

On Centralized and Distributed Mirror Descent: Convergence Analysis Using Quadratic Constraints

Youbang Sun¹, Mahyar Fazlyab² and Shahin Shahrampour¹

Abstract—Mirror descent (MD) is a powerful first-order optimization technique that subsumes several optimization algorithms including gradient descent (GD). In this work, we develop a semi-definite programming (SDP) framework to analyze the convergence rate of MD in centralized and distributed settings under both strongly convex and non-strongly convex assumptions. We view MD with a dynamical system lens and leverage quadratic constraints (QCs) to provide explicit convergence rates based on Lyapunov stability. For centralized MD under strongly convex assumption, we develop a SDP that certifies exponential convergence rates. We prove that the SDP always has a feasible solution that recovers the optimal GD rate as a special case. We complement our analysis by providing the $O(1/k)$ convergence rate for convex problems. Next, we analyze the convergence of distributed MD and characterize the rate using SDP. To the best of our knowledge, the numerical rate of distributed MD has not been previously reported in the literature. We further prove an $O(1/k)$ convergence rate for distributed MD in the convex setting. Our numerical experiments on strongly convex problems indicate that our framework certifies superior convergence rates compared to the existing rates for distributed GD.

I. INTRODUCTION

Over the last two decades, distributed optimization over multi-agent networks has received a lot of attention in control, optimization, machine learning, and signal processing. In distributed optimization, a group of n agents are connected via a graph and can communicate locally with their neighbors. Each agent is assigned a local objective function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, and the agents aim to collectively minimize the global objective function,

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}. \quad (1)$$

The most intuitive gradient-based algorithm to tackle the problem above is *distributed gradient descent* [1], where at each iteration k , each agent i updates its decision variables by a (private) local gradient descent combined with an averaging of its neighbors variables. In the unconstrained case, this update is given by

$$x_i^{(k+1)} = x_i^{(k)} - \eta^{(k)} \nabla f_i(x_i^{(k)}) + \beta \sum_{j \in \mathcal{N}_i} (x_j^{(k)} - x_i^{(k)}),$$

where $\eta^{(k)} > 0$ is the step-size and $\beta > 0$ is the consensus parameter. In the form given above, this update is able to achieve optimal rates for convex problems using a diminishing

step-size sequence. Optimality here refers to matching the centralized convergence rate (iteration complexity) up to some errors related to the network structure. However, when the local functions are smooth, the centralized GD algorithm employs a *constant* step-size sequence for which the above distributed counterpart fails to converge.

The mirror descent (MD) algorithm [2] is a primal-dual method that has been actively studied in recent years. MD can be seen as a generalization of GD, in which the squared Euclidean distance is replaced by Bregman divergence as the regularizer. The freedom in the choice of Bregman divergence makes MD suitable for various problem geometries. MD has been proven to have the same iteration complexity as GD for non-strongly convex problems [3], and it may even scale better with respect to the dimension of the decision variables [4]. In the strongly convex scenario, MD is less studied, and very recently its exponential convergence was established under the Polyak-Łojasiewicz (PL) condition [5]. Inspired by the success of MD in centralized optimization, MD has also been studied in the distributed setting. To the best of authors' knowledge, the convergence rate of distributed MD is not established for strongly convex and smooth problems, and only recently [6] provided a continuous-time analysis suggesting local exponential rate (without explicitly characterizing the rate).

In this paper, we leverage the framework of Quadratic Constraints (QCs) to certify numerical exponential convergence rates for centralized as well as distributed MD for strongly convex and smooth problems using SDP. For merely convex and smooth problems, we also establish an ergodic $O(1/k)$ convergence rate. We first analyze centralized MD, for which we derive linear matrix inequalities (LMIs) as sufficient conditions for convergence of the algorithm at a specified rate (Theorem 2, Theorem 6 and Proposition 3). For the strongly convex case, we prove that these LMIs always have a feasible solution that matches the optimal convergence rate of GD when the Bregman divergence is chosen as the squared Euclidean distance (Proposition 4 and Corollary 5). Next, we analyze the convergence of distributed MD and characterize the rate using LMIs (Theorem 8, 9). To the best of our knowledge, the exponential rate of distributed MD has not been previously established in the literature. Our numerical experiments on strongly convex problems indicate that our framework certifies superior convergence rates compared to the existing rates for distributed GD.

A. Related Literature

1) *Distributed Optimization*: To ensure that distributed GD (or sub-gradient descent) reach consensus, many methods [1], [7], [8] use diminishing step-size (commonly $1/k$).

This work is supported in part by NSF ECCS-2136206 Award.

¹ Y. Sun and S. Shahrampour are with the Department of Mechanical and Industrial Engineering at Northeastern University, Boston, MA 02115, USA. email: {sun.youb, s.shahrampour}@northeastern.edu.

² M. Fazlyab is with the Department of Electrical and Computer Engineering at Johns Hopkins University, Baltimore, MD 21218, USA. email: mahyarfazlyab@jhu.edu.

For distributed MD, similar studies have been conducted for stochastic optimization [9], [10] and online optimization [11], [12]. Doan et al. [13] provide convergence results for both centralized and decentralized MD algorithms. However, convergence rates obtained using diminishing step-size are sub-exponential and sub-optimal under assumptions of strong convexity and smoothness.

To address this issue, a number of recent works introduce an additional variable in the state vectors to track past gradients (see e.g., [14]–[17]). One of the earlier works in this direction is the EXTRA algorithm proposed by Shi et al. [14], which uses the information from past two iterations to perform each update. For smooth problems, EXTRA provably achieves $O(1/k)$ convergence rate under the convexity assumption and exponential convergence rate under the strong convexity assumption, respectively.

A closely relevant literature is the continuous-time distributed GD, where the algorithms are constructed by a set of ordinary differential equations (ODEs). These works are mostly based on the idea of *integral feedback*, which can be thought as the continuous-time analog of gradient tracking. In this case, each agent uses an integration term as a part of the ODE (see e.g., [18]–[21]). In these works, the analysis is carried out by proving the Lyapunov stability for the corresponding continuous-time dynamics, and exponential stability can be obtained in certain cases [20]. For MD, the continuous-time algorithm in [6], [22] and the discrete-time algorithm in [23] both adapt the integral feedback (or gradient tracking) method and propose algorithms that do not suffer from sub-optimal convergence rates. Specifically, Sun et al. [6] propose a continuous-time distributed MD that achieves a “local” exponential rate for strongly convex problems, and Yu et al. [23] provide an $O(1/k)$ convergence rate under the convexity assumption in discrete time. Nevertheless, the exponential rate of (discrete-time) distributed MD for strongly convex and smooth problems remains an open problem, which we target in the current work.

2) *Integral Quadratic Constraints*: Deriving convergence rates for iterative optimization algorithms in the worst-case is an integral part of algorithm design. However, this procedure is not principled, requires a case-by-case analysis, and might lead to conservative rates. To automate convergence analysis and derive sharp convergence rates, several past works have used Integral Quadratic Constraints (IQCs) and semidefinite programming in various settings [24]–[30], pioneered by the work in [24]. IQCs are a tool from robust control to analyze dynamical systems that contain components that are nonlinear, uncertain, or difficult to model [31]. The basic idea is to abstract these troublesome components by constraints on their input and output signals. This approach to algorithm analysis can also guide the search for parameter selection in algorithm design. [32], [33] are of particular relevance to our work. They both provide IQC-based analysis of distributed gradient-based algorithms in strongly convex settings. Compared to these works, our framework focuses on distributed mirror descent in both strongly convex and convex settings.

II. PRELIMINARIES

A. Notations

The identity matrix of dimension n is denoted by I_n and the n -dimensional vector with all entries 1 is represented by $\mathbf{1}_n$. We denote the set of n -dimensional symmetric matrices by S^n . The positive (negative) semi-definiteness of matrix M is denoted as $M \succeq 0$ ($M \preceq 0$). We use \otimes and $\|\cdot\|$ to denote the Kronecker product and spectral norm, respectively. We define the norm of vector v with respect to a positive semi-definite matrix M as $\|v\|_M$. The indicator function of a set $\mathcal{X} \subseteq \mathbb{R}^d$ is defined as $\mathbb{1}_{\mathcal{X}}(x) = 0$ if $x \in \mathcal{X}$ and $\mathbb{1}_{\mathcal{X}}(x) = +\infty$ otherwise.

Definition 1 (Strong convexity). *A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ_f -strongly convex on \mathbb{R}^d if the following inequality is true for all $x, y \in \mathbb{R}^d$.*

$$f(x) + \nabla f(x)^\top (y - x) + \frac{\mu_f}{2} \|y - x\|^2 \leq f(y).$$

Definition 2 (Lipschitz smoothness). *A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L_f -smooth on \mathbb{R}^d if $\frac{L_f}{2} \|x\|^2 - f(x)$ is convex, which implies that for all $x, y \in \mathbb{R}^d$.*

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L_f}{2} \|y - x\|^2.$$

We further denote the condition number of function f by $\kappa_f \triangleq \frac{L_f}{\mu_f} \geq 1$. When $\mu_f = 0$, the function is only convex.

Proposition 1. *Suppose f is μ_f -strongly convex and L_f -smooth on \mathbb{R}^d . Then, the following inequality holds for all $x, y \in \mathbb{R}^d$, and $u = \nabla f(x)$, $v = \nabla f(y)$,*

$$\begin{bmatrix} x - y \\ u - v \end{bmatrix}^\top \begin{bmatrix} \frac{-\mu_f L_f}{\mu_f + L_f} I_d & \frac{1}{2} I_d \\ \frac{1}{2} I_d & \frac{-1}{\mu_f + L_f} I_d \end{bmatrix} \begin{bmatrix} x - y \\ u - v \end{bmatrix} \geq 0. \quad (2)$$

The above QC follows from the combination of strong convexity and Lipschitz smoothness [24], [34].

