




# Compressed Gradient Tracking for Decentralized Optimization Over General Directed Networks

Zhuoqing Song , Lei Shi, Shi Pu , and Ming Yan 

**Abstract**—In this paper, we propose two communication-efficient decentralized optimization algorithms over a general directed multi-agent network. The first algorithm, termed Compressed Push-Pull (CPP), combines the gradient tracking Push-Pull method with communication compression. We show that CPP is applicable to a general class of unbiased compression operators and achieves linear convergence rate for strongly convex and smooth objective functions. The second algorithm is a broadcast-like version of CPP (B-CPP), and it also achieves linear convergence rate under the same conditions on the objective functions. B-CPP can be applied in an asynchronous broadcast setting and further reduce communication costs compared to CPP. Numerical experiments complement the theoretical analysis and confirm the effectiveness of the proposed methods.

**Index Terms**—Decentralized optimization, compression, directed graphs, first-order methods, linear convergence.

## I. INTRODUCTION

IN THIS paper, we focus on solving the decentralized optimization problem, where a system of  $n$  agents, each having access to a private function  $f_i(x)$ , work collaboratively to obtain a consensual solution to the following problem:

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

where  $x$  is the global decision variable. The  $n$  agents are connected through a general directed network and only communicate directly with their immediate neighbors. The problem (1)

Manuscript received August 13, 2021; revised January 1, 2022 and February 14, 2022; accepted March 3, 2022. Date of publication March 17, 2022; date of current version April 20, 2022. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. David I. Shuman. The work of Zhuoqing Song was supported by the Shanghai Science and Technology Program under Grant 21JC1400600. The work of Lei Shi was supported in part by the Shanghai Science and Technology Program under Grants 20JC1412700 and 19JC1420101, and in part by the National Natural Science Foundation of China (NSFC) under Grant 12171093. The work of Shi Pu was supported in part by the Shenzhen Research Institute of Big Data under Grant J00120190011, and in part by the National Natural Science Foundation of China (NSFC) under Grant 62003287. The work of Ming Yan was supported by NSF under Grant DMS-2012439. (Corresponding authors: Shi Pu; Ming Yan.)

Zhuoqing Song is with the Shanghai Center for Mathematical Sciences, Fudan University, Shanghai 200437, China (e-mail: zqsong19@fudan.edu.cn).

Lei Shi is with the School of Mathematical Sciences, Shanghai Key Laboratory for Contemporary Applied Mathematics, Fudan University, Shanghai 200437, China (e-mail: leishi@fudan.edu.cn).

Shi Pu is with the School of Data Science, Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen 518172, China (e-mail: shipu.idda@gmail.com).

Ming Yan is with the Department of Computational Mathematics, Science and Engineering and Department of Mathematics, Michigan State University, East Lansing, MI 48864 USA (e-mail: myan@msu.edu).

Digital Object Identifier 10.1109/TSP.2022.3160238

has received much attention in recent years due to its wide applications in distributed machine learning [2]–[4], multi-agent target seeking [5], [6], and wireless networks [7]–[9], among many others. For example, the rapid development of distributed machine learning involves data whose size is getting increasingly large, and they are usually stored across multiple computing agents that are spatially distributed. Centering large amounts of data is often undesirable due to limited communication resources and/or privacy concerns, and decentralized optimization serves as an important tool to solve such large-scale distributed learning problems due to its scalability, sparse communication, and better protection for data privacy [10].

Many methods have been proposed to solve the problem (1) under various settings on the optimization objectives, network topologies, and communication protocols. The paper [11] developed a decentralized subgradient descent method (DGD) with diminishing stepsizes to reach the optimum for convex objective functions over an undirected network topology. Subsequently, decentralized optimization methods for undirected networks, or more generally, with doubly stochastic mixing matrices, have been extensively studied in the literature; see, e.g., [12]–[17]. Among these works, EXTRA [15] was the first method that achieves linear convergence for strongly convex and smooth objective functions under symmetric stochastic matrices. For directed networks, however, constructing a doubly stochastic mixing matrix usually requires a weight-balancing step, which could be costly when carried out in a distributed manner. Therefore, the push-sum technique [18] was utilized to overcome this issue. Specifically, the push-sum based subgradient method in [19] can be implemented over time-varying directed graphs, and linear convergence rates were achieved in [20], [21] for minimizing strongly convex and smooth objective functions by applying the push-sum technique to EXTRA.

Gradient tracking is an important technique that has been successfully applied in many decentralized optimization algorithms. Specifically, the methods proposed in [13], [22]–[24] employ gradient tracking to achieve linear convergence for strongly convex and smooth objective functions, where the work in [22]–[24] particularly considered combining gradient tracking with the push-sum technique to accommodate directed graphs. The methods can also be applied to time-varying graphs [22] and asynchronous settings [23]. The Push-Pull/AB method introduced in [25], [26] modified the gradient tracking methods to deal with directed network topologies without the push-sum technique. The algorithm uses a row stochastic matrix to mix the local decision variables and a column stochastic matrix

to mix the auxiliary variables that track the average gradients over the network. It can unify different network architectures, including peer-to-peer, master-slave, and leader-follower architectures [25]. For minimizing strongly convex and smooth objectives, the Push-Pull/ $\mathcal{AB}$  method not only enjoys linear convergence over fixed graphs [25], [26], but also works well under time-varying graphs and asynchronous settings [25], [27], [28].

In decentralized optimization, efficient communication is critical for enhancing algorithm performance and system scalability. One major approach to reduce communication costs is considering communication compression, which is essential especially under limited communication bandwidth. Recently, several compression methods have been proposed for distributed and federated learning, including [29]–[41]. Recent works have tried to combine the communication compression methods with decentralized optimization. The existence of compression errors may result in inferior convergence performance compared to uncompressed or centralized algorithms. For example, the methods considered by [42]–[47] can only guarantee to reach a neighborhood of the desired solutions when the compression errors exist. QDGD [48] achieves a vanishing mean solution error with a slower rate than the centralized method. To reduce the error from compression, some works [49]–[51] increase compression accuracy as the iteration grows to guarantee the convergence. However, they still need high communication costs to get highly accurate solutions. Techniques to remedy this increased communication costs include gradient difference compression [35], [52], [53] and error compensation [38], [54], [55], which enjoy better performance than direct compression. In [56], the difference compression (DCD-PSGD) and extrapolation compression (ECD-PSGD) algorithms were proposed to reach the same convergence rate of the centralized schemes with additional requirements on the compression ratio. In [55], a quantized decentralized SGD (CHOCO-SGD) method was proposed and shown to converge to the optimal solution or a stationary point at a comparable rate as the centralized SGD method. The works [54] and [57] also achieve comparable convergence rates with the centralized scheme for solving non-convex problems.

For strongly convex and smooth objective functions, [58] first considered a linearly convergent gradient tracking method based on a specific quantizer. More recently, the paper [53] introduced LEAD that works with a general class of compression operators and still enjoys linear convergence. Some recent developments can be found in [59]–[61], where [60] particularly combined Push-Pull/ $\mathcal{AB}$  with a special quantizer to achieve linear convergence over directed graphs.

In this paper, we consider decentralized optimization over general directed networks and propose a novel Compressed Push-Pull method (CPP) that combines Push-Pull/ $\mathcal{AB}$  with a general class of unbiased compression operators. CPP enjoys large flexibility in both the compression method and the network topology. We show CPP achieves linear convergence rate under strongly convex and smooth objective functions.

Broadcast or gossip-based communication protocols are popular choices for distributed computation due to their low communication costs [62]–[64]. In the second part of this paper, we propose a broadcast-like CPP algorithm (B-CPP) that allows for

asynchronous updates of the agents: at every iteration of the algorithm, only a subset of the agents wake up to perform prescribed updates. Thus, B-CPP is more flexible, and due to its broadcast nature, it can further save communication over CPP in certain scenarios [64]. We show that B-CPP also achieves linear convergence for minimizing strongly convex and smooth objectives.

The main contributions of this paper are summarized as follows.

- We propose CPP – a novel decentralized optimization method with communication compression. The method works under a general class of compression operators and is shown to achieve linear convergence for strongly convex and smooth objective functions over general directed graphs. To the best of our knowledge, CPP is the first method that enjoys linear convergence under such a general setting.
- We consider an asynchronous broadcast version of CPP (B-CPP). B-CPP further reduces the communicated data per iteration and is also provably linearly convergent over directed graphs for minimizing strongly convex and smooth objective functions. Numerical experiments demonstrate the advantages of B-CPP in saving communication costs.

The rest of this paper is organized as follows. We provide necessary notation and assumptions in Section II. CPP is introduced and analyzed in Section III. In Section IV, we consider the algorithm B-CPP. Numerical examples are presented in Section V, and we conclude the paper in Section VI.

## II. NOTATION AND ASSUMPTIONS

Denote the set of agents as  $\mathcal{N} = \{1, 2, \dots, n\}$ . At iteration  $k$ , each agent  $i$  has a local estimation  $\mathbf{x}_i^k \in \mathbb{R}^p$  of the global decision variable and an auxiliary variable  $\mathbf{y}_i^k$ . We use lowercase bold letters to denote the local variables, and their counterpart uppercase bold letters denote the concatenation of these local variables. For instance,  $\mathbf{X}^k$ ,  $\nabla \mathbf{F}(\mathbf{X}^k)$  are the concatenation of  $\{\mathbf{x}_i^k\}_{i \in \mathcal{N}}$ ,  $\{\nabla f_i(\mathbf{x}_i^k)\}_{i \in \mathcal{N}}$ , respectively, and their connections are

$$\mathbf{X}^k = [\mathbf{x}_1^k, \dots, \mathbf{x}_n^k]^\top \in \mathbb{R}^{n \times p},$$

$$\nabla \mathbf{F}(\mathbf{X}^k) = [\nabla f_1(\mathbf{x}_1^k), \dots, \nabla f_n(\mathbf{x}_n^k)]^\top \in \mathbb{R}^{n \times p}.$$

*Assumption 1:* Each  $f_i$  is  $\mu$ -strongly convex ( $\mu > 0$ ) and  $L$ -smooth, i.e., for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p$ ,

$$\langle \mathbf{x}_1 - \mathbf{x}_2, \nabla f_i(\mathbf{x}_1) - \nabla f_i(\mathbf{x}_2) \rangle \geq \mu \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2,$$

$$\|\nabla f_i(\mathbf{x}_1) - \nabla f_i(\mathbf{x}_2)\|_2 \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|_2.$$

Since all  $f_i(\mathbf{x})$  are strongly convex,  $f(\mathbf{x})$  admits a unique minimizer  $\mathbf{x}^*$ . Denote  $\mathbf{x}^* = \mathbf{1} \mathbf{X}^{* \top}$ , where  $\mathbf{1}$  is the all-ones column vector.

Given any nonnegative matrix  $\mathbf{M}$ , we denote by  $\mathcal{G}_{\mathbf{M}}$  the induced graph by  $\mathbf{M}$ . The sets  $\mathcal{N}_{\mathbf{M},i}^- \stackrel{\text{def}}{=} \{j \in \mathcal{N} : \mathbf{M}_{ij} > 0\}$  and  $\mathcal{N}_{\mathbf{M},i}^+ \stackrel{\text{def}}{=} \{j \in \mathcal{N} : \mathbf{M}_{ji} > 0\}$  are called the in-neighbors and out-neighbors of agent  $i$ .

The communication between all the agents is modeled by directed graphs. Given a strongly connected graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  with  $\mathcal{E} \subset \mathcal{N} \times \mathcal{N}$  being the edge set, agent  $i$  can receive information from agent  $j$  if and only if  $(j, i) \in \mathcal{E}$ . There are

two  $n$ -by- $n$  nonnegative matrices  $\mathbf{R}$  and  $\mathbf{C}$ . A spanning tree  $\mathcal{T}$  rooted at some  $i \in \mathcal{N}$  in  $\mathcal{G}_{\mathbf{R}}$  is a subgraph of  $\mathcal{G}_{\mathbf{R}}$  containing  $n - 1$  edges, and each node except  $i$  can be connected to  $i$  by a path in  $\mathcal{T}$ . Let  $\mathcal{R}_{\mathbf{R}}, \mathcal{R}_{\mathbf{C}^\top}$  denote the roots of the spanning trees in  $\mathcal{G}_{\mathbf{R}}$  and  $\mathcal{G}_{\mathbf{C}^\top}$ . We have the following assumption on  $\mathbf{R}$  and  $\mathbf{C}$ .

**Assumption 2:** The matrix  $\mathbf{R}$  is supported by  $\mathcal{G}$ , i.e.,  $\mathcal{E}_{\mathbf{R}} = \{(i, j) \in \mathcal{N} \times \mathcal{N} \mid \mathbf{R}_{ij} > 0\} \subset \mathcal{E}$ , and  $\mathbf{R}$  is a row stochastic matrix, i.e.,  $\mathbf{R}\mathbf{1} = \mathbf{1}$ . The matrix  $\mathbf{C}$  is also supported by  $\mathcal{G}$ , and  $\mathbf{C}$  is a column stochastic matrix, i.e.,  $\mathbf{1}^\top \mathbf{C} = \mathbf{1}^\top$ . In addition,  $\mathcal{R}_{\mathbf{R}} \cap \mathcal{R}_{\mathbf{C}^\top} \neq \emptyset$ .

By [25, Lemma 1], Assumption 2 implies the following facts:  $\mathbf{R}$  has a unique left eigenvector  $\mathbf{s}_{\mathbf{R}}$  with respect to  $\mathbf{1}$  such that  $\mathbf{s}_{\mathbf{R}}^\top \mathbf{1} = n$ .  $\mathbf{C}$  has a unique right eigenvector  $\mathbf{s}_{\mathbf{C}}$  with respect to  $\mathbf{1}$  such that  $\mathbf{s}_{\mathbf{C}}^\top \mathbf{1} = n$ . In addition, the entries of  $\mathbf{s}_{\mathbf{R}}$  and  $\mathbf{s}_{\mathbf{C}}$  are all nonnegative. All nonzero entries of  $\mathbf{s}_{\mathbf{R}}, \mathbf{s}_{\mathbf{C}}$  correspond to  $\mathcal{R}_{\mathbf{R}}, \mathcal{R}_{\mathbf{C}^\top}$ , respectively. Because  $\mathcal{R}_{\mathbf{R}} \cap \mathcal{R}_{\mathbf{C}^\top} \neq \emptyset$ , we have  $\mathbf{s}_{\mathbf{R}}^\top \mathbf{s}_{\mathbf{C}} > 0$ .

