

# Improved Convergence Rate for a Distributed Two-Time-Scale Gradient Method under Random Quantization

Marcos M. Vasconcelos, Think T. Doan and Urbashi Mitra

**Abstract**—We study the so-called distributed two-time-scale gradient method for solving convex optimization problems over a network of agents when the communication bandwidth between the nodes is limited, and so information that is exchanged between the nodes must be quantized. Our main contribution is to provide a novel analysis, resulting to an improved convergence rate of this method as compared to the existing works. In particular, we show that the method converges at a rate  $\mathcal{O}(\log^2(k)/\sqrt{k})$  to the optimal solution, when the underlying objective function is strongly convex and smooth. The key technique in our analysis is to consider a Lyapunov function that simultaneously captures the coupling of the consensus and optimality errors generated by the method.

## I. INTRODUCTION

Next generation cyber-physical systems will seamlessly incorporate machine learning capabilities to enable adaptability and resiliency to guarantee performance in non-stationary environments. In decentralized cyber-physical systems, the application of real-time, collaborative machine learning methods require a collection of agents to *train* a global model efficiently based on local processing and information exchange over a network [1]. Such problems can be formulated using the framework of distributed optimization, which has a long and rich history, see for example [2] and references therein.

One popular class of distributed optimization algorithms consists of each agent performing gradient descent on its local function, obtaining a local estimate of the optimal solution, which is then shared with its neighbors. Such approach essentially *interleaves* average-consensus and gradient descent. In contrast to quantized-consensus which has been extensively studied and is now well understood, *distributed optimization under quantization* has received significantly less attention. Since quantization is a necessary component in any application involving digital communication among agents, it is essential to design algorithms that take quantization into account and quantify the impact of quantization on the performance of such algorithms. As [1] points out that quantization is a big bottleneck in the design and analysis of distributed optimization algorithms. Our work seeks to address this fundamental issue.

M. M. Vasconcelos is with the Commonwealth Cyber Initiative and the Bradley Department of Electrical and Computer Engineering, Virginia Tech. T. T. Doan is with the Bradley Department of Electrical and Computer Engineering, Virginia Tech. U. Mitra is with the Ming Hsieh Department of Electrical Engineering, University of Southern California. E-mails: {marcosv, thinkdoan}@vt.edu, ubli@usc.edu.

In [3], a quantized subgradient method is considered, where the quantization scheme is non-adaptive, i.e., the quantization bins do not change with time. In this case, convergence to a ball centered at the optimal solution is shown. While [4] establishes convergence to the optimal solution with adaptive quantization, a drawback of [4] is the need for periodic communication of the quantization intervals, which violates hard constraints on the fixed number of bits allowed over each link, since these intervals are represented by real numbers. Such requirement is also considered in the recent work in both centralized (there exists a centralized coordinator) and distributed communication frameworks; see for example [5] and the references therein. Finally, the work in [6] provides an adaptive quantization framework to achieve an optimal rate of distributed subgradient methods under quantization, however, it requires a (strict) condition on the capacity of the communication bandwidth; see Assumption 2 in [6].

These challenges can be overcome by two-time scale algorithms from stochastic approximation [7], [8]. In particular, [7] established convergence and a convergence rate of a distributed optimization algorithm under strict restrictions on the number of bits that are transmitted *at all times*. However, that algorithm requires the use of an operator that projects the iterates onto compact sets, which are then quantized. This projection operator creates a *projection error* in addition to the *quantization error*. The interplay between these two quantities introduces significant complexity to the analysis and limits the achieved convergence rate. For example, the convergence rate achieved in [7] is  $\mathcal{O}(1/\sqrt[3]{k})$ .

The main contribution herein is to provide a novel analysis, resulting to an improved convergence rate of the distributed two-time-scale gradient method under random quantization. In our method, the quantization bins increase with time in a deterministic way, while keeping the number of bits in communication constant at all time. This avoids the strong requirement of exchanging real numbers per communication time in the existing literature. By introducing a Lyapunov function that simultaneously captures the coupling between the consensus and optimality errors, we show that our method converges at a rate  $\mathcal{O}(1/\sqrt{k})$  to the optimal solution, which is better than the one in [7] by a factor of  $1/\sqrt[3]{k}$ .

### A. Notation

We adopt the following notation: scalar and vector valued random variables, and corresponding realizations are denoted by lowercase letters, such as  $x$ . Matrices and some constants

are denoted by uppercase letters, such as  $X$ . The all-one vector is denoted by  $\mathbf{1}$ . The expectation of a random variable  $x$  is denoted by  $\mathbf{E}[x]$ . The Euclidean norm of a vector  $x$  is denoted by  $\|x\|$ , and the Frobenius norm of a matrix  $X$  is denoted by  $\|X\|_F$ .