B. Centralized Mirror Descent Algorithm

We start by providing some background on the centralized MD algorithm. For simplicity in the exposition, we study the unconstrained case, but our analysis can also be extended to the constrained case. Let us start with the GD algorithm, whose update is equivalent to the following minimization,

$$x^{(k+1)} = \underset{x}{\operatorname{argmin}} \left\{ f(x^{(k)}) + \nabla f(x^{(k)})^\top (x - x^{(k)}) + \frac{1}{2\eta} \|x - x^{(k)}\|^2 \right\},$$

where $\eta > 0$ is the step-size. In each iteration, the algorithm seeks to minimize a first-order approximation of the function with a Euclidean regularizer. As a generalization of gradient descent, MD replaces the squared Euclidean distance with Bregman divergence, which is defined with respect to a distance generating function (DGF) $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ as follows

$$D_\phi(x, x') \triangleq \phi(x) - \phi(x') - \langle \nabla \phi(x'), x - x' \rangle. \quad (3)$$

Assumption 1. *The distance generating function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ_ϕ -strongly convex and L_ϕ -smooth.*

The centralized (unconstrained) MD algorithm with step-size η is written as

$$x^{(k+1)} = \underset{x}{\operatorname{argmin}} \left\{ f(x^{(k)}) + \nabla f(x^{(k)})^\top (x - x^{(k)}) + \frac{1}{\eta} \mathcal{D}_\phi(x, x^{(k)}) \right\}, \quad (4)$$

where if we choose the Bregman divergence to be the squared Euclidean distance, the update above reduces to GD.

We can also view the MD update through a different lens using the *convex conjugate* of function ϕ . The convex conjugate of function ϕ , denoted by ϕ^* , is defined as $\phi^*(z) \triangleq \sup_{x \in \mathbb{R}^d} \{ \langle x, z \rangle - \phi(x) \}$. Assumption 1 guarantees that ϕ^* is L_ϕ^{-1} -strongly convex and μ_ϕ^{-1} -smooth. We refer the reader to [35] for further details. Correspondingly, the following equivalence can be established,

$$z' = \nabla \phi(x') \iff x' = \nabla \phi^*(z').$$

Then, the centralized MD update can be rewritten in the following form,

$$\begin{aligned} z^{(k+1)} &= z^{(k)} - \eta \nabla f(x^{(k)}) \\ x^{(k+1)} &= \nabla \phi^*(z^{(k+1)}), \end{aligned} \quad (5)$$

or, equivalently, $z^{(k+1)} = z^{(k)} - \eta(\nabla f \circ \nabla \phi^*)(z^{(k)})$, which is reminiscent of GD. We can see that MD is more general than GD in that we can exploit the geometry of the problem using an appropriate choice of ϕ , which makes MD more suitable for problems such as convex clustering, matrix optimization with regularization, etc. [36], [37].

Denoting x^* and z^* as the fixed points of (5), we have

$$z^* = z^* - \eta \nabla f(x^*) \quad x^* = \nabla \phi^*(z^*),$$

which implies that x^* is a minimizer of f .

III. CONVERGENCE ANALYSIS OF CENTRALIZED MIRROR DESCENT

In this section, we provide a convergence analysis of the centralized MD using semidefinite programming. Our starting point is to describe all the nonlinear functions in the algorithm, namely ∇f and $\nabla \phi^*$ by QCs on their input-output pairs, resulting in a *quadratically-constrained linear system*. We then find a suitable ‘‘rate-generating’’ Lyapunov function for this constrained system using semidefinite programming. We derive exponential (respectively, sub-exponential) convergence rate for strongly convex (respectively, convex) functions.

A. Exponential Convergence for Strongly Convex f

In the following theorem, we characterize an LMI that depends on parameters of f (μ_f and L_f), parameters of ϕ (μ_ϕ and L_ϕ), and several decision variables (including the step-size η and the convergence rate $\rho \in (0, 1)$). We prove that if the LMI is satisfied, the iterates converge exponentially fast to the unique fixed point (x^*, z^*) with the rate ρ .

Theorem 2. *Let Assumption 1 hold and assume that f is μ_f -strongly convex and L_f -smooth. Define matrices M_{sc}, M_f, M_ϕ as follows,*

$$\begin{aligned} M_{sc} &= \begin{bmatrix} \frac{1-\rho}{2\mu_\phi} I_d & 0 & 0 \\ 0 & 0 & \frac{-\eta}{2} I_d \\ 0 & \frac{-\eta}{2} I_d & \frac{\eta^2}{2\mu_\phi} I_d \end{bmatrix} \\ M_f &= \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{-\mu_f L_f}{\mu_f + L_f} I_d & \frac{1}{2} I_d \\ 0 & \frac{1}{2} I_d & \frac{-1}{\mu_f + L_f} I_d \end{bmatrix} \\ M_\phi &= \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} I_d & \frac{1}{2} I_d & 0 \\ \frac{1}{2} I_d & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} I_d & 0 \\ 0 & 0 & 0 \end{bmatrix}. \end{aligned} \quad (6)$$

If there exist some $\rho \in (0, 1), \eta > 0, \sigma_f \geq 0, \sigma_\phi \geq 0$, such that the following matrix inequality holds

$$M_{sc} + \sigma_f M_f + \sigma_\phi M_\phi \preceq 0, \quad (7)$$

then the mirror descent algorithm in (5) converges exponentially fast with the rate of ρ . In particular,

$$\|x^{(k)} - x^*\|^2 \leq \frac{2\mathcal{D}_{\phi^*}(z^{(0)}, z^*)}{\mu_\phi} \rho^k.$$

Proof. Denote $u^{(k)} \triangleq \nabla f(x^{(k)})$ and define the stacked vector

$$e^{(k)\top} = \left[(z^{(k)} - z^*)^\top \quad (x^{(k)} - x^*)^\top \quad (u^{(k)} - u^*)^\top \right]. \quad (8)$$

Then, from Proposition 1, we obtain the following quadratic inequalities

$$e^{(k)\top} M_f e^{(k)} \geq 0, \quad e^{(k)\top} M_\phi e^{(k)} \geq 0 \quad \forall k,$$

which are imposed by ∇f and $\nabla \phi$, respectively. Consider the Lyapunov candidate $V^{(k)} = \rho^{-k} \mathcal{D}_{\phi^*}(z^{(k)}, z^*)$. Recall that ϕ^* is L_ϕ^{-1} -strongly convex and μ_ϕ^{-1} -smooth, so the Lyapunov function is indeed non-negative and continuously differentiable. Using Lemma 10 (provided in the appendix of [38]), we can calculate the Lyapunov function difference between two consecutive iterations as

$$V^{(k+1)} - V^{(k)} \leq \rho^{-k-1} e^{(k)\top} M_{sc} e^{(k)}. \quad (9)$$

Utilizing the two quadratic inequalities imposed by the nonlinearities, we can write

$$\begin{aligned} V^{(k+1)} - V^{(k)} &\leq \rho^{-k-1} e^{(k)\top} M_{sc} e^{(k)} \\ &\leq \rho^{-k-1} e^{(k)\top} (M_{sc} + \sigma_f M_f + \sigma_\phi M_\phi) e^{(k)}. \end{aligned}$$

Now if the LMI in (7) holds, the Lyapunov function is non-increasing, which yields

$$\mathcal{D}_{\phi^*}(z^{(k)}, z^*) = \rho^k V^{(k)} \leq \rho^k V^{(0)} = \rho^k \mathcal{D}_{\phi^*}(z^{(0)}, z^*). \quad (10)$$

Observing $\mathcal{D}_{\phi^*}(z^{(k)}, z^*) = \mathcal{D}_\phi(x^*, x^{(k)})$ and

$$\frac{\mu_\phi}{2} \|x^{(k)} - x^*\|^2 \leq \mathcal{D}_\phi(x^*, x^{(k)}),$$

completes the proof. \square

Theorem 2 provides a matrix inequality feasibility problem that establishes the exponential convergence rate of MD for a

given ρ . This matrix inequality is linear in $(\rho, \sigma_f, \sigma_\phi)$ (but not in η), allowing us to find the smallest ρ by the semidefinite program

$$\begin{aligned} & \underset{\rho, \sigma_\phi, \sigma_f}{\text{minimize}} && \rho \\ & \text{subject to} && 0 < \rho \leq 1 \\ & && \eta, \sigma_\phi, \sigma_f \geq 0 \\ & && M_{sc} + \sigma_f M_f + \sigma_\phi M_\phi \preceq 0. \end{aligned} \quad (11)$$

If in addition we want to optimize ρ over the step-size η , we can use Schur Complements to ‘‘convexify’’ the matrix inequality with respect to η . We state this result formally in the next proposition.

Proposition 3. *The optimization problem in (11) is equivalent to the following SDP,*

$$\begin{aligned} & \underset{\eta, \rho, \sigma_\phi, \sigma_f}{\text{minimize}} && \rho \\ & \text{subject to} && 0 < \rho \leq 1 \\ & && \eta, \sigma_\phi, \sigma_f \geq 0 \end{aligned} \quad (12)$$

$$\begin{bmatrix} \frac{\sigma_\phi}{\mu_\phi + L_\phi} + \frac{\rho - 1}{2\mu_\phi} & -\frac{\sigma_\phi}{2} & 0 & 0 \\ -\frac{\sigma_\phi}{2} & \frac{\mu_\phi L_\phi \sigma_\phi}{\mu_\phi + L_\phi} + \frac{\mu_\phi}{2} + \frac{\mu_f L_f \sigma_f}{\mu_f + L_f} & -\frac{\sigma_f}{2} & \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ 0 & -\frac{\sigma_f}{2} & \frac{\sigma_f}{\mu_f + L_f} & \frac{\eta}{\sqrt{2\mu_\phi}} \\ 0 & \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} & \frac{\eta}{\sqrt{2\mu_\phi}} & 1 \end{bmatrix} \succeq 0$$

We refer to the appendix of [38] for the proof of this proposition. We now show that the SDP in (12) has a feasible solution for which we can analytically calculate the convergence rate.