We denote the spectral radius of matrix  $\mathbf{A}$  as  $\rho(\mathbf{A})$ . The inner product of two matrices is defined as  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{tr}(\mathbf{X}^\top \mathbf{Y})$ , and the Frobenius norm is  $\|\mathbf{X}\|_{\text{F}} = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle}$ . Given a vector  $\mathbf{d}$ ,  $\mathbf{d}_{a:b}$  is the subvector of  $\mathbf{d}$  containing the entries indexed from  $a$  to  $b$ . Given a matrix  $\mathbf{A}$ , the notion  $\mathbf{A}_{a:b,c:d}$  denotes the submatrix containing the entries with row index in  $[a, b]$  and column index in  $[c, d]$ . We abbreviate “ $1 : n$ ” by “ $:$ ” and “ $i : i$ ” is abbreviated as “ $i$ ”. Especially in our notations,  $\mathbf{A}_{i,:}$  and  $\mathbf{A}_{:,j}$  denote the  $i$ -th row and  $j$ -th column of  $\mathbf{A}$ , respectively. For vector  $\mathbf{v} \in \mathbb{R}^m$ ,  $\text{Diag}(\mathbf{v})$  is an  $m$ -by- $m$  diagonal matrix with  $\text{Diag}(\mathbf{v})_{ii} = v_i$ .

**Definition 1:** Given a vector norm  $\|\cdot\|_*$ , we define the corresponding matrix norm on an  $n \times p$  matrix  $\mathbf{A}$  as

$$\|\|\mathbf{A}\|\|_* = \left\| \left[ \|\mathbf{A}_{:,1}\|_*, \|\mathbf{A}_{:,2}\|_*, \dots, \|\mathbf{A}_{:,p}\|_* \right] \right\|_2.$$

When  $\|\cdot\|_* = \|\cdot\|_2$ , we have  $\|\|\mathbf{A}\|\|_2 = \|\mathbf{A}\|_{\text{F}}$ , the Frobenius norm.

**Definition 2:** Given a vector norm  $\|\cdot\|_*$ , we define the corresponding induced norm on an  $n \times n$  matrix  $\mathbf{A}$  as  $\|\mathbf{A}\|_* = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|_*}{\|\mathbf{x}\|_*}$ .

**Lemma 1 (Lemma 5 in [25]):** For any matrices  $\mathbf{A} \in \mathbb{R}^{n \times p}$ ,  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , and a vector norm  $\|\cdot\|_*$ , we have  $\|\|\mathbf{W}\mathbf{A}\|\|_* \leq \|\mathbf{W}\|_* \|\|\mathbf{A}\|\|_*$ . For any vectors  $\mathbf{a} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^p$ ,  $\|\|\mathbf{a}\mathbf{b}^\top\|\|_* \leq \|\mathbf{a}\|_* \|\mathbf{b}\|_2$ .

The compression in this paper is denoted by  $\text{COMPRESS}(\cdot)$ . For any matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , we denote  $\text{COMPRESS}(\mathbf{X})$  as an  $n$ -by- $p$  matrix with the  $i$ -th row being  $\text{COMPRESS}(\mathbf{X}_{i,:})$ .

**Assumption 3:** The compression is unbiased, i.e., given  $\mathbf{x} \in \mathbb{R}^p$  and  $\hat{\mathbf{x}} = \text{COMPRESS}(\mathbf{x})$ , there exists  $C_2 > 0$  such that  $\mathbb{E}[\hat{\mathbf{x}}|\mathbf{x}] = \mathbf{x}$ , and  $\mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2 | \mathbf{x}] \leq C_2 \|\mathbf{x}\|_2^2$ . And the random variables generated inside the procedure  $\text{COMPRESS}(\mathbf{x})$  depends on  $\mathbf{x}$  only.

### III. A PUSH-PULL METHOD WITH COMPRESSION

In this section, we propose a Push-Pull method with Compression (CPP) in Algorithm 1.

---

**Algorithm 1:** Compressed Push-Pull (in agent  $i$ 's perspective).

---

**Input:** initial position  $\mathbf{x}_i^0$ , stepsize  $\alpha_i$ , averaging parameters  $\beta, \gamma, \eta \in (0, 1]$ , network information  $\mathbf{R}_{ij} (j \in \mathcal{N}_{\mathbf{R},i}^-)$ ,  $\mathbf{C}_{ij} (j \in \mathcal{N}_{\mathbf{C},i}^-)$  and the sets  $\mathcal{N}_{\mathbf{R},i}^+, \mathcal{N}_{\mathbf{C},i}^+$ , total iteration number  $K$

**Output:**  $\mathbf{x}_i^K$

Initiate  $\mathbf{y}_i^0 = \nabla f_i(\mathbf{x}_i^0)$ ,  $\mathbf{u}_i^0 = \mathbf{u}_{\mathbf{R},i}^0 = \mathbf{0}$ .

**for**  $k = 0$  **to**  $K - 1$  **do**

    Call procedure PULL:

$$\left( \hat{\mathbf{x}}_{\mathbf{R},i}^k, \mathbf{u}_i^{k+1}, \mathbf{u}_{\mathbf{R},i}^{k+1} \right) =$$

$$\text{PULL}(\mathbf{x}_i^k, \mathbf{u}_i^k, \mathbf{u}_{\mathbf{R},i}^k).$$

$$\text{Update } \mathbf{x}_i^{k+1} = (1 - \beta)\mathbf{x}_i^k + \beta\hat{\mathbf{x}}_{\mathbf{R},i}^k - \alpha_i \mathbf{y}_i^k.$$

    Call procedure PUSH:

$$\left( \hat{\mathbf{y}}_i^k, \hat{\mathbf{y}}_{\mathbf{C},i}^k \right) = \text{PUSH}(\mathbf{y}_i^k).$$

$$\text{Compute local gradient } \nabla f_i(\mathbf{x}_i^{k+1}).$$

$$\text{Update } \mathbf{y}_i^{k+1} =$$

$$\mathbf{y}_i^k + \gamma \hat{\mathbf{y}}_{\mathbf{C},i}^k - \gamma \hat{\mathbf{y}}_i^k + \nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k).$$

**end**

**procedure** PULL( $\mathbf{x}_i, \mathbf{u}_i, \mathbf{u}_{\mathbf{R},i}$ )

$$\mathbf{p}_i = \text{COMPRESS}(\mathbf{x}_i - \mathbf{u}_i).$$

$$\mathbf{p}_{\mathbf{R},i} = \sum_{j \in \mathcal{N}_{\mathbf{R},i}^-} \mathbf{R}_{ij} \mathbf{p}_j.$$

$$\hat{\mathbf{x}}_i = \mathbf{u}_i + \mathbf{p}_i.$$

$$\hat{\mathbf{x}}_{\mathbf{R},i} = \mathbf{u}_{\mathbf{R},i} + \mathbf{p}_{\mathbf{R},i}.$$

$$\mathbf{u}_i \leftarrow (1 - \eta) \mathbf{u}_i + \eta \hat{\mathbf{x}}_i.$$

$$\mathbf{u}_{\mathbf{R},i} \leftarrow (1 - \eta) \mathbf{u}_{\mathbf{R},i} + \eta \hat{\mathbf{x}}_{\mathbf{R},i}.$$

**Return:**  $\hat{\mathbf{x}}_{\mathbf{R},i}, \mathbf{u}_i, \mathbf{u}_{\mathbf{R},i}$ .

**end procedure**

**procedure** PUSH( $\mathbf{y}_i$ )

$$\hat{\mathbf{y}}_i = \text{COMPRESS}(\mathbf{y}_i).$$

$$\hat{\mathbf{y}}_{\mathbf{C},i} = \sum_{j \in \mathcal{N}_{\mathbf{C},i}^-} \mathbf{C}_{ij} \hat{\mathbf{y}}_j.$$

**Return:**  $\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_{\mathbf{C},i}$ .

**end procedure**

---

We start from discussing the following scheme of Push-Pull/AB from the viewpoint of agent  $i$  [25], [26]:

$$\begin{cases} \mathbf{x}_i^{k+1} = \sum_{j \in \mathcal{N}_{\mathbf{R},i}^-} \mathbf{R}_{ij} \mathbf{x}_j^k - \alpha_i \mathbf{y}_i^k \\ \mathbf{y}_i^{k+1} = \sum_{j \in \mathcal{N}_{\mathbf{C},i}^-} \mathbf{C}_{ij} \mathbf{y}_j^k + \nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k) \end{cases}$$

and  $\mathbf{y}_i^0 = \nabla f_i(\mathbf{x}_i^0)$ . At each iteration, agent  $i$  computes  $\sum_{j \in \mathcal{N}_{\mathbf{R},i}^-} \mathbf{R}_{ij} \mathbf{x}_j^k$  to average the local decision variables received from its neighbors. As each agent takes such a step, the local decision variables tend to get closer with each other and eventually reach consensus when  $\mathbf{y}_i^k$  goes to zero. In the next “gradient tracking” step, by adding the gradient difference term  $\nabla f_i(\mathbf{x}_i^{k+1}) - \nabla f_i(\mathbf{x}_i^k)$  into  $\mathbf{y}_i^{k+1}$  and considering that  $\mathbf{y}_i^0 = \nabla f_i(\mathbf{x}_i^0)$  and  $\mathbf{1}^\top \mathbf{C} = \mathbf{1}^\top$ , we have  $\sum_{i \in \mathcal{N}} \mathbf{y}_i^k = \sum_{i \in \mathcal{N}} \nabla f_i(\mathbf{x}_i^k)$  by induction, which indicates that  $\{\mathbf{y}_i^k\}_{i \in \mathcal{N}}$  can help track the average gradients over the network.



The Compressed Push-Pull (CPP) method differs from Push-Pull/AB mainly in the averaging step which requires communication. To alleviate the impact of compression errors, we use a “damped” averaging step and replace the exact local variables by their communication-efficient version, i.e.,  $(1 - \beta)x_i^k + \beta \sum_{j \in \mathcal{N}_{R,i}^-} R_{ij} \hat{x}_j^k$ . For  $y_i^k$ , we take a slightly different averaging step so that the relation  $\sum_{i \in \mathcal{N}} y_i^k = \sum_{i \in \mathcal{N}} \nabla f_i(x_i^k)$  is preserved and the compression errors can also be bounded “safely”.

However, sending  $\hat{x}_i^k$  itself can still be expensive. Therefore, we call the procedure PULL to save the communication costs. In this procedure,  $(u_i, u_{R,i})$  is used to track  $(x_i, \sum_{j \in \mathcal{N}_{R,i}^-} R_{ij} x_j)$ . The compression and the communication are applied on the difference  $(x_i - u_i)$  and its compressed version, respectively. It can be derived from the relation  $p_{R,i} = \sum_{j \in \mathcal{N}_{R,i}^-} R_{ij} p_j$  and induction that in each call to PULL,  $u_{R,i} = \sum_{j \in \mathcal{N}_{R,i}^-} R_{ij} u_j$ . And the compression error could be very small if the difference  $x_i - u_i$  is small. With this observation, we have  $\hat{x}_{R,i} = \sum_{j \in \mathcal{N}_{R,i}^-} R_{ij} \hat{x}_j$  with  $\hat{x}_j$  being an unbiased approximation of  $x_j$ , whose variance converges to 0.

In the PUSH procedure, since the variable  $y_i^k$  converges to zero as we will show later, we can simply compress it and estimate  $\sum_{j \in \mathcal{N}_{C,i}^-} C_{ij} y_j^k$  using the compressed values. Similarly, we have that in each call to PUSH,  $\hat{y}_{C,i} = \sum_{j \in \mathcal{N}_{C,i}^-} C_{ij} \hat{y}_j$ , with  $\hat{y}_j$  being the unbiased compression of  $y_j$ , whose variance converges to 0 as the input  $y_i^k$  converges to 0.

Let  $\hat{\alpha} = \max_{1 \leq i \leq n} \alpha_i$ , and define  $\alpha$  as a diagonal matrix with  $\alpha_{ii} = \alpha_i$ . Then, we can rewrite Algorithm 1 into the following matrix form

$$\begin{cases} \hat{X}^k = U^k + \text{COMPRESS}(X^k - U^k), \\ \hat{Y}^k = \text{COMPRESS}(Y^k), \\ x^{k+1} = (1 - \beta)x^k + \beta R \hat{X}^k - \alpha Y^k, \\ Y^{k+1} = Y^k + \gamma(C - I)\hat{Y}^k \\ \quad + \nabla F(X^{k+1}) - \nabla F(X^k), \\ U^{k+1} = (1 - \eta)U^k + \eta \hat{X}^k. \end{cases} \quad (2)$$

Comparing to the standard Push-Pull method in [25], the procedures PUSH and PULL save communication costs at the cost of error in the mixing. Note  $\hat{X}^k$  and  $\hat{Y}^k$  are approximations for  $X^k$  and  $Y^k$ , respectively. The difference between  $X^k$  (or  $Y^k$ ) and  $\hat{X}^k$  (or  $\hat{Y}^k$ ) are induced by the compression errors. Let us denote the approximation error as  $w_i^k = x_i^k - \hat{x}_i^k$  and  $e_i^k = y_i^k - \hat{y}_i^k$ . We further write them into the following matrix form

$$W^k = X^k - \hat{X}^k, E^k = Y^k - \hat{Y}^k. \quad (3)$$

From the PULL procedure, we can see that  $W^k$  is also the compression error for  $X^k - U^k$ . The update of  $U^{k+1} = (1 - \eta)U^k + \eta \hat{X}^k$  indicates that  $U^k$  is tracking the motions of  $X^k$ . As  $U^k$  approaches  $X^k$ , by Assumption 3, the variance of the approximation error  $W^k$  will also tends to 0.

Since  $Y^0 = \nabla F(X^0)$  and  $1^\top C = 1^\top$ , it follows from (5d) and induction that

$$1^\top Y^k = 1^\top \nabla F(X^k). \quad (4)$$

Then, as the local variables become closer to each other,  $y_i^k$  are more and more parallel to  $1^\top \nabla F(X^k)$ . As  $X^k$  approaches the optimal point,  $1^\top \nabla F(X^k)$  tends to  $1^\top \nabla F(X^*) = 0^\top$ . This indicates that  $y^k$  also tends to 0. Then, by Assumption 3, the variance of compression error  $E^k$  will also tend to 0. As the compression errors  $W^k, E^k$  tend to 0, we will have  $\hat{X}^k \approx X^k$  and  $\hat{Y}^k \approx Y^k$ . We remark that if  $\hat{X}^k = X^k, \hat{Y}^k = Y^k$ , Algorithm 1 will reduce to the PUSH-PULL method in [25].

To show the convergence of the proposed compressed algorithm, we define  $R_\beta = (1 - \beta)I + \beta R$ ,  $C_\gamma = (1 - \gamma)I + \gamma C$ , and rewrite (2) as

$$\begin{cases} W^k = X^k - U^k - \text{COMPRESS}(X^k - U^k), \end{cases} \quad (5a)$$

$$\begin{cases} E^k = Y^k - \text{COMPRESS}(Y^k), \end{cases} \quad (5b)$$

$$\begin{cases} X^{k+1} = R_\beta X^k - \alpha Y^k - \beta R W^k, \end{cases} \quad (5c)$$

$$\begin{cases} Y^{k+1} = C_\gamma Y^k + \nabla F(X^{k+1}) - \nabla F(X^k) \\ \quad + \gamma(I - C)E^k, \end{cases} \quad (5d)$$

$$\begin{cases} U^{k+1} = (1 - \eta)U^k + \eta X^k - \eta W^k, \end{cases} \quad (5e)$$

We define the value  $\bar{\alpha} = s_R^\top \alpha s_C / n$ . Under Assumption 2,  $\bar{\alpha} > 0$  can be satisfied by setting  $\alpha_i > 0$  for at least one  $i \in \mathcal{R}_R \cap \mathcal{R}_{C^\top}$ . When  $\bar{\alpha} > 0$ , let us define a constant  $w \geq \hat{\alpha} / \bar{\alpha} \geq \frac{n}{s_R s_C}$ . Note that the analysis in the rest of this section holds true for arbitrary  $w \geq \hat{\alpha} / \bar{\alpha}$ . As we will choose  $w \geq \hat{\alpha} / \bar{\alpha}$  as an auxiliary parameter to help choose proper stepsizes  $\{\alpha_i\}_{i \in \mathcal{N}}$  in Theorem 7 which is the main theorem of this section, we use the same  $w$  in the analysis below.

