

---

# FULLY ADAPTIVE STEPSIZES: WHICH SYSTEM BENEFIT MORE—CENTRALIZED OR DECENTRALIZED?<sup>\*†</sup>

---

A PREPRINT

**Diyako Ghaderyan**

Department of Information and Communications Engineering  
Aalto University  
diyako.ghaderyan@aalto.fi

**Stefan Werner**

Department of Electronic Systems, Norwegian University of Science and Technology (NTNU)  
Department of Information and Communications Engineering, Aalto University  
stefan.werner@ntnu.no

## ABSTRACT

In decentralized optimization, the choice of stepsize plays a critical role in algorithm performance. A common approach is to use a shared stepsize across all agents to ensure convergence. However, selecting an optimal stepsize often requires careful tuning, which can be time-consuming and may lead to slow convergence, especially when there is significant variation in the smoothness ( $L$ -smoothness) of local objective functions across agents. Individually tuning stepsizes per agent is also impractical, particularly in large-scale networks. To address these limitations, we propose AdGT, an adaptive gradient tracking method that enables each agent to adjust its stepsize based on the smoothness of its local objective. We prove that AdGT generates a sequence of iterates that converges to the optimal consensus solution. Through numerical experiments, we compare AdGT with fixed-stepsize gradient tracking methods and demonstrate its superior performance. Additionally, we compare AdGT with adaptive gradient descent (AdGD) in a centralized setting and observe that fully adaptive stepsizes offer greater benefits in decentralized networks than in centralized ones.

**Keywords** Adaptive stepsize, decentralized optimization, gradient tracking, consensus.

## 1 Introduction

Recent advances in artificial intelligence and communication technologies have accelerated the adoption of distributed optimization techniques to solve large-scale problems. These methods allow individual agents to perform local computations and share information with their neighbors, which is particularly advantageous when data is stored across distributed locations or partitioned across multiple servers to improve training efficiency. Key challenges in decentralized optimization include maintaining data privacy, ensuring algorithmic scalability, and achieving robustness to

---

<sup>\*</sup>This work was partially supported by the Research Council of Finland (Grant 354523) and the Research Council of Norway.

<sup>†</sup>This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

network and system-level disturbances. Distributed optimization has found widespread applications in areas such as machine learning [1, 2] and control systems [3]. A common objective in distributed optimization is to solve the consensus problem, which can be formulated as:

$$\mathbf{x}^* \in \underset{\mathbf{x} \in \mathbb{R}^P}{\operatorname{argmin}} f(\mathbf{x}) \triangleq \sum_{i=1}^n f_i(\mathbf{x}), \quad (1)$$

where the objective function  $f$  represents the summation of all individual cost functions  $\{f_i\}_{i=1}^n$ , where  $f_i : \mathbb{R}^P \rightarrow \mathbb{R}$  is the private function of agent  $i$ .

Building on the foundational work of [4], the authors of [5] introduced a distributed (sub)gradient method for solving the consensus problem in (1). When each local function  $f_i$  is convex, the method achieves sublinear convergence. This relatively slow rate is a known limitation of subgradient methods and is largely due to the use of diminishing stepsizes. To improve convergence, [6] proposed an accelerated method that assumes bounded and Lipschitz continuous gradients. By incorporating Nesterov’s acceleration technique, they obtained a convergence rate of  $\mathcal{O}(\log(k)/k^2)$ , which significantly outperforms standard subgradient methods. Under similar assumptions, the EXTRA [7] was developed for smooth convex functions with Lipschitz gradients. This method uses a fixed stepsize and introduces a correction term involving the difference between consecutive gradients. It achieves a convergence rate of  $\mathcal{O}(1/k)$ , which improves to linear convergence when the objective functions are strongly convex. In [8], the authors proposed the *Exact Diffusion* method that was shown to achieve linear convergence under standard assumptions.

In addition to gradient-based methods, a number of distributed algorithms have been developed based on the *alternating direction method of multipliers (ADMM)*. These approaches are known for their strong convergence guarantees and ability to handle more complex and structured optimization problems. Notable examples include [9, 10, 11, 12], which address general consensus optimization settings. Furthermore, several ADMM-based methods have been proposed for solving composite convex objectives [13, 14], as well as for problems involving complicated constraints [15, 16].

Gradient tracking methods (GT) [17, 18, 19, 20, 21] form a widely studied class of algorithms in decentralized optimization. These methods introduce an auxiliary variable to estimate the global gradient by locally tracking the average gradient across the network. This mechanism enables their performance to closely match that of centralized algorithms, which often exhibit linear convergence rates. One of the key strengths of gradient tracking methods lies in their flexibility and broad applicability across a wide range of distributed optimization problems. These techniques have been effectively adapted to a wide range of challenging settings, including directed and time-varying communication graphs, asynchronous computation, composite optimization problems, and even nonconvex objectives; see, for instance, [19, 22, 23, 24, 25, 26, 27, 28, 29, 30]. More recently, [31] proposed a new approach that incorporates gradient tracking in an implicit way by introducing a novel momentum term, specifically designed for use over directed networks. To improve convergence rates and computational complexity, a variety of accelerated decentralized gradient-based methods have been proposed; see, for example, [32, 33, 34, 35, 36, 37, 38, 39, 40, 41].

However, tuning a single fixed stepsize that works well for all agents is a difficult task, particularly in heterogeneous networks. This challenge becomes more severe when some agents hold more complex or ill-conditioned data, often associated with larger smoothness constants. In such cases, the stepsize must be chosen conservatively to ensure stability for all agents. However, this conservative choice typically results in a much smaller stepsize, as the global smoothness constant is dominated by the worst-case agents. As a result, the overall convergence rate is significantly slowed down. Moreover, several works have proposed using uncoordinated, agent-dependent stepsizes in gradient tracking methods over both undirected and directed graphs [21, 42, 43]. However, manually tuning individual stepsizes for each agent is often impractical, especially in large-scale decentralized systems. Therefore, these challenges motivate the following question:

*Is it possible to develop a decentralized optimization method where each agent automatically adjusts its local stepsize at every iteration in an adaptive manner, achieving exact and faster convergence?*

This question has been addressed in the centralized setting when function  $f$  is convex or strongly convex by the work of [44], which introduced an adaptive stepsize scheme for gradient descent based on the local smoothness of the objective function  $f$ , which is often substantially smaller than the global smoothness constant  $L$ . This approach updates the stepsize at each iteration based on the local smoothness of the objective, which often leads to faster and more reliable convergence compared to using a fixed, globally chosen stepsize. Several extensions have been proposed, including approaches that modify the stepsize rule or incorporate adaptive proximal gradient methods. See, for example, [45, 46, 47].

Furthermore, adaptive stepsize techniques have received increasing interest in federated learning. These include client adaptivity, adaptive optimizers at the server aimed at improving convergence rates, see, for example, [48, 49, 50]. In this context, a fully adaptive stepsize scheme was proposed in a recent study by [51], which—building on the ideas of [44]—enables each client to independently select its own stepsize based on the local smoothness of its objective function.

If the answer to the above question is yes, and we are indeed able to develop a fully adaptive decentralized optimization framework inspired by the ideas presented in [44], then a natural and significant follow-up question arises:

*Which benefits more from fully adaptive stepsizes: decentralized systems or centralized systems?*

In this paper, we address both questions above by introducing an adaptive gradient tracking method (AdGT). While we were in the process of establishing the theoretical guarantees of our algorithm, the work of [52] appeared, proposing a decentralized optimization approach that also employs an adaptive stepsize for based on local agent information, following the idea introduced in [44]. Although their method shares similarities with ours, it is developed under a slightly different setting, and their analysis only guarantees convergence to a neighborhood of the optimal solution. In contrast, our method is designed to achieve exact convergence to the optimal solution.

*Contributions:* This paper makes two main contributions.

1. We propose a decentralized algorithm, AdGT, that eliminates the need for manually selecting the stepsize. In most scenarios, AdGT adaptively adjusts the stepsize using only local information available to each agent. This feature facilitates the solution of distributed consensus optimization problems over undirected communication networks. Consequently, it significantly reduces the time and effort typically required for tuning stepsizes through trial and error to achieve fast convergence. We show that the iterates of AdGT with a fully adaptive stepsize converge to the optimal consensus solution.
2. Our algorithm uses adaptive stepsizes, which leads to much faster convergence, especially in situations where a few agents have large  $L$ -smoothness values, whereas most others have much smaller values. In these cases, existing decentralized methods with fixed stepsizes must choose a small value based on the agent with the largest  $L$ -smoothness. This slows down the overall convergence. In contrast, our algorithm allows each agent to adjust its own stepsize based on its local smoothness, which makes the algorithm more efficient and much faster in practice. We see this ability to automatically adapt stepsizes for each agent as a key contribution. This feature is not available in the centralized setting, such as in [44], and has not been properly studied in recent adaptive methods for federated learning [51] or decentralized optimization [52].

The rest of the paper is organized as follows. Section II presents the proposed fully adaptive decentralized gradient tracking (AdGT) algorithm, which enables agents to automatically adjust their local stepsizes based solely on local gradient information. Section III establishes the convergence guarantees of AdGT under both convex and strongly convex settings, characterizing the stability conditions introduced by stepsize adaptivity. Section IV provides comprehensive experimental validation on real and synthetic datasets, demonstrating the benefits of local adaptivity in decentralized systems across various problem types, network topologies, and heterogeneity levels. Finally, Section V concludes the paper and discusses potential extensions toward more general network models and stochastic settings

*Notation:* In this paper, bold letters represent vectors; for example,  $\mathbf{x} \in \mathbb{R}^p$ . Uppercase letters are used for matrices; for instance, given  $X \in \mathbb{R}^{n \times n}$ . Throughout the paper,  $\|\cdot\|$  denotes either the Euclidean norm or the spectral norm, depending on whether the argument is a vector or a matrix. Notation  $\odot$  denote the component-wise multiplication of two column vectors.

## 2 Design of AdGT

Consider the consensus optimization problem described by (1) over a communication network represented as an *undirected* graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . The graph consists of a set of nodes (agents)  $\mathcal{V} = \{1, 2, \dots, n\}$  and a set of edges  $\mathcal{E}$ . Each node  $i \in \mathcal{V}$  has an individual cost function  $f_i : \mathbb{R}^p \rightarrow \mathbb{R}$ , which is only accessible to that node.