Proposition 4. *The following selection*

$$\begin{aligned} \eta &= \sigma_f = \frac{2\mu_\phi}{\mu_f + L_f} \\ \sigma_\phi &= \frac{4\mu_f L_f}{(\mu_f + L_f)^2} \frac{(1 + \kappa_\phi)}{\kappa_\phi(\kappa_\phi - 1)} \\ \rho_{opt} &= 1 - \frac{4\mu_f L_f}{(\mu_f + L_f)^2 \kappa_\phi^2}, \end{aligned} \quad (13)$$

is a feasible solution to the SDP in (12).

The proof of the proposition can be found in the appendix of [38]. Note that ρ_{opt} is an upper bound on the optimal value of (12).

The recent work of [5] also proposed an explicit rate of $1 - \frac{1}{5\kappa_\phi^2 \kappa_f^2}$ for MD under the PL condition. Though PL condition is weaker than strong convexity, ρ_{opt} is strictly smaller than the rate of [5]. Furthermore, in our result we do not make full use of strong convexity: we only require the quadratic inequality (2) to hold for the pair (x, x^*) (x arbitrary and x^* the fixed point of the algorithm), whereas for strongly convex f this inequality holds for all (x, y) . Our rate also recovers the optimal rate of GD as a special case.

Corollary 5. *For $\phi(x) = \frac{1}{2} \|x\|^2$ the optimal rate ρ_{opt} in (13) coincides with the optimal convergence rate of gradient descent.*

Proof. If $\phi(x) = \frac{1}{2} \|x\|^2$, we have that $\phi^*(z) = \frac{1}{2} \|z\|^2$ and (5) is equivalent to GD. In this case, the condition number $\kappa_\phi = \frac{L_\phi}{\mu_\phi} = 1$, and ρ_{opt} reduces to the optimal convergence rate for GD (see Theorem 2.1.15 in [34]). \square

B. $O(1/k)$ Convergence for Convex f

We now propose an LMI which establishes subexponential convergence rate for the MD algorithm when the objective function is convex ($\mu_f = 0$).

Theorem 6. *Let Assumption 1 hold and assume that f is convex ($\mu_f = 0$) and L_f -smooth ($0 < L_f < \infty$), and define the matrix M_c as follows,*

$$M_c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \frac{\epsilon - \eta}{2} I \\ 0 & \frac{\epsilon - \eta}{2} I & \frac{\eta}{2\mu_\phi} I \end{bmatrix}. \quad (14)$$

If there exist some $\eta > 0, \sigma_f \geq 0, \sigma_\phi \geq 0, \epsilon \geq 0$, such that the following matrix inequality holds

$$M_c + \sigma_f M_f + \sigma_\phi M_\phi \preceq 0, \quad (15)$$

then the ergodic mean of function value at iteration K satisfies

$$f(\bar{x}^{(K)}) - f(x^*) \leq \frac{\mathcal{D}_{\phi^*}(z^{(0)}, z^*)}{\epsilon K},$$

where $\bar{x}^{(K)} = \frac{1}{K} \sum_{i=1}^K x^{(i)}$.

We remark that a similar analysis can be applied to Theorem 6 to find the best step-size that maximizes ϵ . The details are omitted due to space limitation.

Remark 1 (Constrained Mirror Descent). *Consider the constrained version of centralized (lazy) MD [39],*

$$\begin{aligned} z^{(k+1)} &= z^{(k)} - \eta \nabla f(x^{(k)}) \\ s^{(k)} &= \nabla \phi^*(z^{(k)}) \\ x^{(k)} &= \arg \min_{x \in \mathcal{X}} \mathcal{D}_\phi(x, s^{(k)}), \end{aligned} \quad (16)$$

where \mathcal{X} is a convex subset of \mathbb{R}^d . By defining $g(x) = \mathbb{1}_{\mathcal{X}}(x)$ as the indicator function of the set \mathcal{X} and denoting its subdifferential by ∂g , the optimality condition that characterizes $x^{(k)}$ is

$$\nabla \phi(x^{(k)}) - z^k \in \partial g(x^{(k)}),$$

Using the fact that the subdifferential ∂g is monotone (since \mathcal{X} is convex), we can rewrite (16) as

$$\begin{aligned} z^{(k+1)} &= z^{(k)} - \eta u^{(k)} \\ u^{(k)} &\triangleq \nabla f(x^{(k)}) \\ v^{(k)} &\triangleq \nabla \phi(x^{(k)}), \end{aligned} \quad (17)$$

subject to the quadratic constraint

$$(v^{(k)} - v^* - (z^{(k)} - z^*))^\top (x^{(k)} - x^*) \geq 0 \quad \forall k,$$

Furthermore, we can write two separate quadratic constraints for the relationships $u^{(k)} = \nabla f(x^{(k)})$ and $v^{(k)} = \nabla \phi(x^{(k)})$. We can therefore employ the same approach and derive an LMI as a sufficient condition to establish exponential and $O(1/k)$ convergence rates for strongly convex and convex problems, respectively.

IV. CONVERGENCE ANALYSIS OF DISTRIBUTED MIRROR DESCENT

In the distributed setup, we have a network of agents, characterized by an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each node in $\mathcal{V} = \{1, \dots, n\}$ represents an agent, and the connection between two agents i and j is captured by the edge $\{i, j\} \in \mathcal{E}$. We use $\mathcal{N}_i \triangleq \{j \in \mathcal{V} : \{i, j\} \in \mathcal{E}\}$ to denote the neighborhood of agent i . The graph Laplacian is represented by $\mathcal{L} \in \mathbb{R}^{n \times n}$.

Assumption 2. *The graph \mathcal{G} is undirected and connected, i.e., there exists a path between any two distinct agents $i, j \in \mathcal{V}$.*

The connectivity assumption implies that \mathcal{L} has a unique null eigenvalue; that is, $\mathcal{L}\mathbf{1}_n = 0$.

A. Distributed Mirror Descent Algorithm

We first introduce the distributed MD update, in which each agent i in the network implements the following iterative algorithm,

$$\begin{aligned} z_i^{(k+1)} &= z_i^{(k)} - \eta_1 (\nabla f_i(x_i^{(k)}) + y_i^{(k)}) - \eta_2 \sum_{j \in \mathcal{N}_i} (z_i^{(k)} - z_j^{(k)}), \\ y_i^{(k+1)} &= y_i^{(k)} + \eta_2 \sum_{j \in \mathcal{N}_i} (z_i^{(k)} - z_j^{(k)}), \\ x_i^{(k+1)} &= \nabla \phi^*(z_i^{(k+1)}). \end{aligned} \quad (18)$$

The first update uses private gradient information as well as the dual variables from the neighbors. It also depends on a variable $y_i^{(k)}$ which acts as an integrator. This algorithm is similar to the discretized version of the distributed MD proposed in [22] using the idea of integral feedback. However, the method differs slightly in the local averaging in that the algorithm in [22] performs local averaging with respect to the primal variable, and here the averaging is done on the dual variable $z_i^{(k)}$.

It is evident that the behavior of this system relies on the network structure through the dependence on the Laplacian of the graph capturing the network. Since $\mathcal{L} \in \mathbb{S}^n$, the LMIs will consist of matrices whose dimensions scale with n , which is not suitable when n is large. Following the idea in [32], [33], we transform the updates such that the dependence on the *full structure* of the network is avoided. Define

$$W \triangleq I_n - \eta_2 \mathcal{L} = \Delta W + \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top,$$

and further denote the spectral norm of ΔW by $\lambda \triangleq \|\Delta W\|$. The quantity $1 - \lambda$ is also known as the spectral gap.

To represent the updates collectively for all the agents, we define the stacked vectors

$$\begin{aligned} z^{(k)} &= [z_1^{(k)\top}, \dots, z_n^{(k)\top}]^\top \\ y^{(k)} &= [y_1^{(k)\top}, \dots, y_n^{(k)\top}]^\top \\ u^{(k)} &= \nabla \mathbf{f}(x^{(k)}) \triangleq [\nabla f_1(x_1^{(k)})^\top, \dots, \nabla f_n(x_n^{(k)})^\top]^\top \\ x^{(k)} &= [\nabla \phi^*(z_1^{(k)})^\top, \dots, \nabla \phi^*(z_n^{(k)})^\top]^\top \\ v^{(k)} &= (\Delta W \otimes I_d) z^{(k)}. \end{aligned} \quad (19)$$

We can now rewrite (18) as

$$\begin{aligned} z^{(k+1)} &= \left(\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d\right) z^{(k)} - \eta_1 (u^{(k)} + y^{(k)}) + v^{(k)} \\ y^{(k+1)} &= y^{(k)} + \left((I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes I_d\right) z^{(k)} - v^{(k)} \\ v^{(k)} &= (\Delta W \otimes I_d) z^{(k)} \\ x^{(k)} &= \nabla \phi^*(z^{(k)}) \\ u^{(k)} &= \nabla \mathbf{f}(x^{(k)}). \end{aligned} \quad (20)$$

To represent (20) in state-space form, we can write

$$\begin{aligned} \begin{bmatrix} z^{(k+1)} \\ y^{(k+1)} \end{bmatrix} &= \begin{bmatrix} \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d & -\eta_1 I_{nd} \\ (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes I_d & I_{nd} \end{bmatrix} \begin{bmatrix} z^{(k)} \\ y^{(k)} \end{bmatrix} \\ &+ \begin{bmatrix} 0 & -\eta_1 I_{nd} & I_{nd} \\ 0 & 0 & -I_{nd} \end{bmatrix} \begin{bmatrix} x^{(k)} \\ u^{(k)} \\ v^{(k)} \end{bmatrix}. \end{aligned} \quad (21)$$

Additionally, we know the following constraints on the updates,

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d \\ 0 & 0 \end{bmatrix} \begin{bmatrix} z^{(k)} \\ y^{(k)} \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d \end{bmatrix} \begin{bmatrix} x^{(k)} \\ u^{(k)} \\ v^{(k)} \end{bmatrix}. \quad (22)$$

We define the state vector $\xi^{(k)\top} \triangleq [z^{(k)\top} \quad y^{(k)\top}]$ as well as the input vector $\zeta^{(k)\top} \triangleq [x^{(k)\top} \quad u^{(k)\top} \quad v^{(k)\top}]$. We can rewrite (21) and (22) as

$$\xi^{(k+1)} = A\xi^{(k)} + B\zeta^{(k)} \quad 0 = F\xi^{(k)} + G\zeta^{(k)}, \quad (23)$$

where A, B, F, G are of appropriate dimensions. For the ease of notation we denote $H \triangleq [F \quad G]$.