Let  $\mathcal{F}_k$  denote the  $\sigma$ -field generated by  $\{E^j, W^j\}_{j=0}^{k-1}$ , and  $\mathcal{F}_k^+$  is the  $\sigma$ -field generated by  $\{E^j, W^j\}_{j=0}^{k-1} \cup \{W^k\}$ . Thus, we have  $\mathcal{F}_1 \subset \mathcal{F}_1^+ \subset \mathcal{F}_2 \subset \mathcal{F}_2^+ \subset \dots \subset \mathcal{F}_k \subset \mathcal{F}_k^+ \subset \dots$ . By Lemma 2 below,  $\mathcal{F}_k = \sigma(X^i, Y^i, U^i : 0 \leq i < k)$  and  $\mathcal{F}_k^+ = \sigma(X^i, U^i, Y^j : 0 \leq i \leq k, 0 \leq j < k)$ .

The following lemma comes from (5c)–(5e) and Assumption 3.

**Lemma 2:** The variables  $X^k, Y^k, U^k, V^k$  are measurable with respect to  $\mathcal{F}_k$ , and  $X^k$  is measurable with respect to  $\mathcal{F}_{k-1}^+$ . Moreover, we have

$$\mathbb{E}[W^k | \mathcal{F}_k] = \mathbb{E}[E^k | \mathcal{F}_k^+] = 0. \quad (6)$$

*Proof:* By expanding (5c)–(5e) recursively,  $X^k, Y^k, U^k$ , and  $V^k$  can be represented by linear combinations of  $X^0, Y^0, U^0, V^0$  and random variables  $\{E^j, W^j\}_{j=0}^{k-1}$ , i.e.,  $X^k, Y^k, U^k, V^k$  are measurable with respect to  $\mathcal{F}_k$ . By (5c),  $X^k$  is a linear combination of  $W^{k-1}$  and  $\mathcal{F}_{k-1}$ -measurable variables  $X^{k-1}, Y^{k-1}$ , thus,  $X^k$  is measurable with respect to  $\mathcal{F}_{k-1}^+$ . And we obtain (6) directly from Assumption 3. ■

#### A. Convergence Analysis for CPP

In this section, we analyze the convergence rates of Algorithm 1. To begin with, we define the averages of  $X^k, Y^k$  as follows,

$$\bar{x}^k = \frac{1}{n} s_R^\top X^k, \bar{y}^k = \frac{1}{n} 1^\top Y^k. \quad (7)$$

We define two matrices  $\Pi_R = I - \frac{1s_R^\top}{n}$ ,  $\Pi_C = I - \frac{s_C 1^\top}{n}$ . It follows from direct calculation that  $\Pi_R R = R \Pi_R = R - \frac{1s_R^\top}{n}$ ,  $\Pi_C C = C \Pi_C = C - \frac{s_C 1^\top}{n}$ ,  $\Pi_R X^k = X^k - 1\bar{x}_k$ ,  $\Pi_C Y^k = Y^k - s_C \bar{Y}_k$ ,  $(R - I)\Pi_R = R - I$ , and  $\Pi_C (I - C) = I - C$ .

In the analysis below,  $\mathbb{E} [\|\bar{x}^k - x^*\|_2^2]$  is employed to measure the closeness to the optimal point;  $\mathbb{E} [\|\Pi_R X^k\|_R^2]$  and  $\mathbb{E} [\|\Pi_C Y^k\|_C^2]$  measure the consensus error and the gradient tracking error, respectively. To bound the compression errors, we use  $\mathbb{E} [\|U^k - X^k\|_2^2]$  to measure the convergence of the momentums. The matrix norms  $\|\cdot\|_R$  and  $\|\cdot\|_C$  are defined in Lemma 3.

Compared with the convergence analysis for Push-Pull, the analysis for CPP additionally requires dealing with the compression errors and establishing the relationship between the error terms of different types. Moreover, another term  $\mathbb{E} [\|Y^k\|_2^2]$  is considered to simplify the proof, and its role will be made clear in the follow-up analysis.

Now, we have the following expansion by multiplying  $\frac{s_R^\top}{n}$  on both sides of (5c)

$$\begin{aligned} \bar{x}^{k+1} &= \frac{1}{n} s_R^\top (R_\beta X^k - \alpha Y^k - \beta R W^k) \\ &= \bar{x}^k - \frac{1}{n} s_R^\top \alpha (Y^k - s_C \bar{Y}^k + s_C \bar{Y}^k) - \frac{\beta}{n} s_R^\top W^k \\ &= \bar{x}^k - \bar{\alpha} \bar{Y}^k - \frac{1}{n} s_R^\top \alpha \Pi_C Y^k - \frac{\beta}{n} s_R^\top W^k \\ &= \bar{x}^k - \bar{\alpha} g^k + \bar{\alpha} (g^k - \bar{Y}^k) - \frac{1}{n} s_R^\top \alpha \Pi_C Y^k \\ &\quad - \frac{\beta}{n} s_R^\top W^k, \end{aligned} \quad (8)$$

where  $g^k = \nabla f(\bar{x}^k)$ .

Multiplying  $\Pi_R$  on both sides of (5c), we have

$$\begin{aligned} \Pi_R X^{k+1} &= \Pi_R R_\beta \Pi_R X^k - \Pi_R \alpha Y^k - \beta \Pi_R R W^k, \end{aligned} \quad (9)$$

Multiplying  $\Pi_C$  on both sides of (5d), we obtain

$$\begin{aligned} \Pi_C Y^{k+1} &= \Pi_C C_\gamma \Pi_C Y^k + \gamma (I - C) E^k \\ &\quad + \Pi_C (\nabla F(X^{k+1}) - \nabla F(X^k)). \end{aligned} \quad (10)$$

**Lemma 3:** There are invertible matrices  $\tilde{R}, \tilde{C} \in \mathbb{R}^{n \times n}$  inducing vector norms  $\|x\|_R \stackrel{\text{def}}{=} \|\tilde{R}x\|_2$ ,  $\|x\|_C \stackrel{\text{def}}{=} \|\tilde{C}x\|_2$  for  $x \in \mathbb{R}^n$ . The corresponding matrix norms  $\|\cdot\|_R, \|\cdot\|_C$  defined by Definition 2 satisfy: for any  $\beta, \gamma \in [0, 1]$ ,

$$\|\Pi_R R_\beta\|_R \leq 1 - \theta_R \beta, \|\Pi_C C_\gamma\|_C \leq 1 - \theta_C \gamma, \quad (11)$$

where  $\theta_R, \theta_C$  are constants in  $(0, 1]$ . Especially,  $\|\Pi_R\|_R = \|\Pi_C\|_C = 1$  and  $\|R\|_R \leq 2, \|R - I\|_R \leq 3, \|C\|_C \leq 2, \|C - I\|_C \leq 3$ . Additionally,  $\|x\|_2 \leq \|x\|_R, \|x\|_2 \leq \|x\|_C, \forall x \in \mathbb{R}^n$ ;  $\|A\|_2 \leq \|A\|_R, \|A\|_2 \leq \|A\|_C, \forall A \in \mathbb{R}^{n \times p}$ . There exist constants  $\delta_{R,2}, \delta_{C,2}$  such that  $\|x\|_R \leq \delta_{R,2} \|x\|_2, \|x\|_C \leq \delta_{C,2} \|x\|_2, \forall x \in \mathbb{R}^n$ . And for any diagonal matrix  $V, \|V\|_R = \|V\|_2, \|V\|_C = \|V\|_2$ .

*Proof:* See Lemma 3 of [1].  $\blacksquare$

In the following,  $\|\cdot\|_R, \|\cdot\|_C, \|\cdot\|_R, \|\cdot\|_C$  and  $\|\cdot\|_C$  refer to the norms defined in the above lemma.

**Lemma 4:** For any vector norm  $\|\cdot\|_*$  induced by the inner product  $\langle \cdot, \cdot \rangle_*$ , and for any  $\theta > 1$ , we have

$$\|A + B\|_*^2 \leq \theta \|A\|_*^2 + \frac{\theta}{\theta - 1} \|B\|_*^2, \forall A, B \in \mathbb{R}^{n \times p}.$$

Especially,  $\|\cdot\|_*$  can be taken as  $\|\cdot\|_2, \|\cdot\|_R, \|\cdot\|_C$ .

*Proof:* By Definition 1, it suffices to prove  $\|a + b\|_*^2 \leq \theta \|a\|_*^2 + \frac{\theta}{\theta - 1} \|b\|_*^2$  for any vectors  $a, b \in \mathbb{R}^n$ . And it can be verified by the elementary inequality  $2\langle a, b \rangle_* \leq (\theta - 1)\|a\|_*^2 + \frac{1}{\theta - 1}\|b\|_*^2$ . The vector norms  $\|\cdot\|_R, \|\cdot\|_C$  are induced by inner products  $\langle \tilde{R}x, \tilde{R}y \rangle_*$ ,  $\langle \tilde{C}x, \tilde{C}y \rangle_*$ , respectively.  $\blacksquare$

The next lemma uses inequalities to derive recursive bounds for the quantities  $\mathbb{E} [\|\bar{x}^k - x^*\|_2^2], \mathbb{E} [\|\Pi_R X^k\|_R^2], \mathbb{E} [\|\Pi_C Y^k\|_C^2]$  and  $\mathbb{E} [\|U^k - X^k\|_2^2]$ . The quantity  $\mathbb{E} [\|Y^k\|_2^2]$  is introduced to help simplify the proof.

**Lemma 5:** For  $k \geq 0$ , take  $d^k \in \mathbb{R}^5$  as

$$d^k = \mathbb{E} \left[ \left( \|\bar{x}^k - x^*\|_2^2, \|\Pi_R X^k\|_R^2, \|\Pi_C Y^k\|_C^2, \|U^k - X^k\|_2^2, \|Y^k\|_2^2 \right) \right]^\top.$$

Define a parameterized matrix  $A(\hat{\alpha}, \beta, \gamma, \eta)$  as follows

$$\begin{aligned} A_{:,1:3} &= \begin{pmatrix} 1 - c_1 \hat{\alpha} & c_2 \hat{\alpha} & c_3 \hat{\alpha} \\ 0 & 1 - \theta_R \beta & 0 \\ 0 & \frac{8nc_7\beta^2}{\gamma} & 1 - \theta_C \gamma \\ 0 & \frac{8n\beta^2}{\eta} & 0 \\ c_{10} & c_9 & 3 \end{pmatrix}, \\ A_{:,4:5} &= \begin{pmatrix} c_4 \beta^2 & 0 \\ c_6 \beta^2 & c_5 \frac{\hat{\alpha}^2}{\beta} \\ \frac{C_2 nc_7 \beta^2}{\gamma} & c_8 \gamma^2 + \frac{2c_7 \hat{\alpha}^2}{\gamma} \\ A_{44} & \frac{2\hat{\alpha}^2}{\eta} \\ 0 & 0 \end{pmatrix}, \end{aligned} \quad (12)$$

where  $A_{44} = 1 - \eta + 2C_2\eta^2 + 2nC_2\beta^2$  and  $c_1 - c_{10}$  are constants given by (34), (36), (39) and (44) in the Supplementary Material of [1]. Then, under Assumptions 1-3 and the condition  $\bar{\alpha} \leq \frac{2}{\mu + L}$ , we have the following inequalities:

$$\mathbb{E} [\|Y^k\|_2^2] \leq A(\hat{\alpha}, \beta, \gamma, \eta)_{5,1:4} d_{1:4}^k, \quad (13a)$$

$$d_{1:4}^{k+1} \leq A(\hat{\alpha}, \beta, \gamma, \eta)_{1:4,:} d^k. \quad (13b)$$

*Proof:* See Lemma 5 of [1].  $\blacksquare$

Now, we proceed to show the spectral radius of  $A(\hat{\alpha}, \beta, \gamma, \eta)$  defined in (12) is less than 1 with properly chosen parameters.

**Lemma 6:** Given any  $\alpha' > 0, \beta' > 0, w \geq \frac{n}{s_R s_C}$  and  $0 < \eta < \frac{1}{2C_2}$  with  $C_2 > 0$  given in Assumption 3, there exists  $\gamma' > 0$  such that for any  $0 < \gamma < \gamma'$ , if we take  $\hat{\alpha} = \alpha' \gamma^3$  and  $\beta = \beta' \gamma^2$ , then

$$\rho(A(\hat{\alpha}, \beta, \gamma, \eta)) < 1.$$

More specifically, we can choose  $\gamma' = \sup_{\{\hat{\mathbf{v}}_i\}_{i=1}^5} \in \mathcal{P} \min \{\xi_1, \xi_2, \xi_3, \xi_4\}$ , where  $\mathcal{P} = \left\{ \{\hat{\mathbf{v}}_i\}_{i=1}^5 \mid \{\hat{\mathbf{v}}_i\}_{i=1}^5 \text{ satisfies (14a), (14b)} \right\}$  and  $\xi_1, \xi_2, \xi_3, \xi_4$  are the minimum positive roots of the polynomials (15a)–(15d).