For each node  $i \in \mathcal{V}$ , its neighbors are defined as the nodes that can communicate directly with  $i$ , i.e.,  $\mathcal{N}_i = \{j \in \mathcal{V} : (j, i) \in \mathcal{E}\} \cup \{i\}$ . Additionally,  $W = [w_{ij}] \in \mathbb{R}^{n \times n}$  is a doubly-stochastic matrix satisfying:

$$w_{ij} = \begin{cases} > 0, & j \in \mathcal{N}_i, \\ 0, & \text{otherwise;} \end{cases} \quad \sum_{j \in \mathcal{V}} w_{ij} = 1, \quad \forall i \in \mathcal{V}. \quad (2)$$

which ensures that the weights sum to one across both rows and columns.

Throughout the paper, we make the following assumptions.

**Assumption 1.**  $\mathcal{G}$  is undirected and connected.

**Remark 1.** Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  denote the eigenvalues of  $W$ , then  $\lambda_1 = 1$ , and for  $i = 2, \dots, n$ , the eigenvalues satisfy  $|\lambda_i| < 1$ . Additionally, let define  $\tilde{W} = W - \frac{\mathbf{1}\mathbf{1}^\top}{n}$ , then  $\|\tilde{W}\| \leq \lambda < 1$ .

**Assumption 2.** For every  $i \in \mathcal{V}$ , the local function  $f_i$  is  $L_i$ -smooth, i.e., it is differentiable with a Lipschitz gradient:

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}')\| \leq L_i \|\mathbf{x} - \mathbf{x}'\|, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^p. \quad (3)$$

**Assumption 3.** Each local objective function  $f_i : \mathbb{R}^p \rightarrow \mathbb{R}$  is convex and coercive, meaning that  $\|\mathbf{x}_i\| \rightarrow \infty \Rightarrow f_i(\mathbf{x}_i) \rightarrow \infty$ .

**Remark 2.** It follows immediately from Assumptions 2 and 3 that the global objective function  $f$  is convex, coercive, and possesses Lipschitz continuous gradients. Moreover, the Lipschitz constant of  $\nabla f$  is given by  $L = \max_i L_i$ , where  $L_i$  is the Lipschitz constant of the local function  $f_i$ .

**Assumption 4.** For all  $i \in \mathcal{V}$ ,  $f_i$  is  $\mu_i$ -strongly convex, i.e.,

$$f_i(\mathbf{y}) \geq f_i(\mathbf{x}) + \nabla f_i(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu_i}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (4)$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ .

**Remark 3.** Under Assumption 4, the optimal solution to (1) is unique, denoted by  $\mathbf{x}^*$ .

**Definition 1.** Define  $\mathbf{x} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p}$  and  $\mathbf{y} \triangleq [\mathbf{y}_1, \dots, \mathbf{y}_n]^\top \in \mathbb{R}^{n \times p}$ , where  $\mathbf{x}_i, \mathbf{y}_i \in \mathbb{R}^p$  are the local variables of agent- $i$  for  $i \in \mathcal{V} \triangleq \{1, \dots, n\}$ , and in an algorithmic framework, their values at iteration  $k \geq 0$  are denoted by  $\mathbf{x}_i^k$  and  $\mathbf{y}_i^k$  for  $i \in \mathcal{V}$ . Let  $f : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$  be a function of local variables  $\{\mathbf{x}_i\}_{i \in \mathcal{V}}$  such that  $f(\mathbf{x}) \triangleq \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i)$  for  $\mathbf{x} \in \mathbb{R}^{n \times p}$  and  $\nabla f(\mathbf{x}) \triangleq [\nabla f_1(\mathbf{x}_1), \dots, \nabla f_n(\mathbf{x}_n)]^\top \in \mathbb{R}^{n \times p}$ , where  $\nabla f_i(\mathbf{x}_i) \in \mathbb{R}^p$  denotes the gradient of  $f_i$  at  $\mathbf{x}_i \in \mathbb{R}^p$ .

We next propose our decentralized optimization algorithm, AdGT, to solve the consensus optimization problem in (1).

Consider the AdGT algorithm presented in Algorithm 1. At each iteration  $k \geq 0$ , every agent  $i \in \mathcal{V}$  updates four variables:  $\mathbf{x}_i^k, \mathbf{y}_i^k \in \mathbb{R}^p$ ,  $\theta_i^k \geq 0$ , and  $\alpha_i^k > 0$ . The following subsection presents a detailed explanation of the stepsize sequence  $\alpha_i^k$ .

**Algorithm 1** AdGT

---

**Input:**  $\mathbf{x}_i^0 \in \mathbb{R}^p, \forall i \in \mathcal{V}$ ,  
1:  $\mathbf{y}_i^0 = \nabla f_i(\mathbf{x}_i^0), \alpha_i^0 > 0, \gamma > 0, \theta_i^0 \geq 0$ , for  $i \in \mathcal{V}$   
2: **for all**  $k = 0, 1, \dots$  **do**  
3:   Each agent  $i \in \mathcal{V}$  independently performs:  
4:    $\mathbf{x}_i^{k+1} = \sum_{j \in \mathcal{N}_i} w_{ij} (\mathbf{x}_j^k - \alpha_i^k \mathbf{y}_j^k)$   
5:    $\mathbf{y}_i^{k+1} = \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{y}_j^k + \nabla f_j(\mathbf{x}_j^{k+1}) - \nabla f_i(\mathbf{x}_i^k)$   
6:    $L_{y_i}^k = \frac{\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|}{\|\mathbf{x}_i^{k+1} - \mathbf{x}_i^k\|}$   
7:    $\alpha_i^{k+1} = \min \left\{ \frac{1}{2\gamma L_{y_i}^k}, \sqrt{1 + \theta_i^k} \alpha_i^k \right\}$   
8:    $\theta_i^{k+1} = \frac{\alpha_i^{k+1}}{\alpha_i^k}$   
9: **end for**

---

**2.1 Designing Adaptive stepsizes  $\alpha_i^k$** 

The stepsize  $\alpha_i^k$ , defined in Step (7) of Algorithm 1, is inspired by the adaptive gradient descent method introduced in [44]. In that setting, the centralized stepsize  $\alpha^k$  is updated according to the following rule:

$$\alpha^{k+1} = \min \left\{ \frac{\|\mathbf{x}^{k+1} - \mathbf{x}^k\|}{2\|\nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k)\|}, \sqrt{1 + \theta^k} \alpha^k \right\}, \quad (5)$$

where  $\theta^k = \alpha^{k+1}/\alpha^k$ . This approach balances stability and adaptivity by adjusting the stepsize based on the local smoothness of the objective function. However, applying the adaptive centralized stepsize (5) to the distributed setting can lead to algorithmic divergence. Recently, in [52], the authors introduced an adaptive decentralized method by combining the GT method with the following adaptive stepsize

$$\alpha_i^{k+1} = \min \left\{ \frac{1}{2L_{f_i}^k}, \frac{1}{\|\mathbf{y}_i^k\|}, \sqrt{2}\alpha_i^k \right\} \quad (6)$$

where

$$L_{f_i}^k = \frac{\|\nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k)\|}{\|\mathbf{x}_i^{k+1} - \mathbf{x}_i^k\|} \quad (7)$$

The method is theoretically proven to converge only to a neighborhood of the optimal solution. However, our numerical results show that the algorithm proposed in [52], when used with the adaptive distributed stepsize (6), often diverges—particularly in settings with a small number of samples or sparse network connectivity.

To ensure both theoretical and practical convergence of the adaptive decentralized method, the adaptive centralized stepsize should be defined in (5) as follows:

$$\alpha_i^{k+1} = \min \left\{ \frac{1}{2\gamma L_{f_i}^k}, \sqrt{1 + \theta_i^k} \alpha_i^k \right\}, \quad (8)$$

where  $L_{f_i}^k$  is given in (7) and  $\gamma > 0$ . However, in our implementation, we replace  $L_{f_i}^k$  with

$$L_{y_i}^k = \frac{\|\mathbf{y}_i^{k+1} - \mathbf{y}_i^k\|}{\|\mathbf{x}_i^{k+1} - \mathbf{x}_i^k\|} \quad (9)$$

leading to the following modified stepsize update:

$$\alpha_i^{k+1} = \min \left\{ \frac{1}{2\gamma L_{y_i}^k}, \sqrt{1 + \theta_i^k} \alpha_i^k \right\}. \quad (10)$$

This modification significantly improves the stability of the algorithm. In practice, we find that in many cases it is sufficient to set  $\gamma = 1$ . The algorithm instead of using the local gradient  $\nabla f_i(\mathbf{x}_i^k)$ , we approximate the global gradient using the local estimator  $\mathbf{y}_i^k$ . As discussed in Remark 8, this estimator satisfies

$$\bar{\mathbf{y}}^k = \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^k)$$

The second term in (10) controls the rate at which  $\alpha_i^k$  is allowed to increase, ensuring a balanced and stable adaptation over the iterations.

**Remark 4.** By including the term  $1/L_{f_i}^k$  in the stepsize definition (10), the updated stepsize  $\alpha_i^{k+1}$  can be expressed as

$$\alpha_i^{k+1} = \min \left\{ \frac{1}{2\gamma L_{f_i}^k}, \frac{1}{2\gamma L_{y_i}^k}, \sqrt{1 + \theta_i^k} \alpha_i^k \right\}. \quad (11)$$

As a result, the conclusion in Remark 7 remains valid and can still be used in theoretical analysis. From a practical standpoint, replacing the stepsize rule in (10) with the form in (11) does not alter the convergence rate of the AdGT method.

We present AdGT stated in Algorithm 1 in a compact form as follows:

$$\mathbf{x}^{k+1} = W \left( \mathbf{x}^k - \alpha^k \odot \mathbf{y}^k \right) \quad (12)$$

$$\mathbf{y}^{k+1} = W \mathbf{y}^k + \nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k) \quad (13)$$

where  $\alpha^k = (\alpha_1^k, \dots, \alpha_n^k)$ .