For the purpose of convergence analysis, we characterize the fixed point of (21). Define $x^* \triangleq \mathbf{1}_n \otimes x_*$, where $x_* \in \mathbb{R}^d$ is a minimizer of (1), and let $z^* \triangleq \nabla \phi(x^*)$, $u^* \triangleq \nabla \mathbf{f}(x^*)$, $y^* \triangleq -\nabla \mathbf{f}(x^*)$ and $v^* = 0$. By letting $z^{(k)}, y^{(k)}, v^{(k)}, x^{(k)}, u^{(k)}$ in (21) take the values of z^*, y^*, v^*, x^*, u^* , it is easy to show that $z^{(k+1)} = z^{(k)}, y^{(k+1)} = y^{(k)}$ using Assumption 2.

B. Exponential Convergence of Distributed Mirror Descent

In the following theorem, we present the main result of this section. We provide two LMIs to characterize the convergence rate of distributed MD. The LMIs are written in terms of several decision variables, including the step-size η_1 and the convergence rate ρ . If we can find a feasible solution for these LMIs, the distributed MD is guaranteed to converge exponentially fast.

Before stating the theorem, we state the following lemma, which will allow us to simplify the resulting SDP.

Lemma 7 (Lemma 6 in [32]). *Suppose that square matrices J_1, J_2 satisfy $J_1^2 = J_1, J_2^2 = J_2, J_1 J_2 = J_2 J_1 = 0$. For square matrices Q_1 and Q_2 , define $Q \triangleq Q_1 \otimes J_1 + Q_2 \otimes J_2$. Then, the following are equivalent: 1) $Q \succeq 0$. 2) $Q_1 \succeq 0, Q_2 \succeq 0$.*

Theorem 8. Let Assumptions 1 and 2 hold and assume all local functions f_i are μ_f -strongly convex and L_f -smooth. Define the following matrices,

$$\begin{aligned} A_1 &= \begin{bmatrix} 0 & -\eta_1 \\ 1 & 1 \end{bmatrix}, B_1 = \begin{bmatrix} 0 & -\eta_1 & 1 \\ 0 & 0 & -1 \end{bmatrix}, \\ A_2 &= \begin{bmatrix} 1 & -\eta_1 \\ 0 & 1 \end{bmatrix}, B_2 = \begin{bmatrix} 0 & -\eta_1 & 1 \\ 0 & 0 & -1 \end{bmatrix}, \\ H_1 &= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, H_2 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}. \end{aligned}$$

Furthermore, define

$$\begin{aligned} M_f &= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{-\mu_f L_f}{\mu_f + L_f} & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & \frac{-1}{\mu_f + L_f} & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ M_\lambda &= \begin{bmatrix} \lambda^2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 \end{bmatrix} \\ M_\phi &= \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}. \end{aligned}$$

If there exists $\rho \in (0, 1)$, $\eta_1 \geq 0$, $P \in \mathbb{S}^2$, $P \succ 0$, $\Sigma_{eq} \in \mathbb{S}^2$, $\sigma_f \geq 0$, $\sigma_\phi \geq 0$, $\sigma_\lambda \geq 0$, such that the following matrix inequalities hold for $i = 1, 2$

$$\begin{aligned} &\begin{bmatrix} A_i^\top P A_i - \rho P & A_i^\top P B_i \\ B_i^\top P A_i & B_i^\top P B_i \end{bmatrix} + \sigma_f M_f + \sigma_\lambda M_\lambda \\ &+ \sigma_\phi M_\phi + H_i^\top \Sigma_{eq} H_i \preceq 0, \end{aligned} \quad (24)$$

then the distributed MD algorithm (18) initialized at $y^{(0)} = 0$ converges exponentially with a rate of ρ as follows

$$\|\xi^{(k)} - \xi^*\|_{P \otimes I_{nd}}^2 \leq \rho^k \|\xi^{(0)} - \xi^*\|_{P \otimes I_{nd}}^2.$$

Proof. Define the vector $e^{(k)\top} = [\xi^{(k)\top} \quad \zeta^{(k)\top}]$. We can establish the following (in)equalities,

$$\begin{aligned} e^{(k)\top} (M_f \otimes I_{nd}) e^{(k)} &\geq 0, \\ e^{(k)\top} (M_\phi \otimes I_{nd}) e^{(k)} &\geq 0, \\ e^{(k)\top} (M_\lambda \otimes I_{nd}) e^{(k)} &\geq 0, \\ e^{(k)\top} H^\top (\Sigma_{eq} \otimes I_{nd}) H e^{(k)} &= 0. \end{aligned}$$

The first two inequalities are derived from Proposition 1, the third inequality is due to the fact that $\lambda = \|\Delta W\|$, and the equality follows from the affine constraint in (22).

Define the Lyapunov function

$$V^{(k)} = \rho^{-k} (\xi^{(k)} - \xi^*)^\top P' (\xi^{(k)} - \xi^*),$$

where $P' \triangleq P \otimes I_{nd}$. Then, using (23) we can write

$$V^{(k+1)} - V^{(k)} = \rho^{-k-1} e^{(k)\top} \begin{bmatrix} A^\top P' A - \rho P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} e^{(k)}.$$

Now, if the following LMI holds

$$\begin{aligned} &\begin{bmatrix} A^\top P' A - \rho P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} + \sigma_f M_f \otimes I_{nd} \\ &+ \sigma_\lambda M_\lambda \otimes I_{nd} + \sigma_\phi M_\phi \otimes I_{nd} + H^\top (\Sigma_{eq} \otimes I_{nd}) H \preceq 0, \end{aligned} \quad (25)$$

then for any $e^{(k)}$, we have that

$$\rho^{-k-1} e^{(k)\top} \begin{bmatrix} A^\top P' A - \rho P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} e^{(k)} \leq 0,$$

or, equivalently,

$$(\xi^{(k)} - \xi^*)^\top P' (\xi^{(k)} - \xi^*) \leq \rho^k (\xi^{(0)} - \xi^*)^\top P' (\xi^{(0)} - \xi^*).$$

In words, the squared norm of system variables decreases exponentially fast to zero.

Next, we simplify the LMI such that the dimension is not dependent on the agent number n . Our approach follows that of [32]. Define J_1, J_2 in Lemma 7 as $J_1 = (I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes I_d$, $J_2 = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d$. It is easy to verify that these matrices satisfy the constraints in Lemma 7. We then have that

$$\begin{aligned} A &= A_1 \otimes J_1 + A_2 \otimes J_2, \\ B &= B_1 \otimes J_1 + B_2 \otimes J_2, \\ H &= H_1 \otimes J_1 + H_2 \otimes J_2. \end{aligned}$$

Since matrices J_1, J_2 satisfy the conditions in Lemma 7, if we consider Q, Q_1, Q_2 as the negative left hand side of (25), (24) respectively, then a feasible set of solutions that satisfy (24) is equivalently a feasible set of solutions for (25), which completes our proof. \square

The theorem provides two LMIs that establish the exponential convergence rate of distributed MD. As we can see the LMIs are more involved compared to the centralized case, and it is challenging to find even a suboptimal analytical rate.

We finally remark that common analysis on distributed MD involves general primal-dual norms [11], whereas QCs are defined with respect to the Euclidean norm. The use of general primal-dual norms in non-strongly convex problems helps with improving the rate up to a multiplicative factor of \sqrt{d} . However, in strongly convex case the rate is exponentially fast, and a more general analysis can only change the iteration complexity by at most logarithmic factors of d , which is an interesting avenue to investigate in the future.

C. $O(1/k)$ Convergence for Convex Functions

In the following theorem, we present the counterpart of Theorem 8 for convex problems.

Theorem 9. Let Assumptions 1 and 2 hold and assume all local functions f_i are convex ($\mu_f = 0$) and L_f -smooth. Recall the definitions of matrices $A_1, A_2, B_1, B_2, H_1, H_2, M_f, M_\lambda, M_\phi$ in Theorem 8 and define the following additional matrices,

$$M_1 = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & L_f & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad M_2 = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

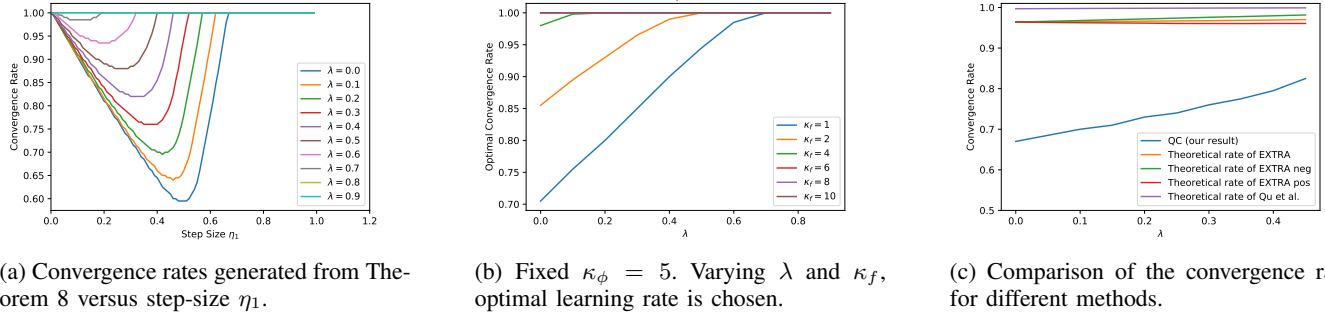


Fig. 1: Optimal convergence rate for Distributed MD obtained by solving LMIs under various assumptions.