*Proof:* To begin with, we choose positive numbers  $\hat{\mathbf{v}}_1, \hat{\mathbf{v}}_2, \hat{\mathbf{v}}_3, \hat{\mathbf{v}}_4, \hat{\mathbf{v}}_5$  satisfying

$$\begin{cases} c_2\hat{\mathbf{v}}_2 + c_3\hat{\mathbf{v}}_3 < c_1\hat{\mathbf{v}}_1, \\ c_{10}\hat{\mathbf{v}}_1 + c_9\hat{\mathbf{v}}_2 + 3\hat{\mathbf{v}}_3 < \hat{\mathbf{v}}_5 \end{cases} \quad (14a)$$

Define  $\hat{\mathbf{v}} = (\hat{\mathbf{v}}_1, \hat{\mathbf{v}}_2, \hat{\mathbf{v}}_3, \hat{\mathbf{v}}_4, \hat{\mathbf{v}}_5)^\top$ , then the five entries  $\{g_i\}_{i=1}^5$  of the vector  $[\mathbf{A}(\alpha'\gamma^3, \beta'\gamma^2, \gamma, \eta)\hat{\mathbf{v}} - \hat{\mathbf{v}}]$  (as functions of  $\gamma$ ) are given by

$$g_1(\gamma) = (-c_1\hat{\mathbf{v}}_1 + c_2\hat{\mathbf{v}}_2 + c_3\hat{\mathbf{v}}_3)\alpha'\gamma^3 + c_4\hat{\mathbf{v}}_4\beta'^2\gamma^4 \quad (15a)$$

$$g_2(\gamma) = -\theta_R\beta'\hat{\mathbf{v}}_2\gamma^2 + \left( c_6\hat{\mathbf{v}}_4\beta'^2 + c_5\hat{\mathbf{v}}_5\frac{\alpha'^2}{\beta'} \right) \gamma^4 \quad (15b)$$

$$g_3(\gamma) = -\theta_C\hat{\mathbf{v}}_3\gamma + c_8\hat{\mathbf{v}}_5\gamma^2 + (8nc_7\hat{\mathbf{v}}_2\beta'^2 + C_2nc_7\hat{\mathbf{v}}_4\beta'^2)\gamma^3 + 2c_7\hat{\mathbf{v}}_5\alpha'^2\gamma^5 \quad (15c)$$

$$g_4(\gamma) = -(\eta - 2C_2\eta^2)\hat{\mathbf{v}}_4 + \left( \frac{8n\beta'^2\hat{\mathbf{v}}_2}{\eta} + 2nC_2\beta'^2\hat{\mathbf{v}}_4 \right) \gamma^4 + \frac{2\alpha'^2\hat{\mathbf{v}}_5}{\eta} \gamma^6 \quad (15d)$$

$$g_5(\gamma) = c_{10}\hat{\mathbf{v}}_1 + c_9\hat{\mathbf{v}}_2 + 3\hat{\mathbf{v}}_3 - \hat{\mathbf{v}}_5. \quad (15e)$$

By (14a), for any  $0 < \gamma < \xi_1 = \frac{\alpha'(c_1\hat{\mathbf{v}}_1 - c_2\hat{\mathbf{v}}_2 - c_3\hat{\mathbf{v}}_3)}{c_4\hat{\mathbf{v}}_4\beta'^2}$ , we have  $g_1(\gamma) < 0$ , where  $\xi_1$  is the minimum positive root of  $g_1(\gamma)$ .

By direct calculation,  $g_2(\gamma) < 0$ , for any  $0 < \gamma < \xi_2 = \sqrt{\frac{\theta_R\beta'\hat{\mathbf{v}}_2}{c_6\hat{\mathbf{v}}_4\beta'^2 + c_5\hat{\mathbf{v}}_5\frac{\alpha'^2}{\beta'}}}$ , where  $\xi_2$  is the minimum positive root of  $g_2(\gamma)$ .

Since  $g_3(0) = 0$ ,  $g'_3(0) = -\theta_C\hat{\mathbf{v}}_3 < 0$  and  $\lim_{\gamma \rightarrow +\infty} g_3(\gamma) = +\infty$ ,  $g_3(\gamma)$  has at least one positive root. By defining  $\xi_3$  as this minimum positive root of  $g_3(\gamma)$ , we have  $g_3(\gamma) < 0$ , for any  $0 < \gamma < \xi_3$ .

As we have assumed  $0 < \eta < \frac{1}{2C_2}$  and  $\hat{\mathbf{v}}_4 > 0$ ,  $g_4(0) = -(\eta - 2C_2\eta^2)\hat{\mathbf{v}}_4 < 0$ . Since  $\lim_{\gamma \rightarrow +\infty} g_4(\gamma) = +\infty$ ,  $g_4(\gamma)$  has a minimum positive root  $\xi_4$ . Then, we have  $g_4(\gamma) < 0$ , for any  $0 < \gamma < \xi_4$ . From (14b),  $g_5(\gamma)$  is a negative constant.

Thus, by the above discussion, for  $0 < \gamma < \min\{\xi_1, \xi_2, \xi_3, \xi_4\}$ , we have  $g_i(\gamma) < 0$  ( $1 \leq i \leq 5$ ), i.e.,  $\mathbf{A}(\hat{\alpha}'\gamma^3, \beta'\gamma^2, \gamma, \eta)\hat{\mathbf{v}} = \mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta)\hat{\mathbf{v}} < \hat{\mathbf{v}}$ . Since  $\hat{\mathbf{v}}$  is a positive vector,  $\mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta)$  is a nonnegative matrix, by [65, Corollary 8.1.29],  $\rho(\mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta)) < 1$ . The above arguments hold true for any positive numbers  $\{\hat{\mathbf{v}}_i\}_{i=1}^5$  satisfying (14a) and (14b). Then, by defining  $\mathcal{P}$  as the set containing these positive number sets  $\{\hat{\mathbf{v}}_i\}_{i=1}^5$ , we can choose  $\gamma'$  as  $\sup_{\{\hat{\mathbf{v}}_i\}_{i=1}^5 \in \mathcal{P}} \min\{\xi_1, \xi_2, \xi_3, \xi_4\}$ . ■

The next theorem shows that Algorithm 1 converges linearly with proper parameters and stepsizes.

*Theorem 7:* Under Assumptions 1-3, given positive numbers  $\alpha', \beta'$  and  $0 < \eta < \frac{1}{2C_2}$ ,  $\eta \leq 1$ ,  $w \geq \frac{n}{s_R^\top s_C}$ , there exists  $\gamma' > 0$  (the value of  $\gamma'$  is given in Lemma 6) such that for any  $\gamma$  satisfying  $0 < \gamma < \gamma'$  and

$$\gamma \leq \min \left\{ 1, \beta'^{-1/2}, \left( \frac{2n}{\alpha'(\mu + L)(s_R^\top s_C)} \right)^{1/3} \right\}, \quad (16)$$

if we set  $\beta = \beta'\gamma^2$  and the stepsizes  $\{\alpha_i\}_{i \in \mathcal{N}}$  satisfy  $\hat{\alpha} = \max_{i \in \mathcal{N}} \alpha_i = \alpha'\gamma^3$ , with  $\hat{\alpha}/\bar{\alpha} \leq w$ , then  $\mathbb{E}[\|\bar{\mathbf{x}}^k - \mathbf{x}^*\|_2^2]$ ,  $\mathbb{E}[\|\Pi_R \mathbf{X}^k\|_R^2]$  converge linearly to 0. More specifically,

$$\mathbb{E}[\|\bar{\mathbf{x}}^k - \mathbf{x}^*\|_2^2] \leq \sigma \bar{\mathbf{v}}_1 \rho(\mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta))^k,$$

$$\mathbb{E}[\|\Pi_R \mathbf{X}^k\|_R^2] \leq \sigma \bar{\mathbf{v}}_2 \rho(\mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta))^k,$$

where  $\sigma, \bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2$  are constants given in the proof.

*Proof:* Denote  $\mathbf{A} = \mathbf{A}(\hat{\alpha}, \beta, \gamma, \eta)$  for simplicity. Since  $\mathbf{A}$  is a regular nonnegative matrix,<sup>1</sup> by the Perron-Frobenius theorem [66],  $\rho(\mathbf{A})$  is an eigenvalue of  $\mathbf{A}$ , and  $\mathbf{A}$  has a unique right positive eigenvector  $\bar{\mathbf{v}}$  with respect to the eigenvalue  $\rho(\mathbf{A})$ . Define the constant  $\sigma = \max_{1 \leq i \leq 4} \frac{d_i^0}{\bar{v}_i}$ . Thus, by the definition of  $\sigma$ ,  $d_{1:4}^0 \leq \sigma \bar{\mathbf{v}}_{1:4}$ .

Next, we prove the linear convergence by induction. If we have proved  $d_{1:4}^k \leq \sigma \rho(\mathbf{A})^k \bar{\mathbf{v}}_{1:4}$ , we will show that it also holds true for  $k+1$ . The requirement (16) guarantees that  $\gamma \leq 1$ ,  $\beta = \beta'\gamma^2 \leq 1$ ,  $\bar{\alpha} \leq \frac{(s_R^\top s_C)\hat{\alpha}}{n} = \frac{(s_R^\top s_C)\alpha'\gamma^3}{n} \leq \frac{2}{\mu+L}$ . According to Lemma 6,  $\rho(\mathbf{A}) < 1$ .

By (13a) and the inductive hypothesis, there holds

$$\mathbb{E}[\|\mathbf{Y}^k\|_2^2] \leq \mathbf{A}_{5,1:4} d_{1:4}^k \leq \sigma \rho(\mathbf{A})^k \mathbf{A}_{5,1:4} \bar{\mathbf{v}}_{1:4}$$

$$= \sigma \rho(\mathbf{A})^k \mathbf{A}_{5,:} \bar{\mathbf{v}} = \sigma \rho(\mathbf{A})^{k+1} \bar{\mathbf{v}}_5 \leq \sigma \rho(\mathbf{A})^k \bar{\mathbf{v}}_5.$$

Combining with the inductive hypothesis, we obtain

$$d^k \leq \sigma \rho(\mathbf{A})^k \bar{\mathbf{v}}. \quad (17)$$

Now, by (13b) and (17),

$$d_{1:4}^{k+1} \leq \mathbf{A}_{1:4,:} d^k \leq \sigma \rho(\mathbf{A})^k \mathbf{A}_{1:4,:} \bar{\mathbf{v}} = \sigma \rho(\mathbf{A})^{k+1} \bar{\mathbf{v}}_{1:4},$$

i.e., the statement holds for  $k+1$ . Therefore, by induction,  $d_{1:4}^k \leq \sigma \rho(\mathbf{A})^k \bar{\mathbf{v}}_{1:4}$ , for any  $k \geq 0$ , which completes the proof. ■

*Remark 8:* In practice, we can take for instance  $\alpha' = \frac{1}{L}$ ,  $\beta' = 1$ ,  $\gamma = 1$ ,  $\eta = \min\{\frac{1}{2C_2}, 1\}$ . Then, we hand-tune  $\gamma$  under the relations  $\alpha_i = \hat{\alpha} = \alpha'\gamma^3$  ( $\forall i \in \mathcal{N}$ ) and  $\beta = \beta'\gamma^2$  to achieve the optimal performance.

*Remark 9:* We can also add momentums for the communication of  $\mathbf{Y}^k$  like what we did for  $\mathbf{X}^k$ . In this case, linear convergence can be proved similarly. However, we see little improvement in the numerical experiments when  $\mathbf{Y}^k$  is equipped with momentums. In addition, more momentums will take more storage space. Therefore, we omit the details of this case here.

<sup>1</sup> A matrix  $\mathbf{A}$  is said to be regular if all the entries of  $\mathbf{A}^k$  are positive for some integer  $k \geq 0$ .

#### IV. A BROADCAST-LIKE GRADIENT TRACKING METHOD WITH COMPRESSION (B-CPP)

In this section, we consider a broadcast-like gradient tracking method with compression (B-CPP). Broadcast or gossip-based communication protocols are popular for distributed computation due to their low communication costs [62]–[64]. In B-CPP, at every iteration  $k$ , one agent  $i_k \in \mathcal{N}$  wakes up with uniform probability. This can be easily realized, for example, if each agent wakes up according to an independent Poisson clock with the same parameter. Hence the probability  $\Pr[i_k = j] = \frac{1}{n}$  for any  $j \in \mathcal{N}$ . In addition,  $\{i_k\}_{k \geq 0}$  are independent with each other.

Briefly speaking, each iteration  $k$  of the B-CPP algorithm consists of the following steps:

- 1) One random agent  $i_k$  wakes up.
- 2) Agent  $i_k$  sends information to all its out-neighbors in  $\mathbf{R}$  ( $\mathcal{N}_{R,i_k}^+$ ) and  $\mathbf{C}$  ( $\mathcal{N}_{C,i_k}^+$ ).
- 3) The agents who received information from  $i_k$  are awoken and update their local variables.

We remark that the number  $k$  is only used for analysis purpose and does not need to be known by the agents. B-CPP can be generalized to the case when each agent wakes up with different but known probabilities. The awakened agents will know whether they are  $i_k$  and whether they are in the set  $\mathcal{N}_{R,i_k}^+$  or in the set  $\mathcal{N}_{C,i_k}^+$ , and will take different actions accordingly. The B-CPP method is illustrated in Algorithm 2. To implement communication compression in a broadcast setting, a naive way is to let each agent  $i$  hold different momentums  $\mathbf{u}_{i,j}$  for each neighbors  $j$ . In this way, at each time agent  $i$  receives information  $\mathbf{p}$  from neighboring agent  $j$ , it can restore the information sent from agent  $j$  directly by summing  $\mathbf{p} + \mathbf{u}_{i,j}$ . However, this procedure would require momentums as many as the total number of edges. Hence it could be impractical when the storage space is limited and the graph is dense. In B-CPP, we overcome this issue so that each agent only uses 2 momentums.

To help analyze the convergence rate of Algorithm 2, for  $0 \leq k \leq K$ , let us define the final state of variable  $\mathbf{x}_j$  before iteration  $k$  as  $\mathbf{x}_j^k$ . These  $\mathbf{x}_j^k$  are written compactly into an  $n$ -by- $p$  matrix  $\mathbf{X}^k$ . And  $\mathbf{Y}^k, \mathbf{U}^k, \mathbf{U}_R^k, \nabla F(\mathbf{X}^k)$  are defined analogously. Let  $\mathcal{D}_k$  be the  $\sigma$ -field generated by  $\{i_j, \mathbf{p}^j, \mathbf{q}^j\}_{0 \leq j \leq k-1}$ . Let  $\mathcal{D}_k^+$  denote the  $\sigma$ -field generated by  $\mathcal{D}_k$  and  $i_k$ . Let  $\mathcal{D}_k^{++}$  denote the  $\sigma$ -field generated by  $\mathcal{D}_k$  and  $\{i_k, \mathbf{p}^k\}$ . For any event  $A$ , let  $\chi_A$  denote the indicator function of  $A$ .

The differences between B-CPP and CPP mainly lie in the averaging step. Taking the averaging step for  $\mathbf{X}^k$  as an example. Firstly, we also have by induction that  $\mathbf{U}_R^k = \mathbf{R}\mathbf{U}^k$ ,  $\forall 0 \leq k \leq K$ . For  $i \in \mathcal{N}$ , denote the in-degree of  $i$  in  $\mathbf{R}$  by  $r_i = |\mathcal{N}_{R,i}^+|$ . Since  $\mathbf{R}$  is row stochastic,  $r_i > 0$  for any  $i \in \mathcal{N}$ . Notice that the update step (18a) will be implemented by agent  $j$  when agent  $i_k$  is an in-neighbor of agent  $j$  (or equivalently,  $j$  is an out-neighbor of  $i_k$ ).

$$\begin{aligned} & \mathbb{E} \left[ \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} n \mathbf{R}_{j,i_k} \mathbf{p}^k \mid \mathcal{D}_k \right] \\ &= \mathbb{E} \left[ \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} n \mathbf{R}_{j,i_k} \mathbb{E} [\mathbf{p}^k \mid i_k] \mid \mathcal{D}_k \right] \end{aligned}$$

---

#### Algorithm 2: B-CPP.

---

**Input:** each agent  $i$  is instilled its initial position  $\mathbf{x}_i$ , parameters  $\beta, \gamma, \eta \in (0, 1]$ , stepsize  $\alpha_i$  and network information  $\mathbf{R}_{ij} (j \in \mathcal{N}_{R,i}^-)$ ,  $\mathbf{C}_{ij} (j \in \mathcal{N}_{C,i}^-)$  and the sets  $\mathcal{N}_{R,i}^+, \mathcal{N}_{C,i}^+$   
**Output:** each agent  $i$  outputs its final local decision variable  $\mathbf{x}_i$

Each agent  $i$  initializes

$$\mathbf{y}_i^0 = \nabla f_i(\mathbf{x}_i^0), \mathbf{u}_i^0 = \mathbf{0}, \mathbf{u}_{R,i}^0 = \mathbf{0}.$$

**for**  $k = 0$  to  $K - 1$  **do**

Agent  $i_k$  awakes.