**Remark 5.** The AdGT method integrates an adaptive stepsize with the GT method framework. It can be adapted to different variants of GT methods by modifying specific steps in Algorithm 1. In particular, replacing Step (5) with its Adapt-Then-Combine (ATC) version, or altering Step (4) to follow the Combine-Then-Adapt (CTA) scheme, results in alternative implementations of the algorithm within the same general framework [18, 53, 19, 20, 24, 25, 54].

**Remark 6.** In this paper, we develop our adaptive stepsize scheme within the framework of gradient tracking over undirected communication graphs. Although our theoretical analysis and numerical experiments are limited to the undirected setting, the adaptive scheme continues to exhibit strong performance when applied to gradient tracking over directed graphs.

### 3 Main Results

In this section, we show that the sequence of iterates generated by the AdGT algorithm, as defined in equations (12) and (13), converges to the optimal solution  $\mathbf{x}^* = 1_n x^*$ . The structure of our proof is inspired by the approach in [21]. Without loss of generality, we consider the case  $p = 1$ .

**Remark 7.** Due to the  $L_i$ -smoothness and  $\mu_i$ -strong convexity of each local function  $f_i$ , the stepsize  $\alpha_i^k$  chosen by each agent at every iteration satisfies the following bounds:

$$\frac{1}{2\gamma L} \leq \alpha_i^k \leq \frac{1}{2\gamma \mu}, \quad (14)$$

where  $L = \max_i l_i$ , and  $\mu = \min_i \mu_i$ ; for more details, see[44]. Let us define  $\alpha_{\max} = \max_{i,k} \alpha_i^k$ . Then, according to (14) we have

$$\alpha_{max} = \frac{1}{2\gamma\mu} \quad (15)$$

**Definition 2.** In iteration  $k$ , the average stepsize is defined as  $\bar{\alpha}^k = \frac{1}{n}\mathbf{1}\mathbf{1}^T \alpha^k$ , and the deviation from this average is given by  $\hat{\alpha}^k = \alpha^k - \bar{\alpha}^k$ . The measure of how different the agents' stepsizes are from each other, and is calculated as:

$$\delta_\alpha^k = \frac{\|\hat{\alpha}^k\|}{\|\bar{\alpha}^k\|}$$

A smaller value of  $\delta_\alpha^k$  indicates that the stepsizes are more similar across agents, whereas a larger value implies greater variability or imbalance in the stepsizes used. We consider  $\delta_\alpha = \max_k \delta_\alpha^k$ <sup>3</sup>.

To support our convergence analysis, we multiply  $\frac{\mathbf{1}\mathbf{1}^T}{n}x^k$  over equation (12):

$$\bar{x}^{k+1} = \bar{x}^k - \overline{\alpha^k \odot y^k} \quad (16)$$

where  $\bar{x}^k = \frac{\mathbf{1}\mathbf{1}^T}{n}x^k$  represents the average of the agents, and  $\overline{\alpha^k \odot y^k} = \frac{\mathbf{1}\mathbf{1}^T}{n}(\alpha^k \odot y^k)$  denotes the average element-wise product of the stepsizes and variable  $y^k$ .

We also define  $\bar{y}^k = \frac{\mathbf{1}\mathbf{1}^T}{n}y^k$  as the average, and  $\hat{y}^k = y^k - \bar{y}^k$  as the consensus error, the term  $\overline{\alpha^k \odot y^k}$  can be decomposed as follows:

$$\overline{\alpha^k \odot y^k} = \bar{\alpha}^k \odot \bar{y}^k + \hat{\alpha}^k \odot \hat{y}^k \quad (17)$$

Taking the norm of both sides and applying the Cauchy-Schwarz inequality, we obtain:

$$\|\hat{\alpha}^k \odot \hat{y}^k\| \leq \frac{1}{\sqrt{n}} \|\hat{\alpha}^k\| \cdot \|\hat{y}^k\| \quad (18)$$

Using this bound, we derive a lower bound on  $\|\overline{\alpha^k \odot y^k}\|$ :

$$\frac{1}{\sqrt{n}} (\|\bar{\alpha}^k\| \cdot \|\bar{y}^k\| - \|\hat{\alpha}^k\| \cdot \|\hat{y}^k\|) \leq \|\overline{\alpha^k \odot y^k}\| \quad (19)$$

**Remark 8.** Let  $y^0 = \nabla f(x^0)$  then, for all  $k \geq 0$ , we have

$$\bar{y}^k = \frac{\mathbf{1}\mathbf{1}^T}{n} \nabla f(x)$$

This ensures that the variable  $\bar{y}^k$  consistently reflects the global average gradient.

We now present two technical results that provide bounds on the consensus error.

**Lemma 1.** The following inequality holds for all  $k \geq 0$  :

$$\|\hat{x}^{k+1}\| \leq \lambda \|\hat{x}^k\| + \lambda \alpha_{\max} (1 + \delta_\alpha) \|\hat{y}^k\| + \lambda \delta_\alpha \|\overline{\alpha^k \odot y^k}\|$$

<sup>3</sup>It can be shown that  $\delta_\alpha \leq \frac{L-\mu}{L+\mu}$ .

*Proof.* Using (12) and (16), we have for all  $k \geq 0$  :

$$\begin{aligned}
\|\hat{\mathbf{x}}^{k+1}\| &= \|\mathbf{x}^{k+1} - \bar{\mathbf{x}}^{k+1}\| \\
&= \|\tilde{W}(\mathbf{x}^k - \bar{\mathbf{x}}^k) - \tilde{W}(\alpha^k \odot \mathbf{y}^k)\| \\
&= \|\tilde{W}(\mathbf{x}^k - \bar{\mathbf{x}}^k) - \tilde{W}(\alpha^k \odot \mathbf{y}^k - \bar{\alpha}^k \odot \bar{\mathbf{y}}^k)\| \\
&= \|\tilde{W}(\mathbf{x}^k - \bar{\mathbf{x}}^k) - \tilde{W}(\alpha^k \odot \hat{\mathbf{y}}^k + \hat{\alpha}^k \odot \bar{\mathbf{y}}^k)\| \\
&\leq \lambda \|\hat{\mathbf{x}}^k\| + \lambda \|\alpha^k \odot \hat{\mathbf{y}}^k\| + \lambda \frac{1}{\sqrt{n}} \|\hat{\alpha}^k\| \|\bar{\mathbf{y}}^k\| \\
&\leq \lambda \|\hat{\mathbf{x}}^k\| + \lambda \alpha_{\max} \|\hat{\mathbf{y}}^k\| + \lambda \frac{\|\hat{\alpha}^k\|}{\|\bar{\alpha}^k\|} \|\alpha^k \odot \bar{\mathbf{y}}^k\| \\
&\leq \lambda \|\hat{\mathbf{x}}^k\| + \lambda \alpha_{\max} (1 + \delta_\alpha) \|\hat{\mathbf{y}}^k\| + \lambda \delta_\alpha \|\overline{\alpha^k \odot \mathbf{y}^k}\|
\end{aligned}$$

In the second equality, we include the term  $\tilde{W}(\bar{\alpha}^k \odot \bar{\mathbf{y}}^k) = 0$ , which does not affect the expression, and in the third equality, we apply

$$\begin{aligned}
\alpha^k \odot \mathbf{y}^k - \bar{\alpha}^k \odot \bar{\mathbf{y}}^k &= \alpha^k \odot \mathbf{y}^k - \alpha^k \odot \bar{\mathbf{y}}^k + \alpha^k \odot \bar{\mathbf{y}}^k - \bar{\alpha}^k \odot \bar{\mathbf{y}}^k \\
&= \alpha^k \odot (\mathbf{y}^k - \bar{\mathbf{y}}^k) + (\alpha^k - \bar{\alpha}^k) \odot \bar{\mathbf{y}}^k \\
&= \alpha^k \odot \tilde{\mathbf{y}}^k + \tilde{\alpha}^k \odot \bar{\mathbf{y}}^k
\end{aligned}$$

Regarding the three inequalities: the first follows from the Cauchy–Schwarz inequality and Remark 1; the second is obtained by noting that  $\frac{1}{\sqrt{n}} \|\bar{\alpha}^k\| \|\bar{\mathbf{y}}^k\| = \|\bar{\alpha}^k \odot \bar{\mathbf{y}}^k\|$ ; and the final inequality follows directly from (19).  $\square$

**Lemma 2.** *Let*

$$\tilde{\lambda} = \lambda + \lambda L \alpha_{\max} (1 + \delta_\alpha) \quad (20)$$

*then, the following inequality holds for all  $k \geq 0$  :*

$$\|\hat{\mathbf{y}}^{k+1}\| \leq \tilde{\lambda} \|\hat{\mathbf{y}}^k\| + L(1 + \lambda) \|\hat{\mathbf{x}}^k\| + L(1 + \lambda \delta_\alpha) \|\overline{\alpha^k \odot \mathbf{y}^k}\|$$

*Proof.* Using (13) and Remark 8, we have for all  $k \geq 0$  :

$$\begin{aligned}
\|\hat{\mathbf{y}}^{k+1}\| &= \|\hat{W} \hat{\mathbf{y}}^k + (I_n - \frac{1_n 1_n^\top}{n}) (\nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k))\| \\
&\leq \lambda \|\hat{\mathbf{y}}^k\| + L \|\mathbf{x}^{k+1} - \mathbf{x}^k\| \\
&= \lambda \|\hat{\mathbf{y}}^k\| + L \|\hat{\mathbf{x}}^{k+1} + \bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k - \hat{\mathbf{x}}^k\| \\
&\leq \lambda \|\hat{\mathbf{y}}^k\| + L (\|\hat{\mathbf{x}}^{k+1}\| + \|\overline{\alpha^k \odot \mathbf{y}^k}\| + \|\hat{\mathbf{x}}^k\|)
\end{aligned} \quad (21)$$

where the first inequality follows from the observation that the spectral norm  $\|(I_n - \frac{1_n 1_n^\top}{n})\| = 1$ , together with Remark 2; the final inequality is a direct consequence of (16). Combining (21) with Lemma 1 subsequently yields

$$\begin{aligned}
\|\hat{\mathbf{y}}^{k+1}\| &\leq \lambda \|\hat{\mathbf{y}}^k\| + L (\lambda \|\hat{\mathbf{x}}^k\| + \lambda \alpha_{\max} (1 + \delta_\alpha) \|\hat{\mathbf{y}}^k\| \\
&\quad + \lambda \delta_\alpha \|\overline{\alpha^k \odot \mathbf{y}^k}\|) + \|\overline{\alpha^k \odot \mathbf{y}^k}\| + \|\hat{\mathbf{x}}^k\| \\
&= \tilde{\lambda} \|\hat{\mathbf{y}}^k\| + L(1 + \lambda) \|\hat{\mathbf{x}}^k\| + L(1 + \lambda \delta_\alpha) \|\overline{\alpha^k \odot \mathbf{y}^k}\|
\end{aligned}$$

where  $\tilde{\lambda}$  is given in (20).  $\square$

**Remark 9.** *It follows directly from (20) that if the step size satisfies*

$$\alpha_{\max} < \frac{1 - \lambda}{L \lambda (1 + \delta_\alpha)}, \quad (22)$$

then one immediately obtains  $\tilde{\lambda} < 1$ .