If there exist $\eta_1 \geq 0, P \in \mathbb{S}^2, P \succ 0, \Sigma_{eq} \in \mathbb{S}^2, \sigma_f \geq 0, \sigma_\phi \geq 0, \sigma_\lambda \geq 0, \epsilon \geq 0$, such that the following matrix inequalities hold for $i = 1, 2$

$$\begin{bmatrix} A_i^\top P A_i - P & A_i^\top P B_i \\ B_i^\top P A_i & B_i^\top P B_i \end{bmatrix} + \sigma_f M_f + \sigma_\lambda M_\lambda + \sigma_\phi M_\phi + \epsilon M_i + H_i^\top \Sigma_{eq} H_i \preceq 0, \quad (26)$$

then, the iterates of the distributed MD algorithm (18) initialized at $y^{(0)} = 0$ satisfy the following inequality,

$$\sum_{i=1}^n \left(f(\bar{x}_i^{(K)}) - f^* \right) \leq \frac{V^{(0)}}{\epsilon K},$$

where $\bar{x}_i^{(K)} \triangleq \frac{1}{K} \sum_{k=0}^{K-1} x_i^{(k)}$.

We refer to the appendix of [38] for the proof of this theorem. Given that $f(\bar{x}_i^{(K)}) - f^*$ is non-negative, it is easy to see that the function evaluated at the ergodic average of each agent iterate converges to minimum with a rate of $O(1/K)$.

D. Evaluating the Tightness of Results

For the distributed MD algorithm, we provide numerical results based on Theorem 8. First, we demonstrate the influence of the network structure, and then we compare the rate recovered by Theorem 8 to existing theoretical rates on distributed GD when it achieves exponential convergence.

1) *Impact of the Network Structure on Convergence Rate:* We calculate the worst-case convergence rate with several choices of λ and plot it with respect to the step-size η_1 . We set the local functions to have condition number $\kappa_f = 2$ and the DGF to have condition number $\kappa_\phi = 2$. Each curve in the plot represents a certain λ and is obtained by scanning feasible values for the decision variables in the LMIs (24). From Fig. 1a, we can see that there exists an optimal step-size to obtain the best convergence rate, and that as λ increases, the best rate becomes worse. Hence, for any given network structure and its corresponding Laplacian matrix, we should select η_2 such that λ is minimized. This is consistent with results on distributed optimization, where having a larger λ deteriorates the performance.

In Fig. 1b, we keep $\kappa_\phi = 5$ constant and study the optimal convergence rate for different λ and κ_f . When the condition number increases, the optimal rate worsens. This behavior aligns with gradient descent, where $\kappa_\phi = 1$.

2) *Comparison with Distributed Gradient Descent:* To the best of our knowledge, there is currently no work that provides an exponential convergence rate for distributed MD algorithm. Hence, we select two previous works on distributed GD, namely [14] and [15], and compare our performance with the theoretical rates provided in these works. In order to provide a fair comparison, we must set $\kappa_\phi = 1$ to ensure that MD reduces to GD. We also set the local functions to have condition number $\kappa_f = 3$.

Of the two related works above, EXTRA [14] is of particular relevance to our algorithm. If the matrix \tilde{W} in EXTRA is set to be $\frac{I_n + W}{2}$, the EXTRA algorithm coincides with our algorithm with the exception of having a coefficient difference of $\frac{1}{2}$ for the tracking term. Note that the theoretical convergence rate of EXTRA relies on the spectral norm of ΔW as well as the *smallest non-zero eigenvalue* λ_n of W . We plot the convergence rate of EXTRA under three different scenarios:

- 1) $\lambda_n = \lambda$, (EXTRA pos)
- 2) $\lambda_n = -\lambda$, (EXTRA neg)
- 3) $\lambda_n \approx 0$, (EXTRA)

From Fig. 1c, we can see that when λ is small, the rate recovered by Theorem 8 significantly outperforms EXTRA. As λ increases, the convergence rate calculated for our method starts increasing. We also include the theoretical convergence results from Qu et al. [15], which is consistently outperformed by EXTRA.

Note that the point of this plot is not to declare a winner among algorithms. The goal is to show that the richness of the Lyapunov function and QC analysis provides a machinery to obtain better convergence rates, especially compared to the rates that are algorithm specific. In this case, our algorithm can coincide with EXTRA, but still our analysis provides better rates. Our observation is in line with empirical results of [32].

V. CONCLUSION

In this paper, we proposed a SDP framework to characterize the convergence rate of the mirror descent algorithm for both centralized and distributed settings, and empirical evaluations were performed under the assumption of strongly convex and smooth local objective functions. For the centralized case, we derived a closed-form feasible solution to the SDP for the convergence rate, which depends on the condition number of the distance generating function. For the decentralized case,

we numerically derived the convergence rates using SDP. These SDPs do not scale with the ambient dimension and the network size. Using the QC framework, we further proved the $O(1/k)$ convergence rate for centralized and distributed MD in the convex and smooth setting. It would be interesting to derive analytical rates for the distributed case. Another important direction is the analysis of the mirror descent algorithm with primal-dual norms. This is a challenging problem as current SDP approaches rely on the Euclidean norm and they do not lend themselves to general primal-dual norms.

REFERENCES

- [1] 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.
- [2] A. S. Nemirovsky and D. B. Yudin, "Problem complexity and method efficiency in optimization." 1983.
- [3] A. Beck and M. Teboulle, "Mirror descent and nonlinear projected subgradient methods for convex optimization," *Operations Research Letters*, vol. 31, no. 3, pp. 167–175, 2003.
- [4] A. Ben-Tal, T. Margalit, and A. Nemirovski, "The ordered subsets mirror descent optimization method with applications to tomography," *SIAM Journal on Optimization*, vol. 12, no. 1, pp. 79–108, 2001.
- [5] A. Radhakrishnan, M. Belkin, and C. Uhler, "Linear convergence and implicit regularization of generalized mirror descent with time-dependent mirrors," *arXiv preprint arXiv:2009.08574*, 2020.
- [6] Y. Sun and S. Shahrampour, "Linear convergence of distributed mirror descent with integral feedback for strongly convex problems," *IEEE Conference on Decision and Control (CDC)*, 2021.
- [7] A. Nedic, A. Ozdaglar, and P. A. Parrilo, "Constrained consensus and optimization in multi-agent networks," *IEEE Transactions on Automatic Control*, vol. 55, no. 4, pp. 922–938, 2010.
- [8] P. Lin, W. Ren, and Y. Song, "Distributed multi-agent optimization subject to nonidentical constraints and communication delays," *Automatica*, vol. 65, pp. 120–131, 2016.
- [9] D. Yuan, Y. Hong, D. W. Ho, and G. Jiang, "Optimal distributed stochastic mirror descent for strongly convex optimization," *Automatica*, vol. 90, pp. 196–203, 2018.
- [10] M. Rabbat, "Multi-agent mirror descent for decentralized stochastic optimization," in *IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, 2015, pp. 517–520.
- [11] S. Shahrampour and A. Jadbabaie, "Distributed online optimization in dynamic environments using mirror descent," *IEEE Transactions on Automatic Control*, vol. 63, no. 3, pp. 714–725, 2018.
- [12] D. Yuan, Y. Hong, D. W. C. Ho, and S. Xu, "Distributed mirror descent for online composite optimization," *IEEE Transactions on Automatic Control*, pp. 1–1, 2020.
- [13] T. T. Doan, S. Bose, D. H. Nguyen, and C. L. Beck, "Convergence of the iterates in mirror descent methods," *IEEE Control Systems Letters*, vol. 3, no. 1, pp. 114–119, 2019.
- [14] 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.
- [15] 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.
- [16] Y. Sun, A. Daneshmand, and G. Scutari, "Convergence rate of distributed optimization algorithms based on gradient tracking," *arXiv preprint arXiv:1905.02637*, 2019.
- [17] S. Pu and A. Nedić, "Distributed stochastic gradient tracking methods," *Mathematical Programming*, pp. 1–49, 2020.
- [18] B. Gharesifard and J. Cortés, "Distributed continuous-time convex optimization on weight-balanced digraphs," *IEEE Transactions on Automatic Control*, vol. 59, no. 3, pp. 781–786, 2013.
- [19] X. Zeng, P. Yi, and Y. Hong, "Distributed continuous-time algorithm for constrained convex optimizations via nonsmooth analysis approach," *IEEE Transactions on Automatic Control*, vol. 62, no. 10, pp. 5227–5233, 2017.
- [20] S. S. Kia, J. Cortés, and S. Martínez, "Distributed convex optimization via continuous-time coordination algorithms with discrete-time communication," *Automatica*, vol. 55, pp. 254–264, 2015.
- [21] S. Yang, Q. Liu, and J. Wang, "A multi-agent system with a proportional-integral protocol for distributed constrained optimization," *IEEE Transactions on Automatic Control*, vol. 62, no. 7, pp. 3461–3467, 2016.
- [22] Y. Sun and S. Shahrampour, "Distributed mirror descent with integral feedback: Asymptotic convergence analysis of continuous-time dynamics," *IEEE Control Systems Letters*, vol. 5, no. 5, pp. 1507–1512, 2020.
- [23] Y. Yu and B. Açikmeşe, "RLC circuits-based distributed mirror descent method," *IEEE Control Systems Letters*, vol. 4, no. 3, pp. 548–553, 2020.
- [24] L. Lessard, B. Recht, and A. Packard, "Analysis and design of optimization algorithms via integral quadratic constraints," *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 57–95, 2016.
- [25] M. Fazlyab, A. Ribeiro, M. Morari, and V. M. Preciado, "Analysis of optimization algorithms via integral quadratic constraints: Nonstrongly convex problems," *SIAM Journal on Optimization*, vol. 28, no. 3, pp. 2654–2689, 2018.
- [26] B. Hu and L. Lessard, "Dissipativity theory for nesterov's accelerated method," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1549–1557.
- [27] A. B. Taylor, J. M. Hendrickx, and F. Glineur, "Smooth strongly convex interpolation and exact worst-case performance of first-order methods," *Mathematical Programming*, vol. 161, no. 1-2, pp. 307–345, 2017.
- [28] N. K. Dhirga, S. Z. Khong, and M. R. Jovanović, "The proximal augmented lagrangian method for nonsmooth composite optimization," *IEEE Transactions on Automatic Control*, vol. 64, no. 7, pp. 2861–2868, 2018.
- [29] N. S. Aybat, A. Fallah, M. Gurbuzbalaban, and A. Ozdaglar, "Robust accelerated gradient methods for smooth strongly convex functions," *SIAM Journal on Optimization*, vol. 30, no. 1, pp. 717–751, 2020.
- [30] C. Scherer and C. Ebenbauer, "Convex synthesis of accelerated gradient algorithms," *arXiv preprint arXiv:2102.06520*, 2021.
- [31] A. Megretski and A. Rantzer, "System analysis via integral quadratic constraints," *IEEE Transactions on Automatic Control*, vol. 42, no. 6, pp. 819–830, 1997.
- [32] A. Sundararajan, B. Hu, and L. Lessard, "Robust convergence analysis of distributed optimization algorithms," in *2017 55th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2017, pp. 1206–1212.
- [33] A. Sundararajan, B. Van Scoy, and L. Lessard, "Analysis and design of first-order distributed optimization algorithms over time-varying graphs," *IEEE Transactions on Control of Network Systems*, vol. 7, no. 4, pp. 1597–1608, 2020.
- [34] Y. Nesterov *et al.*, *Lectures on convex optimization*. Springer, 2018, vol. 137.
- [35] J.-B. Hiriart-Urruty and C. Lemaréchal, *Fundamentals of convex analysis*. Springer Science & Business Media, 2012.
- [36] D. Lashkari and P. Golland, "Convex clustering with exemplar-based models," *Advances in Neural Information Processing Systems*, vol. 20, 2007.
- [37] A. Benfenati, E. Chouzenoux, and J.-C. Pesquet, "Proximal approaches for matrix optimization problems: Application to robust precision matrix estimation," *Signal Processing*, vol. 169, p. 107417, 2020.
- [38] Y. Sun, M. Fazlyab, and S. Shahrampour, "On centralized and distributed mirror descent: Exponential convergence analysis using quadratic constraints," *arXiv preprint arXiv:2105.14385*, 2021.
- [39] E. Hazan, "Introduction to online convex optimization," *Foundations and Trends in Optimization*, vol. 2, no. 3-4, pp. 157–325, 2016.