## II. PROBLEM FORMULATION

Consider a distributed optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \sum_{i=1}^n f_i(x), \quad (1)$$

where each function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex,  $i \in \{1, \dots, n\}$ . Each function  $f_i$  is only available to agent  $i$  and satisfies the following conditions.

*Assumption 1:* The function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  has Lipschitz continuous gradient with constant  $L$ , i.e.,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d. \quad (2)$$

In addition,  $f_i$  is strongly convex with constant  $\mu$ , i.e.,

$$\frac{\mu}{2}\|x - y\|^2 \leq f(x) - f(y) - \nabla f(y)^T(x - y), \quad x, y \in \mathbb{R}^d. \quad (3)$$

We also assume that the gradient of  $f_i$  is bounded.

*Assumption 2:* There exists a constant  $C > 0$  such that

$$\|\nabla f_i(x)\| \leq C, \quad x \in \mathbb{R}^d, \quad \forall i \in \{1, \dots, n\}. \quad (4)$$

*Remark 1:* We note that the boundedness condition of  $\nabla f_i$  can be guaranteed through a projection step of the variable  $x$  to a predefined compact set, as considered in [7], [8]. For simplicity, we will consider only unconstrained problems under this assumption in this paper. However, the proposed methods extend to constrained problems as well as the relaxation of Assumption 2.

Agents are connected by an undirected and static connected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{1, \dots, n\}$  and  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  are the sets of vertices and edges, respectively. Each communication link  $e \in \mathcal{E}$  can support  $b$  bits per dimension per transmission. Therefore, when a vector  $x_i \in \mathbb{R}^d$  is transmitted from agent  $i$  to one of its neighbors, it must first be quantized into a vector  $q_i$  using  $b \times d$  bits. We do not assume any form of analog communication among agents. The goal of the agents is to cooperatively solve the problem in Eq. (1) by exchanging quantized messages over communication links of limited bandwidth.

### A. Random quantization

Here we describe the quantization algorithm employed in our framework, where each component of the local solution estimate  $x_i \in \mathbb{R}^d$  is quantized using  $b$  bits. The following discussion specifies the quantizer for a real number representing a single component, and the same operation is repeated for all the  $d$  dimensions of  $x_i$ .

Let  $x \in [\ell, u]$ . We partition the interval  $[\ell, u]$  into  $B$  equal length bins whose endpoints are denoted by  $\tau_m$ ,  $m \in \{1, \dots, B+1\}$ , such that  $\tau_1 = \ell$ , and  $\tau_{B+1} = u$ . Define the length of each bin by  $\Delta$ , such that:

$$\Delta \stackrel{\text{def}}{=} \frac{u - \ell}{B}. \quad (5)$$

We use  $\{\tau_m\}_{m=1}^{B+1}$  as the representation symbols for the quantizer. As such, each  $\tau_m$  is mapped into a codeword of  $b$  bits. Hence, for a given number of bits  $b$ , the number of bins is given by  $B = 2^b - 1$ , which implies that:

$$\Delta = \frac{u - \ell}{2^b - 1}. \quad (6)$$

Let  $x \in [\tau_m, \tau_{m+1})$ , then we choose either  $\tau_m$  or  $\tau_{m+1}$  as a representation point for  $x$  based on the following stochastic rule,  $\mathcal{Q}$ , defined as:

$$\mathcal{Q}(x) \stackrel{\text{def}}{=} \begin{cases} \tau_m & \text{with prob. } 1 - (x - \tau_m)/\Delta \\ \tau_{m+1} & \text{with prob. } (x - \tau_m)/\Delta, \end{cases} \quad (7)$$

where  $\Delta$  is given by Eq. (6).

The randomized quantizer  $\mathcal{Q}$  described above satisfies three important properties:

$$\mathbf{E}[\mathcal{Q}(x) \mid x] = x, \quad (8)$$

$$\mathbf{E}[(\mathcal{Q}_b(x) - x)^2 \mid x] \leq \frac{\Delta^2}{4}, \quad (9)$$

and

$$\mathbf{P}(|\mathcal{Q}(x) - x| \leq \Delta) = 1. \quad (10)$$

### B. Distributed two-time scale gradient method under random quantization

We are interested in exploiting distributed two-time-scale gradient methods, formally stated in Algorithm 1, for solving problem in Eq. (1) under random quantization. This algorithm has a simple interpretation given as follows.