**Lemma 3.** For all  $k \geq 0$ , the following inequality holds:

$$(\bar{y}^k)^\top (\overline{\alpha^k \circ y^k}) \geq \frac{1}{\alpha_{\max}} \|\overline{\alpha^k \circ y^k}\|^2 - \delta_\alpha \|\hat{y}^k\| \|\overline{\alpha^k \circ y^k}\|$$

*Proof.* By applying the term in (17), we obtain

$$\begin{aligned} (\bar{y}^k)^\top (\overline{\alpha^k \circ y^k}) &= \sqrt{n} \|\bar{\alpha}^k\| (\overline{\alpha^k \circ y^k} - \hat{\alpha}^k \circ \hat{y}^k)^\top (\overline{\alpha^k \circ y^k}) \\ &\geq \sqrt{n} \|\bar{\alpha}^k\| (\|\overline{\alpha^k \circ y^k}\|^2 - \|\hat{\alpha}^k \circ \hat{y}^k\| \cdot \|\overline{\alpha^k \circ y^k}\|) \\ &\geq \sqrt{n} \|\bar{\alpha}^k\| (\|\overline{\alpha^k \circ y^k}\|^2 - \|\hat{\alpha}^k\| \|\hat{y}^k\| \|\overline{\alpha^k \circ y^k}\|) \\ &\geq \frac{1}{\alpha_{\max}} \|\overline{\alpha^k \circ y^k}\|^2 - \delta_\alpha \|\hat{y}^k\| \|\overline{\alpha^k \circ y^k}\| \end{aligned}$$

where the second inequality follows from the inequality (18), and the final bound uses the definition of  $\delta_\alpha$  given in Definition 2.  $\square$

To ensure completeness, we state the following technical result, whose proof is available in [21, Lemma 2].

**Lemma 4.** Let  $\{s^k\}_{k \geq 0}$  and  $\{t^k\}_{k \geq 0}$  be two sequences of positive real numbers satisfying

$$s^{k+1} \leq \lambda s^k + t^k, \quad (23)$$

where  $\lambda \in (0, 1)$ . Define

$$\Upsilon^k = \sqrt{\sum_{i=0}^k \|s^i\|^2}, \quad \Psi^k = \sqrt{\sum_{i=0}^k \|t^i\|^2}$$

Then, for every  $k \geq 0$ , we have

$$\Upsilon^k \leq \theta \Psi^k + \varepsilon, \quad (24)$$

where

$$\theta = \frac{\sqrt{2}}{1-\lambda}, \quad \varepsilon = \frac{\sqrt{2}}{1-\lambda^2} s^0.$$

**Lemma 5.** Consider Algorithm 1, and suppose that Remark 1 and Assumption 2 hold. Define the accumulated energy up to iteration  $k$  as

$$\begin{aligned} \hat{X}^k &= \sqrt{\sum_{i=0}^k \|\hat{x}^i\|^2}, & \hat{Y}^k &= \sqrt{\sum_{i=0}^k \|\hat{y}^i\|^2}, \\ Y^k &= \sqrt{\sum_{i=0}^k \|\overline{\alpha^k \circ y^k}\|^2} \end{aligned}$$

Let  $\alpha_{\max} := \max_{i,k} \alpha_i^k$ , and  $\tilde{\lambda}$  is given in (20). Suppose the following condition is satisfied:

$$\alpha_{\max} < \frac{(1-\lambda)^2}{(1+\delta_\alpha)L\lambda(\lambda+3)} \quad (25)$$

so that both  $\theta_1\theta_2 < 1$  and  $\tilde{\lambda} < 1$ . Then, for all  $k \geq 0$ , we have

$$\hat{X}^k \leq \frac{\theta_1\beta_2 + \beta_1}{1 - \theta_1\theta_2} Y^k + \frac{\varepsilon_1 + \theta_1\varepsilon_2}{1 - \theta_1\theta_2}, \quad (26)$$

$$\hat{Y}^k \leq \frac{\theta_2\beta_1 + \beta_2}{1 - \theta_1\theta_2} Y^k + \frac{\varepsilon_2 + \theta_2\varepsilon_1}{1 - \theta_1\theta_2}, \quad (27)$$

where

$$\begin{aligned} \theta_1 &= \frac{\sqrt{2}\lambda\alpha_{\max}(1 + \delta_\alpha)}{1 - \lambda}, & \beta_1 &= \frac{\sqrt{2}\lambda\delta_\alpha}{1 - \lambda}, & \varepsilon_1 &= \frac{\sqrt{2}\|\hat{x}^0\|}{\sqrt{1 - \tilde{\lambda}^2}} \\ \theta_2 &= \frac{\sqrt{2}L(1 + \lambda)}{1 - \tilde{\lambda}}, & \beta_2 &= \frac{\sqrt{2}L(1 + \lambda\delta_\alpha)}{1 - \tilde{\lambda}}, & \varepsilon_2 &= \frac{\sqrt{2}\|\hat{y}^0\|}{\sqrt{1 - \tilde{\lambda}^2}} \end{aligned}$$

*Proof.* From Remark 9, Lemma 4, and Lemma 1, we obtain the bound

$$\hat{X}^k \leq \theta_1\hat{Y}^k + \beta_1Y^k + \varepsilon_1 \quad (28)$$

In addition, Lemma 2 gives

$$\hat{Y}^k \leq \theta_2\hat{X}^k + \beta_2Y^k + \varepsilon_2 \quad (29)$$

Substituting (29) into (28), we have

$$\hat{X}^k \leq \theta_1\theta_2\hat{X}^k + (\theta_1\beta_2 + \beta_1)Y^k + \theta_1\varepsilon_2 + \varepsilon_1$$

Rearranging terms yields

$$(1 - \theta_1\theta_2)\hat{X}^k \leq (\theta_1\beta_2 + \beta_1)Y^k + \theta_1\varepsilon_2 + \varepsilon_1$$

Now suppose that

$$\alpha_{\max} < \frac{(1 - \lambda)^2}{(1 + \delta_\alpha)L\lambda(\lambda + 3)} \quad (30)$$

so that  $\theta_1\theta_2 < 1$ . Dividing both sides of the above inequality by  $(1 - \theta_1\theta_2)$  then establishes the bound in (26). The inequality in (27) follows by a similar argument, by substituting (28) into (29).  $\square$

**Theorem 1.** Consider the algorithm defined by equations (12) and (13), initialized with  $y^0 = \nabla f(x^0)$ . Suppose that Assumptions 1 to 4 hold. Then,  $\alpha_{max}$  is given in (15), satisfies in (44) such that if each agent's stepsize satisfies  $\alpha_i^k < \alpha_{max}$  for all  $i \in V$  and for all  $k \geq 0$ , the algorithm converges to the global optimum. In particular, sequence  $\{x^k\}$  satisfies

$$\lim_{k \rightarrow \infty} f(x^k) = f(1 \otimes x^*) = f^*,$$

where  $f^*$  is the optimal value of the original problem.

*Proof.* Since the function  $f$  has a Lipschitz continuous gradient by Assumption 2 and Remark 2, we obtain the following estimate:

$$\begin{aligned}
f(\bar{x}^{k+1}) &\leq f(\bar{x}^k) + \nabla f(\bar{x}^k)^T(\bar{x}^{k+1} - \bar{x}^k) + \frac{L}{2}\|\bar{x}^{k+1} - \bar{x}^k\|^2 \\
&\leq f(\bar{x}^k) + \nabla f(x^k)^T(\bar{x}^{k+1} - \bar{x}^k) + \frac{L}{2}\|\bar{x}^{k+1} - \bar{x}^k\|^2 \\
&\quad + (\nabla f(\bar{x}^k) - \nabla f(x^k))^T(\bar{x}^{k+1} - \bar{x}^k) \\
&\leq f(\bar{x}^k) - (\bar{y}^k)^\top(\alpha^k \odot \bar{y}^k) + \frac{L}{2}\|\bar{x}^{k+1} - \bar{x}^k\|^2 \\
&\quad + (\nabla f(\bar{x}^k) - \nabla f(x^k))^T(\bar{x}^{k+1} - \bar{x}^k)
\end{aligned} \tag{31}$$

where the last inequality uses the identity

$$\begin{aligned}
\nabla f(x^k)^T(\bar{x}^{k+1} - \bar{x}^k) &= \left( \frac{11^\top}{n} \nabla f(x) \right)^T (\bar{x}^{k+1} - \bar{x}^k) \\
&= (\bar{y}^k)^T(\bar{x}^{k+1} - \bar{x}^k),
\end{aligned}$$

The first equality follows from the standard identity  $a^\top \bar{b} = \bar{a}^\top b$  for any vectors  $a, b \in \mathbb{R}^n$ , and the second is a consequence of Remark 8.