Agent  $i_k$  compresses

$\mathbf{p}^k = \text{COMPRESS}(\mathbf{x}_{i_k} - \mathbf{u}_{i_k})$ , sends  $\mathbf{p}^k$  to all agents  $j \in \mathcal{N}_{R,i_k}^+$  and wakes up these agents.

Agent  $i_k$  compresses  $\mathbf{q}^k = \text{COMPRESS}(\mathbf{y}_{i_k})$ , sends  $\mathbf{q}^k$  to all agents in  $\mathcal{N}_{C,i_k}^+$  and wakes up these agents.

For all agents  $j \in \{i_k\} \cup \mathcal{N}_{R,i_k}^+ \cup \mathcal{N}_{C,i_k}^+$ , agent  $j$  records  $\mathbf{d}_j \leftarrow \nabla f_j(\mathbf{x}_j)$ .

For  $j \in \mathcal{N}_{R,i_k}^+$ , agent  $j$  updates

$$\mathbf{u}_{R,j} \leftarrow \mathbf{u}_{R,j} + \eta n \mathbf{R}_{j,i_k} \mathbf{p}^k, \quad (18a)$$

$$\begin{aligned} \mathbf{x}_j &\leftarrow \left(1 - \frac{\beta n}{r_j}\right) \mathbf{x}_j \\ &\quad + \frac{\beta n}{r_j} \mathbf{u}_{R,j} + \beta n \mathbf{R}_{j,i_k} \mathbf{p}^k. \end{aligned} \quad (18b)$$

Agent  $i_k$  updates  $\mathbf{u}_{i_k} \leftarrow \mathbf{u}_{i_k} + \eta n \mathbf{p}^k$ .

For all  $j \in \{i_k\} \cup \mathcal{N}_{R,i_k}^+ \cup \mathcal{N}_{C,i_k}^+$ , agent  $j$  updates  $\mathbf{x}_j \leftarrow \mathbf{x}_j - \alpha_j \mathbf{y}_j$ .

For all  $j \in \{i_k\} \cup \mathcal{N}_{R,i_k}^+ \cup \mathcal{N}_{C,i_k}^+$ , agent  $j$  updates  $\mathbf{y}_j \leftarrow \mathbf{y}_j + \nabla f_j(\mathbf{x}_j) - \mathbf{d}_j$ .

Agent  $i_k$  update  $\mathbf{y}_{i_k} \leftarrow \mathbf{y}_{i_k} - \gamma n \mathbf{q}^k$ .

For  $j \in \mathcal{N}_{C,i_k}^+$ , agent  $j$  updates

$$\mathbf{y}_j \leftarrow \mathbf{y}_j + \gamma n \mathbf{C}_{j,i_k} \mathbf{q}^k.$$

**end**

Return the final local decision variable  $\{\mathbf{x}_i\}_{i \in \mathcal{N}}$ .

---

$$= \mathbb{E} [n \mathbf{R}_{j,i_k} (\mathbf{x}_{i_k}^k - \mathbf{u}_{i_k}^k) \mid \mathcal{D}_k]$$

$$= \frac{1}{n} \sum_{i_k=1}^n n \mathbf{R}_{j,i_k} (\mathbf{x}_{i_k}^k - \mathbf{u}_{i_k}^k)$$

$$= \mathbf{R}_{j,:} (\mathbf{X}^k - \mathbf{U}^k).$$

Analogously,

$$\mathbb{E} \left[ \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} \left(1 - \frac{\beta n}{r_j}\right) \mathbf{x}_j^k \mid \mathcal{D}_k \right]$$

$$= \frac{1}{n} \sum_{i_k \in \mathcal{N}, \mathbf{R}_{j,i_k} > 0} \left(1 - \frac{\beta n}{r_j}\right) \mathbf{x}_j^k$$

$$= \frac{r_j}{n} \mathbf{x}_j^k - \beta \mathbf{x}_j^k,$$



and

$$\begin{aligned}
& \mathbb{E} \left[ \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} \left( \frac{n}{r_j} \mathbf{u}_{R,j}^k \right) | \mathcal{D}_k \right] \\
&= \frac{n}{r_j} \mathbf{u}_{R,j}^k \Pr \left[ j \in \mathcal{N}_{R,i_k}^+ | \mathcal{D}_k \right] \\
&= \frac{1}{n} \sum_{i_k \in \mathcal{N}, \mathbf{R}_{j,i_k} > 0} \frac{n}{r_j} \mathbf{u}_{R,j}^k \\
&= \mathbf{u}_{R,j}^k = \mathbf{R}_{j,:} \mathbf{U}^k.
\end{aligned}$$

Thus, taking expectation on the RHS of (18b) yields

$$\begin{aligned}
& \mathbb{E} \left[ \chi_{\{j \notin \mathcal{N}_{R,i_k}^+\}} \mathbf{x}_j^k + \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} \left( \left( 1 - \frac{\beta n}{r_j} \right) \mathbf{x}_j^k \right) \right. \\
& \quad \left. + \chi_{\{j \in \mathcal{N}_{R,i_k}^+\}} \left( \frac{\beta n}{r_j} \mathbf{u}_{R,j}^k + \beta n \mathbf{R}_{j,i_k} \mathbf{p}^k \right) | \mathcal{D}_k \right] \\
&= \left( \frac{n - r_j}{n} + \frac{r_j}{n} - \beta \right) \mathbf{x}_j^k + \beta \mathbf{R}_{j,:} (\mathbf{U}^k + \mathbf{X}^k - \mathbf{U}^k) \\
&= [\mathbf{R}_\beta]_{j,:} \mathbf{X}^k.
\end{aligned}$$

Briefly speaking, after taking conditional expectation, the averaging step of B-CPP reduces to that of CPP (the first term of the RHS of (5c)). By choosing a proper value for  $\beta$ , the variance of the RHS of (18) can be made small enough. In this way, the averaging step in B-CPP could have a similar effect to that of CPP.

To show the linear convergence of B-CPP, for simplicity, we assume  $\alpha = \alpha_1 = \alpha_2 = \dots = \alpha_n$ . We remark that the analysis below can be easily generalized to the case when the stepsizes differ among the agents.

We denote  $\mathbf{1}_i \in \mathbb{R}^{n \times 1}$  as the vector with 1 on the  $i$ -th component and 0 on the others. Define

$$\begin{aligned}
\mathbf{P}^k &= \mathbf{X}^k - \mathbf{U}^k + \mathbf{1}_{i_k} (\mathbf{p}^k - \mathbf{x}_{i_k}^k + \mathbf{u}_{i_k}^k), \\
\mathbf{Q}^k &= \mathbf{Y}^k + \mathbf{1}_{i_k} (\mathbf{q}^k - \mathbf{Y}_{i_k}^k).
\end{aligned}$$

As  $\mathbf{U}^k + \mathbf{P}^k$  and  $\mathbf{Q}^k$  are used to estimate  $\mathbf{X}^k$  and  $\mathbf{y}^k$ , respectively, we define the error matrices induced by the compression as follows

$$\begin{aligned}
\mathbf{W}^k &= \mathbf{X}^k - \mathbf{U}^k - \mathbf{P}^k = \mathbf{1}_{i_k} (\mathbf{x}_{i_k}^k - \mathbf{u}_{i_k}^k - \mathbf{p}^k), \\
\mathbf{E}^k &= \mathbf{Y}^k - \mathbf{Q}^k = \mathbf{1}_{i_k} (\mathbf{Y}_{i_k}^k - \mathbf{q}^k).
\end{aligned}$$

The following lemma is a direct corollary of Assumption 3.

**Lemma 10:** The random variable  $i_k$  is independent with  $\mathcal{D}_k$ . And  $\mathbb{E}[\mathbf{W}^k | \mathcal{D}_k^+] = \mathbf{0}$ ,  $\mathbb{E}[\mathbf{E}^k | \mathcal{D}_k^{++}] = \mathbf{0}$ .

Moreover, by Assumption 3,

$$\begin{aligned}
& \mathbb{E} [\|\mathbf{W}^k\|_2^2 | \mathcal{D}_k^+] = \mathbb{E} [\|\mathbf{x}_{i_k}^k - \mathbf{u}_{i_k}^k - \mathbf{p}^k\|_2^2 | \mathcal{D}_k^+] \\
& \leq C_2 \mathbb{E} [\|\mathbf{x}_{i_k}^k - \mathbf{u}_{i_k}^k\|_2^2 | \mathcal{D}_k] = \frac{C_2}{n} \|\mathbf{X}^k - \mathbf{U}^k\|_2^2, \quad (19)
\end{aligned}$$

$$\begin{aligned}
& \mathbb{E} [\|\mathbf{E}^k\|_2^2 | \mathcal{D}_k^{++}] = \mathbb{E} [\|\mathbf{Y}_{i_k}^k - \mathbf{q}^k\|_2^2 | \mathcal{D}_k] \\
& \leq C_2 \mathbb{E} [\|\mathbf{Y}_{i_k}^k\|_2^2 | \mathcal{D}_k^{++}] = \frac{C_2}{n} \|\mathbf{Y}^k\|_2^2. \quad (20)
\end{aligned}$$

We will also use the following random matrices, which are all measurable with respect to  $i_k$ ,

$$\Lambda^k = n \mathbf{Diag}(\mathbf{1}_{i_k}), \tilde{\Lambda}^k = \sum_{j \in \{i_k\} \cup \mathcal{N}_{R,i_k}^+ \cup \mathcal{N}_{C,i_k}^+} \mathbf{Diag}(\mathbf{1}_j),$$

$$\tilde{\Lambda}_R^k = \sum_{j \in \mathcal{N}_{R,i_k}^+} \frac{n}{r_j} \mathbf{Diag}(\mathbf{1}_j), \bar{\mathbf{R}}^k = \sum_{j \in \mathcal{N}_{R,i_k}^+} \frac{n}{r_j} \mathbf{1}_j \mathbf{R}_{j,:},$$

$$\hat{\mathbf{R}}^k = n \mathbf{R}_{:,i_k} \mathbf{1}_{i_k}^\top, \hat{\mathbf{C}}^k = n \mathbf{C}_{:,i_k} \mathbf{1}_{i_k}^\top.$$

Now, Algorithm 2 can be rewritten into a more compact form:

$$\begin{cases} \mathbf{X}^{k+1} = (\mathbf{I} - \beta \tilde{\Lambda}_R^k) \mathbf{X}^k + \beta \tilde{\Lambda}_R^k \mathbf{R} \mathbf{U}^k \\ \quad + \beta \hat{\mathbf{R}}^k \mathbf{P}^k - \alpha \tilde{\Lambda}^k \mathbf{Y}^k, \end{cases} \quad (21a)$$

$$\begin{cases} \mathbf{Y}^{k+1} = \mathbf{Y}^k + \gamma (\hat{\mathbf{C}}^k - \Lambda^k) \mathbf{Q}^k \\ \quad + \nabla \mathbf{F}(\mathbf{X}^{k+1}) - \nabla \mathbf{F}(\mathbf{X}^k), \end{cases} \quad (21b)$$

$$\begin{cases} \mathbf{U}^{k+1} = \mathbf{U}^k + \eta \Lambda^k \mathbf{P}^k. \end{cases} \quad (21c)$$

Compared to CPP, the stochastic matrices  $\hat{\mathbf{R}}^k, \bar{\mathbf{R}}^k, \Lambda^k, \tilde{\Lambda}^k, \tilde{\Lambda}_R^k$  will induce additional errors that need to be dealt with in the convergence analysis of B-CPP.

To analyze the convergence rate of B-CPP, similar to what we did for CPP, we also use  $\mathbb{E}[\|\mathbf{x}^k - \mathbf{x}^*\|_2^2]$ ,  $\mathbb{E}[\|\Pi_R \mathbf{X}^k\|_2^2]$ ,  $\mathbb{E}[\|\Pi_C \mathbf{Y}^k\|_2^2]$  and  $\mathbb{E}[\|\mathbf{U}^k - \mathbf{X}^k\|_2^2]$  to measure the closeness to the optimal point, consensus error, gradient tracking error and the convergence of the momentums, respectively. And we use  $\mathbb{E}[\|\mathbf{Y}^k\|_2^2]$ ,  $\mathbb{E}[\|\mathbf{X}^{k+1} - \mathbf{X}^k\|_2^2]$  and  $\mathbb{E}[\|\mathbf{V}^k\|_2^2]$  to help simplify the recursive relations between the above four quantities. The definition of  $\mathbf{V}^k$  is given by (47) in the Supplementary Material of [1].

**Lemma 11:** For  $k \geq 0$ , define  $\tilde{\mathbf{d}}^k \in \mathbb{R}^7$  as an extension of  $\mathbf{d}^k$ :

$$\tilde{\mathbf{d}}^k = \left( \mathbf{d}^k, \mathbb{E}[\|\mathbf{V}^k\|_2^2], \mathbb{E}[\|\mathbf{X}^{k+1} - \mathbf{X}^k\|_2^2] \right)^\top.$$

Define a parameterized matrix  $\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)$  as follows

$$\begin{aligned}
\tilde{\mathbf{A}}_{:,1:3} &= \begin{pmatrix} 1 - d_1 \alpha & d_2 \alpha & d_3 \alpha \\ 0 & 1 - \theta_R \beta & 0 \\ 0 & 0 & 1 - \theta_C \gamma \\ 0 & 0 & 0 \\ c_{10} & c_9 & 3 \\ 0 & 3V_1 \beta^2 & 0 \\ 0 & 18\beta^2 & 0 \end{pmatrix}, \\
\tilde{\mathbf{A}}_{:,4:7} &= \begin{pmatrix} 0 & 0 & d_4 & 0 \\ 0 & \frac{\alpha^2}{1 - \theta_C \gamma} + d_7 \gamma^2 & \delta_{R,2}^2 & 0 \\ 0 & \frac{d_5 \gamma^2}{1 - \theta_C \gamma} + d_7 \gamma^2 & 0 & \frac{d_6}{\gamma} \\ \tilde{\mathbf{A}}_{44} & 0 & 0 & \frac{2}{\eta} \\ 0 & 0 & 0 & 0 \\ d_9 \beta^2 & 3V_4 \alpha^2 & 0 & 0 \\ 0 & 2\alpha^2 & 1 & 0 \end{pmatrix} \quad (22)
\end{aligned}$$

where  $\tilde{\mathbf{A}}_{44} = 1 - \eta + \frac{\eta^2(n-1)^2}{1-\eta} + d_8 \frac{\beta^2}{\eta} + 2C_2 n \eta^2 + d_8 \beta^2$  and the constants  $d_1$ - $d_9$  are given by (51), (55), (57) and (59) in the Supplementary Material of [1]. Then, under Assumptions 1-3



and the condition  $\tilde{\alpha} \leq \frac{2}{\mu+L}$ , we have the following inequalities:

$$\mathbb{E} \left[ \left\| \mathbf{Y}^k \right\|_2^2 \right] \leq \tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)_{5,1:4} \tilde{\mathbf{d}}_{1:4}^k, \quad (23a)$$

$$\mathbb{E} \left[ \left\| \mathbf{V}^k \right\|_2^2 \right] \leq \tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)_{6,1:5} \tilde{\mathbf{d}}_{1:5}^k, \quad (23b)$$

$$\mathbb{E} \left[ \left\| \mathbf{X}^{k+1} - \mathbf{X}^k \right\|_2^2 \right] \leq \tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)_{7,1:6} \tilde{\mathbf{d}}_{1:6}^k, \quad (23c)$$

$$\tilde{\mathbf{d}}_{1:4}^{k+1} \leq \tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)_{1:4,:} \tilde{\mathbf{d}}^k. \quad (23d)$$

*Proof:* See Lemma 11 of [1]. ■

Next, we show that with properly chosen parameters and stepsize, the spectral radius of the parameterized matrix  $\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)$  in (22) will be less than 1.

