

Frontiers of Information Technology & Electronic Engineering  
 www.jzus.zju.edu.cn; engineering.cae.cn; www.springerlink.com  
 ISSN 2095-9184 (print); ISSN 2095-9230 (online)  
 E-mail: jzus@zju.edu.cn



# A distributed stochastic optimization algorithm with gradient-tracking and distributed heavy-ball acceleration\*

Bihao SUN<sup>1</sup>, Jinhui HU<sup>1</sup>, Dawen XIA<sup>2</sup>, Huaqing LI<sup>†‡1</sup>

<sup>1</sup>Chongqing Key Laboratory of Nonlinear Circuits and Intelligent Information Processing,

College of Electronic and Information Engineering, Southwest University, Chongqing 400715, China

<sup>2</sup>College of Data Science and Information Engineering, Guizhou Minzu University, Guiyang 550025, China

<sup>†</sup>E-mail: huaqingli@swu.edu.cn

Received Nov. 8, 2020; Revision accepted Feb. 15, 2021; Crosschecked Apr. 1, 2021; Published online July 24, 2021

**Abstract:** Distributed optimization has been well developed in recent years due to its wide applications in machine learning and signal processing. In this paper, we focus on investigating distributed optimization to minimize a global objective. The objective is a sum of smooth and strongly convex local cost functions which are distributed over an undirected network of  $n$  nodes. In contrast to existing works, we apply a distributed heavy-ball term to improve the convergence performance of the proposed algorithm. To accelerate the convergence of existing distributed stochastic first-order gradient methods, a momentum term is combined with a gradient-tracking technique. It is shown that the proposed algorithm has better acceleration ability than GT-SAGA without increasing the complexity. Extensive experiments on real-world datasets verify the effectiveness and correctness of the proposed algorithm.

**Key words:** Distributed optimization; High-performance algorithm; Multi-agent system; Machine-learning problem; Stochastic gradient

<https://doi.org/10.1631/FITEE.2000615>

**CLC number:** TP14

## 1 Introduction

Distributed optimization has attracted much attention in many fields, such as wireless sensor networks (Yin et al., 2010; Cohen et al., 2017), machine learning (Xia and Wang, 2004; McMahan et al., 2017; Liu et al., 2019), and coordinated control (Han et al.,

2015; Cheng B and Li, 2019). It is evident that practical problems based on distributed optimization can be modeled as the minimization of the global objective function, i.e., solving the problem over a connected network consisting of  $n$  agents cooperatively over a common variable  $\mathbf{x}$ :

$$P0 : \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}),$$

where  $f_i$  is only known by agent  $i$  to exchange information with its neighbor nodes over an undirected network.

### 1.1 Literature review

In the literature, distributed optimization has attracted widespread attention based mainly on the

<sup>‡</sup> Corresponding author

\* Project supported by the Open Research Fund Program of Data Recovery Key Laboratory of Sichuan Province, China (No. DRN2001), the National Natural Science Foundation of China (Nos. 61773321 and 61762020), the Science and Technology Top-Notch Talents Support Project of Colleges and Universities in Guizhou Province, China (No. QJHKY2016065), the Science and Technology Foundation of Guizhou Province, China (No. QKHJC20181083), and the Science and Technology Talents Fund for Excellent Young of Guizhou Province, China (No. QKHPTRC20195669)

ORCID: Bihao SUN, <https://orcid.org/0000-0002-2790-2528>; Huaqing LI, <https://orcid.org/0000-0001-6310-8965>

© Zhejiang University Press 2021

Lagrangian, (sub)gradient descent, and Newton's methods. As for Lagrangian methods, Boyd et al. (2011) proposed a classic distributed alternating direction method of multipliers (ADMM) on the basis of the dual-domain over augmented Lagrangian functions, and many other distributed algorithms (Ling and Tian, 2010; Erseghe et al., 2011; Cheng S et al., 2014; Ling and Ribeiro, 2014; Wang B et al., 2018; Zhang et al., 2019; Guan et al., 2020) have been proposed based on this method. Although distributed ADMM with a constant step-size has the advantage of a linear convergence rate for a strongly convex function, it requires a large amount of calculation because each node has to optimize a local problem at each iteration. Unlike the Lagrangian method, Newton's method (Bertsekas and Gafni, 1983; Wei et al., 2013) is a common method to solve unconstrained optimization problems due to the advantage of fast convergence. To simplify the calculation process of Newton's method, the quasi-Newton method (Eisen et al., 2017) approximates the inverse Hessian matrix or the Hessian matrix using a positive definite matrix. The (sub)gradient descent method includes dual average (Duchi et al., 2012; Yuan et al., 2013; Matthews et al., 2016), distributed gradient descent (DGD) (Mateos et al., 2010; Xu et al., 2015; Nedić et al., 2017a, 2017b; Xi and Khan, 2017; Xin et al., 2019c), and stochastic gradient descent (SGD) methods (Zinkevich et al., 2010). In the pioneering work of Tsitsiklis et al. (1986), a framework for analyzing distributed computing models was developed. Nedić and Ozdaglar (2009) applied this method to the distributed convex optimization problem in the network, and achieved convergence of the strongly convex non-smooth structure. The SGD algorithm uses only a set of data from the sample to perform gradient descent, which improves the training speed of samples and reduces the calculation cost. Stochastic average gradient (SAG) (Schmidt et al., 2017; Wang Z and Li, 2020), SAGA (Defazio et al., 2014), and SVRG (Johnson and Zhang, 2013; Tan et al., 2016) methods have been proposed because the SGD algorithm can reach only an optimal interval instead of an optimal value. At each iteration, only one randomly selected gradient of a subfunction is evaluated at a node, and a variance-reduced stochastic averaging gradient technique is applied to approximate the gradient of the local objective function. SAG calculates a random vector as the average value of the

random gradient in the previous iterations, where in the  $k^{\text{th}}$  iteration, the algorithm stores derivatives to achieve exact convergence. In addition, based on SAG, the SAGA algorithm can directly support non-strongly convex problems, and is adaptive to any inherently strong convexity of the problem. Another algorithm, SVRG, which is performed in a loop, uses considerably amount of calculation to reduce the influence of noise, and is better than SGD at achieving more accurate convergence. GT-SAGA and GT-SVRG (Xin and Khan, 2020), which are based on SAGA and SVRG, respectively, achieve accelerated linear convergence by combining distributed stochastic gradient-tracking methods with variance-reduced techniques. In addition to the above methods (Johnson and Zhang, 2013; Defazio et al., 2014; Tan et al., 2016; Schmidt et al., 2017; Wang Z and Li, 2020), the  $\mathcal{AB}$  method based on row and column randomization (Xin and Khan, 2018), its acceleration method, the  $\mathcal{AB}m$  algorithm (Xin et al., 2020), and  $\mathcal{S}\text{-}\mathcal{AB}$  algorithm (Xin et al., 2019b) have also made significant contributions to solving the directed network.

## 1.2 Motivations and contributions

We find that plenty of works are interested in solving large-scale optimization problems with numerous and complex local objective functions. However, in distributed settings, distributed algorithms with exact gradient need massive calculation. Therefore, plenty of methods such as SAG, SAGA, and SVRG have been proposed to reduce the cost of full evaluation and retain the advantage of fast convergence under strongly convex and smooth conditions. Based on these works, GT-SAGA (Xin and Khan, 2020) combined a gradient-tracking technique with a variance-reduced technique that accelerates linear convergence. By introducing the above method to our work, we propose a new algorithm with a faster linear convergence rate and provide the concise proof. We summarize the following contributions in this paper:

1. Aiming to accelerate the convergence of stochastic first-order gradient methods (Zinkevich et al., 2010; Johnson and Zhang, 2013; Defazio et al., 2014; Tan et al., 2016; Schmidt et al., 2017; Lan et al., 2018; Wang Z and Li, 2020), we incorporate the momentum term combining with a gradient-tracking technique to achieve an accelerated convergence rate over undirected networks.

2. Compared with the distributed gradient methods (Ling and Tian, 2010; Boyd et al., 2011; Erseghe et al., 2011; Ling and Ribeiro, 2014; Wang B et al., 2018; Xin and Khan, 2018; Zhang et al., 2019; Xin et al., 2020) with deterministic gradients, the proposed algorithm requires much less computation when facing large-scale dataset optimization problems, because SGD method can update the model parameters with just a single training dataset.

3. In this paper, we present rigorous theoretical analysis for the proposed algorithm in Section 3. In addition, extensive experiments are provided to verify the correctness of theoretical analysis, such as the distributed logistic regression experiment based on a real-world dataset, a signal processing oriented least-square experiment, and distributed quadratic programming.

Notations used in this paper are summarized as follows: lowercase bold letters denote vectors, while uppercase denotes a matrix. The Euclidean norm of a vector is denoted as  $\|\cdot\|$ , and  $\|\|\cdot\|\|$  denotes the spectral norm of a matrix.  $\mathbf{1}_n$  means the  $n$ -dimensional column vector with all ones. For arbitrary matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{A} \otimes \mathbf{B}$  denotes the Kronecker product.  $\rho(\mathbf{A})$  is regarded as the spectral radius of  $\mathbf{A}$ .

## 2 Problem formulation and algorithm development

### 2.1 Distributed optimization problem

We consider  $n$  nodes communicating over an undirected network  $\mathcal{G}$ , so that it is capable of accessing a local cost function at each node  $i$ . The goal is to solve the following optimization problem with  $m$  nodes:

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

where each local cost  $f_i$  is averaged by  $p_i$  constituent functions  $\{f_i^h\}_{h=1}^{p_i}$ . The number of local samples is retained by agent  $i$ ,  $i \in U$ . The results of the proposed algorithm are on the basis of the following assumptions:

**Assumption 1** The global objective function  $f(\mathbf{x})$  is strongly convex with strong convexity parameter  $\mu$ ; i.e.,  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ , we have

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq \mu \|\mathbf{x} - \mathbf{y}\|^2. \quad (1)$$

**Assumption 2** Each local objective function  $f_i^h$  has Lipschitz continuous gradient with Lipschitz constant  $L_f > 0$ ; i.e.,  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ , we have

$$\|\nabla f_i^h(\mathbf{x}) - \nabla f_i^h(\mathbf{y})\| \leq L_f \|\mathbf{x} - \mathbf{y}\|, \quad (2)$$

where  $L_f > \mu > 0$ .

**Assumption 3** The weight matrix  $\mathbf{W}$  is doubly stochastic and associated with the undirected network  $\mathcal{G}$ .