Next, we consider the final term in (31). Based on Assumption 2 and Remark 2, we derive the following bound:

$$\|(\nabla f(\bar{x}^k) - \nabla f(x^k))^T(\bar{x}^{k+1} - \bar{x}^k)\| \leq L\|\hat{x}^k\|\|\bar{x}^{k+1} - \bar{x}^k\| \tag{32}$$

Then, by combining (31), Lemma 3 and (32), we obtain:

$$\begin{aligned}
f(\bar{x}^{k+1}) &\leq f(\bar{x}^k) - \left( \frac{1}{\alpha_{\max}} - \frac{L}{2} \right) \|\bar{x}^{k+1} - \bar{x}^k\|^2 \\
&\quad + (\delta_\alpha \|\hat{y}^k\| + L\|\hat{x}^k\|)\|\bar{x}^{k+1} - \bar{x}^k\|
\end{aligned}$$

By summing the above inequality over  $i$  from 0 to  $k$ , we have:

$$\begin{aligned}
f(\bar{x}^{k+1}) &\leq f(\bar{x}^0) - \left( \frac{1}{\alpha_{\max}} - \frac{L}{2} \right) \sum_{j=0}^k \|\bar{x}^{j+1} - \bar{x}^j\|^2 \\
&\quad + \delta_\alpha \sum_{j=0}^k \|\hat{y}^j\| \|\bar{x}^{j+1} - \bar{x}^j\| + L \sum_{j=0}^k \|\hat{x}^j\| \|\bar{x}^{j+1} - \bar{x}^j\|
\end{aligned}$$

Next, using the Cauchy–Schwarz inequality and (16), we derive the following result.

$$f(\bar{x}(t+1)) \leq f(\bar{x}(0)) - \left( \frac{1}{\alpha_{\max}} - \frac{L}{2} \right) (Y^k)^2 + \delta_\alpha \hat{Y}^k Y^k + L \bar{X}^k Y^k \tag{33}$$

Assume that all conditions in Lemma 5 are satisfied and that the stepsize satisfies in (25) which ensures that  $\theta_1 \theta_2 < 1$  and  $\tilde{\lambda} < 1$ . Then, by applying Lemma 5, we obtain the following bound:

$$f(\bar{x}(t+1)) \leq f(\bar{x}(0)) - \rho_1 (Y^k)^2 + \rho_2 Y^k, \tag{34}$$

where the constants  $\rho_1$  and  $\rho_2$  are given by

$$\rho_1 = \frac{1}{\alpha_{\max}} - \frac{L}{2} - \frac{L(\beta_1 + \theta_1\beta_2) + \delta_\alpha(\beta_2 + \theta_2\beta_1)}{1 - \theta_1\theta_2}, \quad (35)$$

$$\rho_2 = \frac{L(\epsilon_1 + \theta_1\epsilon_2) + \delta_\alpha(\epsilon_2 + \theta_2\epsilon_1)}{1 - \theta_1\theta_2}. \quad (36)$$

Let  $u(k) = f(\bar{x}(k)) - f^*$  represent the optimality gap at iteration  $k$ . Since  $u(k) \geq 0$  for all  $k \geq 0$ , inequality (34) leads to:

$$-\rho_1(Y^k)^2 + \rho_2 Y^k + u(0) \geq 0.$$

One can verify that  $\rho_1 > 0$  holds provided the stepsize satisfies

$$0 < \alpha_{\max} < \frac{-b - \sqrt{b^2 - 4ac}}{2a}, \quad (37)$$

where the constants  $a$ ,  $b$ , and  $c$  are given by:

$$a = L^2\lambda(1 + \delta_\alpha) \left( (\lambda - 1) + (2\sqrt{2} - 4)\lambda\delta_\alpha \right) < 0, \quad (38)$$

$$\begin{aligned} b &= -2L\lambda(\lambda + 3)(1 + \delta_\alpha) - 2\sqrt{2}L(1 - \lambda)\delta_\alpha(\lambda + 1) \\ &\quad - L\lambda\delta_\alpha^2 \left( (4 - 2\sqrt{2})\lambda + (4 + 2\sqrt{2}) \right) - L(1 - \lambda)^2 < 0, \end{aligned} \quad (39)$$

$$c = 2(1 - \lambda)^2 > 0. \quad (40)$$

This implies that

$$\lim_{k \rightarrow \infty} Y^k \leq B := \frac{\rho_2 + \sqrt{\rho_2^2 - 4\rho_1 u(0)}}{2\rho_1} < \infty. \quad (41)$$

Then, by the monotone convergence theorem, it follows that  $\lim_{k \rightarrow \infty} \|\bar{\alpha}^k \circ y^k\| = 0$ . Moreover, by combining (27) with (41), we obtain:

$$\lim_{k \rightarrow \infty} \hat{Y}^k \leq \frac{(\theta_2\beta_1 + \beta_2)B + (\epsilon_2 + \theta_2\epsilon_1)}{1 - \theta_1\theta_2} < \infty,$$

which yields  $\lim_{k \rightarrow \infty} \|\hat{y}^k\| = 0$ . Furthermore, using (19) together with the result above, we derive:

$$\lim_{k \rightarrow \infty} \|\bar{y}^k\| \leq \lim_{k \rightarrow \infty} \left( \frac{\sqrt{n}\|\bar{\alpha}^k \circ y^k\|}{\|\bar{\alpha}^k\|} + \delta_\alpha \|\hat{y}^k\| \right) = 0.$$

Therefore, we conclude that  $\lim_{k \rightarrow \infty} \|\bar{y}^k\| = 0$ . Similarly, combining (27) with (41), we obtain  $\lim_{k \rightarrow \infty} \|\hat{x}^k\| = 0$ . Since  $f$  is convex by Assumption 3 and Remark 2, for any  $\bar{x}^k \in \mathbb{R}^m$ , we have:

$$\begin{aligned} f(\bar{x}^k) - f(x^*) &\leq \nabla f(\bar{x}^k)^T (\bar{x}^k - x^*) \\ &= \nabla f(x^k)^T (\bar{x}^k - x^*) + (\nabla f(\bar{x}^k) - \nabla f(x^k))^T (\bar{x}^k - x^*) \\ &= \left( \frac{11^\top}{n} \nabla f(x) \right)^T (\bar{x}^k - x^*) + (\nabla f(\bar{x}^k) - \nabla f(x^*))^T (\bar{x}^k - x^*) \\ &\leq \|\bar{y}^k\| \cdot \|\bar{x}^k - x^*\| + L \cdot \|\hat{x}^k\| \cdot \|\bar{x}^k - x^*\| \end{aligned}$$

where the last inequality follows from the Lipschitz continuity of  $\nabla f$ , and  $\bar{y}^k = \left( \frac{11^\top}{n} \nabla f(x) \right)^T$ .

From (34) and (41),  $f(\bar{x}^k)$  is bounded. Furthermore, since  $f$  is coercive by Assumption 3, the sequence  $\|\bar{x}^k - x^*\|$  is also bounded. Then, we conclude:

$$\lim_{k \rightarrow \infty} f(\bar{x}^k) = f(x^*) = f(1 \otimes x^*). \quad (42)$$

Additionally, by the mean value theorem, there exists  $0 \leq \xi \leq 1$  such that:

$$f(x^k) = f(\bar{x}^k) + \nabla f(\bar{x}^k + \xi \hat{x}^k)^T \hat{x}^k.$$

Since  $\|\bar{x}^k + \xi \hat{x}^k\|$  remains bounded and  $\nabla f$  is Lipschitz continuous, it follows that  $\|\nabla f(\bar{x}^k + \xi \hat{x}^k)\|$  is also bounded. Therefore,

$$\lim_{k \rightarrow \infty} |f(x^k) - f(\bar{x}^k)| \leq \lim_{k \rightarrow \infty} \|\nabla f(\bar{x}^k + \xi \hat{x}^k)\| \cdot \|\hat{x}^k\| = 0. \quad (43)$$

Combining (42) and (43), we conclude that

$$\lim_{k \rightarrow \infty} f(x^k) = f(1 \otimes x^*) = f^*.$$

□

**Remark 10.** As stated in Remark 7, the upper bound of the stepsize is given by  $\alpha_{\max} = \frac{1}{2\gamma\mu}$ . Based on the results in equations (22), (30), and (37), this value is further bounded by

$$\alpha_{\max} \leq D = \min \left( \frac{(1-\lambda)^2}{(1+\delta_\alpha)L\lambda(\lambda+3)}, \frac{-b - \sqrt{b^2 - 4ac}}{2a} \right) \quad (44)$$

where constants  $a$ ,  $b$ , and  $c$  are defined in (38), (39), and (40), respectively.

This implies that the condition on  $\alpha_{\max}$  holds as long as the stepsize parameter  $\gamma$  satisfies  $\gamma \geq \gamma_{\min} = \frac{1}{2D\mu}$ . Hence, from Remark 7, we can have

$$\frac{\mu}{L}D \leq \alpha \leq D \quad (45)$$

Fig 1 illustrates how  $\gamma_{\min}$  and  $D$  vary with respect to  $\lambda$  and  $\delta_\alpha$  when  $L = 3$ ,  $\mu = 1$ . A larger spectral gap  $(1 - \lambda)$  indicates well-connected network, enabling a larger step size. Conversely, a smaller spectral gap, reflecting weaker connectivity, requires a reduction in step size to ensure stability.

As illustrated in Fig. 1a, a larger  $\gamma$  is theoretically required, particularly when the spectral gap is small. However, in practice,  $\gamma = 1$  is commonly used across many scenarios.

## 4 Numerical Results

In this section, we present numerical results to evaluate the performance of AdGT, with a particular focus on our contributions. To keep the comparison clear and meaningful, we compare AdGT with two methods: the fixed stepsize Gradient Tracking (GT) method and an adaptive decentralized method, referred to as Method-DM, proposed in [52]. All experiments are carried out under undirected communication graphs. We chose these two methods because GT has already been widely studied and compared with other advanced decentralized optimization methods in previous works.

Our experiments are divided into three parts. First, in Section 4.1, we show that AdGT can change stepsizes automatically, removing the need for manual tuning and making the algorithm easier to use. Second, in Section 4.2, we point out the advantage of letting each agent use a different stepsize based on the smoothness of its local function.