**Lemma 12:** Given any positive numbers  $\alpha, \beta, \gamma$  and  $\eta$  satisfying

$$0 < \eta \leq \frac{1}{2}, \quad \eta < \frac{1}{2(n-1)^2 + 2C_2n}, \quad (24)$$

there exists  $\gamma' > 0$  such that for any  $0 < \gamma < \gamma'$ , if we take  $\alpha = \alpha'\gamma^3$  and  $\beta = \beta'\gamma^2$ , then

$$\rho(\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)) < 1.$$

More specifically, we can take  $\gamma' = \sup_{\{\tilde{\mathbf{v}}_i\}_{i=1}^7 \in \mathcal{P}'} \min\{\psi_2, \psi_3, \psi_4, \psi_6, \psi_7, \frac{1}{\theta_C}\}$ , where  $\theta_C$  is given in Lemma 3 and  $\mathcal{P}'$ ,  $\psi_i$  ( $i = 2, 3, 4, 6, 7$ ) are given in the proof.

*Proof:* Given positive numbers  $\{\tilde{\mathbf{v}}_i\}_{i=1}^7$  satisfying

$$\begin{cases} d_2\tilde{\mathbf{v}}_2\alpha' + d_3\tilde{\mathbf{v}}_3\alpha' + d_4\tilde{\mathbf{v}}_6 < d_1\tilde{\mathbf{v}}_1\alpha', & (25a) \\ c_{10}\tilde{\mathbf{v}}_1 + c_9\tilde{\mathbf{v}}_2 + 3\tilde{\mathbf{v}}_3 < \tilde{\mathbf{v}}_5, & (25b) \\ \tilde{\mathbf{v}}_6 < \tilde{\mathbf{v}}_7. & (25c) \end{cases}$$

Define  $\tilde{\mathbf{v}} = (\tilde{\mathbf{v}}_1, \tilde{\mathbf{v}}_2, \tilde{\mathbf{v}}_3, \tilde{\mathbf{v}}_4, \tilde{\mathbf{v}}_5, \gamma^3\tilde{\mathbf{v}}_6, \gamma^3\tilde{\mathbf{v}}_7)^\top$ . The seven entries  $\{g_i\}_{i=1}^7$  of the vector  $[\tilde{\mathbf{v}} - \tilde{\mathbf{A}}(\alpha'\gamma^3, \beta'\gamma^2, \gamma, \eta)\tilde{\mathbf{v}}]$  (as functions of  $\gamma$ ) are given by

$$g_1(\gamma) = \gamma^3(-d_1\tilde{\mathbf{v}}_1\alpha' + d_2\tilde{\mathbf{v}}_2\alpha' + d_3\tilde{\mathbf{v}}_3\alpha' + d_4\tilde{\mathbf{v}}_6)$$

$$g_2(\gamma) = \gamma^2\left(-\theta_R\beta'\tilde{\mathbf{v}}_2 + \frac{\alpha'^2\tilde{\mathbf{v}}_5}{\beta'\theta_R}\gamma^2 + \delta_{R,2}^2\tilde{\mathbf{v}}_6\gamma\right),$$

$$g_3(\gamma) = \gamma\left(-\theta_C\tilde{\mathbf{v}}_3 + \frac{d_5\tilde{\mathbf{v}}_5\gamma}{1-\theta_C\gamma} + d_7\tilde{\mathbf{v}}_5\gamma + d_6\tilde{\mathbf{v}}_7\gamma\right),$$

$$g_4(\gamma) = \left(-\eta + \frac{\eta^2(n-1)^2}{1-\eta} + 2C_2n\eta^2\right)\tilde{\mathbf{v}}_4 + d_8\tilde{\mathbf{v}}_4\frac{\beta'^2}{\eta}\gamma^4 + d_8\tilde{\mathbf{v}}_4\beta'^2\gamma^4 + \frac{2\tilde{\mathbf{v}}_7}{\eta}\gamma^3,$$

$$g_5(\gamma) = c_{10}\tilde{\mathbf{v}}_1 + c_9\tilde{\mathbf{v}}_2 + 3\tilde{\mathbf{v}}_3 - \tilde{\mathbf{v}}_5,$$

$$g_6(\gamma) = \gamma^3(3V_1\tilde{\mathbf{v}}_2\beta'^2\gamma + d_9\tilde{\mathbf{v}}_4\beta'^2\gamma + 3V_4\tilde{\mathbf{v}}_5\alpha'^2\gamma^3 - \tilde{\mathbf{v}}_6),$$

$$g_7(\gamma) = \gamma^3((\tilde{\mathbf{v}}_6 - \tilde{\mathbf{v}}_7) + 18\tilde{\mathbf{v}}_2\beta'^2\gamma + 2\tilde{\mathbf{v}}_5\alpha'^2\gamma^3).$$

Denote  $\tilde{g}_i = g_i/\gamma^3$  ( $i = 1, 6, 7$ ),  $\tilde{g}_2 = g_2/\gamma^2$ ,  $\tilde{g}_3 = g_3/\gamma$ . By (25a),  $\tilde{g}_1 = -d_1\tilde{\mathbf{v}}_1\alpha' + d_2\tilde{\mathbf{v}}_2\alpha' + d_3\tilde{\mathbf{v}}_3\alpha' + d_4\tilde{\mathbf{v}}_6 < 0$ . By (25b),  $\tilde{g}_5 = c_{10}\tilde{\mathbf{v}}_1 + c_9\tilde{\mathbf{v}}_2 + 3\tilde{\mathbf{v}}_3 - \tilde{\mathbf{v}}_5 < 0$ . By (25c),  $\tilde{g}_7(0) = \tilde{\mathbf{v}}_6 - \tilde{\mathbf{v}}_7 < 0$ . By (24),  $\eta \leq 1/2$ , then  $g_4(0) \leq -\eta + 2\eta^2(n-1)^2 + 2C_2n\eta^2 < 0$ . It follows directly that  $\tilde{g}_2(0) = -\theta_R\beta'\tilde{\mathbf{v}}_2 < 0$ ,  $\tilde{g}_3(0) = -\theta_C\tilde{\mathbf{v}}_3 < 0$ ,  $\tilde{g}_6(0) = -\tilde{\mathbf{v}}_6 < 0$ .

If  $0 < \gamma < \frac{1}{\theta_C}$ , we have  $\lim_{\gamma \rightarrow +\infty} g_3(\gamma) = +\infty$ . Combining with the facts that  $\lim_{\gamma \rightarrow +\infty} g_i(\gamma) = +\infty$ , for  $i = 2, 4, 6, 7$ , we

can define  $\psi_i$  to be the minimum positive solution of  $g_i(\gamma)$  ( $i = 2, 3, 4, 6, 7$ ) when  $0 < \gamma < \frac{1}{\theta_C}$ .

Thus, for any  $0 < \gamma < \min\{\psi_2, \psi_3, \psi_4, \psi_6, \psi_7, \frac{1}{\theta_C}\}$ , there holds  $g_i(\gamma) < 0$ , for any  $1 \leq i \leq 7$ , i.e.,  $\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)\tilde{\mathbf{v}} = \tilde{\mathbf{A}}(\alpha'\gamma^3, \beta'\gamma^2, \gamma, \eta)\tilde{\mathbf{v}} < \tilde{\mathbf{v}}$ . Since  $\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)$  is a nonnegative matrix,  $\tilde{\mathbf{v}}$  is a positive vector, by [65, Corollary 8.1.29], we have  $\rho(\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)) < 1$ .

The above arguments hold true for arbitrary  $\{\tilde{\mathbf{v}}_i\}_{i=1}^7$  satisfying (25a), (25b) and (25c). Then, let the set  $\mathcal{P}'$  consist of all these  $\{\tilde{\mathbf{v}}_i\}_{i=1}^7$ , we have  $\gamma' \geq \sup_{\{\tilde{\mathbf{v}}_i\}_{i=1}^7 \in \mathcal{P}'} \min\{\psi_2, \psi_3, \psi_4, \psi_6, \psi_7, \frac{1}{\theta_C}\}$ . ■

The following theorem shows the linear convergence of Algorithm 2 given proper parameters and stepsize.

**Theorem 13:** Under Assumptions 1-3, if we choose arbitrary positive numbers  $\alpha, \beta, \gamma, \eta$  satisfying  $0 < \eta \leq \frac{1}{2}$ ,  $\eta < \frac{1}{2(n-1)^2 + 2C_2n}$ , and choose  $\gamma' > 0$  as in Lemma 12. Then, for any  $\gamma$  satisfying  $0 < \gamma < \gamma'$  and

$$\gamma \leq \min \left\{ 1, \beta'^{-\frac{1}{2}}, \left( \frac{2n}{\alpha'(\mu+L)(\mathbf{s}_R^\top \tilde{\mathbf{A}} \mathbf{s}_C)} \right)^{\frac{1}{3}} \right\}, \quad (26)$$

$\beta = \beta'\gamma^2$ , and  $\alpha = \alpha'\gamma^3$ , there holds, for any  $k \geq 0$ ,

$$\mathbb{E} \left[ \left\| \bar{\mathbf{x}}^k - \mathbf{x}^* \right\|_2^2 \right] \leq \sigma' \bar{\mathbf{v}}_1 \rho(\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta))^k,$$

$$\mathbb{E} \left[ \left\| \Pi_R \mathbf{X}^k \right\|_R^2 \right] \leq \sigma' \bar{\mathbf{v}}_2 \rho(\tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta))^k,$$

where  $\sigma', \bar{\mathbf{v}}_1, \bar{\mathbf{v}}_2$  are constants given in the proof.

*Proof:* To begin with, we denote  $\tilde{\mathbf{A}} = \tilde{\mathbf{A}}(\alpha, \beta, \gamma, \eta)$  for simplicity. Since  $\tilde{\mathbf{A}}$  is a regular nonnegative matrix, by the Perron-Frobenius theorem [66],  $\rho(\tilde{\mathbf{A}})$  is an eigenvalue of  $\tilde{\mathbf{A}}$ , and  $\tilde{\mathbf{A}}$  has a unique positive right eigenvector  $\bar{\mathbf{v}}$  with respect to the eigenvalue  $\rho(\tilde{\mathbf{A}})$ . Define  $\sigma' = \max_{1 \leq i \leq 4} \frac{\tilde{\mathbf{d}}_i^0}{\bar{\mathbf{v}}_i}$ . Then,  $\tilde{\mathbf{d}}_{1:4}^0 \leq \sigma' \bar{\mathbf{v}}_{1:4}$ .

Next, we prove the linear convergence by induction. If we have proved  $\tilde{\mathbf{d}}_{1:4}^k \leq \sigma' \rho(\tilde{\mathbf{A}})^k \bar{\mathbf{v}}_{1:4}$ , we will show that it also holds for  $k+1$ . The requirement (26) guarantees that  $\gamma \leq 1$ ,  $\beta \leq 1$ ,  $\tilde{\alpha} \leq \frac{2}{\mu+L}$ . From Lemma 12,  $\rho(\tilde{\mathbf{A}}) < 1$ .

By (23a) and the inductive hypothesis, there holds

$$\begin{aligned} \mathbb{E} \left[ \left\| \mathbf{Y}^k \right\|_2^2 \right] &\leq \tilde{\mathbf{A}}_{5,1:4} \tilde{\mathbf{d}}_{1:4}^k \leq \sigma' \rho(\tilde{\mathbf{A}})^k \tilde{\mathbf{A}}_{5,1:4} \bar{\mathbf{v}}_{1:4} \\ &= \sigma' \rho(\tilde{\mathbf{A}})^k [\tilde{\mathbf{A}} \bar{\mathbf{v}}]_5 = \sigma' \rho(\tilde{\mathbf{A}})^{k+1} \bar{\mathbf{v}}_5 \leq \sigma' \rho(\tilde{\mathbf{A}})^k \bar{\mathbf{v}}_5. \end{aligned}$$

Then combining with the inductive hypothesis, we have

$$\tilde{\mathbf{d}}_{1:5}^k \leq \sigma' \rho(\tilde{\mathbf{A}})^k \bar{\mathbf{v}}_{1:5}. \quad (27)$$

By (23b) and (27),

$$\begin{aligned} \mathbb{E} \left[ \left\| \mathbf{V}^k \right\|_2^2 \right] &\leq \tilde{\mathbf{A}}_{6,1:5} \tilde{\mathbf{d}}_{1:5}^k \leq \sigma' \rho(\tilde{\mathbf{A}})^k \tilde{\mathbf{A}}_{6,1:5} \bar{\mathbf{v}}_{1:5} \\ &= \sigma' \rho(\tilde{\mathbf{A}})^k [\tilde{\mathbf{A}} \bar{\mathbf{v}}]_6 = \sigma' \rho(\tilde{\mathbf{A}})^{k+1} \bar{\mathbf{v}}_6 \leq \sigma' \rho(\tilde{\mathbf{A}})^k \bar{\mathbf{v}}_6. \end{aligned}$$

Then combining with (27), we have

$$\tilde{\mathbf{d}}_{1:6}^k \leq \sigma' \rho(\tilde{\mathbf{A}})^k \bar{\mathbf{v}}_{1:6}. \quad (28)$$

By (23c) and (28),

$$\begin{aligned} \mathbb{E} \left[ \left\| \mathbf{X}^{k+1} - \mathbf{X}^k \right\|_2^2 \right] &\leq \tilde{\mathbf{A}}_{7,1:6} \tilde{\mathbf{d}}_{1:6}^k \leq \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \tilde{\mathbf{A}}_{7,1:6} \tilde{\mathbf{v}}_{1:6}' \\ &= \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \left[ \tilde{\mathbf{A}} \tilde{\mathbf{v}}' \right]_7 = \sigma' \rho \left( \tilde{\mathbf{A}} \right)^{k+1} \tilde{\mathbf{v}}_7' \leq \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \tilde{\mathbf{v}}_7'. \end{aligned}$$

Then combining with (28), we have obtained

$$\tilde{\mathbf{d}}^k \leq \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \tilde{\mathbf{v}}'. \quad (29)$$

Using (23d) and (29), we have

$$\begin{aligned} \tilde{\mathbf{d}}_{1:4}^{k+1} &\leq \tilde{\mathbf{A}}_{1:4,:} \tilde{\mathbf{d}}^k \leq \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \tilde{\mathbf{A}}_{1:4,:} \tilde{\mathbf{v}}' \\ &= \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \left[ \tilde{\mathbf{A}} \tilde{\mathbf{v}}' \right]_{1:4} = \sigma' \rho \left( \tilde{\mathbf{A}} \right)^{k+1} \tilde{\mathbf{v}}_{1:4}', \end{aligned}$$

i.e., the statement holds for  $k+1$ . Therefore, by induction,  $\tilde{\mathbf{d}}_{1:4}^k \leq \sigma' \rho \left( \tilde{\mathbf{A}} \right)^k \tilde{\mathbf{v}}_{1:4}'$ , for any  $k \geq 0$ , which completes the proof. ■

*Remark 14:* In practice, the parameters  $\beta, \gamma, \eta$  and stepsize  $\alpha$  can be chosen in the same way as in CPP.