Let  $\sigma$  indicate the spectral norm of the matrix  $\mathbf{W} - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T$ . According to Assumption 3, we have  $\sigma < 1$ , i.e.,  $\sigma = \|\mathbf{W} - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T\| < 1$ . In addition, we denote  $P := \max_i \{p_i\}$ ,  $p := \min_i \{p_i\}$ .

**Remark 1** These three assumptions are standard in recent literature.

### 2.2 Algorithm development

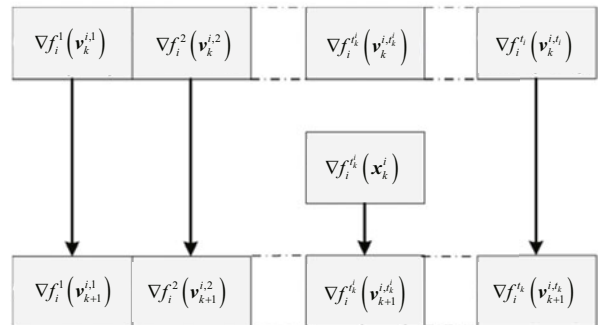
Previous work on the DGD method updates two vectors  $\mathbf{x}_k^i, \mathbf{y}_k^i \in \mathbb{R}^p$  at each node  $i$ , where  $\mathbf{x}_k^i$  is the local estimate of the global minimizer and  $\mathbf{y}_k^i$  is an auxiliary variable:

$$\begin{cases} \mathbf{x}_{k+1}^i = \sum_{j=1}^m w_{ij} \mathbf{x}_k^j - \alpha \mathbf{y}_k^i, \\ \mathbf{y}_{k+1}^i = \sum_{j=1}^m w_{ij} \mathbf{y}_k^j + \nabla f_i(\mathbf{x}_{k+1}^i) - \nabla f_i(\mathbf{x}_k^i). \end{cases} \quad (3)$$

Following Algorithm 1, we assume that the gradient of agent  $i$  at the  $k^{\text{th}}$  iteration is updated as shown in Fig. 1.

When  $\mathbf{g}_k^i$  is updated, the  $\nabla f_i^{t_k^i}(\mathbf{v}_{k+1}^{i,t_k^i})$  entry in the gradient table is replaced with  $\nabla f_i^{t_k^i}(\mathbf{x}_k^i)$ , while others remain unchanged:

$$\begin{aligned} \nabla f_i^j(\mathbf{v}_{k+1}^{i,j}) &\leftarrow \nabla f_i^j(\mathbf{v}_k^{i,j}), \\ j &= \{j \mid j = 1, 2, \dots, t_i, j \neq t_k^i\}. \end{aligned}$$



**Fig. 1** Gradient updating of agent  $i$  in the  $k^{\text{th}}$  iteration

Therefore, one can obtain the updated formula as follows:

$$\sum_{j=1}^{p_i} \nabla f_i^j(\mathbf{v}_{k+1}^{i,j}) = \sum_{j=1}^{p_i} \nabla f_i^j(\mathbf{v}_k^{i,j}) + \nabla f_i^{t_{k+1}^i}(\mathbf{x}_k^i) - \nabla f_i^{t_{k+1}^i}(\mathbf{v}_k^{i,t_{k+1}^i}). \tag{4}$$

**Algorithm 1** The proposed algorithm at each node  $i$

---

**Initialization:**  $\mathbf{x}_i^0; \mathbf{z}_{i,j}^1 = \mathbf{x}_j^0, \forall j \in \{1, 2, \dots, p_i\}; \alpha > 0; \{\tilde{w}_{ij}\}_{j=1}^m; \mathbf{y}_i^0 = \mathbf{g}_i^0 = \nabla f_i(\mathbf{x}_i^0)$   
**for**  $k = 0, 1, 2, \dots$  **do**  
 Choose  $t_{k+1}^i$  uniformly from local sample set  $\{1, 2, \dots, p_i\}$  at random  
 Update the variable  $\mathbf{g}_{k+1}^i$  as  

$$\mathbf{g}_{k+1}^i = \nabla f_i^{t_{k+1}^i}(\mathbf{x}_{k+1}^i) - \nabla f_i^{t_{k+1}^i}(\mathbf{v}_{k+1}^{i,t_{k+1}^i}) + \frac{1}{p_i} \sum_{j=1}^{p_i} \nabla f_i^j(\mathbf{v}_{k+1}^{i,j})$$
  
 Update the variable  $\mathbf{y}_{k+1}^i$  as  

$$\mathbf{y}_{k+1}^i = \sum_{j=1}^m w_{ij} \mathbf{y}_k^j + \mathbf{g}_{k+1}^i - \mathbf{g}_k^i$$
  
**if**  $j = t_{k+1}^i$  **then**  

$$\mathbf{v}_{k+1}^{i,j} = \mathbf{x}_k^{i,j}$$
  
**else**  

$$\mathbf{v}_{k+1}^{i,j} = \mathbf{v}_k^{i,j}$$
  
**end if**  
**end for**

---

### 3 Convergence analysis of the proposed algorithm

We aim to introduce unified analysis for the proposed algorithm that relies on the following dynamical system, with  $\mathbf{g}_0 = \mathbf{y}_0$  and  $\forall k \geq 0$ :

$$\mathbf{x}_{k+1} = \mathcal{W} \mathbf{x}_k - \alpha \mathbf{y}_k + \beta (\mathbf{x}_k - \mathbf{x}_{k-1}), \tag{5a}$$

$$\mathbf{y}_{k+1} = \mathcal{W} \mathbf{y}_k + \mathbf{g}_{k+1} - \mathbf{g}_k, \tag{5b}$$

where  $\mathcal{W} = \mathbf{W} \otimes \mathbf{I}_n$  and

$$\mathbf{x}_k = \begin{bmatrix} \mathbf{x}_k^1 \\ \mathbf{x}_k^2 \\ \vdots \\ \mathbf{x}_k^m \end{bmatrix}, \mathbf{y}_k = \begin{bmatrix} \mathbf{y}_k^1 \\ \mathbf{y}_k^2 \\ \vdots \\ \mathbf{y}_k^m \end{bmatrix}, \mathbf{g}_k = \begin{bmatrix} \mathbf{g}_k^1 \\ \mathbf{g}_k^2 \\ \vdots \\ \mathbf{g}_k^m \end{bmatrix}.$$

#### 3.1 Preliminaries

We denote  $\mathcal{F}_k$  as the history of the dynamical system generated by  $\{t_s^i\}_{i=1,2,\dots,m}^{s \leq k-1}$ . To simplify the subsequent analysis, we use the following notations:

$$\bar{\mathbf{x}}_k = \frac{1}{m} (\mathbf{1}_m^T \otimes \mathbf{I}_n) \mathbf{x}_k,$$

$$\bar{\mathbf{y}}_k = \frac{1}{m} (\mathbf{1}_m^T \otimes \mathbf{I}_n) \mathbf{y}_k,$$

$$\bar{\mathbf{g}}_k = \frac{1}{m} (\mathbf{1}_m^T \otimes \mathbf{I}_n) \mathbf{g}_k,$$

$$\nabla F(\mathbf{x}_k) = [(\nabla f_1(\mathbf{x}_k^1))^T, \dots, (\nabla f_m(\mathbf{x}_k^m))^T]^T,$$

$$\nabla \bar{F}(\mathbf{x}_k) = \frac{1}{m} (\mathbf{1}_m^T \otimes \mathbf{I}_n) \nabla F(\mathbf{x}_k).$$

**Lemma 1**  $\forall k \geq 0$ , it holds that

$$\mathbb{E}[\bar{\mathbf{y}}_k | \mathcal{F}_k] = \mathbb{E}[\bar{\mathbf{g}}_k | \mathcal{F}_k] = \frac{1}{m} (\mathbf{1}_m^T \otimes \mathbf{I}_n) \nabla F(\mathbf{x}_k). \tag{6}$$

**Proof** Multiplying  $\mathbf{1}_m^T$  by both sides of dynamical iterating Eq. (5b) yields

$$\bar{\mathbf{y}}_{k+1} = \bar{\mathbf{y}}_k + \bar{\mathbf{g}}_{k+1} - \bar{\mathbf{g}}_k. \tag{7}$$

Then, recursively updating Eq. (7) reduces to

$$\bar{\mathbf{y}}_k = \bar{\mathbf{g}}_k. \tag{8}$$

Thus, we obtain

$$\mathbb{E}[\bar{\mathbf{y}}_k | \mathcal{F}_k] = \mathbb{E}[\bar{\mathbf{g}}_k | \mathcal{F}_k]. \tag{9}$$

According to Algorithm 1, there exists

$$\begin{aligned} & \mathbb{E}[\mathbf{g}_k^i | \mathcal{F}_k] \\ &= \mathbb{E} \left[ \nabla f_i^{t_{k+1}^i}(\mathbf{x}_k^i) - \nabla f_i^{t_{k+1}^i}(\mathbf{v}_k^i) + \nabla f_i(\mathbf{v}_k^i) \mid \mathcal{F}_k \right] \\ &= \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{x}_k^i) - \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{v}_k^i) + \nabla f_i(\mathbf{v}_k^i) \\ &= \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{x}_k^i) - \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{v}_k^i) + \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{v}_k^i) \\ &= \frac{1}{p_i} \sum_{h=1}^{p_i} \nabla f_i^h(\mathbf{x}_k^i) \\ &= \nabla f_i(\mathbf{x}_k^i). \end{aligned}$$

Therefore, it is straightforward to obtain

$$\mathbb{E}[\bar{\mathbf{g}}_k | \mathcal{F}_k] = \frac{1}{m} (\mathbf{1}_m \otimes \mathbf{I}_n) \nabla f(\mathbf{x}_k). \tag{10}$$

Combining Eqs. (9) and (10) completes the proof.

**Lemma 2** Let Assumption 3 hold.  $\forall \mathbf{x} \in \mathbb{R}^{mn}$ , the inequality holds as follows:

$$\|\mathcal{W} \mathbf{x} - \mathcal{W}_\infty \mathbf{x}\| \leq \sigma \|\mathbf{x} - \mathcal{W}_\infty \mathbf{x}\|, \tag{11}$$

where  $0 < \sigma < 1$  is a constant.

**Proof** According to the definition of  $\mathcal{W}_\infty$ , we know that  $\mathcal{W}_\infty = \mathcal{W}\mathcal{W}_\infty$  and  $\|\mathcal{W} - \mathcal{W}_\infty\| < 1$ . Therefore,

$$\begin{aligned} & \|\mathcal{W}\mathbf{x} - \mathcal{W}_\infty\mathbf{x}\| \\ &= \|(\mathcal{W} - \mathcal{W}_\infty)(\mathbf{x} - \mathcal{W}_\infty\mathbf{x})\| \\ &\leq \|\mathcal{W} - \mathcal{W}_\infty\| \|\mathbf{x} - \mathcal{W}_\infty\mathbf{x}\| = \sigma \|\mathbf{x} - \mathcal{W}_\infty\mathbf{x}\|, \end{aligned}$$

where  $\mathcal{W}_\infty(\mathbf{x} - \mathcal{W}_\infty\mathbf{x}) = \mathbf{0}$  is used in the inequality, and  $\sigma = \|\mathcal{W} - \mathcal{W}_\infty\|$  is used to complete the proof.

