

Retraction-Free Decentralized Non-convex Optimization with Orthogonal Constraints

Youbang Sun¹, Shixiang Chen², Alfredo Garcia³ and Shahin Shahrampour¹

Abstract—In this paper, we investigate decentralized non-convex optimization with orthogonal constraints. Conventional algorithms for this setting require either manifold retractions or other types of projection to ensure feasibility, both of which involve costly linear algebra operations (e.g., SVD or matrix inversion). On the other hand, infeasible methods are able to provide similar performance with higher computational efficiency. Inspired by this, we propose the first decentralized version of the retraction-free landing algorithm, called Decentralized Retraction-Free Gradient Tracking (DRFGT). We theoretically prove that DRFGT enjoys the ergodic convergence rate of $\mathcal{O}(1/K)$, matching the convergence rate of centralized, retraction-based methods. We further establish that under a local Riemannian PL condition, DRFGT achieves the much faster linear convergence rate. Numerical experiments demonstrate that DRFGT performs on par with the state-of-the-art retraction-based methods with substantially reduced computational overhead.

I. INTRODUCTION

This paper focuses on the decentralized non-convex optimization problem under orthogonality constraints:

$$\begin{aligned} \min_{x_1, \dots, x_n \in \mathbb{R}^{d \times r}} \quad & \left\{ \frac{1}{n} \sum_{i=1}^n f_i(x_i) \right\} \\ \text{s.t.} \quad & x_i \in \text{St}(d, r) \triangleq \{x \mid x^\top x = I_r\}, \quad \forall i \in [n], \\ & x_1 = \dots = x_n, \end{aligned} \quad (1)$$

where $\text{St}(d, r)$ denotes the orthogonal constraint, also known as the Stiefel manifold, and each individual function $f_i(x_i)$ is assumed to be smooth and non-convex. Problem (1) seeks to find the minimum of a global objective function (or network function) $f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x)$ while maintaining consensus across the network, which naturally appears in many machine learning applications, e.g., principal component analysis (PCA) [1], canonical correlation analysis (CCA) [2], decentralized spectral analysis [3], low-rank matrix approximation [4] and dictionary learning [5]. More recently, due to the distinctive properties of orthogonal matrices, orthogonality constraints and regularization methods have been used for deep neural networks [6], [7], providing improvements in model robustness and stability [8] and adaptive fine-tuning of large language models [9].

This work is supported in part by NSF ECCS-2240788 and NSF ECCS-2240789 Awards.

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 Shixiang Chen is with the School of Mathematical Sciences, University of Science and Technology of China (USTC), Hefei, China. email: shxchen@ustc.edu.cn.

³ Alfredo Garcia is with the Department of Industrial & Systems Engineering at Texas A & M University, College Station, TX 77845, USA. email: alfredo.garcia@tamu.edu.

Typically, Problem (1) is numerically solved by extending the classical gradient descent (GD) to the Riemannian framework using manifold optimization algorithms. Instead of the Euclidean gradient, the Riemannian gradient is calculated, and a retraction step follows in the direction of the Riemannian gradient [10]. Many retraction-based algorithms, often referred to as “feasible” methods, exhibit similar iteration complexity rates [11], [12] to the Euclidean problems. In many instances, however, computing retractions becomes prohibitively *expensive* or *even infeasible*, especially in high-dimensional problems. For example, performing retractions on the Stiefel manifold $\text{St}(d, r)$ requires $\mathcal{O}(dr^2)$ algebraic operations. Therefore, when r is large, the computation of retractions becomes the predominant factor in terms of the algorithm runtime. Moreover, retractions typically necessitate SVD or QR operations that are not as efficient and GPU-friendly as matrix multiplications; these operations introduce bottlenecks for large-scale problems.

To address these issues, “retraction-free” or “infeasible” methods have been proposed. These methods only require matrix multiplications and are especially desirable when strict feasibility is not enforced during the optimization process. In retraction-free methods, iterates do not always remain on a manifold \mathcal{M} but gradually converge towards it. A prominent example of retraction-free methods is the landing algorithm [13], analyzed in the centralized setting. Similar to the majority of modern machine learning algorithms, when retraction-free algorithms are applied to large-scale problems, scenarios such as distributed datasets or intractable computational complexity in the centralized setting necessitate distributed and/or decentralized implementation across a set of agents in a network. However, compared to the Euclidean setting, this task is extremely challenging partly due to the non-convexity of the Stiefel manifold in Problem (1).

A. Contributions

In this paper, we consider a decentralized, retraction-free algorithm to solve (1), referred to as the **Decentralized Retraction-Free Gradient Tracking (DRFGT)** algorithm. The algorithm is fully decentralized and only requires agents to communicate with their neighbors to ensure convergence with consensus. We list our main contributions as follows:

- The iterates of retraction-free algorithm must stay close enough to the manifold in safety regions. We characterize the safety constraint for iterates of DRFGT to remain in the neighborhood of the Stiefel manifold for the purpose of ensuring eventual feasibility (Proposition IV.1).
- We prove that for general smooth functions, the iteration complexity of obtaining an ϵ -stationary point for

TABLE I

A COMPARISON OF EXISTING ALGORITHMS: DRCS AND DEEPCA ACHIEVE LINEAR RATE ONLY FOR SPECIAL CASES OF PROBLEM (1) (CONSENSUS AND PCA, RESPECTIVELY). OTHER COMPETITORS SOLVE (1) WITH SUBLINEAR RATES. L : SMOOTHNESS FACTOR, κ : CONDITION NUMBER, n : NETWORK SIZE.

	DRGTA [14]	DRCS [15]	DeEPCA [16]	DESTINY [17]	DRFGT (This paper)	DRFGT (This paper)
Retraction	yes	yes	yes ¹	no	no	no
PL Condition	no	N/A	implied	no	no	yes
Communication	multi-step	multi-step	single-step	single-step ²	single-step ²	single-step ²
Optimal Step size	$\mathcal{O}(1/L)$	N/A	N/A	$\min\{\mathcal{O}(L^{-4}), \mathcal{O}(L^{-2}/n)\}$	$\mathcal{O}(1/L)$	$\mathcal{O}(\kappa^{-1/4}/L)$
Convergence Rate	$\mathcal{O}(1/K)$	linear (consensus only)	linear (PCA only)	$\mathcal{O}(1/K)$	$\mathcal{O}(1/K)$	linear

DRFGT is $\mathcal{O}(1/\epsilon)$ (Theorem IV.7). This convergence result matches the centralized version and is the first convergence result for a decentralized version of the landing algorithm.

- We also establish that under a local Riemannian Polyak-Łojasiewicz (PL) condition on $\text{St}(d, r)$, a much faster local linear convergence can be guaranteed (Theorem IV.11). To the best of the authors' knowledge, this is the *first-ever linear convergence result* of any decentralized Riemannian optimization algorithm (see Table I).
- We provide numerical experiments to verify our theoretical results and compare the efficiency of our algorithm with existing retraction-based algorithms.

II. RELATED LITERATURE

A. Optimization on Manifolds

Optimization on manifolds has gained significant attention recently, both in the theoretical and practical directions. With the introduction of *retraction*, a projection mapping from the manifold tangent space back to the manifold, many Euclidean optimization algorithms have been adapted to the Riemannian setting, including gradient descent [12], [18], quasi-Newton methods [19], and even accelerated gradient methods [20].

When applied to the Stiefel manifold, the aforementioned methods produce iterates that always stay on $\text{St}(d, r)$; however, computing retractions are often a lot more expensive than gradient calculations. Instead, many application-oriented algorithms, especially for deep learning, often rely on adding a regularizer so that the iterates only stay relatively close to $\text{St}(d, r)$. These algorithms are easy to implement due to cheaper computation costs, but they can only approximately solve the problem with sub-optimal solutions.

Recently, as an effort for computational efficiency, many works have studied unconstrained surrogate problems for manifold constraints. Although the resulting algorithms do not enforce strict feasibility at every iteration, they can be efficiently implemented with convergence guarantees. [21] introduced ODCGM, which uses a non-smooth penalty function. Alternatively, Fletcher's penalty method [22] introduced a smooth function, yet solving the problem often requires

second-order gradient information or the help of augmented Lagrangian methods [23]. Another approach is to view the manifold as a functional constraint. [24]–[26] focused on optimization tasks in the Euclidean space that are constrained by functional equality or inequality. Most of these approaches require hierarchical optimizations in order to derive convergence. Recently, some works leveraged the structure of manifolds and used a simplified Fletcher's merit function, including PenC [27], CDF [28], and the landing algorithm [13], which is of great relevance to our work.

B. Decentralized Extensions of Riemannian Optimization

Decentralized optimization have been well-studied in the Euclidean domain. The decentralized GD algorithm [29] and its variations [30], [31] perform a local gradient step with a neighborhood averaging term. When the agents are heterogeneous, these methods require a diminishing step size in order to achieve consensus among agents. Since by using a constant step size these methods are only guaranteed to converge to a neighborhood of the optimal solution, gradient tracking style methods [32]–[35] have been proposed to achieve exact convergence.

Nevertheless, these methods cannot be applied to Problem (1). Due to the non-convexity of $\text{St}(d, r)$, an average of points on the manifold is not guaranteed to be on the manifold. As a result, we need to consider the problem under the scope of Riemannian optimization. Most existing Riemannian studies are designed for specific tasks, such as the PCA problem [16] or the consensus problem [15], [36]. The work by [26] addressed a parallel functional inequality constrained problem, yet the convergence is only asymptotic. Recent work by [14], [17], [37]–[40] addressed decentralized Riemannian optimization. DPRGD proposed by [37] is close to DGD [29] in nature and uses a diminishing step size and a projection operator to ensure the feasibility of each iterate. [14] introduced DRGTA, which uses multi-step consensus and retractions. DRCGD in [40] considered a conjugate gradient approach with asymptotic convergence rates. The recent work of [38] focused on non-smooth composite problems based on projection, which entails the same complexity as retraction. Another approach, presented in [17], introduced an infeasible decentralized method named DESTINY, founded on the exact penalty function proposed for smooth problems on $\text{St}(d, r)$ in [41]. This was recently extended to CDADT [42], a double

¹ While the updates in [16] do not use retraction explicitly, the computation complexity of QR-decomposition is the same as a retraction step.

² Although both algorithms use single-step updates, the step size depends on the network size n .

tracking method that addresses the generalized orthogonal constrained problems. In contrast to [41], we employ a distinct merit function, defined in (8), which serves as a pivotal factor for achieving a faster convergence rate and simpler analysis under various assumptions.

III. PRELIMINARIES

A. Notations

We start with definitions and notations. The Frobenius norm of matrix A is denoted as $\|A\|_F$, the spectral norm of matrix A is denoted as $\|A\|_2$, the skew and symmetric parts of a square matrix A are denoted as $\text{skew}(A) = \frac{1}{2}(A - A^\top)$ and $\text{sym}(A) = \frac{1}{2}(A + A^\top)$, respectively. The inner product of matrices A and B is defined as $\langle A, B \rangle \triangleq \text{Tr}(AB^\top)$. We denote the identity matrix of rank r as I_r and denote a vector of all ones as $\mathbf{1}$. The Kronecker product of matrices A and B is denoted as $A \otimes B$. The spectral radius of matrix A is denoted as $\rho(A)$. We use $[n]$ to denote the set $\{1, 2, \dots, n\}$. Given a set \mathcal{S} , we represent $\mathcal{D}(\mathcal{S}, \delta) \triangleq \{x \mid \text{dist}(\mathcal{S}, x) \leq \delta\}$, where $\text{dist}(\mathcal{S}, x)$ denotes $\min_{y \in \mathcal{S}} \|x - y\|_F$. Proj_{St} represents projection on the Stiefel manifold.

B. Optimization on the Stiefel Manifold

Let us consider the problem $\min_{x \in \text{St}(d, r)} f(x)$, which is an optimization on the Stiefel manifold. Gradient descent on the Stiefel manifold requires adapting the standard GD with respect to the manifold's geometric constraints [18]. In each iteration, instead of taking a step using the Euclidean gradient $\nabla f(x)$, we use the Riemannian gradient $\text{grad}f(x)$. Here, we use the canonical definition of the Riemannian gradient, which maps $\nabla f(x)$ onto the tangent space of $\text{St}(d, r)$ at point x . The Riemannian gradient with respect to the canonical metric [43] for the Stiefel manifold is defined as

$$\text{grad}f(x) = \text{skew}(\nabla f(x)x^\top)x, \quad (2)$$

where we drop a multiplicative factor of 2 for convenience following [44].

The next iterate of Riemannian GD is then calculated by a retraction, which defines how a point moves on the manifold. There are many ways to define a retraction on the Stiefel manifold, such as the traditional exponential retraction [43], the projection retraction [10], QR-based retraction [18] and the Cayley retraction [45]. The convergence of the retraction-based gradient methods has been well-studied in the literature. Although retraction-based methods are relatively well-understood, the computational burden and poor scalability of the retraction operation have motivated the development of novel approaches without the use of retractions.

In this paper, our focus is on the landing algorithm proposed in [13], which does not require retractions and updates each iterate with a *landing field* $\Lambda(x)$ as follows

$$x_{k+1} = x_k - \alpha \Lambda(x_k), \quad (3)$$

where $\alpha > 0$ is the step size, and the landing field is calculated as

$$\begin{aligned} \Lambda(x) &\triangleq \text{grad}f(x) + \lambda \nabla p(x), \\ p(x) &\triangleq \frac{1}{4} \|x^\top x - I_r\|_F^2. \end{aligned} \quad (4)$$

Note that when $x \notin \text{St}(d, r)$, we refer to $\text{grad}f(x)$ as the relative gradient. The direction $-\text{grad}f(x)$ aims at minimizing the function $f(x)$ and the direction $-\nabla p(x)$ works towards the feasibility constraint. In fact, these two components are orthogonal to each other, and $\langle \text{grad}f(x), \nabla p(x) \rangle = 0$ [44]. The orthogonality implies that $\Lambda(x) = 0$ if and only if $\nabla p(x) = 0$ and $\text{grad}f(x) = 0$, which coincides with the stationarity condition of x on the manifold.