## V. NUMERICAL EXPERIMENTS

In this section, we compare the numerical performance of CPP and B-CPP with the Push-Pull/ $\mathcal{AB}$  method [25], [26]. In the experiments, we equip CPP and B-CPP with different compression operators and consider different graph topologies.

We consider the following decentralized  $\ell_2$ -regularized logistic regression problem:

$$\min_{\mathbf{x} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-\lambda_i \mathbf{z}_i^\top \mathbf{x}} \right) + \frac{\mu}{2} \|\mathbf{x}\|_2^2,$$

where  $(\mathbf{z}_i, \lambda_i) \in \mathbb{R}^p \times \{-1, +1\}$  is the training example possessed by agent  $i$ . The data is from QSAR biodegradation Data Set<sup>2</sup> [67], where each feature vector is of  $p = 41$  dimension. In our experiments, we set  $n = 20$  and  $\mu = 0.001$ . We construct the directed graphs  $\mathcal{G}_R$  and  $\mathcal{G}_C$  by adding  $d$  directed links to the undirected cycle, respectively. In our experiments below, we will consider the cases  $d = 5, 20, 50$ , representing sparse, normal and dense graphs respectively.

We consider two kinds of compression operators. The first one is the  $b$  bit 2-quantization given by

$$\text{COMPRESS}(\mathbf{x}) = \left( \|\mathbf{x}\|_2 \text{sign}(\mathbf{x}) 2^{1-b} \right) \left\lfloor \frac{2^{b-1} \|\mathbf{x}\|_2}{\|\mathbf{x}\|_2} + \mathbf{v} \right\rfloor,$$

where  $\mathbf{v} \in \mathbb{R}^p$  is a random vector picked uniformly from  $(0, 1)^p$  and we choose  $b = 2, 4, 6$  in the experiments. The second one is Rand- $k$  which is defined by

$$\text{COMPRESS}(\mathbf{x}) = \frac{p}{k} \hat{\mathbf{x}},$$

where  $\hat{\mathbf{x}}$  is obtained by randomly choosing  $k$  entries of  $\mathbf{x}$  and setting the rest entries to be 0. We will set  $k = 5, 10, 20$  respectively. Then, we hand-tune to find the optimal parameters for all the numerical cases below. In the figures, the  $y$ -axis represents the loss  $f(\bar{\mathbf{x}}^i) - f(\mathbf{x}^*)$  and  $x$ -axis is the iteration number or the number of transmitted bits.

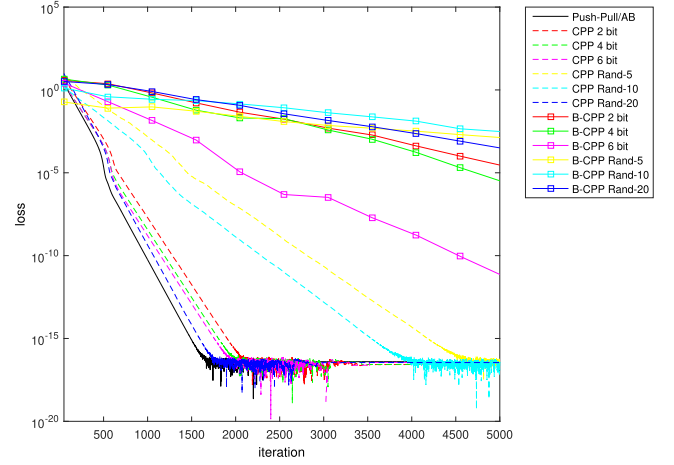


Fig. 1. Linear convergence of Push-Pull/ $\mathcal{AB}$ , CPP, and B-CPP with  $b$  bit quantization ( $b = 2, 4, 6$ ) and Rand- $k$  ( $k = 5, 10, 20$ ) compressors.

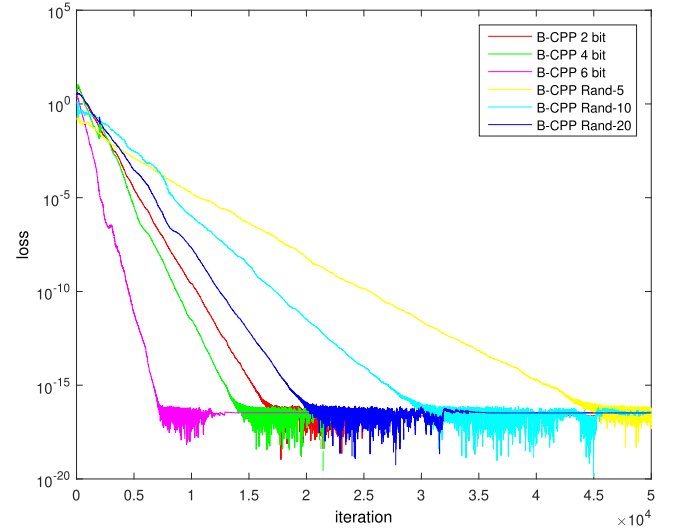


Fig. 2. Linear convergence of B-CPP with  $b$  bit quantization ( $b = 2, 4, 6$ ) and Rand- $k$  ( $k = 5, 10, 20$ ) compressors.

### A. Linear Convergence

The performance of Push-Pull/ $\mathcal{AB}$ , CPP and B-CPP is illustrated in Fig. 1. To see the performance of B-CPP more clearly, we plot the trajectories of B-CPP additionally at a larger scale in Fig. 2. The experiments illustrated in Fig. 1 and Fig. 2 are conducted on graphs with  $d = 20$ . The quantization with  $b = 2, 4, 6$  and Rand- $k$  with  $k = 5, 10, 20$  are considered. It can be seen from Fig. 1 and Fig. 2 that all the trajectories exhibit linear convergence, and the exact Push-Pull/ $\mathcal{AB}$  method is faster than CPP and B-CPP with any compression methods. This is reasonable as the compression operator induces additional errors compared to the exact method, and these additional errors could slow down the convergence. Meanwhile, as the values of  $b$  or  $k$  increases, both CPP and B-CPP speed up since the compression errors decrease. When  $b = 6$  or  $k = 20$ , the trajectories of CPP are very close to that of exact Push-Pull/ $\mathcal{AB}$ , which indicates that when the compression errors are small, they are no longer the bottleneck of convergence.

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/QSAR+biodegradation>

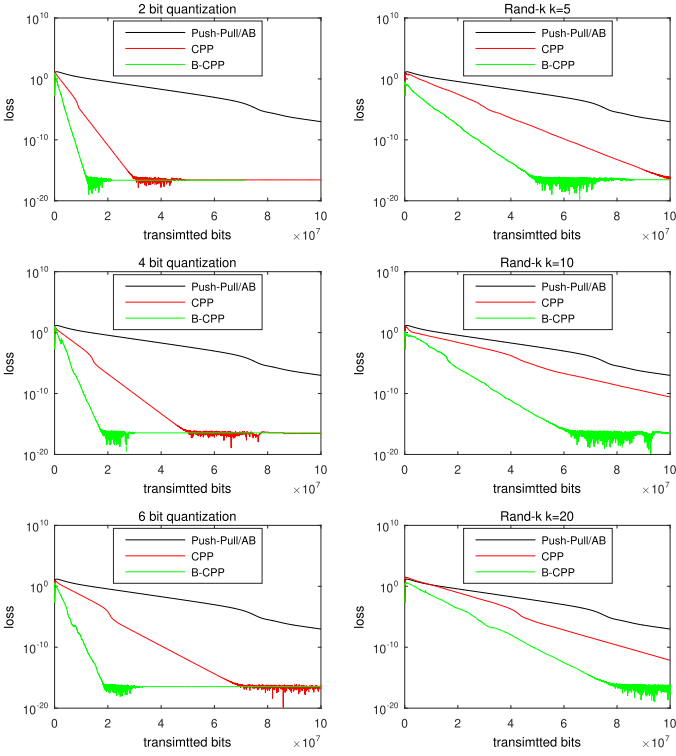


Fig. 3. Performance of Push-Pull/AB, CPP, B-CPP against the number of transmitted bits: the left column shows the results with quantization ( $b = 2, 4, 6$ ) and the right column shows the results with Rand- $k$  ( $k = 5, 10, 20$ ).

Within the same number of iterations, CPP outperforms B-CPP, and the trajectories of CPP are smoother than B-CPP. These results can be expected since CPP updates all the local variables in a single iteration, while in B-CPP, only one node updates with its neighbors. To guarantee the linear convergence, B-CPP requires smaller parameters and stepsizes. In addition, the mixing matrix at each iteration of B-CPP can be regarded as a stochastic matrix, which will induce additional variances.

### B. Communication Efficiency

When we compare the number of transmitted bits to reach certain levels of accuracy, CPP and B-CPP show their superiority. We consider the same settings as in Section V-A. We can see from all of the sub-figures of Fig. 3 that, to reach a high accuracy within about  $10^{-15}$ , the number of transmitted bits required by these methods have the ranking: B-CPP < CPP < Push-Pull/AB.

To see why CPP outperforms Push-Pull/AB, note that the vectors sent in CPP have been compressed, and hence the transmitted bits at each iteration are greatly reduced compared to Push-Pull/AB. Although the additional compression errors slow down the convergence, our design for CPP guarantees that the impact on the convergence rate is relatively small. Therefore, CPP is much more communication-efficient than the exact Push-Pull/AB method. Moreover, in all cases, B-CPP is much more communication-efficient than CPP. This is because when CPP converges, each agent will receive similar, or “overlapping” vectors from different neighbors, which impacts the communication-efficiency. By contrast, for B-CPP, only one

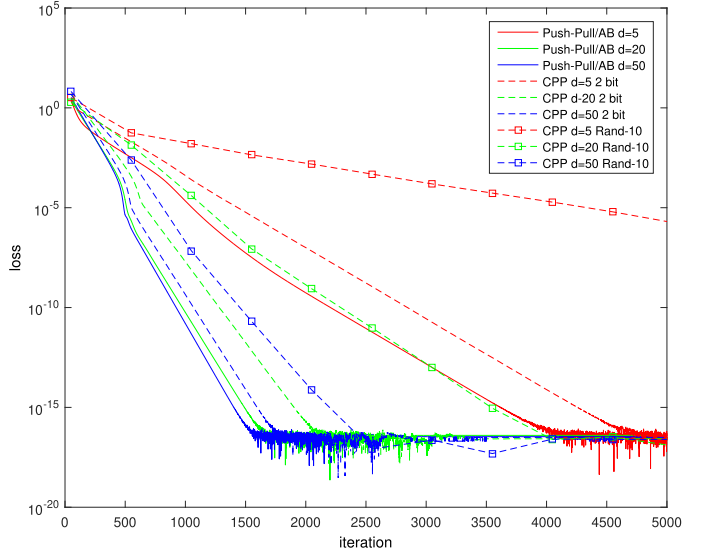


Fig. 4. Performance of CPP and Push-Pull/AB with different communication networks under both quantization and Rand- $k$  compressors.

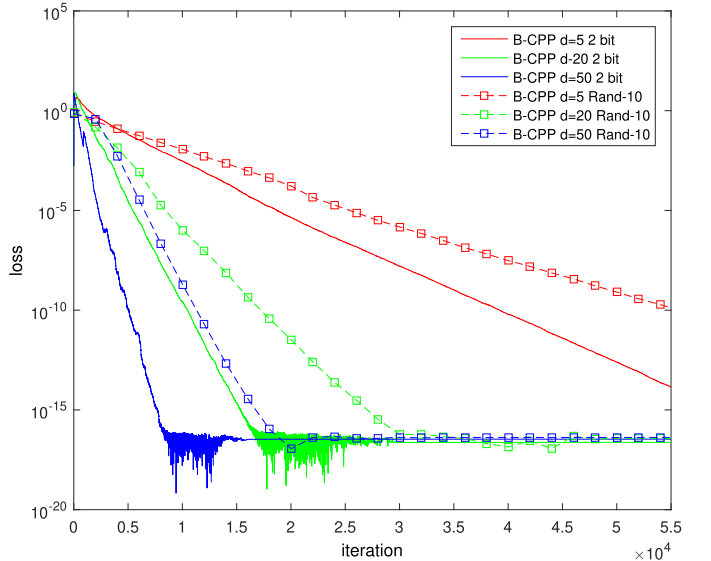


Fig. 5. Performance of B-CPP with different communication networks under both quantization and Rand- $k$  compressors.

agent wakes up in each iteration, and its information can be utilized by its neighbors more efficiently.

It is worth noting that for both CPP and B-CPP, the choices  $b = 2$  for quantization or  $k = 5$  for Rand- $k$  are more communication-efficient than  $b = 4, 6$  or  $k = 10, 20$ . This indicates that as the compression accuracy becomes smaller, its impact exhibits “marginal effects”. In other words, when the compression errors are not the bottleneck for the convergence, sacrificing the communication costs for faster convergence will reduce the communication efficiency.

### C. Different Topologies

We also examine the performance of CPP and B-CPP on communication networks with different levels of connectivity.

We consider  $d = 5, 10, 20$  respectively: as the value of  $d$  increases, the communication network has a better connectivity. As before, both the quantization and the Rand-k compressors are considered, and we show the performance of CPP and B-CPP separately for better clarity.

In Fig. 4 and Fig. 5, we can see that when  $d$  increases, the convergence of both CPP/B-CPP and Push-Pull/ $AB$  speed up. This is expected since better connectivity implies less consensus errors in each iteration, and the algorithms perform closer to the centralized gradient descent algorithm which is faster.