**Lemma 3** Let Assumption 2 hold. Considering dynamical system (5),  $\forall k \geq 0$ , it holds that

$$\|\nabla \bar{F}(\mathbf{x}_k) - \nabla f(\bar{\mathbf{x}}_k)\| \leq \frac{L_f}{\sqrt{m}} \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|. \tag{12}$$

**Proof** Recalling the definition of  $\nabla \bar{F}(\mathbf{x}_k)$ , it holds that

$$\begin{aligned} & \|\nabla \bar{F}(\mathbf{x}_k) - \nabla f(\bar{\mathbf{x}}_k)\| \\ &= \left\| \frac{1}{m} \sum_{i=1}^m \nabla f_i(\mathbf{x}_k^i) - \frac{1}{m} \sum_{i=1}^m \nabla f_i(\bar{\mathbf{x}}_k) \right\| \\ &= \frac{1}{m} \left\| \sum_{i=1}^m \nabla f_i(\mathbf{x}_k^i) - \sum_{i=1}^m \nabla f_i(\bar{\mathbf{x}}_k) \right\| \\ &= \frac{1}{m} \sum_{i=1}^m \|\nabla f_i(\mathbf{x}_k^i) - \nabla f_i(\bar{\mathbf{x}}_k)\| \\ &\leq \frac{L_f}{m} \sum_{i=1}^m \|\mathbf{x}_k^i - \bar{\mathbf{x}}_k\| \\ &= \frac{L_f}{\sqrt{m}} \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|, \end{aligned}$$

where the Lipschitz continuity is used in the inequality. The proof is completed.

**Lemma 4** Suppose that Assumptions 1 and 2 hold.  $\forall \mathbf{x} \in \mathbb{R}^{mn}$ , if  $0 < \alpha < 1/L_f$ , it holds that

$$\|\mathbf{x} - \alpha \nabla f(\mathbf{x}) - \mathbf{x}^*\| \leq (1 - \mu\alpha) \|\mathbf{x} - \mathbf{x}^*\|, \tag{13}$$

where  $0 < \mu \leq L_f$ ,  $f(\mathbf{x})$  is the global objective function, and we define  $\mathbf{x}^*$  as the global optimum,  $\mathbf{x}^* = \mathbf{1}_m \otimes \tilde{\mathbf{x}}^* \in \mathbb{R}^{mn}$ .

### 3.2 Auxiliary results

We analyze the convergence of the general dynamical system using the following four formulas:

- (1)  $\mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \right]$ , the consensus error in the network;
- (2)  $\mathbb{E} \left[ \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right]$ , the optimal gap;

- (3)  $\mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right]$ , the state difference;

- (4)  $\mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty\mathbf{y}_k\|^2 \right]$ , the gradient tracking error.

**Lemma 5** Suppose that Assumption 3 holds. Considering the sequence  $\{\mathbf{x}_k\}$  yielded by dynamical system (5),  $\forall k \geq 0$ , it holds that

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathcal{W}_\infty\mathbf{x}_{k+1}\|^2 \right] \\ &\leq \frac{1-\sigma^2}{2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \right] + \frac{4\alpha^2}{1-\sigma^2} \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty\mathbf{y}_k\|^2 \right] \\ &\quad + \frac{4\beta^2}{1-\sigma^2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \end{aligned} \tag{14}$$

and

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathcal{W}_\infty\mathbf{x}_{k+1}\|^2 \right] \\ &\leq 2\sigma^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \right] + 4\alpha^2 \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty\mathbf{y}_k\|^2 \right] \\ &\quad + 4\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right], \end{aligned} \tag{15}$$

where  $0 < \mu \leq L_f$ , and  $f(\mathbf{x})$  is the global objective function.

**Proof** Notice that  $\mathcal{W}_\infty\mathcal{W} = \mathcal{W}\mathcal{W}_\infty = \mathcal{W}_\infty$  and  $\|\mathbf{I}_{mn} - \mathcal{W}_\infty\| = 1$ . There exists

$$\begin{aligned} & \|\mathbf{x}_{k+1} - \mathcal{W}_\infty\mathbf{x}_{k+1}\|^2 \\ &\leq (1+\eta) \|\mathcal{W}\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \\ &\quad + (1+\eta^{-1}) \left( 2\alpha^2 \|\mathbf{y}_k - \mathcal{W}_\infty\mathbf{y}_k\|^2 + 2\beta^2 \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right), \end{aligned}$$

where Young's inequality and Lemma 2 are used. Setting  $\eta$  as  $\frac{1-\sigma^2}{2\sigma^2}$  and 1 in the above inequality can lead to inequalities (14) and (15), respectively.

**Lemma 6** Let Assumptions 1–3 hold. Considering the sequence  $\{\mathbf{x}_k\}$  yielded by dynamical system (5),  $\forall k \geq 0$ , if  $0 < \alpha \leq 1/L_f$ , we can obtain

$$\begin{aligned} & \mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] \\ &\leq 4\alpha^2 L_f^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \right] + 4\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ &\quad + \frac{\alpha^2}{m} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right] + 2\mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \end{aligned} \tag{16}$$

and

$$\begin{aligned} & \mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] \\ &\leq (1-\mu\alpha) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + \frac{2\alpha L_f^2}{\mu} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty\mathbf{x}_k\|^2 \right] \\ &\quad + \frac{2\beta^2}{\mu\alpha} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + \frac{\alpha^2}{m^2} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right]. \end{aligned} \tag{17}$$

**Proof** According to dynamical iterating Eq. (5a) and recalling the definition of  $\bar{\mathbf{x}}_{k+1}$ , we know that  $\bar{\mathbf{x}}_{k+1} = \bar{\mathbf{x}}_k - \alpha \bar{\mathbf{y}}_k + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})$ :

$$\begin{aligned} & \mathbb{E} \left[ \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] \\ &= \mathbb{E} \left[ \|\bar{\mathbf{x}}_k - \alpha \bar{\mathbf{y}}_k + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1}) - \tilde{\mathbf{x}}^*\|^2 \right] \\ &= \mathbb{E} \left[ \|\alpha (\nabla f(\bar{\mathbf{x}}_k) - \bar{\mathbf{g}}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\|^2 \right] \\ & \quad + 2 \langle \bar{\mathbf{x}}_k - \alpha \nabla f(\bar{\mathbf{x}}_k) - \tilde{\mathbf{x}}^*, \alpha (\nabla f(\bar{\mathbf{x}}_k) - \bar{\mathbf{g}}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1}) \rangle \\ & \quad + \|\bar{\mathbf{x}}_k - \alpha \nabla f(\bar{\mathbf{x}}_k) - \tilde{\mathbf{x}}^*\|^2, \end{aligned}$$

where  $[\bar{\mathbf{y}}_k | \mathcal{F}_k] = \nabla \bar{F}(\mathbf{x}_k)$  is employed. Then we expand  $\mathbb{E} \left[ \|\alpha (\nabla f(\bar{\mathbf{x}}_k) - \bar{\mathbf{g}}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\|^2 \right]$  as follows:

$$\begin{aligned} & \mathbb{E} \left[ \|\alpha (\nabla f(\bar{\mathbf{x}}_k) - \bar{\mathbf{g}}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\|^2 \right] \\ &= \|\alpha [\nabla f(\bar{\mathbf{x}}_k) - \nabla \bar{F}(\mathbf{x}_k)] + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\|^2 \\ & \quad + \alpha^2 \mathbb{E} \left[ \|\nabla \bar{F}(\mathbf{x}_k) - \bar{\mathbf{g}}_k\|^2 \right] \\ &\leq \frac{2\alpha^2 L_f^2}{m} \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 + \frac{2\beta^2}{m} \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \\ & \quad + \frac{\alpha^2}{m^2} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right], \end{aligned}$$

since

$$\langle \alpha \nabla f(\bar{\mathbf{x}}_k) - \alpha \nabla \bar{F}(\mathbf{x}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1}), \alpha \nabla \bar{F}(\mathbf{x}_k) - \alpha \bar{\mathbf{g}}_k \rangle = 0.$$

Next, we simplify

$$\begin{aligned} & 2 \langle \bar{\mathbf{x}}_k - \alpha \nabla f(\bar{\mathbf{x}}_k) - \tilde{\mathbf{x}}^*, \\ & \quad \alpha \nabla f(\bar{\mathbf{x}}_k) - \alpha \bar{F}(\mathbf{x}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1}) \rangle \\ &= 2 \|\bar{\mathbf{x}}_k - \alpha \nabla f(\bar{\mathbf{x}}_k) - \tilde{\mathbf{x}}^*\| \\ & \quad \|\alpha \nabla f(\bar{\mathbf{x}}_k) - \alpha \bar{F}(\mathbf{x}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\| \\ &\leq \frac{1}{\lambda} \|\bar{\mathbf{x}}_k - \alpha \nabla f(\bar{\mathbf{x}}_k) - \tilde{\mathbf{x}}^*\|^2 \\ & \quad + \lambda \|\alpha \nabla f(\bar{\mathbf{x}}_k) - \alpha \bar{F}(\mathbf{x}_k) + \beta (\bar{\mathbf{x}}_k - \bar{\mathbf{x}}_{k-1})\|^2. \end{aligned}$$

Setting  $\lambda$  as  $(1 - \mu\alpha)/(\mu\alpha)$  and 1 can lead to inequalities (16) and (17), respectively, which completes the proof.