Although the iterates of the landing algorithm are not on the manifold, they can stay close to it due to $-\nabla p(x)$. Suppose there exists a uniform upper bound on the Riemannian gradient of f . Then, each iterate of the landing algorithm stays within a relatively close neighborhood of $\text{St}(d, r)$, referred to as the “safety region” and formally defined below.

Definition III.1 (Safety Region [44]). With $\epsilon \in (0, 3/4)$, we define

$$\text{St}(d, r)^\epsilon \triangleq \{x \in \mathbb{R}^{d \times r} \mid \|x^\top x - I_r\|_F \leq \epsilon\}. \quad (5)$$

We use the above definition to address the safety properties of our proposed decentralized algorithm later in Section IV.

C. Technical Assumptions

In the decentralized setting, we consider a network of n agents modeled with a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the agents are represented by nodes $\mathcal{V} = [n]$ and each edge $\{i, j\} \in \mathcal{E}$ denotes a connection between agent i and j . As such, the neighborhood of agent i is defined as $\mathcal{N}_i \triangleq \{j \mid \{i, j\} \in \mathcal{E}\}$. Each agent $i \in [n]$ is associated with the individual objective function $f_i : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}$ in (1). The agents work collectively to find the optimum of the global objective function f . We provide the following assumptions on the network as well as the objective functions.

Assumption III.2 (Network Model). We assume the graph \mathcal{G} is undirected and connected, i.e., there exists a path between any two distinct agents $i, j \in \mathcal{V}$.

Furthermore, we use a symmetric and doubly stochastic matrix $W \in \mathbb{R}^{n \times n}$ to capture the communication among agents. It is easy to show from stochasticity that $W\mathbf{1} = \mathbf{1}$ and that one is an eigenvalue of W . Additionally, the connectivity assumption guarantees that all other eigenvalues of W are strictly less than 1 in magnitude. The second largest singular value of W is denoted as σ_W . In general, $0 \leq \sigma_W < 1$ describes the connectivity of \mathcal{G} , and smaller values of σ_W often imply a better-connected graph. For a fully connected graph with $W = \mathbf{1}\mathbf{1}^\top/n$ we have $\sigma_W = 0$.

We also provide the following assumptions on the objective functions, all of which are considered to be standard in the literature [12].

Assumption III.3 (Lipschitz Smoothness). We assume that all local objective functions $f_i : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}, i \in [n]$ are differentiable and L -smooth on $\mathbb{R}^{d \times r}$, i.e. for all i and $x, y \in \mathbb{R}^{d \times r}$, we have

$$f_i(y) \leq f_i(x) + \langle \nabla f_i(x), y - x \rangle + \frac{L}{2} \|y - x\|_F^2. \quad (6)$$

We also assume $\max_{x \in \text{St}(d, r)^\epsilon} \|\text{grad}f_i(x)\|_F \leq G, \forall i \in [n]$.

The above assumption implies that the global objective f is also L -smooth. We further assume that f is C^2 -smooth.

Assumption III.4 (Local Riemannian PL on $\text{St}(d, r)$). The global objective function $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}$ satisfies the local Riemannian Polyak-Łojasiewicz (PL) condition on the Stiefel manifold with a factor $\mu > 0$ if

$$|f(x) - f_S^*| \leq \frac{1}{2\mu} \|\text{grad} f(x)\|_F^2, \quad (7)$$

for any point $x \in \text{St}(d, r) \cap \mathcal{D}(\mathcal{S}, 2\delta)$, where \mathcal{S} denotes the set of all local minima with a given value f_S^* .

The local PL condition is less restrictive than other commonly studied conditions, such as geodesic strong convexity. Problems satisfying the local PL condition include the PCA problem and generalized quadratic problem [46]. In addition, we note that although for our proposed algorithm iterates do not strictly stay on the manifold (similar to the landing algorithm), we only impose a local Riemannian PL condition on $\text{St}(d, r)$ with no other assumptions on the landing field itself.

D. A Smooth Merit Function and the Optimality Condition

For the convergence analysis, we introduce the following smooth merit function [44], defined with respect to the global objective function:

$$\begin{aligned} \mathcal{L}(x) &= f(x) + h(x) + \gamma p(x), \\ \text{where } h(x) &\triangleq -\frac{1}{2} \langle \text{sym}(x^\top \nabla f(x)), x^\top x - I_r \rangle. \end{aligned} \quad (8)$$

For $\epsilon \in (0, \frac{3}{4})$, γ must satisfy the following relationship,

$$\begin{aligned} \gamma &\geq \frac{2}{3-4\epsilon} \left(L(1-\epsilon) + 3s + \hat{L}^2 \frac{(1+\epsilon)^2}{\lambda(1-\epsilon)} \right), \\ \text{where } s &= \sup_{x \in \text{St}(d, r)^\epsilon} \|\text{sym}(x^\top \nabla f(x))\|_F, \quad \text{and} \\ \hat{L} &= \max(L, \max_{x \in \text{St}(d, r)^\epsilon} \|\nabla f(x)\|_F). \end{aligned} \quad (9)$$

In addition, it can be verified that $\forall x \in \text{St}(d, r)$,

$$\nabla \mathcal{L}(x) = \nabla f(x) - x \text{sym}(x^\top \nabla f(x)). \quad (10)$$

Some existing works [21] have used the non-smooth merit function $\mathcal{L}'(x) = f(x) + \gamma \|x^\top x - I_r\|_F$ to better capture the function value change towards the normal direction in a local neighborhood of $\text{St}(d, r)$. However, we choose to work with the smooth merit function in (8), which introduces the term $h(x)$ to enable a more fine-grained analysis for the orthogonal direction, facilitating the deduction of a new PL inequality.

We now state the following proposition on the analytical properties of the merit function $\mathcal{L}(x)$ and refer to [44] for the detailed proof of these results.

Proposition III.5. *The merit function $\mathcal{L}(x)$ satisfies the following properties.*

- 1) $\mathcal{L}(x)$ is $L_{\mathcal{L}}$ -smooth on $x \in \text{St}(d, r)^\epsilon$, with $L_{\mathcal{L}} \leq L_{f+h} + (2 + 3\epsilon)\gamma$, where L_{f+h} is the smoothness of $f + h$.

- 2) For $\rho = \min\{\frac{1}{2}, \frac{\gamma}{4\lambda(1+\epsilon)}\}$ (γ given in (9)) and $x \in \text{St}(d, r)^\epsilon$, we have

$$\langle \Lambda(x), \nabla \mathcal{L}(x) \rangle \geq \rho \|\Lambda(x)\|_F^2.$$

The two properties of Proposition III.5 suggest that the merit function is indeed Lipschitz smooth, and a landing step can be seen as a descent step on the merit function. Moreover, since $\langle \text{grad} f(x), \nabla p(x) \rangle = 0$, the second property shows that within the neighborhood $\text{St}(d, r)^\epsilon$, if $\nabla \mathcal{L}(x) = 0$, we have $\nabla p(x) = 0$ and $\text{grad} f(x) = 0$.

Here, the landing field is defined as (4) with respect to the global function $f(x)$. Similarly, we can define the local landing field $\Lambda_i(x)$ using the local objective $f_i(x)$, which we later use in the next section.

Additionally, with Assumption III.3, it is easy to show the Lipschitz continuity of $\Lambda(x)$ and $\Lambda_i(x)$, we denote the Lipschitz continuity constant as L_Λ . In an ideal scenario, γ and λ are chosen to be on the same order as L . It can then be seen that $L_\Lambda, L_{\mathcal{L}}$ are also in the same order as L . For the ease of analysis, we define

$$L' \triangleq \max\{\hat{L}, L_{\mathcal{L}}, L_\Lambda\}.$$

IV. MAIN RESULTS

In this section, we first propose a decentralized algorithm that solves Problem (1). We then formulate the update as a dynamical system and study its safety and stability conditions. Next, we establish both global and local convergence results for the algorithm.

A. Decentralized Retraction-Free Gradient Tracking

Existing algorithms such as DRGTA [14] have introduced an auxiliary tracking sequence (similar to the Euclidean case [34]) to estimate the global Riemannian gradient in the manifold setting. However, the DRGTA updates require retractions and complicated consensus computations. The consensus problem on the Stiefel manifold is a highly non-trivial task in itself as shown by [15], and the retractions could be inefficient and computationally expensive.

This paper seeks to develop a *decentralized retraction-free* algorithm, only involving matrix multiplications and working based on single-step consensus, so that its implementation is significantly easier than DRGTA. Our exact update for the **Decentralized Retraction-Free Gradient Tracking** (DRFGT) is provided in Algorithm 1. It is easily verified that the algorithm is fully decentralized, and the communication complexity is no more than the Euclidean counterpart in [34]. We also note that apart from saving on the retraction calculations, DRFGT offers additional benefits compared to [14].

- In Algorithm 1, the agents only require one gradient calculation per iteration, which is more efficient than previous works such as [17].
- Apart from not requiring retraction calculations, assuming availability of the Riemannian gradient via closed-form (2), DRFGT does not need additional projections onto the Riemannian tangent space, which is necessary in [14].

Algorithm 1 Decentralized Retraction-free Gradient Tracking

1: **Input:** initial point $x_0 \in \text{St}(d, r)^\epsilon$, $\alpha > 0$, $\lambda > 0$, $\epsilon \in (0, \frac{3}{4})$.
2: Set $x_{i,0} = x_0$; $y_{i,0} = 0$; $\Lambda_i(x_{i,0}) = 0$ for all agents $i \in [n]$.
3: **for** $k = 0, 1, \dots$ **do**
4: **for** $i \in [n]$ **do**
5: Update x , $\Lambda(x)$ and tracking term y :

$$x_{i,k+1} = \sum_{j \in \mathcal{N}_i} W_{ij} x_{j,k} - \alpha y_{i,k},$$

$$\Lambda_i(x_{i,k+1}) = \text{grad} f_i(x_{i,k+1}) + \lambda x_{i,k+1} (x_{i,k+1}^\top x_{i,k+1} - I_r),$$

$$y_{i,k+1} = \sum_{j \in \mathcal{N}_i} W_{ij} y_{j,k} + \Lambda_i(x_{i,k+1}) - \Lambda_i(x_{i,k}).$$

6: **end for**
7: **end for**

- Also, the consensus component of DRFGT is single-step unlike [14], [15] that work with multi-step consensus. This is a desired property, consistent with the Euclidean gradient tracking algorithms.

As a result, Algorithm 1 is very efficient both in terms of communication and computation, thanks to the unique properties offered by the retraction-free landing update.

B. Linear System Analysis

We now study the proposed algorithm from a dynamical system perspective. For the notation convenience, we define

$$\mathbf{x}_k \triangleq [x_{1,k}^\top, \dots, x_{n,k}^\top]^\top.$$

We also write the stacked version of the communication matrix as $\mathbf{W} \triangleq W \otimes I_d$. In addition, we denote the Euclidean average of the iterate over the network as

$$\bar{x}_k \triangleq \frac{1}{n} \sum_{i \in [n]} x_{i,k}, \quad \bar{\mathbf{x}}_k \triangleq \mathbf{1} \otimes \bar{x}_k.$$

With the above notations, we can write Algorithm 1 in a matrix format as

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{W} \mathbf{x}_k - \alpha \mathbf{y}_k, \\ \mathbf{y}_{k+1} &= \mathbf{W} \mathbf{y}_k + \mathbf{\Lambda}_{k+1} - \mathbf{\Lambda}_k, \end{aligned} \quad (11)$$

where \mathbf{y}_k and $\mathbf{\Lambda}_k$ are stacked matrices, defined similar to \mathbf{x}_k .

For analyzing the dynamics of the distributed system, we need to identify an ϵ -safety region similar to the centralized landing algorithm. In the decentralized setting, we evenly split the margin of safety into two parts: (i) the distance of the average iterate from the manifold, which we define safe as $\bar{x} \in \text{St}(d, r)^{\frac{\epsilon}{2}}$; (ii) the consensus safety, which we define safe when $\|\mathbf{x} - \bar{\mathbf{x}}\|_F \leq \frac{\epsilon}{2}$. Combining the two, we have the following proposition on the system safety over the network.

Proposition IV.1 (Safe Step Size in Networks). *Let $\bar{x}_k \in \text{St}(d, r)^{\frac{\epsilon}{2}}$ and $\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F \leq \frac{\epsilon}{10}$. Given $\|\text{grad} f_i(x)\|_F \leq G$ for $x \in \text{St}(d, r)^\epsilon$ and $\lambda > 0$, if the step size α satisfies*

$$\alpha_{safe} \triangleq \min \left\{ \frac{(1 - \sigma_W)^2 \epsilon}{20\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}, \frac{\lambda\epsilon^2(1 - \sigma_W)^2}{16L'(G + \lambda\epsilon(1 + \epsilon))}, \frac{1}{2\lambda}, \frac{\lambda\epsilon(1 - \epsilon)}{2(G^2 + \lambda^2(1 + \epsilon)\epsilon^2 + \frac{\epsilon^4\lambda^2}{16})} \right\},$$

the next iterate satisfies $\bar{x}_{k+1} \in \text{St}(d, r)^{\frac{\epsilon}{2}}$ and $\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}_{k+1}\|_F \leq \frac{\epsilon}{10}$.

Proposition IV.1 shows that with a sufficiently small step size α , Algorithm 1 will stay within the ϵ -safety region for any network structure and objective function.

With the safety constraint satisfied for Algorithm 1, we next study the stability aspects of the system. Let us first introduce the following lemmas to bound the system consensus errors in both \mathbf{x} and \mathbf{y} .

Lemma IV.2. *Let Assumption III.2 hold with σ_W as the second largest singular value of W . The consensus error on \mathbf{x} satisfies the following inequality for $\alpha > 0$,*

$$\begin{aligned} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 &\leq \frac{1 + \sigma_W^2}{2} \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2. \end{aligned} \quad (12)$$

Lemma IV.3. *Let Assumption III.2 hold with σ_W as the second largest singular value of W . The consensus error on \mathbf{y} satisfies the following inequality for $\alpha > 0$,*

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &\leq \left(\frac{1 + \sigma_W^2}{2} + 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\quad + 8L'^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \|\bar{\mathbf{y}}_{k-1}\|_F^2. \end{aligned} \quad (13)$$

We use the previous two lemmas to characterize an LTI system with a state related to consensus errors. The following theorem identifies the stability condition for such LTI system, allowing us to analyze the convergence of Algorithm 1.

