spreads across 100 classes with 600 images per class. The 50,000 images are used for training and 10,000 for testing. However, due to its larger number of categories and the fine-grained nature of many classes, the "CIFAR-100" dataset is regarded as the most challenging dataset within the CIFAR series.

Mushrooms. The "Mushrooms" dataset is a classic benchmark dataset from the UCI Machine Learning Repository (Tutuncu et al., 2022). It contains 8,124 instances of gilled mushrooms, each described by 22 categorical attributes, such as cap shape, surface, and color. The prediction task is to classify each mushroom as either edible or poisonous. In this paper, we focus on l_2 -logistic regression on the "Mushrooms" dataset, as the task naturally fits into a binary classification problem.

C.2 EXPERIMENTAL SETUPS

Convolutional neural network (CNN) training. For the "MNIST" dataset, we trained a two-layer CNN. The first convolutional layer has 64 output channels with 3×3 kernels, stride 1, and padding 1, followed by batch normalization, LeakyReLU activation, and 2×2 max pooling. The second convolutional layer has 128 output channels with the same kernel configuration. The feature maps are then passed through adaptive average pooling to a 1×1 representation, flattened, and fed into a fully connected layer to produce the output classes. The model was trained with a batch size of 128 using the cross-entropy loss.

For the "CIFAR-10" dataset, we trained a four-layer CNN consisting of four convolutional layers with progressively increasing channel sizes of 32, 64, 128, and 256. Each convolution uses a 3×3 kernel with stride 1 and padding 1. To stabilize training and reduce spatial resolution, we employed batch normalization, a LeakyReLU activation, and 2×2 max pooling after every convolutional block. The resulting feature maps are aggregated by adaptive average pooling to a 1×1 representation, which is then flattened and passed to a fully connected layer to produce the final class predictions. The model was trained with a batch size of 128 using the cross-entropy loss.

For the "CIFAR-100" dataset, we trained a five-layer CNN with residual connections to enhance feature extraction. The network begins with a 32-channel convolutional layer (3×3 kernels, stride 1, padding 1), followed by batch normalization, LeakyReLU activation, and 2×2 max pooling. The subsequent convolutional blocks progressively increase the channels to 64, 128, 256, and 512. To enhance feature extraction, we introduced residual paths: one from the raw input through a 2×2

convolution with stride 2, another from the second block via a 2×2 convolution, and a direct path from the raw input via an 8×8 convolution. The model was trained with a batch size of 128 using the cross-entropy loss.

Logistic regression. For the logistic regression task using the "Mushrooms" dataset, we employed a single-layer linear model, which directly maps the 22 input features to two output logits corresponding to the classes. Training was conducted using the loss function given in equation 82.

C.3 ADDITIONAL EXPERIMENTAL RESULTS

In this section, we provide two additional experimental results: 1) the performance evaluation of Algorithm 1 on logistic regression with strongly convex and smooth loss functions; and 2) the comparison of Algorithm 1 and distributed ADAM in Nazari et al. (2022).

Logistic regression using the "Mushrooms" dataset. We evaluate the effectiveness of Algorithm 1 by using an l_2 -logistic regression classification problem on the "Mushrooms" dataset (Tutuncu et al., 2022). To ensure heterogeneous data distribution, we spread data samples among five agents according to their target values. Specifically, agents 1, 2, and 3 have samples with the target value of 0, while agents 4 and 5 have samples with the target value of 1. All agents cooperatively learn an optimal model parameter x^* to problem 1, in which the loss function of agent i is given by

$$l(x,\xi_i) = \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \left(-(1-b_{ij}) \ln \left(\frac{e^{x_1 a_{ij}}}{e^{x_1 a_{ij}} + e^{x_2 a_{ij}}} \right) - b_{ij} \ln \left(\frac{e^{x_2 a_{ij}}}{e^{x_1 a_{ij}} + e^{x_2 a_{ij}}} \right) + \frac{L_2}{2} ||x||^2 \right),$$
(82)

where $|\mathcal{B}|$ represents the number of sampled data points per iteration. In this experiment, we used a full batch setting, i.e., $|\mathcal{B}| = |\mathcal{D}_i|$ with \mathcal{D}_i denoting the local dataset of agent i. Here, $x = [x_1, x_2]^{\top}$ is the model parameter and the positive constant L_2 is a regularization parameter. It is clear that the loss function in equation 82 is strongly convex and smooth.

In this experiment, we compared the test accuracies of Algorithm 1 with existing distributed optimization algorithms, including distributed GD in Nedic & Ozdaglar (2009) and deterministic GT in Nedić et al. (2017). The stepsizes for distributed GD and deterministic GT are the same as those employed in our "MNIST" experiment in the main text (i.e., $\eta_i^t = \frac{0.1}{(1+t)^{0.5}}$ for distributed GD and $\eta_i = 0.1$ for deterministic GT). The training process spanned 250 iterations.

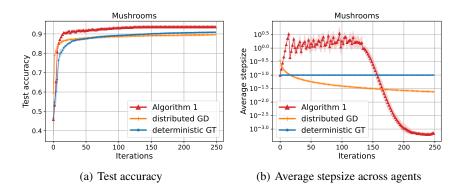


Figure 4: Test-accuracy and average-stepsize (across five agents) evolutions of Algorithm 1, distributed GD in Nedic & Ozdaglar (2009), and deterministic GT in Nedić et al. (2017). The 95% confidence intervals were computed from three independent runs with seeds 42, 1010, and 2024.

Fig. 4(a) shows that Algorithm 1 achieves the highest test accuracy and convergence speed compared with distributed GD and deterministic GT. This is because larger stepsizes is allowed in the early stages of Algorithm 1 than distributed GD and deterministic GT (as shown in Fig. 4(b)). Furthermore, Fig. 4 shows that Algorithm 1 exhibits stable convergence accuracy after 200 iterations. This result implies a clear stopping criterion for our algorithm, that is, by setting $\tau=10^{-3}$, each agent i can stop training once $|\eta_i^t|<\tau$.

Comparison of Algorithm 1 and distributed ADAM in Nazari et al. (2022). To compare the convergence accuracy of Algorithm 1 with the existing adaptive stepsize approach for distributed (online) learning, i.e., distributed ADAM in Nazari et al. (2022), we conducted additional experiments by comparing their test accuracies on image classification using the "CIFAR-10" dataset.

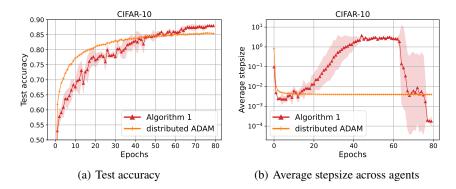


Figure 5: Test-accuracy and average-stepsize (across five agents) evolutions of Algorithm 1 and distributed ADAM (Nazari et al., 2022). The 95% confidence intervals were computed from three independent runs with random seeds 42, 1010, and 2024.

Fig. 5(a) shows that our Algorithm 1 outperforms distributed ADAM in terms of both test accuracy and steady-state performance. Furthermore, Fig. 5(b) indicates that the stepsize in distributed ADAM decays rapidly, which leads to a low convergence speed in the later stages of the algorithm.