**Lemma 7** Considering the sequences  $\{\mathbf{x}_k\}_{k \geq 0}$  and  $\{\mathbf{y}_k\}_{k \geq 0}$  yielded by dynamical system (5),  $\forall k \geq 0$ , if  $0 < \alpha \leq 1/L_f$ , the inequality holds as follows:

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \right] \\ &\leq (16 + 16\alpha^2 L_f^2) \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \end{aligned}$$

$$\begin{aligned} & + 16\alpha^2 L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + 16\alpha^2 \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right] \\ & + 16\alpha^2 \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] + 2\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right]. \end{aligned} \tag{18}$$

**Proof** From dynamical iterating Eq. (5a), we have

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \right] \\ &\leq \mathbb{E} \left[ \|\mathcal{W} - \mathbf{I}_{mn}\| (\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k) - \alpha \mathbf{y}_k + \beta (\mathbf{x}_k - \mathbf{x}_{k-1}) \|^2 \right] \\ &\leq 2 \mathbb{E} \left[ \|\mathcal{W} - \mathbf{I}_{mn}\| (\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k) - \alpha \mathbf{y}_k \|^2 \right] \\ & \quad + 2\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ &\leq 16 \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] + 4\alpha^2 \mathbb{E} \left[ \|\mathbf{y}_k\|^2 \right] \\ & \quad + 2\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right]. \end{aligned} \tag{19}$$

Let Assumption 2 hold. Then one can obtain

$$\begin{aligned} & \|\mathbf{y}_k\| \\ &\leq \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\| + \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\| + L_f \|\mathbf{x}_k - \mathbf{x}^*\| \\ &\leq \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\| + \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\| \\ & \quad + L_f \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\| + \sqrt{m} L_f \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|, \end{aligned}$$

where we used  $\bar{\mathbf{y}}_k = \bar{\mathbf{g}}_k, \forall k \geq 0$ , for squaring the above inequality to obtain

$$\begin{aligned} \|\mathbf{y}_k\|^2 &\leq 4 \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 + 4 \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \\ & \quad + 4L_f^2 \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 + 4L_f^2 m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2. \end{aligned} \tag{20}$$

Therefore, inequality (18) can be proved by plugging inequality (20) into inequality (19).

**Lemma 8** Let Assumptions 2 and 3 hold. Regarding the sequence  $\{\mathbf{y}_k\}$  yielded by dynamical iterating Eq. (5b),  $\forall k \geq 0$ , if  $0 < \alpha \leq 1/(4\sqrt{2}L_f)$ , it holds that

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{y}_{k+1} - \mathcal{W}_\infty \mathbf{y}_{k+1}\|^2 \right] \\ &\leq \left( \frac{1 + \sigma^2}{2} + \frac{32\alpha^2 L_f^2}{1 - \sigma^2} \right) \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] \\ & \quad + \frac{33L_f^2}{1 - \sigma^2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{4L_f^2 \beta^2}{1 - \sigma^2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ & \quad + \frac{4}{1 - \sigma^2} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_{k+1}) - \mathbf{g}_{k+1}\|^2 \right] \\ & \quad + \frac{L_f^2}{1 - \sigma^2} \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \quad + \frac{4 + 32\alpha^2 L_f^2}{1 - \sigma^2} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right]. \end{aligned} \tag{21}$$

**Proof** From inequality (21), we can obtain

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{y}_{k+1} - \mathcal{W}_\infty \mathbf{y}_{k+1}\|^2 \right] \\ = & \mathbb{E} \left[ \|(\mathcal{W} - \mathcal{W}_\infty)(\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k) \right. \\ & \left. + (\mathbf{I}_{mn} - \mathcal{W}_\infty)(\mathbf{g}_{k+1} + \mathbf{g}_k)\|^2 \right]. \end{aligned}$$

Then we employ  $\|\mathbf{I}_{mn} - \mathcal{W}_\infty\| = 1$  and Young's inequality to set  $\eta$  as  $(1 - \sigma^2)/(2\sigma^2)$  in the above equality, and obtain

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{y}_{k+1} - \mathcal{W}_\infty \mathbf{y}_{k+1}\|^2 \right] \\ \leq & \frac{1 + \sigma^2}{2} \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] \quad (22) \\ & + \frac{2}{1 - \sigma^2} \mathbb{E} \left[ \|\mathbf{g}_{k+1} - \mathbf{g}_k\|^2 \right]. \end{aligned}$$

Next we expand  $\mathbb{E} \left[ \|\mathbf{g}_{k+1} - \mathbf{g}_k\|^2 \right]$  as

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{g}_{k+1} - \mathbf{g}_k\|^2 \right] \\ \leq & 2\mathbb{E} \left[ \|\mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1})\|^2 \right] + 2\mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right] \\ & + 2\mathbb{E} \left[ \|\nabla F(\mathbf{x}_{k+1}) - \nabla F(\mathbf{x}_k)\|^2 \right] \\ & + 4\mathbb{E} \left[ \langle \mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1}), \nabla F(\mathbf{x}_k) - \mathbf{g}_k \rangle \right] \\ \leq & 2L_f^2 \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \right] + 2\mathbb{E} \left[ \|\mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1})\|^2 \right] \\ & + 2\mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right], \quad (23) \end{aligned}$$

where in inequality (23), we use the fact that

$$\mathbb{E} \left[ \langle \mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1}), \nabla F(\mathbf{x}_k) - \mathbf{g}_k \rangle \right] = 0.$$

Using the bound of inequality (23) in inequality (22), we obtain inequality (21), which completes the proof.

**Lemma 9** Considering  $\{r_k\}$ , the following holds:

$$\begin{aligned} \mathbb{E} [r_{k+1}] \leq & \left( 1 - \frac{1}{P} \right) \mathbb{E} [r_k] + \frac{2}{p} \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ & + \frac{2}{p} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right]. \quad (24) \end{aligned}$$

**Proof** From Algorithm 1, if  $j = t_{k+1}^i$ , then  $\mathbf{v}_{k+1}^{i,j} = \mathbf{x}_{k+1}^{i,j}$  and  $\mathbf{v}_{k+1}^{i,j} = \mathbf{v}_k^{i,j}$ . Then we can obtain

$$\begin{aligned} & \mathbb{E} [r_{k+1}] \\ = & \frac{1}{p_i} \sum_{j=1}^{p_i} \mathbb{E} \left[ \|\mathbf{v}_{k+1}^{i,j} - \tilde{\mathbf{x}}^*\|^2 \right] \\ = & \frac{1}{p_i} \sum_{j=1}^{p_i} \mathbb{E} \left[ \left( 1 - \frac{1}{p_i} \right) \|\mathbf{v}_{k+1}^{i,j} - \tilde{\mathbf{x}}^*\|^2 + \frac{1}{p_i} \|\mathbf{x}_k^i - \tilde{\mathbf{x}}^*\|^2 \right] \end{aligned}$$

$$\begin{aligned} & = \left( 1 - \frac{1}{p_i} \right) r_k^i + \frac{1}{p_i} \|\mathbf{x}_k^i - \tilde{\mathbf{x}}^*\|^2 \\ \leq & \left( 1 - \frac{1}{P} \right) r_k^i + \frac{2}{p} \|\mathbf{x}_k^i - \bar{\mathbf{x}}_k\|^2 + \frac{2}{p} \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2. \end{aligned}$$

In the next lemma, we derive an upper bound on  $\mathbb{E} \left[ \|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\|^2 \right]$ , and then define  $r_k^i = \frac{1}{p_i} \sum_{j=1}^{p_i} \|\mathbf{v}_k^{i,j} - \tilde{\mathbf{x}}^*\|^2$  and  $r_k = \sum_i^m r_k^i$ , recalling that  $P = \max_i \{p_i\}$  and  $p = \min_i \{p_i\}$ .

**Lemma 10** Let Assumption 2 hold. Considering the sequence  $\{\mathbf{g}_k\}$  yielded by dynamical system (5),  $\forall k \geq 0$ , the inequality is obtained as follows:

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\|^2 \right] \\ \leq & 4L_f^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \quad (25) \\ & + 4L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + 2L_f^2 \mathbb{E} [r_k]. \end{aligned}$$

**Proof** Recalling Algorithm 1, we obtain

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{g}_k^i - \nabla f_i(\mathbf{x}_k^i)\|^2 \right] \\ = & \mathbb{E} \left[ \left\| \nabla f_i^{q_k^i}(\mathbf{x}_k^i) - \nabla f_i^{q_k^i}(\mathbf{v}_k^i) - \nabla f_i(\mathbf{x}_k^i) \right. \right. \\ & \left. \left. + \frac{1}{p_i} \sum_{j=1}^{p_i} \nabla f_i^j(\mathbf{v}_k^{i,j}) \right\|^2 \right] \\ \leq & \mathbb{E} \left[ \left\| \nabla f_i^{t_k^i}(\mathbf{x}_k^i) - \nabla f_i^{t_k^i}(\mathbf{v}_k^i) \right\|^2 \right] \\ = & \frac{1}{p_i} \sum_{j=1}^{p_i} \left\| \nabla f_i^j(\mathbf{x}_k^i) - \nabla f_i^j(\tilde{\mathbf{x}}^*) + \nabla f_i^j(\tilde{\mathbf{x}}^*) - \nabla f_i^j(\mathbf{v}_k^i) \right\|^2 \\ \leq & 2 \|\nabla f_i(\mathbf{x}_k^i) - \nabla f_i(\tilde{\mathbf{x}}^*)\|^2 + 2 \|\nabla f_i^j(\tilde{\mathbf{x}}^*) - \nabla f_i^j(\mathbf{v}_k^i)\|^2 \\ \leq & 2L_f^2 \|\mathbf{x}_k^i - \tilde{\mathbf{x}}^*\|^2 + 2L_f^2 \|\mathbf{v}_k^i - \tilde{\mathbf{x}}^*\|^2 \\ \leq & 4L_f^2 \|\mathbf{x}_k^i - \bar{\mathbf{x}}_k\|^2 + 4L_f^2 \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 + 2L_f^2 r_k^i. \end{aligned}$$

**Corollary 1** Let Assumptions 2 and 3 hold. Considering the sequence  $\{\mathbf{g}_k\}_{k \geq 0}$  yielded by dynamical system (5),  $\forall k \geq 0$ , we can obtain

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1})\|^2 \right] \\ \leq & \left( 8L_f^2 \sigma^2 + \frac{4L_f^2}{p} + 16L_f^2 \alpha^2 L_f^2 \right) \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & + 16L_f^2 \alpha^2 \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] + 32L_f^2 \beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ & + \left( 2 + \frac{1}{p} \right) 4L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + \left( 1 - \frac{1}{P} \right) 2L_f^2 \mathbb{E} [r_k] \\ & + \frac{4L_f^2 \alpha^2}{m} \mathbb{E} \left[ \|\nabla F(\mathbf{x}_k) - \mathbf{g}_k\|^2 \right]. \quad (26) \end{aligned}$$

**Proof** From Lemma 10, it is straightforward to obtain

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1})\|^2 \right] \\ & \leq 4L_f^2 \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathcal{W}_\infty \mathbf{x}_{k+1}\|^2 \right] \\ & \quad + 4L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] + 2L_f^2 \mathbb{E} [r_{k+1}]. \end{aligned}$$

Then, we can expand the above inequality by plugging inequalities (15), (16), and (24), from Lemmas 5, 6, and 9, respectively, which completes the proof.