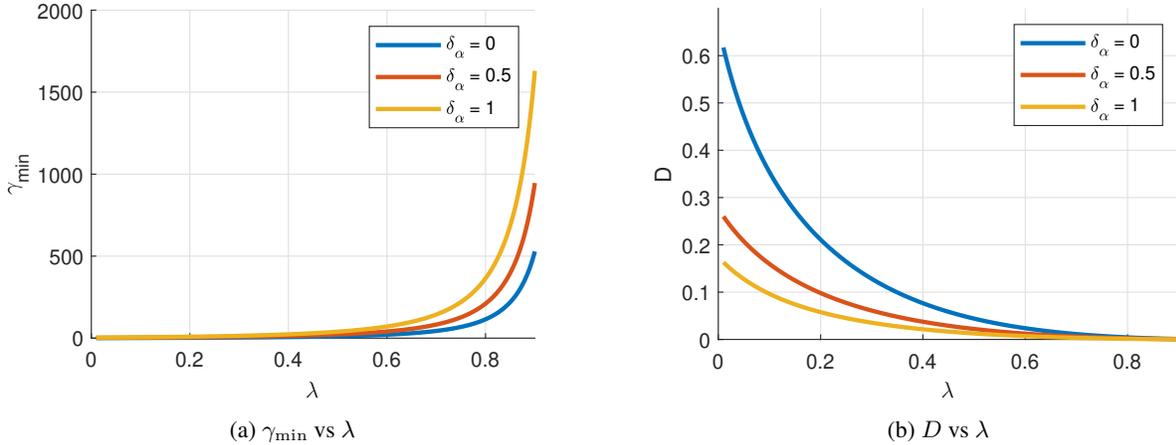


Figure 1: Variation of  $\gamma_{\min}$  and  $D$  vs.  $\lambda$  and  $\delta_\alpha$ . The smoothness parameter is set to  $L = 3$  and the strong convexity parameter to  $\mu = 1$ .

In this case, where agents have different smoothness levels, we show that the gap between AdGT and GT is quite noticeable compared to the gap between Adaptive Gradient Descent (AdGD) [44] and regular Gradient Descent (GD) in a centralized setting. Third, in Section 4.3, we study how the number of nodes affects the performance of AdGT and GT in two types of networks: random graphs and cycle graphs.

#### 4.1 Auto-Tuned stepsizes

In this experiment, our objective is to estimate  $\tilde{x}$  using the optimal solution  $x^*$  of the decentralized logistic regression problem, which is formulated as:

$$x^* \in \underset{x \in \mathbb{R}^p}{\operatorname{argmin}} f(x) = \sum_{i=1}^n f_i(x), \quad (46)$$

where each local function  $f_i(x)$  is defined as

$$f_i(x) = \sum_{j=1}^{m_i} \log \left( 1 + \exp(-y_i^j \cdot M_i^j x) \right) + \frac{\rho}{2} \|x\|^2.$$

We solve the optimization problem in (46) across a variety of known network topologies—specifically, Star, Cycle, Line, Ladder—as well as two randomly generated graphs with connectivity ratios of 0.2 and 0.35. The dataset used is the w8a dataset [55], which contains 49,749 examples, each with a 300-dimensional feature vector and a corresponding binary label. To include an intercept term in the model, we set  $p = 301$ .

The network consists of  $n = 16$  nodes (or agents). Each agent  $i \in \mathcal{V}$  randomly selects  $m_i = 25$  data points from the training set without replacement. For each agent, the local data matrix  $M_i^j \in \mathbb{R}^p$ ,  $j = 1, \dots, m_i$  is formed by standardizing the original  $p - 1$  feature vectors using Z-score normalization, and then adding a column of ones to include the intercept term, and the regularizer is  $\rho = 0.01$ . To generate the random graphs, we used Ggen code<sup>4</sup>.

In all simulations presented in Sections 4.1 to 4.3, we fix the random seed to 42, and set the initial point as  $x^0 = 0$ , and the convergence behavior of different methods is illustrated through the residual sequence  $r(k)_{k \geq 1}$ , which is defined

<sup>4</sup>The Ggen code is written by Dr. W. Shi; see <https://sites.google.com/view/wilburshi/home/research/software-a-z-order/graph>

as

$$r(k) \triangleq \frac{\|\mathbf{x}^k - \mathbf{x}^*\|}{\|\mathbf{x}^0 - \mathbf{x}^*\|}$$

In Fig. 2, we plot the residual sequence  $\{r(k)\}_{k \geq 0}$  for all methods. To improve the practical convergence rate,

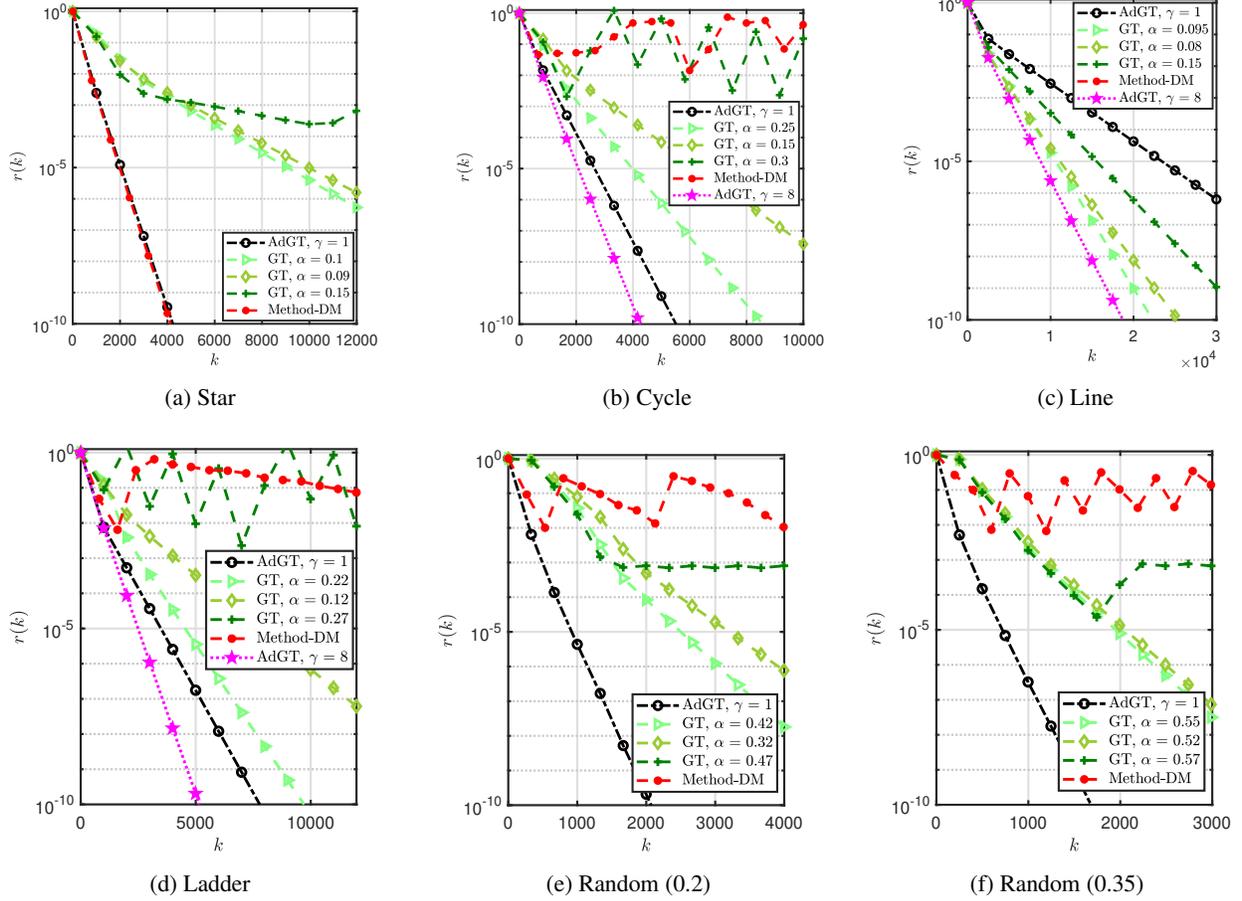


Figure 2: Decentralized logistic regression with  $n = 16$ . (a)  $\{r(k)\}_k$  for Star, (b)  $\{r(k)\}_k$  for Cycle, (c)  $\{r(k)\}_k$  for Line, (d)  $\{r(k)\}_k$  for Ladder, (e)  $\{r(k)\}_k$  for Random graph with connectivity ratio 0.2, (f)  $\{r(k)\}_k$  for Random graph with connectivity ratio 0.35.

the stepsize in the GT algorithm was tuned via a grid search. For the AdGT algorithm, we set  $\gamma = 1$  for all graph topologies, which generally resulted in improved performance over GT. An exception occurred with line graphs, where using  $\gamma = 8$  led to faster convergence. Additionally, we applied  $\gamma = 8$  to the cycle and ladder graphs, which yielded slightly better results than the default setting. We also include results for the GT algorithm using two additional stepsize values: one larger and one smaller than the empirically optimal value. This is to illustrate the sensitivity of GT to stepsize selection and highlight the significant effort required to properly tune it. In contrast, using an adaptive stepsize, as in AdGT, substantially reduces this burden. As shown in the results, AdGT not only requires little to no tuning in most cases, but it also often achieves better convergence rates compared to the constant-stepsize version. In some scenarios, such as random or star graph topologies, the improvement is significant.

As noted previously, Method-DM fails to converge in all graph settings except for the star graph.

## 4.2 Adaptive Stepsizes via Disparity in $L_i$ -Smoothness

In this section, we study the benefits of using different stepsizes, especially in situations where some agents have larger smoothness constants (higher  $L$ -smoothness). Such differences in smoothness can have a strong impact on the convergence speed. We show that AdGT, which adjusts stepsizes automatically, achieves much faster convergence than methods with fixed stepsizes. To explore this, we consider a decentralized quadratic optimization problem in which each agent  $i \in \{1, \dots, n\}$  has a local quadratic cost function of the form

$$f_i(x) := \frac{1}{2}x^\top A_i x + b_i^\top x,$$

where  $A_i \in \mathbb{S}_{++}^p$  is a positive definite diagonal matrix, and  $b_i \in \mathbb{R}^p$  is a vector. The overall objective is to minimize the sum of all local functions:

$$f(x) := \sum_{i=1}^n \frac{1}{2}x^\top A_i x + b_i^\top x. \quad (47)$$

To model differences in smoothness across the local objectives, we generate matrices  $A_i$  with varying condition numbers, which directly affect the Lipschitz constants  $L_i \approx \|A_i\| = \lambda_{\max}(A_i)$ . Following the setup in [56], we construct each diagonal matrix  $A_i$  by sampling its entries as follows: the first half of the diagonal (i.e., the first  $p/2$  entries) is drawn uniformly at random from the set  $\{1, 10^{-1}, \dots, 10^{-\tau}\}$ , and the second half from  $\{1, 10^1, \dots, 10^\tau\}$ . This ensures that the eigenvalues of each matrix  $A_i$  lie within the interval  $[10^{-\tau}, 10^\tau]$ . The linear terms  $b_i^\top x$  are added to ensure that the local objective functions have different minimizers. The vectors  $b_i$  are generated uniformly at random from the box  $[0, 1]^p$ .