APPENDIX

A. Preliminary Lemmas for Proof of Theorems

In this section we provide a few lemmas used in the proof of main theorems later.

Lemma 10. *Let Assumption 1 hold and consider the Lyapunov function $V^{(k)} = \rho^{-k} \mathcal{D}_{\phi^*}(z^{(k)}, z^*)$. Then, the following inequality,*

$$V^{(k+1)} - V^{(k)} \leq \rho^{-k-1} e^{(k)\top} M_{sc} e^{(k)},$$

is satisfied, where M_{sc} is given in Theorem 2, and $e^{(k)}$ is defined in (8).

Proof. From the definition of Lyapunov function and Bregman divergence, we have that □

$$\begin{aligned}
& V^{(k+1)} - V^{(k)} \\
&= \rho^{-k-1} \mathcal{D}_{\phi^*}(z^{(k+1)}, z^*) - \rho^{-k} \mathcal{D}_{\phi^*}(z^{(k)}, z^*) \\
&= \rho^{-k-1} (\phi^*(z^{(k+1)}) - \phi^*(z^*) - \langle \nabla \phi^*(z^*), z^{(k+1)} - z^* \rangle) - \\
&\quad \rho^{-k} (\phi^*(z^{(k)}) - \phi^*(z^*) - \langle \nabla \phi^*(z^*), z^{(k)} - z^* \rangle) \\
&= \rho^{-k-1} \phi^*(z^{(k+1)}) - \rho^{-k-1} \langle x^*, z^{(k)} - z^* - \eta u^{(k)} \rangle \\
&\quad + \rho^{-k} \langle x^*, z^{(k)} - z^* \rangle - (\rho^{-k-1} - \rho^{-k}) \phi^*(z^*) - \rho^{-k} \phi^*(z^{(k)}).
\end{aligned}$$

Since ϕ^* is μ_ϕ^{-1} -smooth, we get

$$\begin{aligned}
& V^{(k+1)} - V^{(k)} \\
&\leq \rho^{-k-1} [\phi^*(z^{(k)}) + \langle x^{(k)}, -\eta u^{(k)} \rangle + \frac{\eta^2}{2\mu_\phi} \|u^{(k)}\|^2] \\
&\quad - \rho^{-k-1} \langle x^*, z^{(k)} - z^* - \eta u^{(k)} \rangle + \rho^{-k} \langle x^*, z^{(k)} - z^* \rangle \\
&\quad - (\rho^{-k-1} - \rho^{-k}) \phi^*(z^*) - \rho^{-k} \phi^*(z^{(k)}) \\
&= (\rho^{-k-1} - \rho^{-k}) (\phi^*(z^{(k)}) - \phi^*(z^*)) + \frac{\rho^{-k-1} \eta^2}{2\mu_\phi} \|u^{(k)}\|^2 \\
&\quad - (\rho^{-k-1} - \rho^{-k}) \langle x^*, z^{(k)} - z^* \rangle - \rho^{-k-1} \langle x^{(k)} - x^*, \eta u^{(k)} \rangle.
\end{aligned}$$

Applying smoothness again, we can bound $V^{(k+1)} - V^{(k)}$ by

$$\begin{aligned}
& (\rho^{-k-1} - \rho^{-k}) (\nabla \phi^*(z^*)^\top (z^{(k)} - z^*) + \frac{1}{2\mu_\phi} \|z^{(k)} - z^*\|^2) \\
&\quad + \rho^{-k-1} \langle x^{(k)} - x^*, -\eta u^{(k)} \rangle + \frac{\rho^{-k-1} \eta^2}{2\mu_\phi} \|u^{(k)}\|^2 \\
&\quad - (\rho^{-k-1} - \rho^{-k}) \langle x^*, z^{(k)} - z^* \rangle \\
&= \rho^{-k-1} (\frac{1-\rho}{2\mu_\phi} \|z^{(k)} - z^*\|^2 - \eta \langle x^{(k)} - x^*, u^{(k)} \rangle) + \frac{\eta^2}{2\mu_\phi} \|u^{(k)}\|^2 \\
&= \rho^{-k-1} e^{(k)\top} M_{sc} e^{(k)},
\end{aligned}$$

and observing $u^* = 0$ finishes the proof. □

Lemma 11. *Let Assumption 1 hold and consider the Lyapunov function $V^{(k)} = \epsilon \sum_{i=0}^{k-1} (f(x^{(i)}) - f(x^*)) + \mathcal{D}_{\phi^*}(z^{(k)}, z^*)$, defined for $\epsilon > 0$. Then, when f is convex, the following inequality holds*

$$V^{(k+1)} - V^{(k)} \leq e^{(k)\top} M_c e^{(k)},$$

where M_c is given in Theorem 6, and $e^{(k)}$ is defined in (8).

Proof. Following the proof of Lemma 10 and by setting $\rho = 1$, we know that

$$\mathcal{D}_{\phi^*}(z^{(k+1)}, z^*) - \mathcal{D}_{\phi^*}(z^{(k)}, z^*) \leq -\eta \langle x^{(k)} - x^*, u^{(k)} \rangle + \frac{\eta^2}{2\mu_\phi} \|u^{(k)}\|^2.$$

Therefore, we can bound $V^{(k+1)} - V^{(k)}$ using the convexity of f and observing $u^* = 0$, as follows

$$\begin{aligned}
& -\eta \langle x^{(k)} - x^*, u^{(k)} \rangle + \frac{\eta^2}{2\mu_\phi} \|u^{(k)}\|^2 + \epsilon (f(x^{(k)}) - f(x^*)) \\
&\leq -\eta \langle x^{(k)} - x^*, u^{(k)} \rangle + \frac{\eta^2}{2\mu_\phi} \|u^{(k)}\|^2 + \epsilon \langle u^{(k)} - u^*, x^{(k)} - x^* \rangle \\
&= e^{(k)\top} M_c e^{(k)}.
\end{aligned}$$

Lemma 12. *Assume all local functions f_i are convex ($\mu_f = 0$) and L_f -smooth. Then, the following inequality holds for the distributed algorithm in (20)*

$$\sum_{i=1}^n (f(x_i^{(k)}) - f^*) \leq e^{(k)\top} M e^{(k)},$$

where $f^* \triangleq \min_x f(x)$ and $M \in \mathbb{R}^{5nd \times 5nd}$ is defined as

$$M \triangleq \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & L_f(I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes I_d & \frac{1}{2n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d & 0 \\ 0 & 0 & \frac{1}{2n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Proof. Recall that we denote an optimal solution of the function in (1) as $x_* \in \mathbb{R}^d$. From the definition of f , we know that $\sum_{i=1}^n \nabla f_i(x_*) = 0$. We note that the dimension of x_* differs from that of the stationary point of the distributed system $x^* \in \mathbb{R}^{nd}$. Specifically, we have $x^* \triangleq \mathbf{1}_n \otimes x_*$.