### 3.3 Main results

To summarize, the bound of the gradient variance is obtained to process the inequality, which is obtained from dynamical system (5), and then we start to derive the convergence rate of the proposed algorithm. First, the upper bounds on  $\mathbb{E} \left[ \|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\|^2 \right]$  and  $\mathbb{E} \left[ \|\mathbf{g}_{k+1} - \nabla F(\mathbf{x}_{k+1})\|^2 \right]$  are obtained from Lemma 10 and Corollary 1, respectively, and then  $\forall k \geq 0$ , we have

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{y}_{k+1} - \mathcal{W}_\infty \mathbf{y}_{k+1}\|^2 \right] \\ & \leq \left( 32\sigma^2 + \frac{16}{p} + 55 + \frac{2}{m} \right) \frac{L_f^2}{1-\sigma^2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{132L_f^2\beta^2}{1-\sigma^2} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ & \quad + \frac{L_f^2}{1-\sigma^2} \left( 18 - \frac{8}{P} + \frac{1}{m} \right) \mathbb{E} [r_k] \\ & \quad + \left( \frac{1+\sigma^2}{2} + \frac{96\alpha^2 L_f^2}{1-\sigma^2} \right) \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] \\ & \quad + \frac{L_f^2}{1-\sigma^2} \left( 53 + \frac{16}{p} + \frac{2}{m} \right) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right]. \end{aligned} \tag{27}$$

Second, we refine the following inequality by applying the upper bound on  $\mathbb{E} \left[ \|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\|^2 \right]$  from Lemma 10:  $\forall k \geq 0$ ,

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \right] \\ & \leq (16 + 80\alpha^2 L_f^2) \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + 2\beta^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + 16\alpha^2 \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] \\ & \quad + 80\alpha^2 L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + 32\alpha^2 L_f^2 \mathbb{E} [r_k]. \end{aligned} \tag{28}$$

Next, combining inequalities (17) and (25) yields

$$\mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right]$$

$$\begin{aligned} & \leq (1 - \mu\alpha) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \quad + \frac{2\alpha L_f^2}{\mu} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{2\beta^2}{\mu\alpha} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + \frac{\alpha^2}{m} 4L_f^2 \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + 4\frac{\alpha^2}{m} L_f^2 \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] + 2\frac{\alpha^2}{m} L_f^2 \mathbb{E} [r_k] \\ & \leq \left( 1 - \mu\alpha + \frac{4L_f^2\alpha^2}{m} \right) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \quad + \alpha L_f^2 \left( \frac{2}{\mu} + \frac{4\alpha}{m} \right) \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{2\beta^2}{\mu\alpha} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + 2\frac{\alpha^2}{m} L_f^2 \mathbb{E} [r_k] \\ & \leq \left( 1 - \mu\alpha + \frac{4L_f^2\alpha^2}{m} \right) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \quad + \alpha L_f^2 \left( \frac{2}{\mu} + \frac{4\alpha}{m} \right) \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{2\beta^2}{\mu\alpha} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + \frac{2L_f^2\alpha^2}{m} \mathbb{E} [r_k]. \end{aligned}$$

If  $0 < \alpha \leq 1/(2\mu)$ , then  $2/\mu + 4\alpha/m \leq 4/\mu$ ; if  $0 < \alpha \leq \mu m/(8L^2)$ , then  $1 - \mu\alpha + 4L_f^2\alpha^2/m \leq 1 - \mu\alpha/2$ . Therefore, we can obtain the following inequality:

$$\begin{aligned} & \mathbb{E} \left[ m \|\bar{\mathbf{x}}_{k+1} - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \leq \left( 1 - \frac{\mu\alpha}{2} \right) \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ & \quad + \frac{4\alpha L_f^2}{\mu} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ & \quad + \frac{2\beta^2}{\mu\alpha} \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] + \frac{2L_f^2\alpha^2}{m} \mathbb{E} [r_k]. \end{aligned} \tag{29}$$

Now we combine inequalities (14), (24), (27), and (28) with inequality (29) as a linear matrix to prepare for the subsequent derivation.

**Proposition 1** Let Assumptions 1–3 hold.  $\forall k \geq 1$ , the following inequality holds entry-wise:

$$\mathbf{J}_{k+1} \leq \mathbf{G}_\alpha \mathbf{J}_k, \tag{30}$$

where  $\mathbf{J}_k \in \mathbb{R}^5$  is given by

$$\mathbf{J}_k = \begin{bmatrix} \mathbb{E} \left[ \|\mathbf{x}_k - \mathcal{W}_\infty \mathbf{x}_k\|^2 \right] \\ \mathbb{E} \left[ \|\mathbf{x}_k - \mathbf{x}_{k-1}\|^2 \right] \\ \mathbb{E} \left[ \|\mathbf{y}_k - \mathcal{W}_\infty \mathbf{y}_k\|^2 \right] \\ \mathbb{E} \left[ m \|\bar{\mathbf{x}}_k - \tilde{\mathbf{x}}^*\|^2 \right] \\ \mathbb{E} [r_k] \end{bmatrix},$$



and  $\mathbf{G}_\alpha \in \mathbb{R}^{5 \times 5}$  is given at the bottom of this page, where  $a_1 = 16/p + 55 + 2/m, a_2 = 16/p + 53 + 2/m, a_3 = 18 - 8/P + 1/m$ . Obviously, it is feasible to obtain the range of  $\alpha$  to satisfy  $\rho(\mathbf{G}_\alpha) < 1$ . To achieve this, we go to the next lemma.

**Lemma 11** Let  $\mathbf{H} \in \mathbb{R}^{d \times d}$  and  $\mathbf{x} \in \mathbb{R}^d$  be non-negative matrix and positive vector, respectively. For  $\gamma > 0$ , if  $\mathbf{H}\mathbf{x} \leq \gamma\mathbf{x}$ , then  $\rho(\mathbf{H}) \leq \|\mathbf{H}\|_\infty \leq \gamma$ .

**Theorem 1** Let Assumptions 1–3 hold. If step-size  $\alpha$  satisfies  $0 < \alpha < \frac{p(1-\sigma^2)^2}{P \cdot 549QL_f}$ , the proposed algorithm is linearly convergent.

**Proof** Let positive vector  $\boldsymbol{\delta} = [\delta_1, \delta_2, \delta_3, \delta_4, \delta_5]^T$ . The linear matrix inequality holds as follows:

$$\mathbf{G}_\alpha \boldsymbol{\delta} < \boldsymbol{\delta},$$

which can be equivalently written as the following inequalities:

$$0 < \frac{1-\sigma^2}{2} - \frac{4\beta^2}{1-\sigma^2} \frac{\delta_2}{\delta_1} - \frac{4\alpha^2}{1-\sigma^2} \frac{\delta_3}{\delta_1}, \quad (31)$$

$$0 < (1-2\beta^2) - (16+80\alpha^2L_f^2) \frac{\delta_1}{\delta_2} - 16\alpha^2 \frac{\delta_3}{\delta_2} - 80\alpha^2L_f^2 \frac{\delta_4}{\delta_2} - 32\alpha^2L_f^2 \frac{\delta_5}{\delta_2}, \quad (32)$$

$$0 < \left( \frac{1-\sigma^2}{2} - \frac{96\alpha^2L_f^2}{1-\sigma^2} \right) \delta_3 - \frac{L_f^2}{1-\sigma^2} (32\sigma^2 + a_1) \delta_1 - \frac{L_f^2}{1-\sigma^2} a_2 \delta_4 - \frac{L_f^2}{1-\sigma^2} a_3 \delta_5, \quad (33)$$

$$0 < \frac{\mu\alpha}{2} - \frac{4\alpha L_f^2}{\mu} \frac{\delta_1}{\delta_4} - \frac{2\beta^2}{\mu\alpha} \frac{\delta_2}{\delta_4} - \frac{2L_f^2\alpha^2}{m} \frac{\delta_5}{\delta_4}, \quad (34)$$

$$0 < \frac{1}{P} - \frac{2}{p} \frac{\delta_1}{\delta_5} - \frac{2}{p} \frac{\delta_4}{\delta_5}. \quad (35)$$

It is evident that for the right-hand sides (RHSs) of inequalities (31)–(35) to be positive, we can bound the range of  $\alpha$ . First, we are supposed to fix vector  $\boldsymbol{\delta}$ , which is positive and independent of  $\alpha$  and  $\beta$ . We

can set  $\delta_1 = 1, \delta_2 = 17$ , and  $Q = L_f/\mu$ , and the following inequality is obtained from inequality (34):

$$\frac{2\beta^2}{\mu\alpha} \frac{\delta_2}{\delta_4} < \frac{\mu\alpha}{2} - \frac{4\alpha L_f^2}{\mu} \frac{\delta_1}{\delta_4} - \frac{2L_f^2\alpha^2}{m} \frac{\delta_5}{\delta_4}. \quad (36)$$

Because the RHS of inequality (36) is positive, it is straightforward that

$$\delta_4 > \frac{8L_f^2}{\mu^2} \delta_1 = 8Q^2. \quad (37)$$

Therefore, we set  $\delta_4 = 80Q^2$ , and obtain the following inequality because the RHS of inequality (35) is positive:

$$\delta_5 > \frac{2P}{p} (\delta_1 + \delta_4) = \frac{2P}{p} + \frac{160P}{p} Q^2. \quad (38)$$

We denote  $\delta_5 = 164PQ^2/p$ . Then, from inequality (33), we obtain

$$\frac{132L_f^2\beta^2}{1-\sigma^2} \delta_2 + \frac{96\alpha^2L_f^2}{1-\sigma^2} \delta_3 < \frac{1-\sigma^2}{2} \delta_3 - \left( \frac{105L_f^2}{1-\sigma^2} \delta_1 + \frac{71L_f^2}{1-\sigma^2} \delta_4 + \frac{19L_f^2}{1-\sigma^2} \delta_5 \right). \quad (39)$$

Because the RHS of inequality (39) is positive, the following inequality is obtained:

$$\delta_3 > \frac{2L_f^2}{(1-\sigma^2)^2} \left( 105 + 5680Q^2 + \frac{3116PQ^2}{p} \right). \quad (40)$$