## VI. CONCLUSION

In this paper, we proposed two communication-efficient algorithms for decentralized optimization over a multi-agent network with general directed topology. First, we consider a novel communication-efficient gradient tracking based method, termed CPP, that combines the Push-Pull method with communication compression. CPP can be applied to a general class of unbiased compression operators and achieves linear convergence for strongly convex and smooth objective functions. Second, we consider a broadcast-like version of CPP (B-CPP) which also achieves linear convergence rate for strongly convex and smooth objective functions. B-CPP can be applied in an asynchronous broadcast setting and further reduce communication costs compared to CPP.

## ACKNOWLEDGMENT

Most work was done while the first author was visiting the School of Data Science, The Chinese University of Hong Kong, Shenzhen.

## REFERENCES

- [1] Z. Song, L. Shi, S. Pu, and M. Yan, "Compressed gradient tracking for decentralized optimization over general directed networks," 2021, *arXiv:2106.07243*.
- [2] K. Cohen, A. Nedić, and R. Srikant, "On projected stochastic gradient descent algorithm with weighted averaging for least squares regression," *IEEE Trans. Autom. Control*, vol. 62, no. 11, pp. 5974–5981, Nov. 2017.
- [3] A. I. Forrester, A. Söbester, and A. J. Keane, "Multi-fidelity optimization via surrogate modelling," in *Proc. Royal Soc. London A: Math. Physical Eng. Sci.*, 2007, pp. 3251–3269.
- [4] A. Nedić, A. Olshevsky, and C. A. Uribe, "Fast convergence rates for distributed non-Bayesian learning," *IEEE Trans. Autom. Control*, vol. 62, no. 11, pp. 5538–5553, Nov. 2017.
- [5] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Trans. Signal Process.*, vol. 60, no. 8, pp. 4289–4305, Aug. 2012.
- [6] S. Pu, A. Garcia, and Z. Lin, "Noise reduction by swarming in social foraging," *IEEE Trans. Autom. Control*, vol. 61, no. 12, pp. 4007–4013, Dec. 2016.
- [7] B. Baingana, G. Mateos, and G. B. Giannakis, "Proximal-gradient algorithms for tracking cascades over social networks," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 4, pp. 563–575, Aug. 2014.
- [8] K. Cohen, A. Nedić, and R. Srikant, "Distributed learning algorithms for spectrum sharing in spatial random access wireless networks," *IEEE Trans. Autom. Control*, vol. 62, no. 6, pp. 2854–2869, Jun. 2017.
- [9] G. Mateos and G. B. Giannakis, "Distributed recursive least-squares: Stability and performance analysis," *IEEE Trans. Signal Process.*, vol. 60, no. 7, pp. 3740–3754, Jul. 2012.
- [10] A. Nedić and J. Liu, "Distributed optimization for control," *Annu. Rev. Control, Robot., Auton. Syst.*, vol. 1, pp. 77–103, 2018.
- [11] A. Nedić and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.
- [12] Z. Li, W. Shi, and M. Yan, "A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates," *IEEE Trans. Signal Process.*, vol. 67, no. 17, pp. 4494–4506, Sep. 2019.
- [13] G. Qu and N. Li, "Harnessing smoothness to accelerate distributed optimization," *IEEE Trans. Control Netw. Syst.*, vol. 5, no. 3, pp. 1245–1260, Sep. 2018.
- [14] K. Scaman, F. Bach, S. Bubeck, L. Massoulié, and Y. T. Lee, "Optimal algorithms for non-smooth distributed optimization in networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 2745–2754.
- [15] W. Shi, Q. Ling, G. Wu, and W. Yin, "EXTRA: An exact first-order algorithm for decentralized consensus optimization," *SIAM J. Optim.*, vol. 25, no. 2, pp. 944–966, 2015.
- [16] C. A. Uribe, S. Lee, A. Gasnikov, A. Gasnikov, and A. Nedić, "A dual approach for optimal algorithms in distributed optimization over networks," *Optim. Methods Softw.*, vol. 36, no. 1, pp. 171–210, 2021, doi: [10.1080/10556788.2020.1750013](https://doi.org/10.1080/10556788.2020.1750013).
- [17] J. Xu, S. Zhu, Y. C. Soh, and L. Xie, "Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes," in *Proc. 54th IEEE Conf. Decis. Control*, 2015, pp. 2055–2060.
- [18] D. Kempe, A. Dobra, and J. Gehrke, "Gossip-based computation of aggregate information," in *Proc. 44th Annu. IEEE Symp. Foundations Comput. Sci.*, 2003, pp. 482–491.
- [19] A. Nedić and A. Olshevsky, "Distributed optimization over time-varying directed graphs," *IEEE Trans. Autom. Control*, vol. 60, no. 3, pp. 601–615, Mar. 2015.
- [20] C. Xi and U. A. Khan, "DEXTRA: A fast algorithm for optimization over directed graphs," *IEEE Trans. Autom. Control*, vol. 62, no. 10, pp. 4980–4993, Oct. 2017.
- [21] J. Zeng and W. Yin, "ExtraPush for convex smooth decentralized optimization over directed networks," *J. Comput. Math.*, vol. 35, no. 4, pp. 383–396, 2017.
- [22] A. Nedić, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," *SIAM J. Optim.*, vol. 27, no. 4, pp. 2597–2633, 2017.
- [23] Y. Tian, Y. Sun, and G. Scutari, "ASY-SONATA: Achieving linear convergence in distributed asynchronous multiagent optimization," in *Proc. 56th Annu. Allerton Conf. Commun., Control, Comput.*, 2018, pp. 543–551.
- [24] C. Xi, R. Xin, and U. A. Khan, "ADD-OPT: Accelerated distributed directed optimization," *IEEE Trans. Autom. Control*, vol. 63, no. 5, pp. 1329–1339, May 2018.
- [25] S. Pu, W. Shi, J. Xu, and A. Nedić, "Push-pull gradient methods for distributed optimization in networks," *IEEE Trans. Autom. Control*, vol. 66, no. 1, pp. 1–16, Jan. 2021.
- [26] R. Xin and U. A. Khan, "A linear algorithm for optimization over directed graphs with geometric convergence," *IEEE Contr. Syst. Lett.*, vol. 2, no. 3, pp. 315–320, Jul. 2018.
- [27] F. Saadatniaki, R. Xin, and U. A. Khan, "Decentralized optimization over time-varying directed graphs with row and column-stochastic matrices," *IEEE Trans. Autom. Control*, vol. 65, no. 11, pp. 4769–4780, Nov. 2020.
- [28] J. Zhang and K. You, "Fully asynchronous distributed optimization with linear convergence in directed networks," 2021, *arXiv:1901.08215*.
- [29] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "QSGD: Communication-efficient SGD via gradient quantization and encoding," *Adv. Neural Inf. Process. Syst.*, pp. 1709–1720, 2017.
- [30] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed optimisation for non-convex problems," in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 560–569.
- [31] A. Beznosikov, S. Horváth, P. Richtárik, and M. Safaryan, "On biased compression for distributed learning," 2020, *arXiv:2002.12410*.
- [32] S. P. Karimireddy, Q. Rehg, S. Stich, and M. Jaggi, "Error feedback fixes signsgd and other gradient compression schemes," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 3252–3261.
- [33] J. Liu et al., "Distributed learning systems with first-order methods," *Foundations Trends Databases*, vol. 9, no. 1, pp. 1–100, 2020.
- [34] X. Liu, Y. Li, J. Tang, and M. Yan, "A double residual compression algorithm for efficient distributed learning," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2020, pp. 133–143.
- [35] K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik, "Distributed learning with compressed gradient differences," 2019, *arXiv:1901.09269*.
- [36] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs," in *Proc. 15th Annu. Conf. Int. Speech Commun. Assoc.*, 2014, pp. 1058–1062.
- [37] S. U. Stich, "On communication compression for distributed optimization on heterogeneous data," 2020, *arXiv:2009.02388*.



- [38] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified SGD with memory," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 4447–4458.
- [39] H. Tang, C. Yu, X. Lian, T. Zhang, and J. Liu, "Doublesqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 6155–6165.
- [40] H. Xu *et al.*, "Compressed communication for distributed deep learning: Survey and quantitative evaluation," King Abdullah University of Science and Technology, Tech. Rep., 2020.
- [41] Y. Lin, S. Han, H. Mao, Y. Wang, and B. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," in *Proc. Int. Conf. Learn. Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=SkhQHMW0W>
- [42] L. Xiao, S. Boyd, and S. Lall, "A scheme for robust distributed sensor fusion based on average consensus," in *Proc. 4th Int. Symp. Inf. Process. Sensor Netw.*, 2005, pp. 63–70.
- [43] R. Carli, F. Fagnani, P. Frasca, T. Taylor, and S. Zampieri, "Average consensus on networks with transmission noise or quantization," in *Proc. Eur. Control Conf.*, 2007, pp. 1852–1857.
- [44] A. Nedic, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis, "Distributed subgradient methods and quantization effects," in *Proc. 47th IEEE Conf. Decis. Control*, 2008, pp. 4177–4184.
- [45] T. C. Aysal, M. J. Coates, and M. G. Rabbat, "Distributed average consensus with dithered quantization," *IEEE Trans. Signal Process.*, vol. 56, no. 10, pp. 4905–4918, Oct. 2008.
- [46] R. Carli, F. Fagnani, P. Frasca, and S. Zampieri, "Gossip consensus algorithms via quantized communication," *Automatica*, vol. 46, no. 1, pp. 70–80, 2010.
- [47] D. Yuan, S. Xu, H. Zhao, and L. Rong, "Distributed dual averaging method for multi-agent optimization with quantized communication," *Syst. Control Lett.*, vol. 61, no. 11, pp. 1053–1061, 2012.
- [48] A. Reiszadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "An exact quantized decentralized gradient descent algorithm," *IEEE Trans. Signal Process.*, vol. 67, no. 19, pp. 4934–4947, Oct. 2019.
- [49] R. Carli, F. Bullo, and S. Zampieri, "Quantized average consensus via dynamic coding/decoding schemes," *Int. J. Robust Nonlinear Control: IFAC-Affiliated J.*, vol. 20, no. 2, pp. 156–175, 2010.
- [50] T. T. Doan, S. T. Maguluri, and J. Romberg, "Fast convergence rates of distributed subgradient methods with adaptive quantization," *IEEE Trans. Autom. Control*, vol. 66, no. 5, pp. 2191–2205, 2020.
- [51] A. S. Berahas, C. Iakovidou, and E. Wei, "Nested distributed gradient methods with adaptive quantized communication," in *Proc. IEEE 58th Conf. Decis. Control*, 2019, pp. 1519–1525.
- [52] Z. Li, D. Kovalev, X. Qian, and P. Richtárik, "Acceleration for compressed gradient descent in distributed and federated optimization," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 5895–5904.
- [53] X. Liu, Y. Li, R. Wang, J. Tang, and M. Yan, "Linear convergent decentralized optimization with compression," in *Proc. Int. Conf. Learn. Representations*, 2021. [Online]. Available: <https://openreview.net/forum?id=84gjULz1t5>
- [54] A. Koloskova, S. Stich, and M. Jaggi, "Decentralized stochastic optimization and gossip algorithms with compressed communication," in *Proc. Int. Conf. Mach. Learn.*, 2019, pp. 3478–3487.
- [55] A. Koloskova, T. Lin, S. U. Stich, and M. Jaggi, "Decentralized deep learning with arbitrary communication compression," in *Proc. Int. Conf. Learn. Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=SkGCKrKvH>
- [56] H. Tang, S. Gan, C. Zhang, T. Zhang, and J. Liu, "Communication compression for decentralized training," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 7663–7673.
- [57] H. Tang *et al.*, "Deepsqueeze: Decentralization meets error-compensated compression," 2019, *arXiv:1907.07346*.
- [58] Y. Kajiyama, N. Hayashi, and S. Takai, "Linear convergence of consensus-based quantized optimization for smooth and strongly convex cost functions," *IEEE Trans. Autom. Control*, vol. 66, no. 3, pp. 1254–1261, Mar. 2020.
- [59] Z. Li, Y. Liao, K. Huang, and S. Pu, "Compressed gradient tracking for decentralized optimization with linear convergence," 2021, *arXiv:2103.13748*.
- [60] Y. Xiong, L. Wu, K. You, and L. Xie, "Quantized distributed gradient tracking algorithm with linear convergence in directed networks," 2021, *arXiv:2104.03649*.
- [61] J. Zhang, K. You, and L. Xie, "Innovation compression for communication-efficient distributed optimization with linear convergence," 2021, *arXiv:2105.06697*.
- [62] T. C. Aysal, M. E. Yildiz, A. D. Sarwate, and A. Scaglione, "Broadcast gossip algorithms for consensus," *IEEE Trans. Signal Process.*, vol. 57, no. 7, pp. 2748–2761, Jul. 2009.
- [63] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Trans. Inf. Theory*, vol. 52, no. 6, pp. 2508–2530, Jun. 2006.
- [64] S. Pu and A. Nedić, "Distributed stochastic gradient tracking methods," *Math. Program.*, vol. 187, no. 1, pp. 409–457, 2021.
- [65] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, U.K.: Cambridge Univ. Press, 2012.
- [66] A. Berman and R. J. Plemmons, *Nonnegative Matrices in the Mathematical Sciences*. Philadelphia, PA, USA: SIAM, 1994.
- [67] K. Mansouri, T. Ringsted, D. Ballabio, R. Todeschini, and V. Consonni, "Quantitative structure-activity relationship models for ready biodegradability of chemicals," *J. Chem. Inf. Model.*, vol. 53, no. 4, pp. 867–878, 2013.



**Zhuoqing Song** received the B.S. degree in mathematics and applied mathematics in 2019 from Fudan University, Shanghai, China, where he is currently working toward the Ph.D. degree in applied mathematics. His research interests include decentralized optimization, graph algorithms, and machine learning.



**Lei Shi** received the Ph.D. degree in mathematics from the City University of Hong Kong, Hong Kong, in 2010. He is currently a Professor with the School of Mathematical Sciences, Fudan University, Shanghai, China. His research interests include learning theory, approximation theory and optimization. He was jointly awarded by University of Science and Technology of China in 2010.



**Shi Pu** received the B.S. degree from Peking University, Beijing, China, in 2012, and the Ph.D. degree in systems engineering from the University of Virginia, Charlottesville, VA, USA, in 2016. He is currently an Assistant Professor with the School of Data Science, The Chinese University of Hong Kong, Shenzhen, China. He is also affiliated with Shenzhen Research Institute of Big Data. He was a Postdoctoral Associate with the University of Florida, Gainesville, FL, USA, Arizona State University, Tempe, AZ, USA, and Boston University, Boston, MA, USA, from 2016

to 2019. His research interests include distributed optimization, network science, machine learning, and game theory.



**Ming Yan** received the B.S. and M.S. degrees from the University of Science and Technology of China, Hefei, China, and the Ph.D. degree from the University of California, Los Angeles, Los Angeles, CA, USA, in 2012. He is currently an Associate Professor with the Department of Computational Mathematics, Science and Engineering and the Department of Mathematics, Michigan State University, East Lansing, MI, USA. His research interests include optimization methods and their applications in sparse recovery and regularized inverse problems, variational

methods for image processing, and parallel and distributed algorithms for solving big data problems.