We consider a network of  $n = 100$  nodes, where the dimension of the local variable  $x$  is  $p = 20$ . We evaluate the performance under four different scenarios based on the choice of the condition number parameter  $\tau$ :

- i) All agents are assigned a condition number parameter of  $\tau = 3$ ;
- ii) Half of the agents (i.e., 50 nodes) have  $\tau = 3$ , while the remaining half have  $\tau = 1$ ;
- iii) Ten percent of the agents (i.e., 10 nodes) are assigned  $\tau = 3$ , and the remaining 90 nodes have  $\tau = 1$ ;
- iv) Only 3 agents are assigned  $\tau = 3$ , with all others using  $\tau = 1$ .

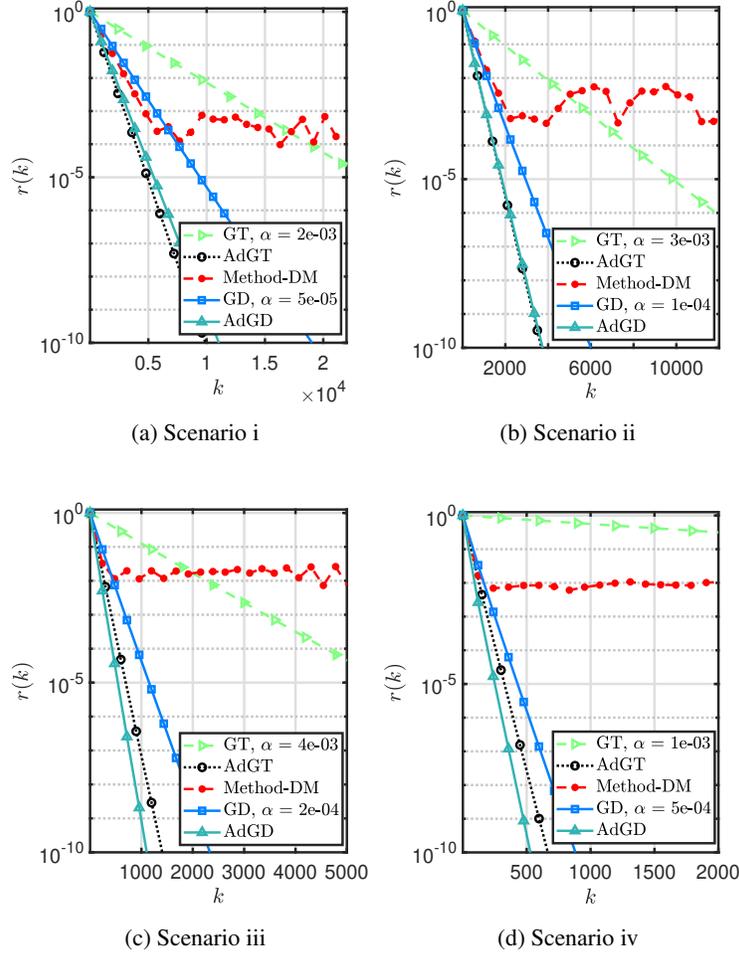
We included AdGD and GD, along with Method-DM and GT, to compare against AdGT. The corresponding numerical results for these four methods are illustrated in Fig. 3. In all simulations, we set  $\gamma = 1$ . For GD and GT, the stepsizes are tuned to achieve the best possible convergence rate. As observed, when the number of nodes with high  $L$ -smoothness decreases, the convergence rate of AdGT improves significantly compared to GT. In contrast, the gap between the convergence rates of AdGD and GD remains relatively unchanged. This highlights the effectiveness of using adaptive stepsizes, particularly in scenarios where the data distributed across nodes exhibit varying levels of smoothness.

## 4.3 Adaptive Stepsizes via number of nodes and connectivity

We consider ridge regression as a specific instance of problem (1), where the local objective function for each agent  $i$  is defined as

$$f_i(x) = \|A_i x_i - b_i\|^2 + \rho \|x_i\|^2,$$

with  $A_i \in \mathbb{R}^{m \times p}$ ,  $b_i \in \mathbb{R}^m$ , and the regularization parameter set to  $\rho = 0.1$ . The data used in this experiment corresponds to the *w8a* dataset, as described in Section 4.1.

Figure 3: Decentralized quadratic problem with  $n = 100$ .

To evaluate performance, we conducted numerical experiments on both random and cycle networks. For the random networks, we generated graphs with a connectivity ratio of 0.35, considering two cases with  $n = 25$  and  $n = 100$  nodes. For the cycle graph, we tested it with  $n = 10$  and  $n = 25$  nodes. In all four experimental settings, the full dataset was evenly partitioned among the  $n$  agents. Each agent was randomly assigned a local dataset of size  $m$ , drawn from the global dataset without replacement. We fix  $\gamma = 1$ , and the stepsizes for GD and GT are carefully tuned to obtain their optimal convergence performance.

The residual sequences  $r(k)_{k \geq 1}$  for all methods are presented in Fig. 4. In the case of the random graph with a moderate connectivity ratio of 0.35, we observe that increasing the number of nodes generally leads to faster convergence for AdGT compared to GT. In contrast, for the cycle graph, AdGT performs better when the number of agents is small. Moreover, Method-DM fails to converge in all scenarios, except for the cycle graph with  $n = 10$  agents. Similar to the results presented in Section 4.2, the difference in convergence rate between AdGT and GT is noticeably larger than that between AdGD and GD, particularly in random graphs with  $n = 100$  nodes.

## 5 Conclusion

In this work, we propose AdGT, a decentralized optimization algorithm that adaptively selects stepsizes based on local information available at each agent. By eliminating the need for manual stepsize tuning, AdGT simplifies prac-

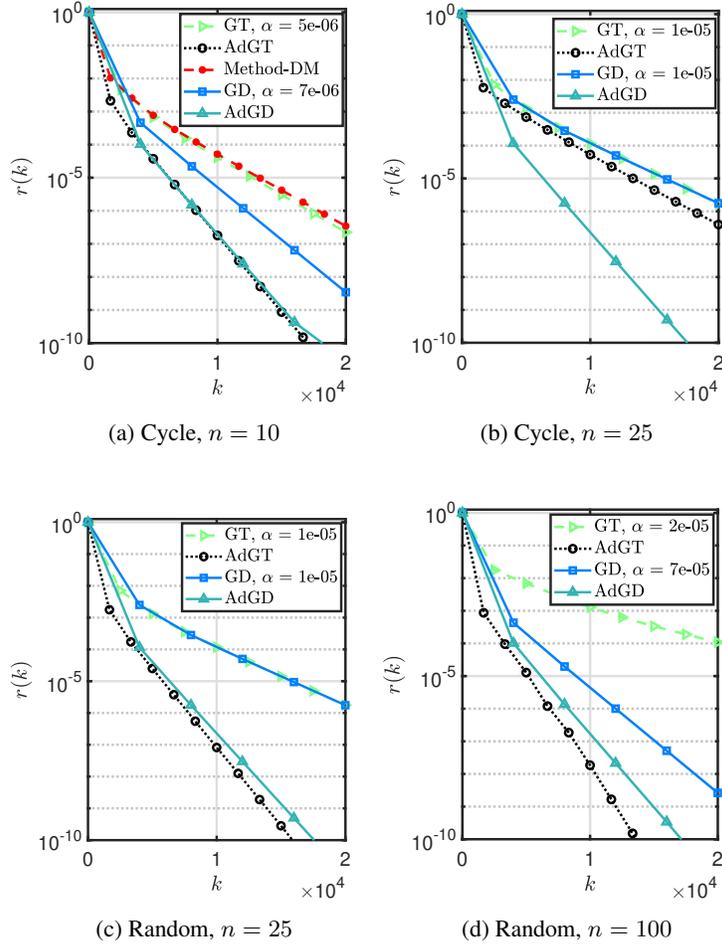


Figure 4: Performance of ridge regression over cycle and random graphs.

tical implementation and improves convergence efficiency, especially in heterogeneous networks with varying local smoothness. Our theoretical analysis establishes that AdGT guarantees consensus convergence. Furthermore, empirical results demonstrate that adaptive stepsizes significantly outperform fixed-stepsize methods in both convergence speed and robustness. These findings underscore the value of local adaptivity in decentralized optimization and suggest promising directions for future research on scalable, communication-efficient algorithms over general network topologies.

## References

- [1] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, “Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent,” *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [2] H. Tang, X. Lian, M. Yan, C. Zhang, and J. Liu, “ $D^2$ : Decentralized training over decentralized data,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 4848–4856.
- [3] A. Nedić and J. Liu, “Distributed optimization for control,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, no. 1, pp. 77–103, 2018.