Because  $105 + 5680Q^2 + 3116PQ^2/p < 8901PQ^2/p$ , we can set  $\delta_3 = \frac{18050L_f^2}{(1-\sigma^2)^2} \frac{PQ^2}{p}$ . Until now, we have fixed values of  $\delta_1$ – $\delta_5$  as follows:

$$\delta_1 = 1, \delta_2 = 17, \delta_3 = \frac{18050L_f^2}{(1-\sigma^2)^2} \frac{PQ^2}{p}, \quad (41)$$

$$\delta_4 = 80Q^2, \delta_5 = \frac{164PQ^2}{p}.$$

The range of  $\alpha$  that satisfies forward inequalities can

$$\mathbf{G}_\alpha = \begin{bmatrix} \frac{1-\sigma^2}{2} & \frac{4\beta^2}{1-\sigma^2} & \frac{4\alpha^2}{1-\sigma^2} & 0 & 0 \\ 16 + 80\alpha^2L_f^2 & 2\beta^2 & 16\alpha^2 & 80\alpha^2L_f^2 & 32\alpha^2L_f^2 \\ \frac{L_f^2}{1-\sigma^2} (32\sigma^2 + a_1) & \frac{132L_f^2\beta^2}{1-\sigma^2} & \frac{1+\sigma^2}{2} + \frac{96\alpha^2}{1-\sigma^2} & \frac{L_f^2}{1-\sigma^2} a_2 & \frac{L_f^2}{1-\sigma^2} a_3 \\ \frac{4\alpha L_f^2}{\mu} & \frac{2\beta^2}{\mu\alpha} & 0 & 1 - \frac{\mu\alpha}{2} & \frac{2L_f^2\alpha^2}{m} \\ \frac{\mu}{2} & 0 & 0 & \frac{2}{p} & 1 - \frac{1}{P} \end{bmatrix}.$$

thus be found. From inequality (31), we obtain

$$\begin{cases} \frac{4\beta^2}{1-\sigma^2}\delta_2 < \frac{1-\sigma^2}{2}\delta_1 - \frac{4\alpha^2}{1-\sigma^2}\delta_3, \\ \alpha^2 < \frac{(1-\sigma^2)^2}{8}\frac{\delta_1}{\delta_3}. \end{cases}$$

Then we have

$$\begin{aligned} \alpha &< \sqrt{\frac{(1-\sigma^2)^2}{8}\frac{\delta_1}{\delta_3}} = \sqrt{\frac{(1-\sigma^2)^2}{8}\frac{(1-\sigma^2)^2}{18050L_f^2}} \\ &= \frac{(1-\sigma^2)^2}{380}\frac{\sqrt{p}}{\sqrt{PQ}}. \end{aligned} \tag{42}$$

From inequality (32), we obtain

$$\begin{aligned} \alpha &< \sqrt{\frac{1}{80L_f^2\delta_1 + 16\delta_3 + 80L_f^2\delta_4 + 32L_f^2\delta_5}} \\ \Leftrightarrow \alpha &< \sqrt{\frac{1}{\left(80 + 16\frac{18050PQ^2}{(1-\sigma^2)^2p} + 6400Q^2 + \frac{5248PQ^2}{p}\right)L_f^2}} \\ \Leftrightarrow \alpha &< \sqrt{\frac{(1-\sigma^2)^2p}{300528PQ^2L_f^2}} \Leftrightarrow \alpha < \frac{1-\sigma^2}{12\sqrt{2087}}\frac{\sqrt{p}}{\sqrt{PQL_f}}. \end{aligned} \tag{43}$$

From inequality (33), we can obtain

$$\begin{aligned} \frac{96\alpha^2L_f^2}{1-\sigma^2}\delta_3 &< \frac{1-\sigma^2}{2}\delta_3 - \frac{L_f^2}{1-\sigma^2}(105\delta_1 + 71\delta_4 + 19\delta_5) \\ \Leftrightarrow \frac{96\alpha^2L_f^2}{1-\sigma^2}\delta_3 &< \frac{1-\sigma^2}{2}\delta_3 - \frac{L_f^2}{1-\sigma^2}\frac{8901PQ^2}{p} \\ \Leftrightarrow \frac{96\alpha^2L_f^2}{1-\sigma^2} &< \frac{1-\sigma^2}{2} - \frac{8901(1-\sigma^2)}{18050} = \frac{62(1-\sigma^2)}{9025} \\ \Leftrightarrow \alpha &< \frac{\sqrt{31}(1-\sigma^2)}{380\sqrt{3}L_f}. \end{aligned} \tag{44}$$

Finally, we can obtain the following inequality by applying inequality (34):

$$\frac{4L_f^2\alpha}{\mu m}\delta_5 < \delta_4 - \frac{8L_f^2}{\mu^2}\delta_1 - \frac{4\beta^2}{\mu^2\alpha^2}\delta_2.$$

Then we can obtain inequalities as follows:

$$\begin{cases} \alpha < \frac{\mu m}{4L_f^2}\frac{\delta_4}{\delta_5} - \frac{2m}{\mu}\frac{\delta_1}{\delta_5}, \\ \alpha < \frac{m}{4QL_f}\frac{80p}{164P} - \frac{2m}{\mu}\frac{p}{164PQ^2}, \\ \alpha < \frac{9}{82}\frac{p}{P}\frac{m}{QL_f}. \end{cases} \tag{45}$$

Therefore,  $\alpha$  satisfies

$$\begin{aligned} \alpha &< \bar{\alpha} \triangleq \min \left\{ \frac{(1-\sigma^2)^2}{380}\frac{\sqrt{p}}{\sqrt{PQ}}, \frac{1-\sigma^2}{12\sqrt{2087}}\frac{\sqrt{p}}{\sqrt{PQL_f}}, \right. \\ &\quad \left. \frac{\sqrt{31}(1-\sigma^2)}{380\sqrt{3}L_f}, \frac{m}{QL_f}\frac{9p}{82P} \right\} \\ &\Leftrightarrow \alpha < \frac{p}{P}\frac{(1-\sigma^2)^2}{549QL_f}. \end{aligned}$$

Next, we set  $\alpha = \frac{p}{P}\frac{(1-\sigma^2)^2}{560QL_f}$  according to the previous range of  $\alpha$  to obtain  $\beta$ . First, from inequality (31), we have

$$\begin{aligned} \frac{4\beta^2}{1-\sigma^2}\delta_2 &< \frac{1-\sigma^2}{2} - \frac{4}{1-\sigma^2}\frac{18050}{(1-\sigma^2)^2}\frac{p}{P}\frac{(1-\sigma^2)^4}{313600} \\ &= \frac{1-\sigma^2}{2} - \frac{p}{P}\frac{361(1-\sigma^2)}{1568} \\ \Leftrightarrow \frac{4\beta^2}{1-\sigma^2}\delta_2 &< \frac{9 \times 47(1-\sigma^2)}{1568}\frac{p}{P} \\ \Leftrightarrow \beta &< \frac{3\sqrt{47}(1-\sigma^2)}{56\sqrt{34}}\frac{\sqrt{p}}{\sqrt{P}}. \end{aligned} \tag{46}$$

Second, from inequality (32), we have

$$\begin{aligned} 2\beta^2\delta_2 &< \delta_2 - 16\delta_1 - 80L_f^2\delta_1\alpha^2 - 16\delta_3\alpha^2 \\ &\quad - 80L_f^2\delta_4\alpha^2 - 32L_f^2\delta_5\alpha^2 \\ \Leftrightarrow 2\beta^2\delta_2 &< 1 - \frac{p^2}{P^2}\frac{(1-\sigma^2)^4}{3920Q^2} - \frac{p}{P}\frac{361(1-\sigma^2)^2}{392} \\ &\quad - \frac{p^2}{P^2}\frac{47(1-\sigma^2)^4}{1568} - \frac{p}{P}\frac{82(1-\sigma^2)^4}{4900} \\ \Leftrightarrow \beta &< \frac{\sqrt{1259}(1-\sigma^2)}{280\sqrt{17}}\frac{\sqrt{p}}{\sqrt{P}}. \end{aligned} \tag{47}$$

Then from inequality (33), we have

$$\begin{aligned} \frac{132L_f^2\beta^2}{1-\sigma^2}\delta_2 &< \frac{1-\sigma^2}{2}\delta_3 - \frac{105L_f^2}{1-\sigma^2}\delta_1 - \frac{96\alpha^2L_f^2}{1-\sigma^2}\delta_3 \\ &\quad - \frac{71L_f^2}{1-\sigma^2}\delta_4 - \frac{19L_f^2}{1-\sigma^2}\delta_5 \\ \Leftrightarrow \frac{132L_f^2\beta^2}{1-\sigma^2}\delta_2 &< \frac{9025L_f^2PQ^2}{1-\sigma^2}\frac{p}{p} - \frac{361PQ^2}{196}\frac{p}{p}\frac{8904L_f^2}{1-\sigma^2} \\ \Leftrightarrow \beta &< \frac{\sqrt{59}}{\sqrt{1122}}\frac{\sqrt{PQ}}{\sqrt{p}}. \end{aligned} \tag{48}$$

Finally, from inequality (34), we have

$$\begin{cases} \frac{4\beta^2}{\mu^2\alpha^2}\delta_2 < \delta_4 - \frac{8L_f^2}{\mu^2}\delta_1 - \frac{4L_f^2\alpha}{\mu m}\delta_5, \\ \beta < \frac{\sqrt{69}(1-\sigma^2)^2}{560\sqrt{68}}\frac{p}{PQ}. \end{cases} \tag{49}$$

Therefore, we obtain the range of  $\beta$  from inequalities (46)–(49) as follows:

$$\beta < \bar{\beta} \triangleq \min \left\{ \frac{3\sqrt{47p}(1-\sigma^2)}{56\sqrt{34P}}, \frac{\sqrt{1259p}(1-\sigma^2)}{280\sqrt{17P}}, \frac{\sqrt{59}}{\sqrt{1122}} \frac{\sqrt{PQ}}{\sqrt{p}}, \frac{\sqrt{69}(1-\sigma^2)^2}{560\sqrt{68}} \frac{p}{PQ} \right\},$$

which completes the proof.

**Remark 2** It is worth emphasizing that this study does not establish explicit expression of the accelerated convergence rate from a theoretical point of view; the acceleration of the proposed algorithm can be demonstrated only from the practical point of view. The acceleration in this study can be observed through a series of experiments because this is still an open problem in recent literature (Xin et al., 2019a; Li HQ et al., 2020; Lü et al., 2020; Xin and Khan, 2020; Hu et al., 2021). Therefore, it would be worthwhile to analyze the acceleration from a theoretical point of view in future work.

## 4 Experimental results

In this section, we compare the proposed algorithm with other methods by solving some optimization problems, such as the distributed logistic regression problem, distributed least-squares regression problem, and distributed quadratic programming problem, to prove the proposed algorithm's practicability. The performances of all the tested algorithms are plotted by  $\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}_k^i - \tilde{\mathbf{x}}^*\|$ .