Theorem IV.4 (Stability Conditions). *Let Assumption III.2 hold and the step size $\alpha \in (0, \alpha_{safe}]$ be safe. We can form a linear system with state $\tilde{\xi}$, transition matrix \tilde{G} and input signal \tilde{u} ,*

$$\tilde{\xi}_k \leq \tilde{G} \tilde{\xi}_{k-1} + \tilde{u}_{k-1}, \quad (14)$$

where “ \leq ” here denotes element-wise inequality and

$$\begin{aligned} \tilde{\xi}_k &\triangleq \begin{bmatrix} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 / L' \\ L' \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 \end{bmatrix}, \\ \tilde{G} &\triangleq \begin{bmatrix} \frac{1 + \sigma_W^2}{2} + 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} & 8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \\ \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 L'^2 & \frac{1 + \sigma_W^2}{2} \end{bmatrix}, \\ \tilde{u}_k &\triangleq \begin{bmatrix} 4L' \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \|\bar{\mathbf{y}}_k\|_F^2 \\ 0 \end{bmatrix}. \end{aligned}$$

In addition, the linear system is stable, i.e., $\rho(\tilde{G}) < 1$ if

$$\alpha < \frac{(1 - \sigma_W^2)^2}{1 + \sigma_W^2} \frac{1}{16L'}.$$

Theorem IV.4 shows that with a sufficiently small step size, the linear system describing consensus errors on \mathbf{x} and \mathbf{y} is stable.

To study the convergence of our algorithm, we need to analyze the merit function introduced in Section III-D. Note that even though we focus on the decentralized Problem (1), we still work with the merit function based on the network function $f(x)$. We provide the following lemma on the merit function, which describes the change in the value of the merit function with respect to the average of iterates.

Lemma IV.5. *Let Assumptions III.2 and III.3 hold and the step size $\alpha \in (0, \alpha_{safe}]$ be safe. Then, the merit function (8) satisfies the following inequality*

$$\begin{aligned} \mathcal{L}(\bar{x}_k) - \mathcal{L}(\bar{x}_{k-1}) &\leq -\frac{\alpha\rho}{2} \|\Lambda(\bar{x}_{k-1})\|_F^2 \\ &+ \frac{\alpha C^2 L'^2}{2\rho n} \|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 + \frac{\alpha^2 L'}{2} \|\bar{\mathbf{y}}_{k-1}\|_F^2, \end{aligned} \quad (15)$$

for ρ given in Proposition III.5 and $C \triangleq \frac{3L'}{\lambda(1-\epsilon)} + 2$.

Lemma IV.5 shows that a decrease in the merit function value can be achieved for DRFGT with a suitable step size.

C. Global Convergence of DRFGT

In this section, we discuss the global convergence of DRFGT, building on the safety and stability guarantees discussed in the previous sections. We start by deriving the following bound on the accumulated consensus errors along the algorithm when the system (14) is stable. The following is a corollary of Theorem IV.4.

Corollary IV.6. *Given a stable step size $\alpha < \frac{(1-\sigma_W^2)^2}{1+\sigma_W^2} \frac{1}{16L'}$, the sum of consensus errors satisfies*

$$\sum_k \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 \leq \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2 \sum_k \|\bar{\mathbf{y}}_{k-1}\|_F^2.$$

Using Lemma IV.5 and Corollary IV.6, we establish one of our main results, which provides ergodic convergence results for the landing field as well as the consensus error.

Theorem IV.7 (Global Convergence). *Let Assumptions III.2 and III.3 hold and the step size $\alpha \in (0, \alpha_{safe}]$ be safe. If the step size additionally satisfies*

$$\alpha < \min \left\{ \frac{(1 - \sigma_W^2)^2}{1 + \sigma_W^2} \frac{1}{16L'}, \sqrt[3]{\frac{\rho(1 - \sigma_W^2)^4}{(1 + \sigma_W^2)^2 C^2}} \frac{1}{4L'}, \frac{\rho}{8L'} \right\},$$

we get the following ergodic convergence results

$$\begin{aligned} \frac{\sum_k \|\Lambda(\bar{x}_{k-1})\|_F^2}{K} &\leq \frac{1}{K} \frac{4}{\alpha\rho} (\mathcal{L}(\bar{x}_0) - \mathcal{L}(\bar{x}_K)), \\ \frac{\sum_k \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2}{K} &\leq \frac{1}{K} \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} \frac{512n\alpha^3 L'^2}{\rho} (\mathcal{L}(\bar{x}_0) - \mathcal{L}(\bar{x}_K)). \end{aligned}$$

The decay of the landing field in Theorem IV.7 shows that DRFGT converges to a stationary point of Problem (1) on

the Stiefel manifold with asymptotically zero consensus error. The constraint on the step size depends on multiple factors, but the step size still scales with $\mathcal{O}(1/L')$ (when $\lambda = \mathcal{O}(L')$) across these terms. As a result, the convergence rate matches the retraction-based algorithms such as [14] and the Euclidean gradient tracking in [34].

D. Local Linear Convergence of DRFGT Under Local PL Condition

In this section, we study the local convergence of DRFGT under the assumption that the network function satisfies the local PL condition (Assumption III.4). Given the global convergence result in the previous section, one can guarantee that the algorithm does converge to a stationary point, but if the stationary point is a local minimizer, the local convergence could be much faster. Therefore, for local analysis, we assume the iterates are already close to a local minimizer, or equivalently, we initialize the algorithm close enough to a local minimizer. This assumption is not surprising; even some classical methods (e.g., Newton's method) achieve faster local rates when initialization is close enough to a local minimizer.

We study the convergence of the algorithm using a dynamical system perspective again. Consider the state vector,

$$\xi_k \triangleq [\|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 / L', L' \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2, n(\mathcal{L}(\bar{x}_k) - \mathcal{L}_S^*)]^\top, \quad (16)$$

where \mathcal{L}_S^* is the merit function evaluated at the local minimizer. It is obvious that $\mathcal{L}_S^* = f_S^*$ based on (8). To establish the local convergence, we need to identify the dynamical system describing the state ξ_k . To that end, having Lemma IV.2 already in place, we further derive Lemmas IV.9 and IV.10 to characterize the dynamical system.

First, let us analyze the safety in the sense of the local PL neighborhood and present the following lemma based on the relation between the PL condition and the quadratic growth (QG) condition [47].

Lemma IV.8 (PL-QG). *Let Assumptions III.3 and III.4 hold. For $x \in \text{St}(d, r)^\epsilon \cap \mathcal{D}(\mathcal{S}, \delta)$, the merit function $\mathcal{L}(x)$ satisfies*

$$\mathcal{L}(x) - \mathcal{L}_S^* \geq \frac{\mu' \rho^2}{4} \text{dist}(\mathcal{S}, x)^2.$$

The above lemma is used to characterize an upper bound on the closeness of the iterates from the local minima set, which is used to guarantee that the iterates never leave a neighborhood of the local minima under appropriate initialization. We next present the following lemma to provide a sufficient descent inequality over $\mathcal{L}(\bar{x}_k)$, which characterizes the behavior of the last entry in (16).

Lemma IV.9. *Given Assumptions III.3 and III.4, if $\bar{x}_{k-1} \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^\epsilon$ and the safe step size satisfies $\alpha \leq \frac{\rho}{2L'}$, the merit function evaluated at \bar{x}_k can be upper bounded using the following inequality,*

$$\begin{aligned} \mathcal{L}(\bar{x}_k) - \mathcal{L}_S^* &\leq (1 - \frac{\rho\alpha\mu'}{4}) (\mathcal{L}(\bar{x}_{k-1}) - \mathcal{L}_S^*) \\ &+ (L'^2 \alpha^2 + \frac{\alpha L' C^2}{\rho}) \frac{L'}{n} \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2. \end{aligned}$$

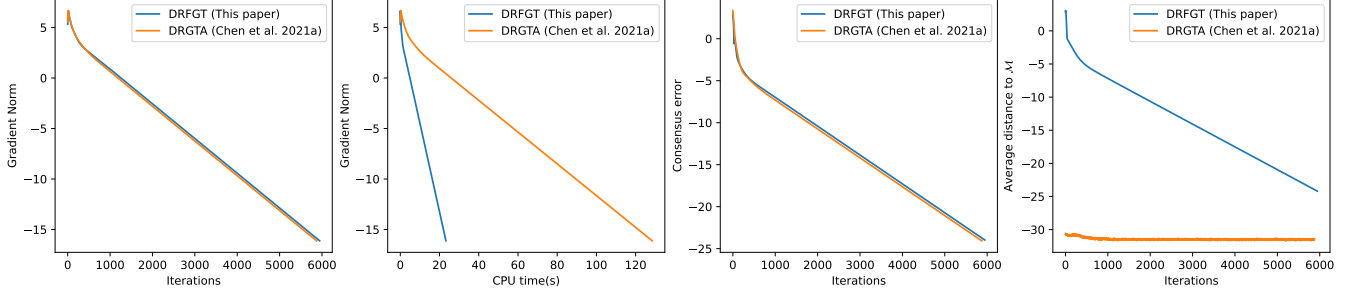


Fig. 1. Convergence of DRFGT (our work) compared to DRGTA in [14].

In addition, since the iterates stay close enough to a local minimizer, we can revisit the result of Lemma IV.3 to get the following relationship.

Lemma IV.10. *Let Assumptions III.2 and III.3 hold, assume $\delta \leq 1$, and select a safe step size $\alpha \in (0, \alpha_{safe}]$. Let also $\bar{x}_{k-1} \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^\epsilon$. Then, the consensus error on \mathbf{y} satisfies the following inequality,*

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &\leq \left(\frac{1 + \sigma_W^2}{2} + 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\quad + \frac{96n\alpha^2 L'^3}{\rho^2} \frac{1 + \sigma_W^2}{1 - \sigma_W^2} (\mathcal{L}(\bar{x}_{k-1}) - \mathcal{L}_S^*) \\ &\quad + 8L'^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} (1 + \alpha^2 L'^2) \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2. \end{aligned}$$

After these lemmas, we present another main result of this paper, which establishes the local linear rate of DRFGT, the first for any decentralized Riemannian algorithm.

Theorem IV.11 (Local Linear Convergence). *Suppose that Assumptions III.2, III.3 and III.4 hold. Additionally, let the step size $\alpha \in (0, \alpha_{safe}]$ be safe and satisfy*

$$\alpha \leq \min \left\{ \frac{\rho}{2L'}, \frac{1 - \sigma_W^2}{\rho\mu'}, \frac{\sqrt{1 - \sigma_W^2}}{4L'\sqrt{\Theta}^4 \sqrt{1 + \frac{12\Phi}{\rho^2}}}, \frac{1 - \sigma_W^2}{16L'\Theta} \right\}, \quad (17)$$

with $\Theta \triangleq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}$ and $\Phi \triangleq \frac{4L'}{\rho\mu'} + \frac{8L'C^2}{\rho^2\mu'}$. Then, the system state vector ξ_k in (16) satisfies $\xi_k \leq M\xi_{k-1}$ element-wise for any iteration k if $\bar{x}_0 \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^{\frac{\delta}{2}}$ and $\xi_0 \leq \frac{n\mu'\rho^2\delta^2}{8}\mathbf{v}$ for $\delta \leq 1$, where M is defined as

$$M \triangleq \begin{bmatrix} \frac{1 + \sigma_W^2}{2} + 4L'^2 \alpha^2 \Theta & 8(1 + \alpha^2 L'^2) \Theta & \frac{96\alpha^2 L'^2 \Theta}{\rho^2} \\ \alpha^2 L'^2 \Theta & \frac{1 + \sigma_W^2}{2} & 0 \\ 0 & \alpha^2 L'^2 + \frac{\alpha L' C^2}{\rho} & 1 - \frac{\alpha \rho \mu'}{4} \end{bmatrix},$$

and \mathbf{v} is the eigenvector corresponding to $\rho(M)$. Furthermore, the system converges with the linear rate $\rho(M) \leq 1 - \frac{\alpha \rho \mu'}{8}$.

We refer to the Appendix for the proof and analysis. Apart from the safety constraints on α , to ensure that the system converges linearly, α also depends on several factors including the PL condition factor μ' and Lipschitz smoothness factor L' . The rate also depends on the graph connectivity, where the convergence becomes slower as σ_W tends to one. The optimal

step size in Euclidean GD is known to be $\alpha = \mathcal{O}(1/L)$. This is mostly true in our analysis as the step size is always $\mathcal{O}(1/L')$; however, if the dominating factor in (17) is the term depending on Φ , we will have α scaling as $\mathcal{O}(\kappa^{-1/4}/L')$, which also recovers the best previously known result on Euclidean distributed non-convex optimization [48]. On the other hand, our safety step size scales as $\alpha_{safe} = \mathcal{O}(1/\sqrt{n}L')$, an improvement compared to DESTINY in [17], where $\alpha_{safe} \leq \min\{\mathcal{O}(L^{-4}), \mathcal{O}(L^{-2}/n)\}$. Our tighter bounds on the step size are partly due to our choice of the landing field and the merit function, between which various relationships were identified (see Lemmas VI.8, VI.9, VI.10).

V. NUMERICAL RESULTS

For the experiments, we evaluate DRFGT on a traditional PCA task with both synthetic and real-world data.

A. PCA Experiment with Synthetic Data

We study the decentralized PCA problem, also evaluated by [14]. The problem is defined as

$$\begin{aligned} \min_{x_1, \dots, x_n \in \mathbb{R}^{d \times r}} & \left\{ \frac{1}{n} \sum_{i \in [n]} \langle A_i^\top A_i x_i, x_i \rangle \right\} \\ \text{s.t. } & x_1 = \dots = x_n, \quad x_i \in \text{St}(d, r), \quad \forall i \in [n], \end{aligned} \quad (18)$$

where A_i is the individual data matrix for agent i , and D is a diagonal matrix with $[D]_{11} > \dots > [D]_{rr} > 0$. The dimension of the A_i matrix is set to $m_i \times d$, where d is the dimension of data and m_i is the number of data samples assigned to agent i . Each entry of A_i is first generated independently from a normal distribution, then the eigenvalues of $A_i^\top A_i$ are manually adjusted to get a larger condition number to increase the problem difficulty. The network has 10 agents, a ring structure, and W has diagonal elements $W_{ii} = 0.8$ and off-diagonal elements $W_{ij} = 0.1$. Each agent i stores $m_i = m = 1000$ total data points with a data dimension of $d = 100$. We set the rank number $r = 10$ and the step size $\alpha = \frac{1}{10m} = 10^{-4}$. Lastly, λ is set to be $0.1/\alpha$.