Each agent  $i \in \mathcal{V}$  keeps a local estimate of the solution to Eq. (1). At time step  $k$ , the  $i$ -th agent's local estimate is denoted by  $x_k^i$ . Each agent has access to a local subgradient of its local function. At time step  $k$ , the  $i$ -th agent's local gradient at  $x_k^i$  is denoted by  $g_k^i$ , i.e.,

$$g_k^i \stackrel{\text{def}}{=} \nabla f_i(x_k^i). \quad (11)$$

At time  $k$  the  $i$ -th agent received the quantized versions of the states from every agent  $j \in \mathcal{N}_i$ . Let

$$q_k^j \stackrel{\text{def}}{=} \mathcal{Q}(x_k^j) \quad (12)$$

denote the quantized version of the local estimate of the  $j$ -th agent at time  $k$ . Each agent  $i$  then iteratively updates its variable as

$$x_{k+1}^i \stackrel{\text{def}}{=} (1 - \beta_k)x_k^i + \beta_k \sum_{j \in \mathcal{N}_i} a_{ij}q_k^j - \alpha_k g_k^i, \quad i \in \mathcal{V}, \quad (13)$$

where  $\{\alpha_k\}_{k=0}^{\infty}$  and  $\{\beta_k\}_{k=0}^{\infty}$  are sequences of diminishing step-sizes and  $a_{ij}$  are averaging coefficients corresponding to the entries of a doubly stochastic matrix  $A$ . The iterates generated by Eq. (13) will be the centerpiece of this paper, for which we will establish the convergence rate under appropriate assumptions on the structure of the problem to be specified in the sequel.

---

**Algorithm 1:** Distributed Two-Time-Scale Gradient Methods Under Random Quantization
 

---

- 1) **Initialize:** Each node  $i$  initializes  $x_0^i = 0$ , and two sequences of stepsizes  $\{\alpha_k, \beta_k\}_{k \in \mathbb{N}}$ .
- 2) **Iteration:** For  $k = 0, 1, \dots$ , node  $i \in \mathcal{V}$  implements:
  - a. Compute random quantization  $q_k^i = \mathcal{Q}(x_k^i)$
  - b. Send  $q_k^i$  to node  $j \in \mathcal{N}_i$
  - c. Receive  $q_k^j$  from node  $j \in \mathcal{N}_i$  and update

$$x_{k+1}^i = (1 - \beta_k)x_k^i + \beta_k \sum_{j \in \mathcal{N}_i} a_{ij}q_k^j - \alpha_k g_k^i.$$

- d. Update the output  $z_{k+1}^i$

$$z_{k+1}^i = \frac{\sum_{t=0}^k (t+1)x_t^i}{\sum_{t=0}^k (t+1)}.$$


---

### C. Quantization with increasing bin size

What enables the analysis of the algorithm in Eq. (13) is the following proposition, which guarantees that there exists a randomized quantization scheme with increasing bin sizes that has a maximal error that grows logarithmically with  $k$ . The purpose of Proposition 1 is to allow each agent to know at any given  $k$  which interval to quantize and decode. By using the construction given in the proof, we only need to transmit  $b \times d$  bits every iteration but not a real number as done in [4], [5].

*Proposition 1:* Let  $\{\alpha_k\}$  be a sequence of diminishing step-sizes in algorithm Eq. (13). There exists a randomized quantization scheme such that:

$$\|x_k^i - q_k^i\| \leq \Delta_k \stackrel{\text{def}}{=} \frac{2C}{2^b - 1} \sum_{t=0}^{k-1} \alpha_t \quad (14)$$

*Proof:* Let  $x_i^0 = 0$ , then since  $|g_i^0| \leq C$  for all  $i$ , we have:

$$|x_1^i| \leq \alpha_0 C. \quad (15)$$

At  $k = 1$ , let agent  $i$  quantize  $x_1^i \in [-\alpha_0 C, \alpha_0 C]$  to obtain  $q_1^i$ . Then for all  $i$

$$\mathbf{E}[q_1^i | x_1^i] = x_1^i \quad \text{and} \quad |x_1^i - q_1^i| \leq \Delta_1, \quad (16)$$

where

$$\Delta_1 = \frac{2\alpha_1 C}{2^b - 1}. \quad (17)$$

Then, at  $k = 1$ , we compute  $x_2^i$  using Eq. (13) as:

$$x_2^i \stackrel{\text{def}}{=} (1 - \beta_1)x_1^i + \beta_1 \sum_{j \in \mathcal{N}_i} a_{ij}q_1^j - \alpha_1 g_1^i \quad (18)$$