### 4.1 Distributed logistic regression

This subsection illustrates the performance comparison between different algorithms. We set  $N = \sum_{i=1}^m q_i$  as the number of samples and allocate  $q_i$  to  $m$  agents on average, i.e.,  $q_i = N/m$ . The undirected network to distributed logistic regression problem is given as follows:

$$\tilde{\mathbf{x}}^* = \arg \min_{\tilde{\mathbf{x}} \in \mathbb{R}^n} \left( \frac{\gamma}{2} \|\tilde{\mathbf{x}}\|^2 + \sum_{i=1}^m \sum_{j=1}^{q_i} \log(1 + \exp(-b_{ij} \mathbf{c}_{ij}^T \tilde{\mathbf{x}})) \right), \quad (50)$$

where  $\mathbf{c}_{ij} \in \mathbb{R}^n$  is the  $j^{\text{th}}$  sample and  $b_{ij} \in \{1, -1\}$  is the corresponding binary label. To avoid over-fitting, we add a regularization term. According to Eq. (50), the local objective function  $f_i$  corresponding to P0

is written as follows:

$$f_i(\tilde{\mathbf{x}}) = \frac{\gamma}{2} \|\tilde{\mathbf{x}}\|^2 + \sum_{j=1}^{q_i} \log(1 + \exp(-b_{ij} \mathbf{c}_{ij}^T \tilde{\mathbf{x}})). \quad (51)$$

We aim to solve the logistic regression problem by choosing the Wisconsin breast cancer (diagnostic) dataset in the UCI Machine Learning Repository (Dua and Graff, 2017). These samples are from a digitized image of a fine needle aspirate of a breast mass. By judging the average of distances from the center to points on the silhouette and the severity of concave portions of the contour, a patient's condition is predicted as malignant or benign. We randomly choose  $N = 500$  samples from the total of 683 dataset samples to train the discriminator  $\tilde{\mathbf{x}}$  and suppose that samples are equally distributed to each agent; the rest of the samples are used for testing. Fig. 2 presents the performance comparison of the proposed algorithm, GT-SAGA, and EXTRA over an undirected network. All step-sizes are set as  $\alpha = 0.05$  and the heavy-ball momentum acceleration term  $\beta = 0.05$  in all the algorithms is the same. The accuracy, in the vertical axis of Fig. 3, is the ratio of the number of correct predictions to the total number of samples in the testing set.

### 4.2 Least-squares method of distributed signal processing

This subsection studies a least-squares problem provided in Li Z et al. (2019), which considers an undirected network of  $m = 10$  agents for an unknown signal  $\tilde{\mathbf{x}} \in \mathbb{R}^n$  for the optimization problem:

$$\min_{\tilde{\mathbf{x}} \in \mathbb{R}^n} \tilde{f}(\tilde{\mathbf{x}}) \triangleq \sum_{i=1}^m \frac{1}{2} \|\mathbf{C}_i \tilde{\mathbf{x}} - \mathbf{d}_i\|^2, \quad (52)$$

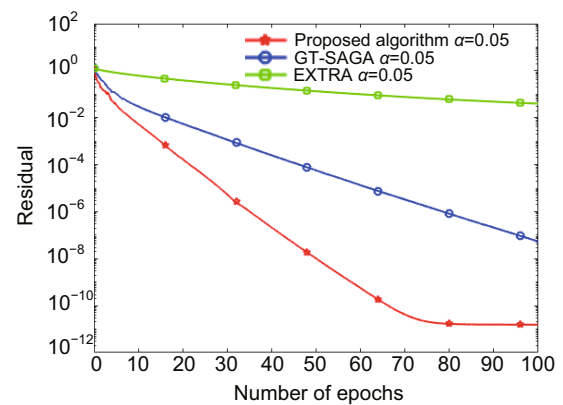


Fig. 2 Logistic regression over an undirected network

where the local objective function maintained by agent  $i$  ( $i = 1, 2, \dots, m$ ) is given by

$$f_i(\tilde{\mathbf{x}}) = \frac{1}{q_i} \sum_{h=1}^{q_i} \|\mathbf{C}_i^h \tilde{\mathbf{x}} - \mathbf{d}_i^h\|^2. \quad (53)$$

From problem (52), we can see that  $\mathbf{C}_i \in \mathbb{R}^{q_i \times n}$  is the sensing matrix and that  $\mathbf{d}_i = \mathbf{C}_i \tilde{\mathbf{x}} + \mathbf{e}_i$  is the measurement, where  $\mathbf{e}_i \in \mathbb{R}^{q_i}$  is the independent and identically distributed noise. When  $q_i = 600$  and  $m = 10$ , Fig. 4 illustrates the proposed algorithm, GT-SAGA, and EXTRA algorithms over an undirected network to verify the acceleration of the proposed algorithm.

### 4.3 Distributed quadratic programming

In this subsection, the problem of quadratic programming can be resolved by an undirected network of  $m = 10$  agents as follows:

$$\min_{\tilde{\mathbf{x}} \in \mathbb{R}^n} \tilde{f}(\tilde{\mathbf{x}}) = \sum_{i=1}^m \tilde{\mathbf{x}}^T \mathbf{G}_i \tilde{\mathbf{x}} + \mathbf{c}_i^T \tilde{\mathbf{x}}, \quad (54)$$

where matrix  $\mathbf{G}_i \in \mathbb{R}^{n \times n}$  is diagonal and positive definite and  $\mathbf{c}_i \in \mathbb{R}^n$  is a randomly generated vector.

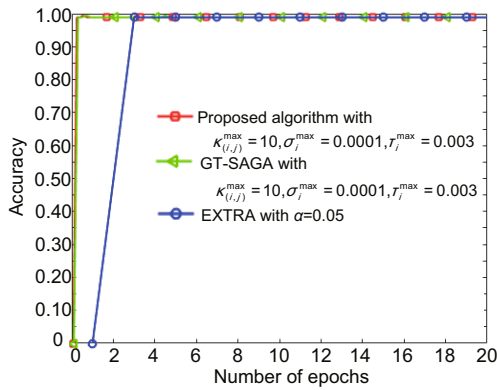


Fig. 3 Comparison of accuracy performance

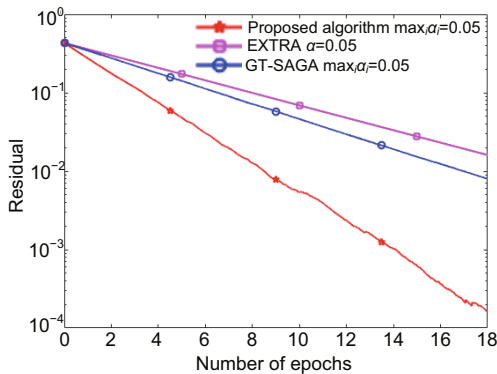


Fig. 4 Performance comparison over the least-squares method of distributed signal processing

We assume the dimension  $n = 10$ . Then we set the step-size  $\alpha = 0.0008$  in all the algorithms and the momentum parameter  $\beta = 0.0008$ . Fig. 5 shows the acceleration performance of the proposed algorithm compare to those of GT-SAGA and EXTRA over an undirected network.

**Remark 3** In view of Tables 1–3, one can clearly see that in different experiments, there are fewer epochs of the proposed algorithm than those of the GT-SAGA and EXTRA algorithms under the same accuracy. This demonstrates that the proposed algorithm accelerates convergence. So, the acceleration of the proposed algorithm can be verified.

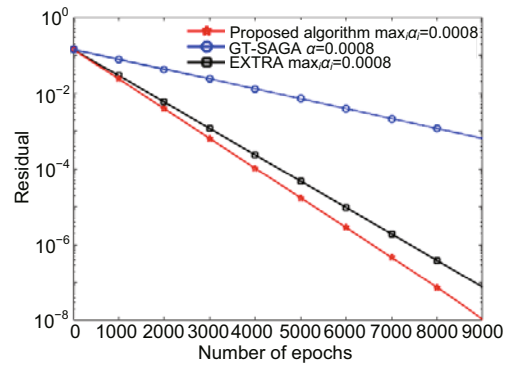


Fig. 5 Performance comparison over an undirected network when condition number  $Q = 500$

Table 1 Convergence performance comparison over logistic regression

Algorithm	Number of epochs		
	Accuracy= $10^{-2}$	$10^{-4}$	$10^{-6}$
Proposed algorithm	8	22	35
GT-SAGA	16	46	79
EXTRA	193	549	935

Table 2 Convergence performance comparison over least-squares

Algorithm	Number of epochs	
	Accuracy= $10^{-1}$	$10^{-2}$
Proposed algorithm	3	18
GT-SAGA	6	8
EXTRA	7	20

Table 3 Convergence performance comparison over the distributed quadratic programming

Algorithm	Number of epochs	
	Accuracy= $10^{-1}$	$10^{-2}$
Proposed algorithm	1472	2747
GT-SAGA	4412	8237
EXTRA	1657	3097

## 5 Conclusions

We presented a distributed stochastic algorithm which is capable of solving large-scale optimization problems over undirected networks. We showed that the proposed algorithm achieves accelerated linear convergence with a constant step-size  $\alpha$ . Extensive experiments on real-world datasets illustrated that the performance of the proposed algorithm is superior to those of other comparable algorithms.

### Contributors

Bihao SUN designed the research, processed the data, and drafted the manuscript. Jinhui HU, Dawen XIA, and Huaqing LI helped organize the manuscript and process the data. Bihao SUN and Huaqing LI revised and finalized the paper.

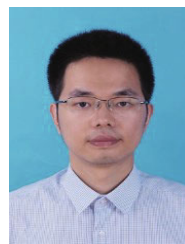
### Compliance with ethics guidelines

Bihao SUN, Jinhui HU, Dawen XIA, and Huaqing LI declare that they have no conflict of interest.