We compare the performance of Algorithm 1 with DRGTA, the retraction-based algorithm in [14] with the same initial point x_0 , graph structure W , and step size α . We plot the evolution of gradient norm $\|\sum_{i \in [n]} \text{grad} f_i(x_i)\|_F$, consensus error $\sum_{i \in [n]} \|x_i - \bar{x}\|_F$, as well as the average distance

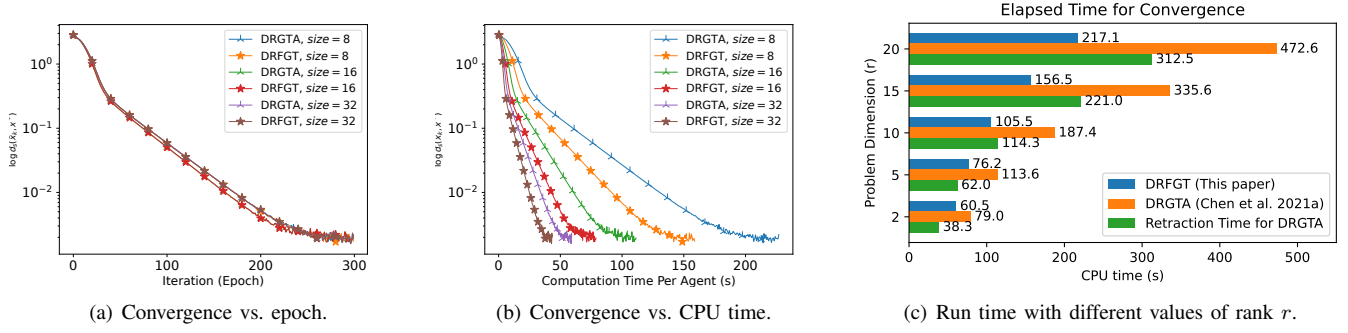


Fig. 2. Experiment with MNIST data.

to $\text{St}(d, r)$, $\frac{1}{n} \sum_{i \in [n]} \|x_i - x'_i\|_F$, with respect to different iterations as well as CPU time. The results are provided on log-scale in Fig. 1: (i) In terms of iterations, the DRFGT algorithm closely matches the performance of DRGTA in gradient norm and consensus error. (ii) The DRFGT algorithm converges much faster than DRGTA in CPU time and is more computationally efficient. (iii) Though DRFGT does not strictly satisfy the feasibility constraints of Problem (1), it eventually converges to a feasible critical point on $\text{St}(d, r)$ and catches up with DRGTA.

B. PCA Experiment with Real-World Data

We also consider the PCA task defined in (18) on the MNIST dataset [49] to show the effectiveness of DRFGT. The MNIST dataset has dimension $d = 784$ and 60000 data samples are evenly distributed across $n = \{8, 16, 32\}$ agents. We set $r = 5$ and the graph structure is set as a ring with the Metropolis constant weight matrix following [14]. In this task, we specifically highlight the effect of network size n in the DRFGT algorithm through a multi-processing experimental setting (mpi4py).

We plot the distance to optimality $\log(d(\bar{x}, x^*))$ versus the epoch number and computation time. The results are provided in Fig. 2(a) and 2(b). For each specific choice of n , DRFGT converges significantly faster than DRGTA in CPU time. The performance with respect to epoch is similar across different n , which is expected. In addition, by increasing the number of agents in the network, computation time (per agent) can be accelerated by a nearly linear ratio, known as the linear speed-up effect [50]. We also plot the run-time comparison for different values of r in Fig. 2(c), explicitly showing the computation time saved by DRFGT by avoiding retractions.

VI. CONCLUSION

We proposed DRFGT to solve distributed non-convex optimization problems under orthogonal constraints and provided a safety step size guarantee to ensure DRFGT remains in the neighborhood of the Stiefel manifold. We proved that the convergence rate for DRFGT is $\mathcal{O}(1/K)$, the first convergence result for the decentralized landing algorithm. When the network function further satisfies a local Riemannian PL condition, we established a local linear convergence, the first-ever linear result among distributed Riemannian optimization

algorithms. Numerical results were provided as validation of our analysis. An interesting future direction is to extend this algorithm beyond Stiefel manifolds, similar to [28].

APPENDIX

A. Useful Lemmas and Inequalities

In this section, we present some fundamental inequalities and lemmas, which are useful for the technical analysis. We first gather the following facts in one place.

Lemma VI.1 ([44]). *For any given $x \in \text{St}(d, r)^\epsilon$, the singular values of x are between $\sqrt{1 - \epsilon}$ and $\sqrt{1 + \epsilon}$.*

In our proofs, we often times use the looser bounds $1 - \epsilon \leq \|x\|_2 \leq 1 + \epsilon$ to simplify the analysis.

Lemma VI.2. *The landing field $\Lambda(x)$, defined in (4), is Lipschitz continuous with a factor $L_\Lambda \leq L'$ and has orthogonal components, i.e., $\langle \text{grad} f(x), \nabla p(x) \rangle = 0$. The same hold for the local landing field $\Lambda_i(x)$ defined with respect to local Riemannian gradient $\text{grad} f_i(x)$.*

The proof is standard and can be found in [44] for $\Lambda(x)$, and the extension to local landing field is trivial.

Lemma VI.3. *For any given $x \in \text{St}(d, r)^\epsilon$, $x' = \text{Proj}_{\text{St}}(x)$, we have the following relationship*

$$\langle \nabla \mathcal{L}(x'), x - x' \rangle = 0.$$

Proof. If $x = USV^\top$ is the SVD decomposition for x , then $x' = UV^\top$. By observing the closed-form formula for $\nabla \mathcal{L}(x')$ in (10) the rest of the proof is simple algebra, omitted due to space limitation. \square

We then state the following lemma, which helps divide the gap between optimality and feasibility at each iteration.

Lemma VI.4. *For any given $x \in \text{St}(d, r)^\epsilon$, $x' = \text{Proj}_{\text{St}}(x)$, we have the following inequality*

$$\|x - x'\|_F \leq \|x^\top x - I_r\|_F. \quad (19)$$

The proof can be found in Lemma 4 of [51]. We then present the following lemma, a well-known result from [52].

Lemma VI.5 (Matrix Spectral Radius [52]). *Let $\mathbf{X} \in \mathbb{R}^{d \times d}$ be a non-negative matrix and $\mathbf{x} \in \mathbb{R}^d$ be a positive vector. If*

$\mathbf{X}\mathbf{x} < \mathbf{x}$, then $\rho(\mathbf{X}) < 1$. Furthermore, if $\mathbf{X}\mathbf{x} \leq z\mathbf{x}$ for some $z > 0$, then $\rho(\mathbf{X}) \leq z$. □

Another useful lemma is the following.

Lemma VI.6. For any point $x \in \text{St}(d, r)$, we have that

$$\|\nabla \mathcal{L}(x)\|_F \leq 2\|\Lambda(x)\|_F. \quad (20)$$

Proof. We know from (10) that

$$\nabla \mathcal{L}(x) = (I_d - \frac{1}{2}xx^\top)\nabla f(x) - \frac{1}{2}x\nabla f(x)^\top x,$$

which implies that

$$\begin{aligned} \|\nabla \mathcal{L}(x)\|_F^2 &= \text{Tr}\left(\nabla f(x)^\top (I_d - \frac{3}{4}xx^\top)\nabla f(x) \right. \\ &\quad + \frac{1}{4}x^\top \nabla f(x)\nabla f(x)^\top x - \frac{1}{4}\nabla f(x)^\top x\nabla f(x)^\top x \\ &\quad \left. - \frac{1}{4}x^\top \nabla f(x)x^\top \nabla f(x)\right). \end{aligned}$$

And for the landing field, we have

$$\Lambda(x) = \text{grad}f(x) = \frac{1}{2}\nabla f(x) - \frac{1}{2}x\nabla f(x)^\top x,$$

which gives

$$\begin{aligned} \|\Lambda(x)\|_F^2 &= \text{Tr}\left(\frac{1}{4}\nabla f(x)^\top \nabla f(x) + \frac{1}{4}x^\top \nabla f(x)\nabla f(x)^\top x \right. \\ &\quad \left. - \frac{1}{4}\nabla f(x)^\top x\nabla f(x)^\top x - \frac{1}{4}x^\top \nabla f(x)x^\top \nabla f(x)\right). \end{aligned}$$

We then have the result as

$$\begin{aligned} 4\|\Lambda(x)\|_F^2 - \|\nabla \mathcal{L}(x)\|_F^2 &= \text{Tr}\left(\frac{3}{4}\nabla f(x)^\top xx^\top \nabla f(x) + \frac{3}{4}x^\top \nabla f(x)\nabla f(x)^\top x \right. \\ &\quad \left. - \frac{3}{4}\nabla f(x)^\top x\nabla f(x)^\top x - \frac{3}{4}x^\top \nabla f(x)x^\top \nabla f(x)\right) \geq 0, \end{aligned}$$

where the inequality is due to the fact that $\text{Tr}(AA^\top) \geq \text{Tr}(AA)$ for a square matrix A . □

Lemma VI.7 (Lipschitz Continuity of $\text{grad}f$). For $x \in \text{St}(d, r)^\epsilon$ and $x' = \text{Proj}_{\text{St}}(x)$, it follows that

$$\|\text{grad}f(x) - \text{grad}f(x')\|_F \leq (3 + 2\epsilon)\hat{L}\|x - x'\|_F, \quad (21)$$

where \hat{L} is defined in (9).

Proof. Based on the definition in (2), we have

$$\begin{aligned} &\|\text{grad}f(x) - \text{grad}f(x')\|_F \\ &= \|\text{skew}(\nabla f(x)x^\top) - \text{skew}(\nabla f(x')x'^\top)\|_F \\ &= \|\text{skew}(\nabla f(x)x^\top) - \text{skew}(\nabla f(x')x'^\top) \\ &\quad + \text{skew}(\nabla f(x')x'^\top) - \text{skew}(\nabla f(x')x'^\top)x'\|_F \\ &\leq \|\text{skew}(\nabla f(x)x^\top) - \text{skew}(\nabla f(x')x'^\top)\|_F \|x\|_2 \\ &\quad + \|\text{skew}(\nabla f(x')x'^\top)\|_2 \|x - x'\|_F \\ &\leq (1 + \epsilon)\|\nabla f(x)x^\top - \nabla f(x')x'^\top\|_F \\ &\quad + (\max_{x' \in \text{St}(d, r)} \|\nabla f(x')\|_F) \|x - x'\|_F \\ &\leq (1 + \epsilon)(\|\nabla f(x)x^\top - \nabla f(x')x'^\top\|_F \\ &\quad + \|\nabla f(x)x'^\top - \nabla f(x')x'^\top\|_F) + \hat{L}\|x - x'\|_F \\ &\leq (3 + 2\epsilon)\hat{L}\|x - x'\|_F. \end{aligned} \quad (22)$$

The next three lemmas describe some useful relationships between the landing field $\Lambda(x)$ and merit function $\mathcal{L}(x)$. First, we show a property similar to the gradient domination on the Riemannian manifold, described as follows.

Lemma VI.8 (Pseudo Gradient Domination). Let $x \in \text{St}(d, r)^\epsilon \cap \mathcal{D}(\mathcal{S}, \delta)$, and without the loss of generality, assume $f_S^* = 0$. Then under Assumptions III.3 and III.4, we have

$$\mathcal{L}(x) \leq \frac{1}{\mu'} \|\Lambda(x)\|_F^2, \quad (23)$$

where $\frac{1}{\mu'} = \max\left\{\frac{1}{\mu}, \frac{2(3+2\epsilon)^2\hat{L}^2 + \mu L'}{2\mu\lambda^2(1-\epsilon)^2}\right\}$.

Proof. We first write the following inequality based on the smoothness of \mathcal{L} for $x' = \text{Proj}_{\text{St}}(x)$,

$$\begin{aligned} \mathcal{L}(x) &\leq \mathcal{L}(x') + \langle \nabla \mathcal{L}(x'), x - x' \rangle + \frac{L'}{2} \|x - x'\|_F^2 \\ &= f(x') + \frac{L'}{2} \|x - x'\|_F^2 \\ &\leq f(x') + \frac{L'}{2} \|x^\top x - I_r\|_F^2. \end{aligned} \quad (24)$$

where the equality is due to Lemma VI.3, and the last inequality follows from Lemma VI.4.

Since $\mathcal{S} \subseteq \text{St}(d, r)$, $\text{dist}(x', x) \leq \text{dist}(\mathcal{S}, x)$, and with triangle inequality, $\text{dist}(\mathcal{S}, x') \leq \text{dist}(\mathcal{S}, x) + \text{dist}(x', x) \leq 2\delta$. Therefore, $x' \in \mathcal{D}(\mathcal{S}, 2\delta)$, and we can apply Assumption III.4 on point x' .