For any $x_j^{(k)}$ at agent j , we have that

$$\begin{aligned}
n(f(x_j^{(k)}) - f^*) &= \sum_{i=1}^n (f_i(x_j^{(k)}) - f_i(x_*)) \\
&= \sum_{i=1}^n (f_i(x_j^{(k)}) - f_i(x_i^{(k)}) + f_i(x_i^{(k)}) - f_i(x_*)) \\
&\leq \sum_{i=1}^n (\nabla f_i(x_i^{(k)})^\top (x_j^{(k)} - x_i^{(k)}) + \frac{L_f}{2} \|x_j^{(k)} - x_i^{(k)}\|^2 + f_i(x_i^{(k)}) - f_i(x_*)) \\
&\leq \sum_{i=1}^n (\nabla f_i(x_i^{(k)})^\top (x_j^{(k)} - x_i^{(k)}) + \frac{L_f}{2} \|x_j^{(k)} - x_i^{(k)}\|^2 + \nabla f_i(x_i^{(k)})^\top (x_i^{(k)} - x_*)) \\
&= \sum_{i=1}^n (\nabla f_i(x_i^{(k)})^\top (x_j^{(k)} - x_i^{(k)} + x_i^{(k)} - x_*) + \frac{L_f}{2} \|x_j^{(k)} - x_i^{(k)}\|^2) \\
&= \sum_{i=1}^n ((\nabla f_i(x_i^{(k)}) - \nabla f_i(x_*))^\top (x_j^{(k)} - x_*) + \frac{L_f}{2} \|x_j^{(k)} - x_i^{(k)}\|^2),
\end{aligned}$$

where the two inequalities are induced by the Lipschitz-smoothness and convexity of f_i , respectively. Since x_* is a global optimal solution, we also have

$$\sum_{i=1}^n (\nabla f_i(x_*)^\top (x_j^{(k)} - x_*)) = (\sum_{i=1}^n \nabla f_i(x_*))^\top (x_j^{(k)} - x_*) = 0.$$

Summing over j , we get

$$\begin{aligned}
& n \sum_{j=1}^n (f(x_j^{(k)}) - f^*) \\
&\leq \sum_{j=1}^n \sum_{i=1}^n (\nabla f_i(x_i^{(k)}) - \nabla f_i(x_*))^\top (x_j^{(k)} - x_*) \\
&\quad + \sum_{j=1}^n \sum_{i=1}^n \frac{L_f}{2} \|x_j^{(k)} - x_i^{(k)}\|^2.
\end{aligned}$$

Writing above in matrix form and dividing by n , we derive

$$\begin{aligned}
& \sum_{j=1}^n (f(x_j^{(k)}) - f^*) \\
&= (u^{(k)} - u^*)^\top (\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes I_d) (x^{(k)} - x^*) \\
&\quad + L_f (x^{(k)} - x^*)^\top \left((I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \otimes I_d \right) (x^{(k)} - x^*) \\
&= e^{(k)\top} M e^{(k)}.
\end{aligned}$$

- We can then remove I inside the block matrix elements in the equation above and apply Lemma 13 to get

$$\begin{aligned}
& \begin{bmatrix} \frac{\rho-1}{2\mu_\phi} & 0 & 0 \\ 0 & \frac{\mu_\phi}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{\mu_f L_f \sigma_f}{\mu_f + L_f} & \frac{-\sigma_f}{2} \\ 0 & \frac{-\sigma_f}{2} & \frac{\sigma_f}{\mu_f + L_f} \end{bmatrix} \\
& + \begin{bmatrix} \frac{\sigma_\phi}{\mu_\phi + L_\phi} & \frac{-\sigma_\phi}{2} & 0 \\ \frac{-\sigma_\phi}{2} & \frac{\mu_\phi L_\phi \sigma_\phi}{\mu_\phi + L_\phi} & 0 \\ 0 & 0 & 0 \end{bmatrix} - \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{-\sigma_f}{\sqrt{2\mu_\phi}} \end{bmatrix} \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{\sigma_f}{\sqrt{2\mu_\phi}} \end{bmatrix}^\top \succeq 0 \\
\Rightarrow & \begin{bmatrix} \frac{\sigma_\phi}{\mu_\phi + L_\phi} + \frac{\rho-1}{2\mu_\phi} & \frac{-\sigma_\phi}{2} & 0 \\ \frac{-\sigma_\phi}{2} & \frac{\mu_\phi L_\phi \sigma_\phi}{\mu_\phi + L_\phi} + \frac{\mu_\phi}{2} + \frac{\mu_f L_f \sigma_f}{\mu_f + L_f} & \frac{-\sigma_f}{2} \\ 0 & \frac{-\sigma_f}{2} & \frac{\sigma_f}{\mu_f + L_f} \end{bmatrix} \\
& - \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{-\sigma_f}{\sqrt{2\mu_\phi}} \end{bmatrix} \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{\sigma_f}{\sqrt{2\mu_\phi}} \end{bmatrix}^\top \succeq 0 \\
\Rightarrow & \begin{bmatrix} \frac{\sigma_\phi}{\mu_\phi + L_\phi} + \frac{\rho-1}{2\mu_\phi} & \frac{-\sigma_\phi}{2} & 0 & 0 \\ \frac{-\sigma_\phi}{2} & \frac{\mu_\phi L_\phi \sigma_\phi}{\mu_\phi + L_\phi} + \frac{\mu_\phi}{2} + \frac{\mu_f L_f \sigma_f}{\mu_f + L_f} & \frac{-\sigma_f}{2} & \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ 0 & \frac{-\sigma_f}{2} & \frac{\sigma_f}{\mu_f + L_f} & \frac{\eta}{\sqrt{2\mu_\phi}} \\ 0 & \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} & \frac{\eta}{\sqrt{2\mu_\phi}} & 1 \end{bmatrix} \succeq 0,
\end{aligned}$$

thereby completing the proof. □

B. Proof of Proposition 3

We start with the following lemma, which helps with turning the non-affine constraint to an affine constraint in the SDP.

Lemma 13. *If matrix $M \in \mathbb{R}^{n \times n}$ can be decomposed as $M = N + SS^\top$, where $S \in \mathbb{R}^{n \times m}$, then a negative semi-definite constraint on M can be equivalently represented by an affine constraint on N and S .*

Proof. Consider the following matrix $M' \in \mathbb{R}^{(n+m) \times (n+m)}$

$$M' = \begin{bmatrix} -N & S \\ S^\top & I_m \end{bmatrix}.$$

By properties of Schur complement, we have that

$$M' \succeq 0 \iff -N - SS^\top \succeq 0 \iff M \preceq 0.$$

Therefore, we can equivalently use $M' \succeq 0$ as the constraint (in lieu of $M \preceq 0$). This constraint is affine with respect to both N and S . □

We now provide the proof for Proposition 3.

Proof. For brevity, in this proof we use $I = I_d$. Given the matrices defined in Theorem 2, we can write the last LMI in (11) as

$$\begin{aligned}
& \begin{bmatrix} \frac{1-\rho}{2\mu_\phi} I & 0 & 0 \\ 0 & 0 & \frac{-\eta}{2} I \\ 0 & \frac{-\eta}{2} I & \frac{\eta}{2\mu_\phi} I \end{bmatrix} + \sigma_f \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{-\mu_f L_f}{\mu_f + L_f} I & \frac{I}{2} \\ 0 & \frac{I}{2} & \frac{-1}{\mu_f + L_f} I \end{bmatrix} \\
& + \sigma_\phi \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} I & \frac{I}{2} & 0 \\ \frac{I}{2} & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} I & 0 \\ 0 & 0 & 0 \end{bmatrix} \succeq 0,
\end{aligned}$$

which implies

$$\begin{aligned}
& \begin{bmatrix} \frac{1-\rho}{2\mu_\phi} I & 0 & 0 \\ 0 & \frac{-\mu_\phi}{2} I & 0 \\ 0 & 0 & 0 \end{bmatrix} + \sigma_f \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{-\mu_f L_f}{\mu_f + L_f} I & \frac{I}{2} \\ 0 & \frac{I}{2} & \frac{-1}{\mu_f + L_f} I \end{bmatrix} \\
& + \sigma_\phi \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} I & \frac{I}{2} & 0 \\ \frac{I}{2} & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} I & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{\eta}{\sqrt{2\mu_\phi}} \end{bmatrix} \begin{bmatrix} 0 \\ \frac{-\sqrt{\mu_\phi}}{\sqrt{2}} \\ \frac{\eta}{\sqrt{2\mu_\phi}} \end{bmatrix}^\top \succeq 0
\end{aligned}$$