### References

- Bertsekas D, Gafni E, 1983. Projected Newton methods and optimization of multicommodity flows. *IEEE Trans Autom Contr*, 28(12):1090-1096. <https://doi.org/10.1109/TAC.1983.1103183>
- Boyd S, Parikh N, Chu E, et al., 2011. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found Trends<sup>®</sup> Mach Learn*, 3(1):1-122. <https://doi.org/10.1561/22000000016>
- Cheng B, Li ZK, 2019. Coordinated tracking control with asynchronous edge-based event-triggered communications. *IEEE Trans Autom Contr*, 64(10):4321-4328. <https://doi.org/10.1109/TAC.2019.2895927>
- Cheng S, Chen MY, Wai RJ, et al., 2014. Optimal placement of distributed generation units in distribution systems via an enhanced multi-objective particle swarm optimization algorithm. *J Zhejiang Univ-Sci C (Comput & Electron)*, 15(4):300-311. <https://doi.org/10.1631/jzus.C1300250>
- Cohen K, Nedić A, Srikant R, 2017. Distributed learning algorithms for spectrum sharing in spatial random access wireless networks. *IEEE Trans Autom Contr*, 62(6):2854-2869. <https://doi.org/10.1109/TAC.2016.2626578>
- Defazio A, Bach F, Lacoste-Julien S, 2014. SAGA: a fast incremental gradient method with support for non-strongly convex composite objectives. Proc 27<sup>th</sup> Int Conf on Neural Information Processing Systems, p.1646-1654.
- Dua D, Graff C, 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- Duchi JC, Agarwal A, Wainwright MJ, 2012. Dual averaging for distributed optimization: convergence analysis and network scaling. *IEEE Trans Autom Contr*, 57(3):592-606. <https://doi.org/10.1109/TAC.2011.2161027>
- Eisen M, Mokhtari A, Ribeiro A, 2017. Decentralized quasi-Newton methods. *IEEE Trans Signal Process*, 65(10):2613-2628. <https://doi.org/10.1109/TSP.2017.2666776>
- Erseghe T, Zennaro D, Dall'Anese E, et al., 2011. Fast consensus by the alternating direction multipliers method. *IEEE Trans Signal Process*, 59(11):5523-5537. <https://doi.org/10.1109/TSP.2011.2162831>
- Guan L, Sun T, Qiao LB, et al., 2020. An efficient parallel and distributed solution to nonconvex penalized linear SVMs. *Front Inform Technol Electron Eng*, 21(4):587-603. <https://doi.org/10.1631/FITEE.1800566>
- Han ZM, Lin ZY, Fu MY, et al., 2015. Distributed coordination in multi-agent systems: a graph Laplacian perspective. *Front Inform Technol Electron Eng*, 16(6):429-448. <https://doi.org/10.1631/FITEE.1500118>
- Hu JH, Yan Y, Li HQ, et al., 2021. Convergence of an accelerated distributed optimisation algorithm over time-varying directed networks. *IET Contr Theory Appl*, 15(1):24-39. <https://doi.org/10.1049/cth2.12022>
- Johnson R, Zhang T, 2013. Accelerating stochastic gradient descent using predictive variance reduction. Proc 26<sup>th</sup> Int Conf on Neural Information Processing Systems, p.315-323.
- Lan Q, Qiao LB, Wang YJ, 2018. Stochastic extra-gradient based alternating direction methods for graph-guided regularized minimization. *Front Inform Technol Electron Eng*, 19(6):755-762. <https://doi.org/10.1631/FITEE.1601771>
- Li HQ, Cheng HQ, Wang Z, et al., 2020. Distributed Nesterov gradient and heavy-ball double accelerated asynchronous optimization. *IEEE Trans Neur Netw Learn Syst*, in press. <https://doi.org/10.1109/TNNLS.2020.3027381>
- Li Z, Shi W, Yan M, 2019. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates. *IEEE Trans Signal Process*, 67(17):4494-4506. <https://doi.org/10.1109/TSP.2019.2926022>
- Ling Q, Ribeiro A, 2014. Decentralized dynamic optimization through the alternating direction method of multipliers. *IEEE Trans Signal Process*, 62(5):1185-1197. <https://doi.org/10.1109/TSP.2013.2295055>
- Ling Q, Tian Z, 2010. Decentralized sparse signal recovery for compressive sleeping wireless sensor networks. *IEEE Trans Signal Process*, 58(7):3816-3827. <https://doi.org/10.1109/TSP.2010.2047721>
- Liu R, Sun WC, Hou T, et al., 2019. Block coordinate descent with time perturbation for nonconvex nonsmooth problems in real-world studies. *Front Inform Technol Electron Eng*, 20(10):1390-1403. <https://doi.org/10.1631/FITEE.1900341>
- Lü QG, Liao XF, Li HQ, et al., 2020. A Nesterov-like gradient tracking algorithm for distributed optimization over directed networks. *IEEE Trans Syst Man Cybern*, in press. <https://doi.org/10.1109/TSMC.2019.2960770>
- Mateos G, Bazerque JA, Giannakis GB, 2010. Distributed sparse linear regression. *IEEE Trans Signal Process*, 58(10):5262-5276. <https://doi.org/10.1109/TSP.2010.2055862>

- Matthews TP, Wang K, Li CP, et al., 2016. Nonlinear waveform inversion by use of the regularized dual averaging method for ultrasound computed tomography. *Progress in Electromagnetic Research Symp*, p.3948. <https://doi.org/10.1109/PIERS.2016.7735487>
- McMahan B, Moore E, Ramage D, et al., 2017. Communication-efficient learning of deep networks from decentralized data. *Proc 20<sup>th</sup> Int Conf on Artificial Intelligence and Statistics*, p.1273-1282.
- Nedić A, Ozdaglar A, 2009. Distributed subgradient methods for multi-agent optimization. *IEEE Trans Autom Contr*, 54(1):48-61. <https://doi.org/10.1109/TAC.2008.2009515>
- Nedić A, Olshevsky A, Shi W, 2017a. Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM J Optim*, 27(4):2597-2633. <https://doi.org/10.1137/16M1084316>
- Nedić A, Olshevsky A, Shi W, et al., 2017b. Geometrically convergent distributed optimization with uncoordinated step-sizes. *American Control Conf*, p.3950-3955. <https://doi.org/10.23919/ACC.2017.7963560>
- Schmidt M, Le Roux N, Bach F, 2017. Minimizing finite sums with the stochastic average gradient. *Math Program*, 162(1-2):83-112. <https://doi.org/10.1007/s10107-016-1030-6>
- Tan CH, Ma SQ, Dai YH, et al., 2016. Barzilai-Borwein step size for stochastic gradient descent. *Proc 30<sup>th</sup> Int Conf on Neural Information Processing Systems*, p.685-693.
- Tsitsiklis J, Bertsekas D, Athans M, 1986. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE Trans Autom Contr*, 31(9):803-812. <https://doi.org/10.1109/TAC.1986.1104412>
- Wang B, Jiang HY, Fang J, et al., 2018. A proximal ADMM for decentralized composite optimization. *IEEE Signal Process Lett*, 25(8):1121-1125. <https://doi.org/10.1109/LSP.2018.2841648>
- Wang Z, Li HQ, 2020. Edge-based stochastic gradient algorithm for distributed optimization. *IEEE Trans Netw Sci Eng*, 7(3):1421-1430. <https://doi.org/10.1109/TNSE.2019.2933177>
- Wei EM, Ozdaglar A, Jadbabaie A, 2013. A distributed Newton method for network utility maximization—I: algorithm. *IEEE Trans Autom Contr*, 58(9):2162-2175. <https://doi.org/10.1109/TAC.2013.2253218>
- Xi CG, Khan UA, 2017. DEXTRA: a fast algorithm for optimization over directed graphs. *IEEE Trans Autom Contr*, 62(10):4980-4993. <https://doi.org/10.1109/TAC.2017.2672698>
- Xia YS, Wang J, 2004. A one-layer recurrent neural network for support vector machine learning. *IEEE Trans Syst Man Cybern B*, 34(2):1261-1269. <https://doi.org/10.1109/TSMCB.2003.822955>
- Xin R, Khan UA, 2018. A linear algorithm for optimization over directed graphs with geometric convergence. *IEEE Contr Syst Lett*, 2(3):315-320. <https://doi.org/10.1109/LCSYS.2018.2834316>
- Xin R, Khan UA, 2020. Distributed heavy-ball: a generalization and acceleration of first-order methods with gradient tracking. *IEEE Trans Autom Contr*, 65(6):2627-2633. <https://doi.org/10.1109/TAC.2019.2942513>
- Xin R, Jakovetić D, Khan UA, 2019a. Distributed Nesterov gradient methods over arbitrary graphs. *IEEE Signal Process Lett*, 26(8):1247-1251. <https://doi.org/10.1109/LSP.2019.2925537>
- Xin R, Sahu AK, Khan UA, et al., 2019b. Distributed stochastic optimization with gradient tracking over strongly-connected networks. *Proc IEEE 58<sup>th</sup> Conf on Decision and Control*, p.8353-8358. <https://doi.org/10.1109/CDC40024.2019.9029217>
- Xin R, Xi CG, Khan UA, 2019c. FROST—fast row-stochastic optimization with uncoordinated step-sizes. *EURASIP J Adv Signal Process*, 2019(1):1. <https://doi.org/10.1186/s13634-018-0596-y>
- Xin R, Khan UA, Kar S, 2020. Variance-reduced decentralized stochastic optimization with accelerated convergence. *IEEE Trans Signal Process*, 68:6255-6271. <https://doi.org/10.1109/TSP.2020.3031071>
- Xu JM, Zhu SY, Soh YC, et al., 2015. Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes. *Proc 54<sup>th</sup> IEEE Conf on Decision and Control*, p.2055-2060. <https://doi.org/10.1109/CDC.2015.7402509>
- Yin R, Zhang Y, Yu GD, et al., 2010. Centralized and distributed resource allocation in OFDM based multi-relay system. *J Zhejiang Univ-Sci C (Comput & Electron)*, 11(6):450-464. <https://doi.org/10.1631/jzus.C0910405>
- Yuan DM, Ma Q, Wang Z, 2013. Distributed dual averaging method for solving saddle-point problems over multi-agent networks. *Proc 32<sup>nd</sup> Chinese Control Conf*, p.6868-6872.
- Zhang CL, Ahmad M, Wang YQ, 2019. ADMM based privacy-preserving decentralized optimization. *IEEE Trans Inform Forens Secur*, 14(3):565-580. <https://doi.org/10.1109/TIFS.2018.2855169>
- Zinkevich MA, Weimer M, Smola A, et al., 2010. Parallelized stochastic gradient descent. *Proc 23<sup>rd</sup> Int Conf on Neural Information Processing Systems*, p.2595-2603.



Huaqing LI received his BS degree in information and computing science in 2009 from Chongqing University of Posts and Telecommunications, Chongqing, China and his PhD degree in computer science and technology in 2013 from Chongqing University. From Sept. 2014 to Sept. 2015, he was a postdoctoral researcher at the School of Electrical and Information Engineering, The University of Sydney, Australia. From Nov. 2015 to Nov. 2016, he was a postdoctoral researcher at the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. He is currently a professor at the College of Electronic and Information Engineering, Southwest University, Chongqing, China. His main research interests include nonlinear dynamics and control, multi-agent system, and distributed optimization. Prof. LI currently serves as a regional editor for *Neur Comput Appl*, an editorial board member for *IEEE Access*, and a corresponding expert for *Front Inform Technol Electron Eng*.