Using the PL inequality for x' and Lipschitz continuity of $\text{grad}f(x)$ in Lemma VI.7, we have

$$\begin{aligned} f(x') &\leq \frac{1}{2\mu} \|\text{grad}f(x')\|_F^2 \\ &= \frac{1}{2\mu} \|\text{grad}f(x) + \text{grad}f(x') - \text{grad}f(x)\|_F^2 \\ &\leq \frac{1}{\mu} \|\text{grad}f(x)\|_F^2 + \frac{1}{\mu} \|\text{grad}f(x') - \text{grad}f(x)\|_F^2 \\ &\leq \frac{1}{\mu} \|\text{grad}f(x)\|_F^2 + \frac{(3 + 2\epsilon)^2\hat{L}^2}{\mu} \|x' - x\|_F^2 \\ &\leq \frac{1}{\mu} \|\text{grad}f(x)\|_F^2 + \frac{(3 + 2\epsilon)^2\hat{L}^2}{\mu} \|x^\top x - I_r\|_F^2, \end{aligned} \quad (25)$$

where the last line is due to Lemma VI.4. Note that

$$\|\Lambda(x)\|_F^2 \geq \|\text{grad}f(x)\|_F^2 + \lambda^2(1 - \epsilon)^2 \|x^\top x - I_r\|_F^2. \quad (26)$$

Therefore,

$$\begin{aligned} \mathcal{L}(x) &\leq \frac{1}{\mu} \|\text{grad}f(x)\|_F^2 + \left(\hat{L}^2 \frac{(3 + 2\epsilon)^2}{\mu} + \frac{L'}{2}\right) \|x^\top x - I_r\|_F^2 \\ &\leq \frac{1}{\mu'} \|\Lambda(x)\|_F^2, \end{aligned}$$

where

$$\frac{1}{\mu'} = \max\left\{\frac{1}{\mu}, \frac{2(3 + 2\epsilon)^2\hat{L}^2 + \mu L'}{2\mu\lambda^2(1 - \epsilon)^2}\right\}.$$

□

Lemma VI.9. *Let Assumption III.3 hold. For any $x \in \text{St}(d, r)^\epsilon$, we have $\|\nabla \mathcal{L}(x)\|_F \leq C\|\Lambda(x)\|_F$, where $C = 3L'\lambda^{-1}(1 - \epsilon)^{-1} + 2$.*

Proof. We denote the projection of x onto the Stiefel manifold as $x' = \text{Proj}_{\text{St}}(x)$. Given Lemma VI.6, we know that for any point x' on the manifold, we have $\|\nabla \mathcal{L}(x')\|_F \leq 2\|\Lambda(x')\|_F$. Hence,

$$\begin{aligned} \|\nabla \mathcal{L}(x)\|_F &= \|\nabla \mathcal{L}(x) - \nabla \mathcal{L}(x') + \nabla \mathcal{L}(x')\|_F \\ &\leq \|\nabla \mathcal{L}(x) - \nabla \mathcal{L}(x')\|_F + 2\|\Lambda(x')\|_F \\ &\leq \|\nabla \mathcal{L}(x) - \nabla \mathcal{L}(x')\|_F + 2\|\Lambda(x') - \Lambda(x) + \Lambda(x)\|_F \\ &\leq L'\|x - x'\|_F + 2L_\Lambda\|x - x'\|_F + 2\|\Lambda(x)\|_F. \end{aligned}$$

Since by definition, $\Lambda(x) = \text{grad}f(x) + \lambda x(x^\top x - I_r)$, we have $\|\Lambda(x)\|_F \geq \|\lambda x(x^\top x - I_r)\|_F \geq \lambda(1 - \epsilon)\|x^\top x - I_r\|_F$. Then, using Lemma VI.4, we get

$$\begin{aligned} \|\nabla \mathcal{L}(x)\|_F &\leq 3L'\|x - x'\|_F + 2\|\Lambda(x)\|_F \\ &\leq 3L'\|x^\top x - I_r\|_F + 2\|\Lambda(x)\|_F \\ &\leq (3L'\lambda^{-1}(1 - \epsilon)^{-1} + 2)\|\Lambda(x)\|_F. \end{aligned}$$

□

Lemma VI.10. *Let Assumption III.3 hold. For any $x \in \text{St}(d, r)^\epsilon \cap \mathcal{D}(\mathcal{S}, \delta)$ and $\delta \leq 1$, we have the following relationship,*

$$\|\Lambda(x)\|_F^2 \leq \frac{12L'}{\rho^2}(\mathcal{L}(x) - \mathcal{L}_S^*).$$

Proof. Consider $z = x - \alpha\Lambda(x)$ where $\alpha = \frac{\rho}{6L'}$ and recall that $\Lambda(x)$ is Lipschitz continuous with a factor L' , implying

$$\|\Lambda(x)\|_F = \|\Lambda(x) - \Lambda(x^*)\|_F \leq L'\|x - x^*\|_F \leq L'\delta \leq L',$$

for a local minimizer $x^* \in \mathcal{S}$. We then have

$$\begin{aligned} \|z^\top z - I_r\|_F &= \|(x - \alpha\Lambda(x))^\top (x - \alpha\Lambda(x)) - I_r\|_F \\ &\leq \|x^\top x - I_r\|_F + 2\alpha\|x^\top \Lambda(x)\|_F + \alpha^2\|\Lambda(x)\|_F^2 \\ &\leq \epsilon + 2\alpha(1 + \epsilon)\|\Lambda(x)\|_F + \alpha^2\|\Lambda(x)\|_F^2 \\ &\leq \epsilon + 2\alpha(1 + \epsilon)L' + \alpha^2L'^2. \end{aligned}$$

Since $\alpha \leq \frac{1}{L'}$, we have $\|z^\top z - I_r\|_F \leq 3 + 3\epsilon$, and $z \in \text{St}(d, r)^{3+3\epsilon}$. Following the analysis of Proposition 8 in [44], it is easy to show that the function $\mathcal{L}(z)$ is Lipschitz smooth with a factor $L + (11 + 9\epsilon)\gamma \leq 6L'$ when $z \in \text{St}(d, r)^{3+3\epsilon}$ (see Eq. 88 of [44]). Using the second condition in Proposition III.5 with $\alpha = \frac{\rho}{6L'}$, we get

$$\begin{aligned} \mathcal{L}(z) &\leq \mathcal{L}(x) - \alpha\rho\|\Lambda(x)\|_F^2 + 3\alpha^2L'\|\Lambda(x)\|_F^2 \\ &\leq \mathcal{L}(x) - \frac{\rho^2}{12L'}\|\Lambda(x)\|_F^2. \end{aligned}$$

Since $z \in \mathcal{D}(\mathcal{S}, 2\delta)$, we have $\mathcal{L}_S^* \leq \mathcal{L}(z)$; therefore,

$$\|\Lambda(x)\|_F^2 \leq \frac{12L'}{\rho^2}(\mathcal{L}(x) - \mathcal{L}_S^*).$$

□

B. Proofs in Section IV-B

Proof of Proposition IV.1: We start the proof by observing that for any set of matrices $\mathbf{Z}^\top = [Z_1^\top, \dots, Z_n^\top]$ and $\bar{\mathbf{Z}}^\top = [\bar{Z}_1^\top, \dots, \bar{Z}_n^\top]$, we have

$$\|\mathbf{Z} - \bar{\mathbf{Z}}\|_F \leq \|\mathbf{I}_n - \mathbf{1}\mathbf{1}^\top/n\|_2\|\mathbf{Z}\|_F \leq \|\mathbf{Z}\|_F. \quad (27)$$

Next, we show that given the assumption on \bar{x}_k and $\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F$, we get $x_{i,k} \in \text{St}(d, r)^\epsilon$ since

$$\begin{aligned} \|x_{i,k}^\top x_{i,k} - I_r\|_F &= \|(x_{i,k} - \bar{x}_k + \bar{x}_k)^\top (x_{i,k} - \bar{x}_k + \bar{x}_k) - I_r\|_F \\ &\leq \|\bar{x}_k^\top \bar{x}_k - I_r\|_F + 2\|\bar{x}_k^\top (x_{i,k} - \bar{x}_k)\|_F + \|x_{i,k} - \bar{x}_k\|_F^2 \\ &\leq \frac{\epsilon}{2} + 2(1 + \epsilon)\frac{\epsilon}{10} + \left(\frac{\epsilon}{10}\right)^2 \leq \epsilon. \end{aligned} \quad (28)$$

Applying (27) to $\mathbf{Z} = \mathbf{\Lambda}_k - \mathbf{\Lambda}_{k-1}$ in updates of \mathbf{y} , we get

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F &\leq \|\mathbf{W}\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F + \|\mathbf{\Lambda}_k - \mathbf{\Lambda}_{k-1}\|_F \\ &\leq \sigma_W\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F + \|\mathbf{\Lambda}_k\|_F + \|\mathbf{\Lambda}_{k-1}\|_F \\ &\leq \sigma_W\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F + 2\sqrt{n}(G + \lambda\epsilon(1 + \epsilon)). \end{aligned}$$

The last inequality is derived from $x_{i,k} \in \text{St}(d, r)^\epsilon$, and that $\|\text{grad}f_i(x_{i,k})\|_F \leq G$. Therefore

$$\|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F \leq \frac{2\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}{1 - \sigma_W}.$$

We also know that $\mathbf{x}_{k+1} = \mathbf{W}\mathbf{x}_k - \alpha\mathbf{y}_k$; then, $\mathbf{x}_{k+1} - \bar{\mathbf{x}}_{k+1} = \mathbf{W}(\mathbf{x}_k - \bar{\mathbf{x}}_k) - \alpha(\mathbf{y}_k - \bar{\mathbf{y}}_k)$, and we get,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \bar{\mathbf{x}}_{k+1}\|_F &= \|\mathbf{W}(\mathbf{x}_k - \bar{\mathbf{x}}_k) - \alpha(\mathbf{y}_k - \bar{\mathbf{y}}_k)\|_F \\ &\leq \|\mathbf{W}(\mathbf{x}_k - \bar{\mathbf{x}}_k)\|_F + \alpha\|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F \\ &\leq \sigma_W\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F + \alpha\|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F \\ &\leq \sigma_W\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F + \frac{2\alpha\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}{1 - \sigma_W}. \end{aligned}$$

If we choose

$$\alpha \leq \frac{(1 - \sigma_W)^2\epsilon}{20\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}, \quad (29)$$

we have $\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}_{k+1}\|_F \leq \frac{2\sqrt{n}\alpha(G + \lambda\epsilon(1 + \epsilon))}{(1 - \sigma_W)^2} \leq \frac{\epsilon}{10}$.

Next, we bound the distance of \bar{x}_{k+1} from $\text{St}(d, r)$. Since $\Lambda(\bar{x}_k) = \text{grad}f(\bar{x}_k) + \lambda\nabla p(\bar{x}_k)$, from the algorithm update we have

$$\bar{x}_{k+1} = \bar{x}_k - \alpha(\text{grad}f(\bar{x}_k) + \lambda\nabla p(\bar{x}_k) + \bar{y}_k - \Lambda(\bar{x}_k)).$$

For the ease of notation, we denote $\Delta_k \triangleq \bar{x}_k^\top \bar{x}_k - I_r$ and $D_k \triangleq \bar{y}_k - \Lambda(\bar{x}_k)$, and we then have

$$\begin{aligned} \Delta_{k+1} &= \bar{x}_{k+1}^\top \bar{x}_{k+1} - I_r \\ &= (\bar{x}_k - \alpha(\Lambda(\bar{x}_k) + D_k))^\top (\bar{x}_k - \alpha(\Lambda(\bar{x}_k) + D_k)) - I_r \\ &= \bar{x}_k^\top \bar{x}_k - I_r + \alpha^2(\Lambda(\bar{x}_k) + D_k)^\top (\Lambda(\bar{x}_k) + D_k) \\ &\quad - \alpha(\bar{x}_k^\top (\lambda\nabla p(\bar{x}_k) + D_k) + (\lambda\nabla p(\bar{x}_k) + D_k)^\top \bar{x}_k) \\ &\quad - \alpha(\bar{x}_k^\top \text{grad}f(\bar{x}_k) + \text{grad}f(\bar{x}_k)^\top \bar{x}_k). \end{aligned}$$

For the last line, given the definition of $\text{grad}f(x)$,

$$\bar{x}_k^\top \text{grad}f(\bar{x}_k) + \text{grad}f(\bar{x}_k)^\top \bar{x}_k = 0.$$

Since $\nabla p(\bar{x}_k) = \bar{x}_k \Delta_k$ we can write

$$\begin{aligned} \Delta_{k+1} = & \Delta_k - \alpha(2\lambda\Delta_k(I + \Delta_k) + \bar{x}_k^\top D_k + D_k^\top \bar{x}_k) \\ & + \alpha^2(\Lambda(\bar{x}_k) + D_k)^\top (\Lambda(\bar{x}_k) + D_k), \end{aligned}$$

which implies

$$\begin{aligned} \|\Delta_{k+1}\|_F \leq & (1 - 2\alpha\lambda)\|\Delta_k\|_F + 2\alpha\lambda\|\Delta_k\|_F^2 + 4\alpha\|D_k\|_F \\ & + 2\alpha^2(\|\Lambda(\bar{x}_k)\|_F^2 + \|D_k\|_F^2). \end{aligned}$$

Given $\alpha \leq \frac{1}{2\lambda}$, $1 - 2\alpha\lambda \geq 0$, and $\|\Lambda(\bar{x}_k)\|_F^2 \leq G^2 + \lambda^2(1 + \epsilon)\epsilon^2$ for $\bar{x}_k \in \text{St}(d, r)^{\frac{\epsilon}{2}}$, we have

$$\begin{aligned} \|\Delta_{k+1}\|_F \leq & (1 - 2\alpha\lambda)\frac{\epsilon}{2} + 2\alpha\lambda\frac{\epsilon^2}{4} + 4\alpha\|D_k\|_F \\ & + 2\alpha^2(G^2 + \lambda^2(1 + \epsilon)\epsilon^2 + \|D_k\|_F^2). \end{aligned}$$

Observe that due to the initialization

$$\begin{aligned} \|\bar{y}_k - \Lambda(\bar{x}_k)\|_F &= \left\| \frac{1}{n} \sum_{i=1}^n \Lambda_i(x_{i,k}) - \Lambda(\bar{x}_k) \right\|_F \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \Lambda_i(x_{i,k}) - \frac{1}{n} \sum_{i=1}^n \Lambda_i(\bar{x}_k) \right\|_F \\ &\leq \frac{L'}{n} \sum_{i=1}^n \|x_{i,k} - \bar{x}_k\|_F \\ &\leq \frac{L'}{\sqrt{n}} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F. \end{aligned} \quad (30)$$

We need $\|\Delta_{k+1}\|_F \leq \frac{\epsilon}{2}$ and we know from above that

$$\|D_k\|_F = \|\bar{y}_k - \Lambda(\bar{x}_k)\|_F \leq \frac{L'}{\sqrt{n}} \frac{2\alpha\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}{(1 - \sigma_W)^2}.$$