Thus, since  $\{|x_1^i|, |q_1^i|\}_{i=1}^n \leq \alpha_0 C$ , we have

$$|x_2^i| \stackrel{(a)}{\leq} (1 - \beta_1)|x_1^i| + \beta_1 \sum_{j \in \mathcal{N}_i} a_{ij}|q_1^j| + \alpha_1 |g_1^i| \quad (19)$$

$$\stackrel{(b)}{\leq} C(\alpha_0 + \alpha_1), \quad (20)$$

where (a) follows from the triangle inequality, and (b) follows from Eq. (16). Therefore, quantizing  $x_2^i \in [-(\alpha_0 + \alpha_1)C, (\alpha_0 + \alpha_1)C]$ , we have

$$|x_2^i - q_2^i| \leq \Delta_2 = \frac{2C}{2^b - 1}(\alpha_0 + \alpha_1). \quad (21)$$

By induction, repeating the steps above for the update and quantization, we obtain:

$$|x_k^i - q_k^i| \leq \Delta_k = \frac{2C}{2^b - 1} \sum_{t=0}^{k-1} \alpha_t. \quad (22)$$

*Remark 2:* For strongly convex functions, we will later choose  $\alpha_k \propto 1/(k+1)$ , thus the quantization error will satisfy:

$$\Delta_k \approx \frac{2C}{2^b - 1} \ln(k) \quad (23)$$

Proposition 1 is proven for scalars, but the same analysis can be carried out for vectors, for which the error bound is simply multiplied by the appropriate dimension  $d$ , i.e.,

$$\Delta_k = \frac{2dC}{2^b - 1} \sum_{t=0}^{k-1} \alpha_t. \quad (24)$$

The main consequence of Proposition 1 is that we do not need to use an additional projection step as in [9]. Unlike [9], the quantization error now depends on time step  $k$  instead of being constant. This error scales with  $\ln(k)$ , and fortunately, will not impact the asymptotic convergence rate.

### D. Preliminary definitions and alternative representations

Throughout the analysis of the algorithm, one quantity will be paramount: the *quantization error* process. For the  $i$ -th agent let the quantization error be defined as:

$$e_k^i \stackrel{\text{def}}{=} x_k^i - q_k^i. \quad (25)$$

To facilitate the analysis, we will use  $X, G(X) \in \mathbb{R}^{n \times d}$  to denote the following matrices:

$$X = \begin{bmatrix} -x_1^T - \\ \dots \\ -x_n^T - \end{bmatrix} \quad \text{and} \quad G(X) = \begin{bmatrix} -g_1^T(x_1) - \\ \dots \\ -g_n^T(x_n) - \end{bmatrix}. \quad (26)$$

Due to the random quantization operator, the sample-paths  $x_k^i$  originating from the recursion in Eq. (13) are stochastic processes. Let  $\mathcal{F}_k$  be the filtration [10] containing all of the history generated by Eq. (13) up to time  $k$ :

$$\mathcal{F}_k \stackrel{\text{def}}{=} \{X_0, Q_0, X_1, Q_1, \dots, X_k, Q_k\}. \quad (27)$$

The conditional expectation operator with respect to the filtration  $\mathcal{F}_k$  is denoted by  $\mathbf{E}_{\mathcal{F}_k}$ .

Let the average of  $x_k^i$ 's be denoted by  $\bar{x}_k$ , such that

$$\bar{x}_k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n x_k^i = X_k^T \mathbf{1}. \quad (28)$$

Finally, define the following stochastic processes:

$$r_k \stackrel{\text{def}}{=} \|\bar{x}_k - x^*\| \quad (29)$$

and

$$Y_k \stackrel{\text{def}}{=} X_k - \mathbf{1}\bar{x}_k^T = WX_k, \quad (30)$$

where  $W = I - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ .

In our analysis of algorithm Eq. (12) we will also make use of the two following alternative representations:

$$\bar{x}_{k+1} = (1 - \beta_k)\bar{x}_k + \beta_k\bar{q}_k - \alpha_k\bar{g}_k, \quad (31)$$

and

$$X_{k+1} = (1 - \beta_k)X_k + \beta_kAQ_k - \alpha_kG_k. \quad (32)$$

### III. MAIN RESULT

Our main contribution is to show that under the strict quantized communication constraints, we can improve upon the previously achievable convergence rates.