C. Proof of Proposition 4

If $\eta = \sigma_f = \frac{2\mu_\phi}{\mu_f + L_f}$, the LMI in (11) becomes

$$\begin{aligned}
& \begin{bmatrix} \frac{(1-\rho)}{2\mu_\phi} I_d & 0 & 0 \\ 0 & 0 & \frac{-\mu_\phi}{\mu_f + L_f} I_d \\ 0 & \frac{-\mu_\phi}{\mu_f + L_f} I_d & \frac{2\mu_\phi}{(\mu_f + L_f)^2} I_d \end{bmatrix} \\
& + \sigma_\phi \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} I_d & \frac{1}{2} I_d & 0 \\ \frac{1}{2} I_d & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} I_d & 0 \\ 0 & 0 & 0 \end{bmatrix} \\
& + \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{2\mu_\phi}{\mu_f + L_f} \frac{-\mu_f L_f}{\mu_f + L_f} I_d & \frac{\mu_\phi}{\mu_f + L_f} I_d \\ 0 & \frac{\mu_\phi}{\mu_f + L_f} I_d & \frac{2\mu_\phi}{\mu_f + L_f} \frac{-1}{\mu_f + L_f} I_d \end{bmatrix} \succeq 0,
\end{aligned}$$

which, after removing I_d , simplifies to

$$\begin{bmatrix} \frac{(1-\rho)}{2\mu_\phi} & 0 & 0 \\ 0 & \frac{-2\mu_\phi \mu_f L_f}{(\mu_f + L_f)^2} & 0 \\ 0 & 0 & 0 \end{bmatrix} + \sigma_\phi \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} & 0 \\ 0 & 0 & 0 \end{bmatrix} \succeq 0,$$

and we get

$$\begin{aligned}
& \begin{bmatrix} \frac{(1-\rho)}{2\mu_\phi} & 0 \\ 0 & \frac{-2\mu_\phi \mu_f L_f}{(\mu_f + L_f)^2} \end{bmatrix} + \sigma_\phi \begin{bmatrix} \frac{-1}{\mu_\phi + L_\phi} & \frac{1}{2} \\ \frac{1}{2} & \frac{-\mu_\phi L_\phi}{\mu_\phi + L_\phi} \end{bmatrix} \succeq 0 \\
\iff & \begin{bmatrix} \frac{(1-\rho)}{2\mu_\phi} - \sigma_\phi \frac{1}{\mu_\phi + L_\phi} & \frac{\sigma_\phi}{2} \\ \frac{\sigma_\phi}{2} & \frac{-2\mu_\phi \mu_f L_f}{(\mu_f + L_f)^2} - \sigma_\phi \frac{\mu_\phi L_\phi}{\mu_\phi + L_\phi} \end{bmatrix} \succeq 0
\end{aligned}$$

This is equivalent to the following constraints on the principal minors of the matrix:

$$\begin{aligned}
1) \quad & -\frac{(1-\rho)}{2\mu_\phi} + \sigma_\phi \frac{1}{\mu_\phi + L_\phi} \geq 0 \\
2) \quad & \frac{2\mu_\phi\mu_f L_f}{(\mu_f + L_f)^2} + \sigma_\phi \frac{\mu_\phi L_\phi}{\mu_\phi + L_\phi} \geq 0 \\
3) \quad & \left(-\frac{(1-\rho)}{2\mu_\phi} + \sigma_\phi \frac{1}{\mu_\phi + L_\phi} \right) \left(\frac{2\mu_\phi\mu_f L_f}{(\mu_f + L_f)^2} \right. \\
& \left. + \sigma_\phi \frac{\mu_\phi L_\phi}{\mu_\phi + L_\phi} \right) - \frac{\sigma_\phi^2}{4} \geq 0
\end{aligned}$$

The last constraint is the most strict of all constraints. Hence, we will focus on the last constraint, where we can alternatively write

$$\rho \geq 1 - \frac{2\sigma_\phi}{1 + \kappa_\phi} + \frac{\sigma_\phi^2}{2} \left(\frac{2\mu_f L_f}{(\mu_f + L_f)^2} + \sigma_\phi \frac{\kappa_\phi}{1 + \kappa_\phi} \right)^{-1}.$$

The right-hand side can be seen as a function of σ_ϕ ; it takes its minimum when derivative of σ_ϕ is zero. We denote the optimal σ_ϕ by σ_ϕ^* . Therefore,

$$\frac{d}{d\sigma_\phi} \left(1 - \frac{2\sigma_\phi}{1 + \kappa_\phi} + \frac{\sigma_\phi^2}{2} \left(\frac{2\mu_f L_f}{(\mu_f + L_f)^2} + \sigma_\phi \frac{\kappa_\phi}{1 + \kappa_\phi} \right)^{-1} \right) = 0$$

The positive solution for the equation above is

$$\sigma_\phi^* = \frac{4\mu_f L_f}{(\mu_f + L_f)^2} \frac{(1 + \kappa_\phi)}{\kappa_\phi(\kappa_\phi - 1)},$$

and the corresponding solution for ρ is

$$\begin{aligned}
\rho_{opt} &= 1 - \frac{2\sigma_\phi^*}{1 + \kappa_\phi} + \frac{\sigma_\phi^{*2}}{2} \left(\frac{2\mu_f L_f}{(\mu_f + L_f)^2} + \sigma_\phi^* \frac{\kappa_\phi}{1 + \kappa_\phi} \right)^{-1} \\
&= 1 - \frac{4\mu_f L_f}{(\mu_f + L_f)^2 \kappa_\phi^2},
\end{aligned}$$

thereby completing the proof. \square

D. Proof of Theorem 6

Proof. We consider the following Lyapunov candidate

$$V^{(k)} = \epsilon \sum_{i=0}^{k-1} (f(x^{(i)}) - f(x^*)) + \mathcal{D}_{\phi^*}(z^{(k)}, z^*).$$

Using Lemma 11, we can calculate an upper bound for the following term

$$V^{(k+1)} - V^{(k)} \leq e^{(k)\top} M_c e^{(k)}. \quad (27)$$

Combined with the two QCs, the above implies that

$$\begin{aligned}
V^{(k+1)} - V^{(k)} &\leq e^{(k)\top} M_c e^{(k)} \\
&\leq e^{(k)\top} M_c e^{(k)} + \sigma_f e^{(k)\top} M_f e^{(k)} + \sigma_\phi e^{(k)\top} M_\phi e^{(k)} \quad (28) \\
&= e^{(k)\top} (M_c + \sigma_f M_f + \sigma_\phi M_\phi) e^{(k)}.
\end{aligned}$$

If the LMI in (15) is feasible, then the Lyapunov function satisfies $V^{(k+1)} \leq V^{(k)}$, which is equivalent to

$$\mathcal{D}_{\phi^*}(z^{(k+1)}, z^*) - \mathcal{D}_{\phi^*}(z^{(k)}, z^*) \leq -\epsilon(f(x^{(k)}) - f(x^*)). \quad (29)$$

Summing up both sides and rearranging terms, we obtain

$$\frac{\sum_{i=1}^K (f(x^{(i)}) - f(x^*))}{K} \leq \frac{\mathcal{D}_{\phi^*}(z^{(0)}, z^*)}{\epsilon K}.$$

The left hand side is again lower bounded by $f(\bar{x}^{(K)}) - f(x^*)$ due to the convexity of f , which completes the proof. \square

E. Proof of Theorem 9

Proof. Recalling Proposition 1, based on the assumptions, we have that

$$\begin{aligned}
e^{(k)\top} (M_f \otimes I_{nd}) e^{(k)} &\geq 0, \\
e^{(k)\top} (M_\phi \otimes I_{nd}) e^{(k)} &\geq 0.
\end{aligned}$$

Note that for the mapping $z \mapsto \Delta W z$, given that $\lambda = \|\Delta W\|$, we can write

$$e^{(k)\top} (M_\lambda \otimes I_{nd}) e^{(k)} \geq 0.$$

Using Lemma 12, we know that

$$\sum_{i=1}^n (f(x_i^{(k)}) - f^*) \leq e^{(k)\top} M e^{(k)},$$

where $M \in \mathbb{R}^{5nd \times 5nd}$ is defined as

$$M \triangleq \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & L_f(I_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top) \otimes I_d & \frac{1}{2n}\mathbf{1}_n\mathbf{1}_n^\top \otimes I_d & 0 \\ 0 & 0 & \frac{1}{2n}\mathbf{1}_n\mathbf{1}_n^\top \otimes I_d & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

Also, from (22), we have the following equality for any $\Sigma_{eq} \in \mathbb{S}^2$,

$$e^{(k)\top} H^\top (\Sigma_{eq} \otimes I_{nd}) H e^{(k)} = 0.$$

Now, let us define the Lyapunov function

$$V^{(k)} = (\xi^{(k)} - \xi^*)^\top P' (\xi^{(k)} - \xi^*),$$

where $P' = P \otimes I_{nd}$. Then, using (23) we can derive

$$V^{(k+1)} - V^{(k)} = e^{(k)\top} \begin{bmatrix} A^\top P' A - P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} e^{(k)}.$$

If the following LMI holds

$$\begin{aligned}
&\begin{bmatrix} A^\top P' A - P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} + H^\top (\Sigma_{eq} \otimes I_{nd}) H \\
&+ \epsilon M + (\sigma_f M_f + \sigma_\lambda M_\lambda + \sigma_\phi M_\phi) \otimes I_{nd} \preceq 0, \quad (30)
\end{aligned}$$

then for any $e^{(k)}$, we have that

$$e^{(k)\top} \left(\begin{bmatrix} A^\top P' A - P' & A^\top P' B \\ B^\top P' A & B^\top P' B \end{bmatrix} + \epsilon M \right) e^{(k)} \leq 0.$$

This inequality implies that

$$V^{(k+1)} - V^{(k)} + \epsilon \sum_{i=1}^n (f(x_i^{(k)}) - f^*) \leq 0,$$

due to Lemma 12. By summing up both sides from $k = 0$ to $K - 1$, applying convexity of f and rearranging, we have

$$\sum_{i=1}^n \left(f(\bar{x}_i^{(K)}) - f^* \right) \leq \frac{V^{(0)}}{\epsilon K},$$

where $\bar{x}_i^{(K)} \triangleq \frac{1}{K} \sum_{k=0}^{K-1} x_i^{(k)}$. Again, the LMI in (30) can be simplified by defining J_1, J_2 in Lemma 7 similar to the proof of Theorem 8, which completes the proof. \square