If we set $\alpha \leq \frac{\lambda\epsilon^2(1 - \sigma_W)^2}{16L'(G + \lambda\epsilon(1 + \epsilon))}$ such that $\|D_k\|_F \leq \frac{\lambda\epsilon^2}{8}$, a sufficient condition for α such that $\|\Delta_{k+1}\|_F \leq \frac{\epsilon}{2}$ is obtained as follows

$$\begin{aligned} 2\alpha(G^2 + \lambda^2(1 + \epsilon)\epsilon^2 + \frac{\epsilon^4\lambda^2}{16}) - \lambda\epsilon + \lambda\epsilon^2 &\leq 0 \\ \Rightarrow \alpha &\leq \frac{\lambda\epsilon(1 - \epsilon)}{2(G^2 + \lambda^2(1 + \epsilon)\epsilon^2 + \frac{\epsilon^4\lambda^2}{16})}. \end{aligned}$$

Combining all the requirements above on α , we get

$$\alpha_{safe} = \min \left\{ \frac{(1 - \sigma_W)^2\epsilon}{20\sqrt{n}(G + \lambda\epsilon(1 + \epsilon))}, \frac{\lambda\epsilon^2(1 - \sigma_W)^2}{16L'(G + \lambda\epsilon(1 + \epsilon))}, \frac{1}{2\lambda}, \frac{\lambda\epsilon(1 - \epsilon)}{2(G^2 + \lambda^2(1 + \epsilon)\epsilon^2 + \frac{\epsilon^4\lambda^2}{16})} \right\}.$$

Proof of Lemma IV.2: Since σ_W is the second largest singular value of W , from Assumption III.2 the following holds for all $\kappa > 0$ by AM-GM inequality

$$\begin{aligned} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 &= \|\mathbf{W}(\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}) - \alpha(\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1})\|_F^2 \\ &\leq (1 + \kappa)\|\mathbf{W}(\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1})\|_F^2 \\ &\quad + \frac{1 + \kappa}{\kappa}\alpha^2\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\leq (1 + \kappa)\sigma_W^2\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + \frac{1 + \kappa}{\kappa}\alpha^2\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2. \end{aligned}$$

We let $\kappa = \frac{1 - \sigma_W^2}{2\sigma_W^2}$ and get:

$$\begin{aligned} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 &\leq \frac{1 + \sigma_W^2}{2}\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + \frac{1 + \sigma_W^2}{1 - \sigma_W^2}\alpha^2\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2, \end{aligned}$$

which completes the proof. \square

Proof of Lemma IV.3: By definition of the update, we have

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &= \|\mathbf{W}(\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}) + (\mathbf{\Lambda}_k - \bar{\mathbf{\Lambda}}_k - \mathbf{\Lambda}_{k-1} + \bar{\mathbf{\Lambda}}_{k-1})\|_F^2. \end{aligned}$$

Using the same technique as Lemma IV.2, we get

$$\begin{aligned} &\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{\Lambda}_k - \bar{\mathbf{\Lambda}}_k - \mathbf{\Lambda}_{k-1} + \bar{\mathbf{\Lambda}}_{k-1}\|_F^2 \\ &\leq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{\Lambda}_k - \mathbf{\Lambda}_{k-1}\|_F^2 \\ &\leq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}L'^2\|\mathbf{x}_k - \mathbf{x}_{k-1}\|_F^2 \\ &= \frac{1 + \sigma_W^2}{1 - \sigma_W^2}L'^2\|(\mathbf{W} - \mathbf{I})\mathbf{x}_{k-1} - \alpha\mathbf{y}_{k-1}\|_F^2 \\ &\leq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}L'^2(2\|(\mathbf{W} - \mathbf{I})\mathbf{x}_{k-1}\|_F^2 + 2\|\alpha\mathbf{y}_{k-1}\|_F^2) \\ &\leq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}L'^2(2\|(\mathbf{W} - \mathbf{I})(\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1})\|_F^2 + 2\|\alpha\mathbf{y}_{k-1}\|_F^2) \\ &\leq 8L'^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 + 2L'^2\alpha^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{y}_{k-1}\|_F^2. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &\leq 8L'^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + \frac{1 + \sigma_W^2}{2}\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 + 2L'^2\alpha^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{y}_{k-1}\|_F^2. \end{aligned} \quad (31)$$

Now, for $\|\mathbf{y}_{k-1}\|_F^2$, we provide the following decomposition:

$$\begin{aligned} \|\mathbf{y}_{k-1}\|_F^2 &= \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1} + \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\leq 2\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 + 2\|\bar{\mathbf{y}}_{k-1}\|_F^2. \end{aligned}$$

Plugging the above inequality into (31), we get

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &\leq \left(\frac{1 + \sigma_W^2}{2} + 4L'^2\alpha^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\quad + 8L'^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 + 4L'^2\alpha^2\frac{1 + \sigma_W^2}{1 - \sigma_W^2}\|\bar{\mathbf{y}}_{k-1}\|_F^2. \end{aligned}$$

\square

Proof of Theorem IV.4: Combining the results of Lemma IV.2 and Lemma IV.3, we can easily verify the linear system relationship. Then, we need to ensure that $\rho(\tilde{G}) < 1$. Based

on Lemma VI.5, $\rho(\tilde{G}) < 1$ holds if there exists a solution $s_1, s_2 > 0$ for the following inequalities

$$\begin{aligned} \left(\frac{1 + \sigma_W^2}{2} + 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) s_1 + 8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} s_2 &< s_1, \\ \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 L'^2 s_1 + \frac{1 + \sigma_W^2}{2} s_2 &< s_2. \end{aligned}$$

Simplifying these equations, we get

$$\begin{aligned} 2 \frac{1 + \sigma_W^2}{(1 - \sigma_W^2)^2} \alpha^2 L'^2 s_1 &< s_2, \\ 8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} s_2 &< \left(\frac{1 - \sigma_W^2}{2} - 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) s_1. \end{aligned}$$

After removing s_2 , we get a sufficient condition for α such that

$$\alpha < \frac{(1 - \sigma_W^2)^2}{1 + \sigma_W^2} \frac{1}{16L'}.$$

□

Proof of Lemma IV.5: We first start with the Lipschitz smoothness of \mathcal{L} and the updates of the algorithm to get

$$\begin{aligned} &\mathcal{L}(\bar{x}_k) - \mathcal{L}(\bar{x}_{k-1}) \\ &\leq \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \bar{x}_k - \bar{x}_{k-1} \rangle + \frac{L'}{2} \|\bar{x}_k - \bar{x}_{k-1}\|_F^2 \\ &= -\alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \bar{y}_{k-1} \rangle + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2 \\ &= -\alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \bar{y}_{k-1} - \Lambda(\bar{x}_{k-1}) + \Lambda(\bar{x}_{k-1}) \rangle \\ &\quad + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2 \\ &\leq -\alpha \rho \|\Lambda(\bar{x}_{k-1})\|_F^2 - \alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \bar{y}_{k-1} - \Lambda(\bar{x}_{k-1}) \rangle \\ &\quad + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2 \\ &\leq -\alpha \rho \|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha \rho}{2\kappa} \|\nabla \mathcal{L}(\bar{x}_{k-1})\|_F^2 \\ &\quad + \frac{\alpha \kappa}{2} \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2, \end{aligned}$$

where the second inequality used the results from Proposition III.5 and the last inequality holds for any $\kappa > 0$ due to AM-GM inequality. Now, with the help of Lemma VI.9 and setting $\kappa = \frac{C^2}{\rho}$ in above, we get

$$\begin{aligned} &\mathcal{L}(\bar{x}_k) - \mathcal{L}(\bar{x}_{k-1}) \\ &\leq -\alpha \rho \|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha \rho}{2C^2} \|\nabla \mathcal{L}(\bar{x}_{k-1})\|_F^2 \\ &\quad + \frac{\alpha C^2}{2\rho} \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2 \\ &\leq -\frac{\alpha \rho}{2} \|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha C^2}{2\rho} \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2 \\ &\leq -\frac{\alpha \rho}{2} \|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha C^2 L'^2}{2\rho n} \|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 + \frac{\alpha^2 L'}{2} \|\bar{y}_{k-1}\|_F^2, \end{aligned}$$

where the last inequality is due to (30). The proof is completed. □

C. Proofs in Section IV-C

Proof of Corollary IV.6: We first calculate the determinant of $I - \tilde{G}$. Given that

$$I - \tilde{G} = \begin{bmatrix} 1 - \frac{1 + \sigma_W^2}{2} - 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} & -8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \\ -\frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 L'^2 & 1 - \frac{1 + \sigma_W^2}{2} \end{bmatrix},$$

we can explicitly write

$$\det(I - \tilde{G}) = \left(\frac{1 - \sigma_W^2}{2} \right)^2 - 2\alpha^2 L'^2 \left((1 + \sigma_W^2) + 4 \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^2} \right).$$

Then, in order to have $\det(I - \tilde{G}) \geq \frac{(1 - \sigma_W^2)^2}{8}$, we need

$$2\alpha^2 L'^2 \left((1 + \sigma_W^2) + 4 \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^2} \right) \leq \frac{(1 - \sigma_W^2)^2}{8},$$

which is satisfied when $\alpha < \frac{(1 - \sigma_W^2)^2}{1 + \sigma_W^2} \frac{1}{16L'}$. Next, we note that

$$\begin{aligned} (I - \tilde{G})^{-1} &= \frac{(I - \tilde{G})^*}{\det(I - \tilde{G})} \\ &\leq \frac{8}{(1 - \sigma_W^2)^2} \begin{bmatrix} \frac{1 - \sigma_W^2}{2} - 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} & -8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \\ -\frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 L'^2 & \frac{1 - \sigma_W^2}{2} \end{bmatrix}^* \\ &= \frac{8}{(1 - \sigma_W^2)^2} \begin{bmatrix} \frac{1 - \sigma_W^2}{2} & 8 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \\ \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \alpha^2 L'^2 & \frac{1 - \sigma_W^2}{2} - 4L'^2 \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{4}{1 - \sigma_W^2} & 64 \frac{1 + \sigma_W^2}{(1 - \sigma_W^2)^3} \\ \frac{1 + \sigma_W^2}{(1 - \sigma_W^2)^3} 8\alpha^2 L'^2 & \frac{4}{1 - \sigma_W^2} - 32L'^2 \alpha^2 \frac{1 + \sigma_W^2}{(1 - \sigma_W^2)^3} \end{bmatrix}. \end{aligned}$$

The above inequality is element-wise, with each entry non-negative.

Now, given that $\tilde{\xi}_k \leq \tilde{G} \tilde{\xi}_{k-1} + \tilde{u}_{k-1}$ and that $\tilde{\xi}_0 = 0$, we have

$$\tilde{\xi}_k \leq \sum_{t=0}^{k-1} \tilde{G}^t \tilde{u}_{k-t-1} \Rightarrow \sum_{k=1}^K \tilde{\xi}_k \leq (I - \tilde{G})^{-1} \sum_{k=0}^{K-1} \tilde{u}_k.$$

Given the inequality above, we know that

$$\sum_k L' \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 \leq \frac{1 + \sigma_W^2}{(1 - \sigma_W^2)^3} 8\alpha^2 L'^2 \sum_k 4L' \alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \|\bar{\mathbf{y}}_{k-1}\|_F^2.$$

Therefore, we can write

$$\sum_k \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 \leq \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2 \sum_k \|\bar{\mathbf{y}}_{k-1}\|_F^2,$$

which completes the proof. □

Proof of Theorem IV.7: We first state that $\alpha < \frac{(1 - \sigma_W^2)^2}{1 + \sigma_W^2} \frac{1}{16L'}$ guarantees the stability of the system (14) and enables the use of Corollary IV.6. We have that

$$\begin{aligned} \|\bar{y}_{k-1}\|_F^2 &= \|\Lambda(\bar{x}_{k-1}) + (\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1}))\|_F^2 \\ &\leq 2\|\Lambda(\bar{x}_{k-1})\|_F^2 + 2\|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2. \end{aligned} \quad (32)$$

Combining this with Lemma IV.5, we get

$$\begin{aligned}
& \mathcal{L}(\bar{x}_k) - \mathcal{L}(\bar{x}_{k-1}) \\
& \leq -\frac{\alpha\rho}{2}\|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha C^2 L'^2}{2\rho n}\|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 + \frac{\alpha^2 L'}{2}\|\bar{y}_{k-1}\|_F^2 \\
& \leq -\frac{\alpha\rho}{2}\|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha C^2 L'^2}{2\rho n}\|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 - \frac{\alpha^2 L'}{2}\|\bar{y}_{k-1}\|_F^2 \\
& \quad + 2\alpha^2 L'\|\Lambda(\bar{x}_{k-1})\|_F^2 + 2\alpha^2 L'\|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 \\
& \leq -\left(\frac{\alpha\rho}{2} - 2\alpha^2 L'\right)\|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + \left(\frac{\alpha C^2 L'^2}{2\rho n} + \frac{2\alpha^2 L'^3}{n}\right)\|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 - \frac{\alpha^2 L'}{2}\|\bar{y}_{k-1}\|_F^2.
\end{aligned}$$

Note that $\alpha \leq \frac{\rho}{8L'}$, and in addition we know from Proposition III.5 that $\rho \leq \frac{1}{2}$. Therefore,

$$\begin{aligned}
\mathcal{L}(\bar{x}_k) - \mathcal{L}(\bar{x}_{k-1}) & \leq -\frac{\alpha\rho}{4}\|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + \frac{\alpha L'^2 C^2}{\rho n}\|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 - \frac{\alpha^2 L'}{2}\|\bar{y}_{k-1}\|_F^2.
\end{aligned}$$