- [4] J. Tsitsiklis, D. Bertsekas, and M. Athans, “Distributed asynchronous deterministic and stochastic gradient optimization algorithms,” *IEEE Transactions on Automatic Control*, vol. 31, no. 9, pp. 803–812, 1986.
- [5] A. Nedic and A. Ozdaglar, “Distributed subgradient methods for multi-agent optimization,” *IEEE Transactions on Automatic Control*, vol. 54, no. 1, pp. 48–61, 2009.
- [6] D. Jakovetić, J. Xavier, and J. M. Moura, “Fast distributed gradient methods,” *IEEE Transactions on Automatic Control*, vol. 59, no. 5, pp. 1131–1146, 2014.
- [7] W. Shi, Q. Ling, G. Wu, and W. Yin, “Extra: An exact first-order algorithm for decentralized consensus optimization,” *SIAM Journal on Optimization*, vol. 25, no. 2, pp. 944–966, 2015.
- [8] K. Yuan, B. Ying, X. Zhao, and A. H. Sayed, “Exact diffusion for distributed optimization and learning—Part I: Algorithm development,” *IEEE Transactions on Signal Processing*, vol. 67, no. 3, pp. 708–723, 2018.
- [9] E. Wei and A. Ozdaglar, “On the  $\mathcal{O}(1/k)$  convergence of asynchronous distributed alternating direction method of multipliers,” in *IEEE Global Conference on Signal and Information Processing*, 2013, pp. 551–554.
- [10] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin, “On the linear convergence of the ADMM in decentralized consensus optimization,” *IEEE Transactions on Signal Processing*, vol. 62, no. 7, pp. 1750–1761, 2014.
- [11] A. Mokhtari, W. Shi, Q. Ling, and A. Ribeiro, “DQM: Decentralized quadratically approximated alternating direction method of multipliers,” *IEEE Transactions on Signal Processing*, vol. 64, no. 19, pp. 5158–5173, 2016.
- [12] —, “A decentralized second-order method with exact linear convergence rate for consensus optimization,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 4, pp. 507–522, 2016.
- [13] N. Aybat, Z. Wang, and G. Iyengar, “An asynchronous distributed proximal gradient method for composite convex optimization,” in *International Conference on Machine Learning*. PMLR, 2015, pp. 2454–2462.
- [14] N. S. Aybat, Z. Wang, T. Lin, and S. Ma, “Distributed linearized alternating direction method of multipliers for composite convex consensus optimization,” *IEEE Transactions on Automatic Control*, vol. 63, no. 1, pp. 5–20, 2017.
- [15] N. S. Aybat and E. Yazdandoost Hamedani, “A primal-dual method for conic constrained distributed optimization problems,” in *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, Eds., vol. 29. Curran Associates, Inc., 2016.
- [16] N. S. Aybat and E. Y. Hamedani, “A distributed ADMM-like method for resource sharing over time-varying networks,” *SIAM Journal on Optimization*, vol. 29, no. 4, pp. 3036–3068, 2019.
- [17] P. Di Lorenzo and G. Scutari, “Distributed nonconvex optimization over networks,” in *IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, 2015, pp. 229–232.
- [18] —, “NEXT: In-network nonconvex optimization,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.
- [19] A. Nedic, A. Olshevsky, and W. Shi, “Achieving geometric convergence for distributed optimization over time-varying graphs,” *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2597–2633, 2017.
- [20] G. Qu and N. Li, “Harnessing smoothness to accelerate distributed optimization,” *IEEE Transactions on Control of Network Systems*, vol. 5, no. 3, pp. 1245–1260, 2017.
- [21] J. Xu, S. Zhu, Y. C. Soh, and L. Xie, “Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes,” in *54th IEEE Conference on Decision and Control (CDC)*, 2015, pp. 2055–2060.
- [22] C. Xi, R. Xin, and U. A. Khan, “ADD-OPT: Accelerated distributed directed optimization,” *IEEE Transactions on Automatic Control*, vol. 63, no. 5, pp. 1329–1339, 2017.
- [23] J. Xu, S. Zhu, Y. C. Soh, and L. Xie, “Convergence of asynchronous distributed gradient methods over stochastic networks,” *IEEE Transactions on Automatic Control*, vol. 63, no. 2, pp. 434–448, 2017.

- [24] S. Pu and A. Nedić, “Distributed stochastic gradient tracking methods,” *Mathematical Programming*, vol. 187, no. 1, pp. 409–457, 2021.
- [25] A. Koloskova, T. Lin, and S. U. Stich, “An improved analysis of gradient tracking for decentralized machine learning,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 11 422–11 435, 2021.
- [26] Y. Sun, G. Scutari, and A. Daneshmand, “Distributed optimization based on gradient tracking revisited: Enhancing convergence rate via surrogation,” *SIAM Journal on Optimization*, vol. 32, no. 2, pp. 354–385, 2022.
- [27] Y. Tian, Y. Sun, and G. Scutari, “Achieving linear convergence in distributed asynchronous multiagent optimization,” *IEEE Transactions on Automatic Control*, vol. 65, no. 12, pp. 5264–5279, 2020.
- [28] R. Xin and U. A. Khan, “A linear algorithm for optimization over directed graphs with geometric convergence,” *IEEE Control Systems Letters*, vol. 2, no. 3, pp. 315–320, 2018.
- [29] Y. Lu and C. De Sa, “Optimal complexity in decentralized training,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 7111–7123.
- [30] Y. Liu, T. Lin, A. Koloskova, and S. U. Stich, “Decentralized gradient tracking with local steps,” *Optimization Methods and Software*, pp. 1–28, 2024.
- [31] D. Ghaderyan, N. S. Aybat, A. P. Aguiar, and F. L. Pereira, “A fast row-stochastic decentralized method for distributed optimization over directed graphs,” *IEEE Transactions on Automatic Control*, vol. 69, no. 1, pp. 275–289, 2023.
- [32] G. Qu and N. Li, “Accelerated distributed nesterov gradient descent,” *IEEE Transactions on Automatic Control*, vol. 65, no. 6, pp. 2566–2581, 2019.
- [33] D. Kovalev, E. Gasanov, A. Gasnikov, and P. Richtarik, “Lower bounds and optimal algorithms for smooth and strongly convex decentralized optimization over time-varying networks,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 22 325–22 335, 2021.
- [34] A. Rogozin, C. A. Uribe, A. V. Gasnikov, N. Malkovsky, and A. Nedić, “Optimal distributed convex optimization on slowly time-varying graphs,” *IEEE Transactions on Control of Network Systems*, vol. 7, no. 2, pp. 829–841, 2019.
- [35] H. Ye, L. Luo, Z. Zhou, and T. Zhang, “Multi-consensus decentralized accelerated gradient descent,” *Journal of Machine Learning Research*, vol. 24, no. 306, pp. 1–50, 2023.
- [36] H. Ye, Z. Zhou, L. Luo, and T. Zhang, “Decentralized accelerated proximal gradient descent,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 18 308–18 317, 2020.
- [37] H. Li, C. Fang, W. Yin, and Z. Lin, “Decentralized accelerated gradient methods with increasing penalty parameters,” *IEEE Transactions on Signal Processing*, vol. 68, pp. 4855–4870, 2020.
- [38] J. Xu, Y. Tian, Y. Sun, and G. Scutari, “Accelerated primal-dual algorithms for distributed smooth convex optimization over networks,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2381–2391.
- [39] H. Li and Z. Lin, “Revisiting extra for smooth distributed optimization,” *SIAM Journal on Optimization*, vol. 30, no. 3, pp. 1795–1821, 2020.
- [40] —, “Accelerated gradient tracking over time-varying graphs for decentralized optimization,” *Journal of Machine Learning Research*, vol. 25, no. 274, pp. 1–52, 2024.
- [41] K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié, “Optimal algorithms for smooth and strongly convex distributed optimization in networks,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 3027–3036.
- [42] A. Nedić, A. Olshevsky, W. Shi, and C. A. Uribe, “Geometrically convergent distributed optimization with uncoordinated step-sizes,” in *2017 American Control Conference (ACC)*. IEEE, 2017, pp. 3950–3955.

- [43] R. Xin, C. Xi, and U. A. Khan, “Frost—fast row-stochastic optimization with uncoordinated step-sizes,” *EURASIP Journal on Advances in Signal Processing*, vol. 2019, pp. 1–14, 2019.
- [44] Y. Malitsky and K. Mishchenko, “Adaptive gradient descent without descent,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 6702–6712.
- [45] P. Latafat, A. Themelis, L. Stella, and P. Patrinos, “Adaptive proximal algorithms for convex optimization under local Lipschitz continuity of the gradient,” *Mathematical Programming*, pp. 1–39, 2024.
- [46] Y. Malitsky and K. Mishchenko, “Adaptive proximal gradient method for convex optimization,” *Advances in Neural Information Processing Systems*, vol. 37, pp. 100 670–100 697, 2024.
- [47] D. Zhou, S. Ma, and J. Yang, “Adabb: Adaptive barzilai-borwein method for convex optimization,” *Mathematics of Operations Research*, 2025.
- [48] J. Wang, Z. Xu, Z. Garrett, Z. Charles, L. Liu, and G. Joshi, “Local adaptivity in federated learning: Convergence and consistency,” *arXiv preprint arXiv:2106.02305*, 2021.
- [49] Y. Wang, L. Lin, and J. Chen, “Communication-efficient adaptive federated learning,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 22 802–22 838.
- [50] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, “Adaptive federated optimization,” *arXiv preprint arXiv:2003.00295*, 2020.
- [51] J. L. Kim, T. Toghiani, C. A. Uribe, and A. Kyrillidis, “Adaptive federated learning with auto-tuned clients via local smoothness,” in *Federated Learning and Analytics in Practice: Algorithms, Systems, Applications, and Opportunities*, 2023.
- [52] Z. Chen and Y. Wang, “Distributed optimization and learning with automated stepsizes,” in *IEEE 63rd Conference on Decision and Control (CDC)*, 2024, pp. 3121–3126.
- [53] G. Scutari and Y. Sun, “Distributed nonconvex constrained optimization over time-varying digraphs,” *Mathematical Programming*, vol. 176, pp. 497–544, 2019.
- [54] D. Jakovetić, “A unification and generalization of exact distributed first-order methods,” *IEEE Transactions on Signal and Information Processing over Networks*, vol. 5, no. 1, pp. 31–46, 2018.
- [55] C.-C. Chang and C.-J. Lin, “LIBSVM: A library for support vector machines,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 2, no. 3, pp. 1–27, 2011.
- [56] A. Mokhtari, Q. Ling, and A. Ribeiro, “Network newton distributed optimization methods,” *IEEE Transactions on Signal Processing*, vol. 65, no. 1, pp. 146–161, 2016.

This figure "fig1.png" is available in "png" format from:

<http://arxiv.org/ps/2504.15196v1>