Rearranging and summing both sides over $k = 1, \dots, K$ and using Corollary IV.6, we have

$$\begin{aligned}
& \sum_k \frac{\alpha\rho}{4}\|\Lambda(\bar{x}_{k-1})\|_F^2 - (\mathcal{L}(\bar{x}_0) - \mathcal{L}(\bar{x}_K)) \\
& \leq -\frac{\alpha^2 L'}{2} \sum_k \|\bar{y}_{k-1}\|_F^2 + \frac{\alpha L'^2 C^2}{\rho n} \sum_k \|\bar{\mathbf{x}}_{k-1} - \mathbf{x}_{k-1}\|_F^2 \\
& \leq -\frac{\alpha^2 L'}{2} \sum_k \|\bar{y}_{k-1}\|_F^2 \\
& \quad + \frac{\alpha L'^2 C^2}{\rho n} \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2 \sum_k \|\bar{\mathbf{y}}_{k-1}\|_F^2 \\
& = -\left(\frac{\alpha^2 L'}{2} - \frac{\alpha L'^2 C^2}{\rho} \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2\right) \sum_k \|\bar{y}_{k-1}\|_F^2.
\end{aligned}$$

Selecting $\alpha \leq \frac{1}{4L'} \sqrt{\frac{\rho(1 - \sigma_W^2)^4}{(1 + \sigma_W^2)^2 C^2}}$ such that $\frac{\alpha^2 L'}{2} - \frac{\alpha L'^2 C^2}{\rho} \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2 \geq 0$, we get the desired result

$$\frac{\sum_k \|\Lambda(\bar{x}_{k-1})\|_F^2}{K} \leq \frac{1}{K} \frac{4}{\alpha\rho} (\mathcal{L}(\bar{x}_0) - \mathcal{L}(\bar{x}_K)).$$

In addition, the consensus error can be bounded with

$$\begin{aligned}
\sum_k \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 & \leq \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 32\alpha^4 L'^2 \sum_k \|\bar{\mathbf{y}}_{k-1}\|_F^2 \\
& \leq \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 64\alpha^4 L'^2 n \sum_k \|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 64\alpha^4 L'^4 \sum_k \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2,
\end{aligned}$$

where we used (30) and (32) again. Since we have

$$\alpha \leq \frac{1}{4L'} \sqrt{\frac{\rho(1 - \sigma_W^2)^4}{(1 + \sigma_W^2)^2 C^2}},$$

and $\rho \leq \frac{1}{2}, C > 1$, it is easy to show that $\frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 64\alpha^4 L'^4 \leq \frac{1}{2}$; hence, we can write

$$\sum_k \|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2 \leq \frac{(1 + \sigma_W^2)^2}{(1 - \sigma_W^2)^4} 128n\alpha^4 L'^2 \sum_k \|\Lambda(\bar{x}_{k-1})\|_F^2,$$

which completes the proof. \square

D. Proofs in Section IV-D

Proof of Lemma IV.8: Since $x \in \text{St}(d, r)^\epsilon$, given Proposition III.5, it is clear that,

$$\|\nabla \mathcal{L}(x)\|_F \geq \rho \|\Lambda(x)\|_F.$$

Combined with Lemma VI.8, we have,

$$\mathcal{L}(x) - \mathcal{L}_S^* \leq \frac{1}{\mu'} \|\Lambda(x)\|_F^2 \leq \frac{1}{\mu' \rho^2} \|\nabla \mathcal{L}(x)\|_F^2.$$

Therefore, we know that the local Euclidean PL condition for $\mathcal{L}(x)$ holds on $\forall x \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^\epsilon$. Given Proposition 2.2 in [47], the local Euclidean PL condition on $\mathcal{L}(x)$ implies the quadratic growth relationship below,

$$\mathcal{L}(x) - \mathcal{L}_S^* \geq \frac{\mu' \rho^2}{4} \text{dist}(\mathcal{S}, x)^2.$$

\square

Proof of Lemma IV.9: Without loss of generality and for the ease of notation, in this proof, we consider the case where $\mathcal{L}_S^* = 0$. By the smoothness of \mathcal{L} and the algorithm update, we have

$$\begin{aligned}
\mathcal{L}(\bar{x}_k) & \leq \mathcal{L}(\bar{x}_{k-1}) + \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \bar{x}_k - \bar{x}_{k-1} \rangle + \frac{L'}{2} \|\bar{x}_k - \bar{x}_{k-1}\|_F^2 \\
& = \mathcal{L}(\bar{x}_{k-1}) - \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \alpha \bar{y}_{k-1} \rangle + \frac{L' \alpha^2}{2} \|\bar{y}_{k-1}\|_F^2 \\
& = \mathcal{L}(\bar{x}_{k-1}) - \alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \Lambda(\bar{x}_{k-1}) \rangle \\
& \quad + \alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \Lambda(\bar{x}_{k-1}) - \bar{y}_{k-1} \rangle + \frac{L' \alpha^2}{2} \|\bar{y}_{k-1}\|_F^2.
\end{aligned}$$

Applying Proposition III.5, we derive

$$\begin{aligned}
\mathcal{L}(\bar{x}_k) & \leq \mathcal{L}(\bar{x}_{k-1}) - \alpha \rho \|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + \alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \Lambda(\bar{x}_{k-1}) - \bar{y}_{k-1} \rangle + \frac{L' \alpha^2}{2} \|\bar{y}_{k-1}\|_F^2 \\
& = \mathcal{L}(\bar{x}_{k-1}) - \alpha \rho \|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + \alpha \langle \nabla \mathcal{L}(\bar{x}_{k-1}), \Lambda(\bar{x}_{k-1}) - \bar{y}_{k-1} \rangle \\
& \quad + \frac{L' \alpha^2}{2} \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1}) + \Lambda(\bar{x}_{k-1})\|_F^2.
\end{aligned}$$

Applying the AM-GM inequality on the last term as well as the inner product term (with a factor $\eta > 0$), we get

$$\begin{aligned}
\mathcal{L}(\bar{x}_k) & \leq \mathcal{L}(\bar{x}_{k-1}) - (\alpha\rho - \alpha^2 L') \|\Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\eta\alpha}{2} \|\nabla \mathcal{L}(\bar{x}_{k-1})\|_F^2 \\
& \quad + L' \alpha^2 \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 + \frac{\alpha}{2\eta} \|\Lambda(\bar{x}_{k-1}) - \bar{y}_{k-1}\|_F^2 \\
& \leq \mathcal{L}(\bar{x}_{k-1}) - (\alpha\rho - \alpha^2 L' - \eta\alpha C^2/2) \|\Lambda(\bar{x}_{k-1})\|_F^2 \\
& \quad + (L' \alpha^2 + \frac{\alpha}{2\eta}) \|\bar{y}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2,
\end{aligned}$$

where the last inequality is due to Lemma VI.9. We then apply Lemma VI.8 on the term $\|\Lambda(\bar{x}_{k-1})\|_F^2$, set $\eta = \frac{\rho}{2C^2}$, consider $\alpha \leq \frac{\rho}{2L'}$, and use (30) on the last term to get the final result. \square

Proof of Lemma IV.10: We first recall (31), showing that the consensus error on \mathbf{y}_k can be bounded using terms related to previous iterations, including $\|\mathbf{y}_{k-1}\|_F^2$ for which we provide the following decomposition,

$$\|\mathbf{y}_{k-1}\|_F^2 \leq 2\|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 + 2n\|\bar{\mathbf{y}}_{k-1}\|_F^2,$$

and observe that

$$\begin{aligned} 2n\|\bar{\mathbf{y}}_{k-1}\|_F^2 &\leq 4n\|\bar{\mathbf{y}}_{k-1} - \Lambda(\bar{x}_{k-1})\|_F^2 + 4n\|\Lambda(\bar{x}_{k-1})\|_F^2 \\ &\leq 4L'^2\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 + 4n\|\Lambda(\bar{x}_{k-1})\|_F^2 \\ &\leq 4L'^2\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 + 48nL'\rho^{-2}\mathcal{L}(\bar{x}_{k-1}), \end{aligned}$$

where last inequality is derived from Lemma VI.10 assuming $\mathcal{L}_S^* = 0$ for notation convenience. Plugging the inequalities above into (31), we get

$$\begin{aligned} \|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2 &\leq \left(\frac{1 + \sigma_W^2}{2} + 4L'^2\alpha^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \right) \|\mathbf{y}_{k-1} - \bar{\mathbf{y}}_{k-1}\|_F^2 \\ &\quad + 8L'^2 \frac{1 + \sigma_W^2}{1 - \sigma_W^2} (1 + \alpha^2 L'^2) \|\mathbf{x}_{k-1} - \bar{\mathbf{x}}_{k-1}\|_F^2 \\ &\quad + \frac{1 + \sigma_W^2}{1 - \sigma_W^2} \frac{96n\alpha^2 L'^3}{\rho^2} \mathcal{L}(\bar{x}_{k-1}). \end{aligned}$$

\square

Proof of Theorem IV.11: The matrix M is constructed by recalling the definition of the state vector in (16),

$$\xi_k \triangleq [\|\mathbf{y}_k - \bar{\mathbf{y}}_k\|_F^2/L', L'\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F^2, n(\mathcal{L}(\bar{x}_k) - \mathcal{L}_S^*)]^\top,$$

and applying Lemmas IV.2, IV.9, IV.10 on the three terms. For the ease of notation, we define $\Theta \triangleq \frac{1 + \sigma_W^2}{1 - \sigma_W^2}$ to get

$$M \triangleq \begin{bmatrix} \frac{1 + \sigma_W^2}{2} + 4L'^2\alpha^2\Theta & 8\alpha^2 L'^2\Theta + 8\Theta & 96L'^2\alpha^2\Theta \frac{1}{\rho^2} \\ \alpha^2 L'^2\Theta & \frac{1 + \sigma_W^2}{2} & 0 \\ 0 & \alpha^2 L'^2 + \frac{\alpha L' C^2}{\rho} & 1 - \frac{\alpha \rho \mu'}{4} \end{bmatrix}.$$

We then prove that the system converges with a linear rate $1 - \frac{\alpha \rho \mu'}{8}$. Suppose that we have a positive vector $\delta = [\delta_1, \delta_2, \delta_3]^\top$, such that

$$M\delta \leq (1 - \frac{\alpha \rho \mu'}{8})\delta,$$

with element-wise inequality. Then, in lieu of Lemma VI.5, the spectral radius of M is upper bounded by $1 - \frac{\alpha \rho \mu'}{8}$.

In the next few lines of the proof, we solve these three inequalities to obtain sufficient conditions on α and find such δ . We start from the third row to get

$$\left(\alpha^2 L'^2 + \frac{\alpha L' C^2}{\rho} \right) \delta_2 + \left(1 - \frac{\alpha \rho \mu'}{4} \right) \delta_3 \leq \left(1 - \frac{\alpha \rho \mu'}{8} \right) \delta_3,$$

for which, given $\alpha \leq \frac{1}{2L'}$, a sufficient condition is

$$\left(\frac{4L'}{\rho \mu'} + \frac{8L' C^2}{\rho^2 \mu'} \right) \delta_2 \leq \delta_3.$$

We define $\Phi \triangleq \frac{4L'}{\rho \mu'} + \frac{8L' C^2}{\rho^2 \mu'}$, and fix $\delta_2 = 1, \delta_3 = \Phi$. Next, we look at the first equation

$$\begin{aligned} \left(\frac{1 + \sigma_W^2}{2} + 4L'^2\alpha^2\Theta \right) \delta_1 + \left(8\alpha^2 L'^2\Theta + 8\Theta \right) \delta_2 \\ + \left(96L'^2\alpha^2\Theta \frac{1}{\rho^2} \right) \delta_3 \leq \left(1 - \frac{\alpha \rho \mu'}{8} \right) \delta_1. \end{aligned}$$

We choose $\alpha > 0$ such that $4L'^2\alpha^2\Theta \leq \frac{1 - \sigma_W^2}{8}$ and $\frac{\alpha \rho \mu'}{8} \leq \frac{1 - \sigma_W^2}{8}$. This puts constraint on α such that

$$\alpha \leq \min \left\{ \sqrt{\frac{1 - \sigma_W^2}{32L'^2\Theta}}, \frac{1 - \sigma_W^2}{\rho \mu'} \right\}. \quad (33)$$

Then, with $\delta_2 = 1, \delta_3 = \Phi$, a sufficient condition for the inequality is

$$(8\alpha^2 L'^2\Theta + 8\Theta) + \left(96L'^2\alpha^2\Theta \frac{\Phi}{\rho^2} \right) \leq \frac{1 - \sigma_W^2}{4} \delta_1.$$

Therefore, we fix

$$\delta_1 = \frac{32}{1 - \sigma_W^2} \left(\alpha^2 L'^2\Theta + \Theta + 12L'^2\alpha^2\Theta \frac{\Phi}{\rho^2} \right).$$

Finally, we look at the second equation using our choices of δ_1 and δ_2 to get,

$$\alpha^4 L'^4 \Theta^2 \left(1 + \frac{12\Phi}{\rho^2} \right) + \alpha^2 L'^2 \Theta^2 \leq \frac{(1 - \sigma_W^2)^2}{128}.$$

A sufficient condition for the above inequality to hold is that

$$\begin{aligned} \alpha^4 L'^4 \Theta^2 \left(1 + \frac{12\Phi}{\rho^2} \right) &\leq \frac{(1 - \sigma_W^2)^2}{256} \\ \alpha^2 L'^2 \Theta^2 &\leq \frac{(1 - \sigma_W^2)^2}{256}. \end{aligned}$$

This can be achieved when

$$\alpha \leq \sqrt[4]{\frac{(1 - \sigma_W^2)^2}{256L'^4\Theta^2(1 + \frac{12\Phi}{\rho^2})}} = \frac{\sqrt{1 - \sigma_W^2}}{4L'\sqrt{\Theta} \sqrt[4]{1 + \frac{12\Phi}{\rho^2}}},$$

and

$$\alpha \leq \frac{1 - \sigma_W^2}{16L'\Theta}.$$

Combining the above constraints on α with (33) and simplifying the sufficient conditions, we obtain

$$\alpha \leq \min \left\{ \frac{1 - \sigma_W^2}{\rho \mu'}, \frac{\sqrt{1 - \sigma_W^2}}{4L'\sqrt{\Theta} \sqrt[4]{1 + \frac{12\Phi}{\rho^2}}}, \frac{1 - \sigma_W^2}{16L'\Theta} \right\}.$$

We further want $\alpha \leq \frac{\rho}{2L'}$ to apply Lemma IV.9, which together with above completes the proof for step size requirement.

We now need to verify that $\bar{x}_k \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^{\frac{5}{2}}$ for every iteration k . We know $\bar{x}_k \in \text{St}(d, r)^{\frac{5}{2}}$ is already satisfied due to the safe step size constraint. Recall that \mathbf{v} is the eigenvector corresponding to the spectral radius of M , which has all positive elements since M is non-negative and irreducible. We prove by induction that if $\bar{x}_k \in \mathcal{D}(\mathcal{S}, \delta)$ and $\xi_k \leq \frac{1}{8} n \mu' \rho^2 \delta^2 \mathbf{v}$,

then $\bar{x}_{k+1} \in \mathcal{D}(\mathcal{S}, \delta)$ and $\xi_{k+1} \leq \frac{1}{8}n\mu'\rho^2\delta^2\mathbf{v}$, where the inequality is element-wise and $\|\mathbf{v}\|_1 = 1$. We have

$$\begin{aligned} \text{dist}(\mathcal{S}, \bar{x}_{k+1}) &\leq \text{dist}(\mathcal{S}, \bar{x}_k) + \|\bar{x}_k - \bar{x}_{k+1}\|_F \\ &= \text{dist}(\mathcal{S}, \bar{x}_k) + \alpha\|\bar{y}_k\|_F \\ &\leq \text{dist}(\mathcal{S}, \bar{x}_k) + \alpha\|\bar{y}_k - \Lambda(\bar{x}_k)\|_F + \alpha\|\Lambda(\bar{x}_k)\|_F \\ &\leq (1 + \alpha L')\text{dist}(\mathcal{S}, \bar{x}_k) + \alpha L' n^{-1/2}\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F \\ &\leq 1.25 \text{dist}(\mathcal{S}, \bar{x}_k) + 0.25n^{-1/2}\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F, \end{aligned}$$

where we used $\alpha L' \leq \rho/2 \leq 0.25$, Lipschitz continuity of Λ , and (30) in the last two lines. Applying Lemma IV.8, we know that

$$\text{dist}(\mathcal{S}, \bar{x}_k)^2 \leq \frac{4(\mathcal{L}(\bar{x}_k) - \mathcal{L}_S^*)}{\mu'\rho^2} \leq \frac{\delta^2}{2},$$

where the last inequality is due to the induction assumption of $\xi_k \leq \frac{1}{8}n\mu'\rho^2\delta^2\mathbf{v}$, which also guarantees that $n^{-1/2}\|\mathbf{x}_k - \bar{\mathbf{x}}_k\|_F \leq \frac{\delta}{2\sqrt{2}}$. Therefore, we have $\text{dist}(\mathcal{S}, \bar{x}_{k+1}) \leq \delta$, i.e., $\bar{x}_{k+1} \in \mathcal{D}(\mathcal{S}, \delta)$. It is also immediate that

$$\xi_{k+1} \leq M\xi_k \leq \frac{1}{8}n\mu'\rho^2\delta^2M\mathbf{v} < \frac{1}{8}n\mu'\rho^2\delta^2\mathbf{v}.$$

Therefore, given the initial conditions on x_0 and ξ_0 , by induction $\bar{x}_k \in \mathcal{D}(\mathcal{S}, \delta) \cap \text{St}(d, r)^{\frac{\epsilon}{2}}$ and $\xi_k \leq \frac{1}{8}n\mu'\rho^2\delta^2\mathbf{v}$ are guaranteed to be satisfied for all k . \square

REFERENCES

- [1] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of educational psychology*, vol. 24, no. 6, p. 417, 1933.
- [2] —, "Canonical correlation analysis (cca)," *Journal of Educational Psychology*, vol. 10, pp. 12913–2, 1935.
- [3] D. Kempe and F. McSherry, "A decentralized algorithm for spectral analysis," *Journal of Computer and System Sciences*, vol. 74, no. 1, pp. 70–83, 2008.
- [4] N. Kishore Kumar and J. Schneider, "Literature survey on low rank approximation of matrices," *Linear and Multilinear Algebra*, vol. 65, no. 11, pp. 2212–2244, 2017.
- [5] H. Raja and W. U. Bajwa, "Cloud k-svd: A collaborative dictionary learning algorithm for big, distributed data," *IEEE Transactions on Signal Processing*, vol. 64, no. 1, pp. 173–188, 2015.
- [6] M. Arjovsky, A. Shah, and Y. Bengio, "Unitary evolution recurrent neural networks," in *International Conference on Machine Learning*. PMLR, 2016, pp. 1120–1128.
- [7] E. Vorontsov, C. Trabelsi, S. Kadoury, and C. Pal, "On orthogonality and learning recurrent networks with long term dependencies," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3570–3578.
- [8] A. Trockman and J. Z. Kolter, "Orthogonalizing convolutional layers with the cayley transform," *arXiv preprint arXiv:2104.07167*, 2021.
- [9] Q. Zhang, M. Chen, A. Bukharin, P. He, Y. Cheng, W. Chen, and T. Zhao, "Adaptive budget allocation for parameter-efficient fine-tuning," *arXiv preprint arXiv:2303.10512*, 2023.
- [10] P.-A. Absil and J. Malick, "Projection-like retractions on matrix manifolds," *SIAM Journal on Optimization*, vol. 22, no. 1, pp. 135–158, 2012.
- [11] N. Boumal, P.-A. Absil, and C. Cartis, "Global rates of convergence for nonconvex optimization on manifolds," *IMA Journal of Numerical Analysis*, vol. 39, no. 1, pp. 1–33, 2019.
- [12] H. Zhang and S. Sra, "First-order methods for geodesically convex optimization," in *Conference on Learning Theory*. PMLR, 2016, pp. 1617–1638.
- [13] P. Ablin and G. Peyré, "Fast and accurate optimization on the orthogonal manifold without retraction," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 5636–5657.
- [14] S. Chen, A. Garcia, M. Hong, and S. Shahrampour, "Decentralized riemannian gradient descent on the stiefel manifold," in *International Conference on Machine Learning*. PMLR, 2021, pp. 1594–1605.
- [15] —, "On the local linear rate of consensus on the stiefel manifold," *IEEE Transactions on Automatic Control*, pp. 1–16, 2023.
- [16] H. Ye and T. Zhang, "Deepca: Decentralized exact pca with linear convergence rate," *The Journal of Machine Learning Research*, vol. 22, no. 1, pp. 10777–10803, 2021.
- [17] L. Wang and X. Liu, "Decentralized optimization over the stiefel manifold by an approximate augmented lagrangian function," *IEEE Transactions on Signal Processing*, vol. 70, pp. 3029–3041, 2022.
- [18] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- [19] C. Qi, K. A. Gallivan, and P.-A. Absil, "Riemannian bfgs algorithm with applications," in *Recent Advances in Optimization and its Applications in Engineering: The 14th Belgian-French-German Conference on Optimization*. Springer, 2010, pp. 183–192.
- [20] Y. Liu, F. Shang, J. Cheng, H. Cheng, and L. Jiao, "Accelerated first-order methods for geodesically convex optimization on riemannian manifolds," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [21] S. Schechtman, D. Tiapkin, M. Muehlebach, and E. Moulines, "Orthogonal directions constrained gradient method: from non-linear equality constraints to stiefel manifold," *arXiv preprint arXiv:2303.09261*, 2023.
- [22] R. Fletcher and S. Leyffer, "Nonlinear programming without a penalty function," *Mathematical programming*, vol. 91, pp. 239–269, 2002.
- [23] B. Gao, X. Liu, and Y.-x. Yuan, "Parallelizable algorithms for optimization problems with orthogonality constraints," *SIAM Journal on Scientific Computing*, vol. 41, no. 3, pp. A1949–A1983, 2019.
- [24] Q. Lin, R. Ma, and Y. Xu, "Complexity of an inexact proximal-point penalty method for constrained smooth non-convex optimization," *Computational optimization and applications*, vol. 82, no. 1, pp. 175–224, 2022.
- [25] D. Boob, Q. Deng, and G. Lan, "Level constrained first order methods for function constrained optimization," *Mathematical Programming*, pp. 1–61, 2024.
- [26] G. Scutari, F. Facchinei, and L. Lampariello, "Parallel and distributed methods for constrained nonconvex optimization—part i: Theory," *IEEE Transactions on Signal Processing*, vol. 65, no. 8, pp. 1929–1944, 2016.
- [27] N. Xiao, X. Liu, and Y.-x. Yuan, "A class of smooth exact penalty function methods for optimization problems with orthogonality constraints," *Optimization Methods and Software*, vol. 37, no. 4, pp. 1205–1241, 2022.
- [28] N. Xiao, X. Liu, and K.-C. Toh, "Dissolving constraints for riemannian optimization," *Mathematics of Operations Research*, vol. 49, no. 1, pp. 366–397, 2024.
- [29] 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.
- [30] 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.
- [31] S. Chen, A. Garcia, and S. Shahrampour, "On distributed nonconvex optimization: Projected subgradient method for weakly convex problems in networks," *IEEE Transactions on Automatic Control*, vol. 67, no. 2, pp. 662–675, 2021.
- [32] 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.
- [33] P. Di Lorenzo and G. Scutari, "Next: In-network nonconvex optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.
- [34] 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.
- [35] Y. Sun, M. Fazlyab, and S. Shahrampour, "On centralized and distributed mirror descent: Convergence analysis using quadratic constraints," *IEEE Transactions on Automatic Control*, 2022.
- [36] A. Sarlette and R. Sepulchre, "Consensus optimization on manifolds," *SIAM journal on Control and Optimization*, vol. 48, no. 1, pp. 56–76, 2009.
- [37] K. Deng and J. Hu, "Decentralized projected riemannian gradient method for smooth optimization on compact submanifolds," *arXiv preprint arXiv:2304.08241*, 2023.
- [38] L. Wang, L. Bao, and X. Liu, "A decentralized proximal gradient tracking algorithm for composite optimization on riemannian manifolds," *arXiv preprint arXiv:2401.11573*, 2024.
- [39] X. Wu, Z. Hu, and H. Huang, "Decentralized riemannian algorithm for nonconvex minimax problems," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 9, 2023, pp. 10370–10378.

- [40] J. Chen, H. Ye, M. Wang, T. Huang, G. Dai, I. Tsang, and Y. Liu, "Decentralized riemannian conjugate gradient method on the stiefel manifold," in *The Twelfth International Conference on Learning Representations*, 2023.
- [41] N. Xiao and X. Liu, "Solving optimization problems over the stiefel manifold by smooth exact penalty function," *arXiv preprint arXiv:2110.08986*, 2021.
- [42] L. Wang, N. Xiao, and X. Liu, "A double tracking method for optimization with decentralized generalized orthogonality constraints," *arXiv preprint arXiv:2409.04998*, 2024.
- [43] A. Edelman, T. A. Arias, and S. T. Smith, "The geometry of algorithms with orthogonality constraints," *SIAM journal on Matrix Analysis and Applications*, vol. 20, no. 2, pp. 303–353, 1998.
- [44] P. Ablin, S. Vary, B. Gao, and P.-A. Absil, "Infeasible deterministic, stochastic, and variance-reduction algorithms for optimization under orthogonality constraints," *arXiv preprint arXiv:2303.16510*, 2023.
- [45] Z. Wen and W. Yin, "A feasible method for optimization with orthogonality constraints," *Mathematical Programming*, vol. 142, no. 1-2, pp. 397–434, 2013.
- [46] H. Liu, A. M.-C. So, and W. Wu, "Quadratic optimization with orthogonality constraint: explicit tojasiewicz exponent and linear convergence of retraction-based line-search and stochastic variance-reduced gradient methods," *Mathematical Programming*, vol. 178, pp. 215–262, 2019.
- [47] Q. Rebjock and N. Boumal, "Fast convergence to non-isolated minima: four equivalent conditions for c^2 functions," *arXiv preprint arXiv:2303.00096*, 2023.
- [48] R. Xin, U. A. Khan, and S. Kar, "An improved convergence analysis for decentralized online stochastic non-convex optimization," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1842–1858, 2021.
- [49] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [50] 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.
- [51] X. Chen, Y. He, and Z. Zhang, "Tight error bounds for the sign-constrained stiefel manifold," 2024. [Online]. Available: <https://arxiv.org/abs/2210.05164>
- [52] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.



Youbang Sun is currently a Ph.D. candidate in the Department of Mechanical and Industrial Engineering at Northeastern University. His research interests include areas in machine learning and optimization with emphasis on distributed and multi-agent systems. He is interested in topics such as distributed optimization, Riemannian optimization, federated learning, and multi-agent reinforcement learning.



Shixiang Chen received the Ph.D. degree in Systems Engineering and Engineering Management from The Chinese University of Hong Kong in July, 2019. He is an assistant professor in the School of Mathematical Sciences, University of Science and Technology of China. He was a postdoctoral associate in the Department of Industrial & Systems Engineering at Texas A&M University. His current research interests include design and analysis of optimization algorithms, and their applications in machine learning and signal processing.



Alfredo Garcia received the Degree in electrical engineering from the Universidad de los Andes, Bogotá, Colombia, in 1991, the Diplôme d'Etudes Approfondies in automatic control from the Université Paul Sabatier, Toulouse, France, in 1992, and the Ph.D. degree in industrial and operations engineering from the University of Michigan, Ann Arbor, MI, USA, in 1997. From 1997 to 2001, he was a consultant to government agencies and private utilities in the electric power industry. From 2001 to 2015, he was a Faculty with the Department of Systems and Information Engineering, University of Virginia, Charlottesville, VA, USA. From 2015 to 2017, he was a Professor with the Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL, USA. In 2018, he joined the Department of Industrial and System Engineering, Texas A&M University, College Station, TX, USA. His research interests include game theory and dynamic optimization with applications in communications and energy networks.



Shahin Shahrampour received the Ph.D. degree in Electrical and Systems Engineering, the M.A. degree in Statistics (The Wharton School), and the M.S.E. degree in Electrical Engineering, all from the University of Pennsylvania, in 2015, 2014, and 2012, respectively. He is currently an Assistant Professor in the Department of Mechanical and Industrial Engineering at Northeastern University. His research interests include machine learning, optimization, sequential decision-making, and distributed learning, with a focus on developing computationally efficient methods for data analytics. He is a Senior Member of the IEEE.