*Theorem 1:* Suppose that Assumptions 1 and 2 hold. Let  $\{x_i^k\}$  be generated by Algorithm 1, and  $\{\alpha_k, \beta_k\}$  satisfy

$$\alpha_k = \frac{4/\mu}{k+1}, \quad \beta_k = \frac{4/(1-\sigma_2)}{(k+1)^{3/4}}, \quad (33)$$

where  $\sigma_2$  is the second largest eigenvalue of the adopted averaging matrix  $A$ . Then the output  $z_i^k$  of Algorithm 1 at each agent  $i$  satisfies

$$\begin{aligned} \mathbf{E} [f(z_i^k)] - f^* &\leq \\ \mathcal{O} \left( \frac{\mathbf{E}\|x_0 - x^*\|^2}{k^2} + \frac{n^2(\ln k)^2}{(1-\sigma_2)^2k^{3/4}} + \frac{(\ln k)^2}{(1-\sigma_2)\sqrt{k}} \right). \end{aligned} \quad (34)$$

*Remark 3:* An exact expression for the rate in (34) is given in (83). Here we provide key observations on the result presented in Theorem 1.

We note that the convergence rate of Algorithm 1 in (34) is better than those in [4], [7] (where we ignore the log factor). In particular, a convergence rate  $\mathcal{O}(1/k^{(1-\gamma)/2})$  for some  $\gamma \in (0, 1)$  is shown in [4]. However, this parameter  $\delta$  impacts the speed of the algorithm as mentioned by the authors (see the discussion after Theorem 1 in [4]). Specifically, as  $\delta$  goes to zero, they achieve a rate close to  $1/\sqrt{k}$ , at the cost of spending more iterations due to some constants in their rates getting larger. It is also worth to recall that this work requires the communication of real numbers between agents at every iteration. On the other hand, a convergence rate  $\mathcal{O}(1/k^{1/3})$  is provided in [7].

We note that for unquantized (infinite bandwidth) problems with strongly convex objectives, distributed consensus-based gradient methods achieve  $\mathcal{O}(1/k)$  and linear convergence rates for non-smooth and smooth functions, respectively, see for example [11]. Thus, quantization does slow down the rate of convergence of this method.

### IV. ANALYSIS

In this section, we state and prove two technical Lemmas that will be used to prove Theorem 1.

*Lemma 1:* Let the sequence  $\{x_i^k\}$  be generated by Eq. (13),  $i \in \mathcal{V}$ . Let  $\{\alpha_k, \beta_k\}$  be two sequences of non-negative and non-increasing step sizes. Let  $(1 - \sigma_2)$  be the

spectral gap of the network connectivity. If  $f$  is  $L$ -Lipschitz continuous, then

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} \|Y_{k+1}\|_F^2 &\leq (1 - (1 - \sigma_2)\beta_k) \|Y_k\|_F^2 \\ &\quad + (1 + (1 - \sigma_2)\beta_0) \beta_k^2 \sigma_2^2 n \Delta_k^2 \\ &\quad + \left( \frac{(1 - \sigma_2)\beta_0 + 1}{1 - \sigma_2} \right) L^2 \frac{\alpha_k^2}{\beta_k}. \end{aligned} \quad (35)$$

*Proof:* Beginning with the definition of  $Y_k$ , we obtain the following sequence of identities:

$$\begin{aligned} \|Y_{k+1}\|_F^2 &= \|WX_{k+1}\|_F^2 \\ &\stackrel{(a)}{=} \|W((1 - \beta_k)X_k + \beta_kAQ_k - \alpha_kG_k)\|_F^2 \\ &\stackrel{(b)}{=} \|(1 - \beta_k)Y_k + \beta_kAWQ_k - \alpha_kWG_k\|_F^2 \\ &\stackrel{(c)}{=} (1 - \eta) \underbrace{\|(1 - \beta_k)Y_k + \beta_kAWQ_k\|_F^2}_{\stackrel{\text{def}}{=}(\diamond)} \\ &\quad + \left(1 + \frac{1}{\eta}\right) \alpha_k^2 \|WG_k\|_F^2 \end{aligned}$$

where (a) follows from Eq. (32), (b) follows from the fact that  $WA = AW$ , and (c) follows from Proposition 2 in Appendix A, where  $\eta$  is an arbitrary positive constant that will be specified later in the proof. Let us analyse the term  $(\diamond)$ , starting with the following identity:

$$\begin{aligned} (\diamond) &= \|(1 - \beta_k)Y_k + \beta_kAW(X_k + E_k)\|_F^2 \\ &= \|(1 - \beta_k)I + \beta_kA\|_F^2 \|Y_k\|_F^2 + \beta_k^2 \|AW E_k\|_F^2 \\ &\quad + 2\langle ((1 - \beta_k)I + \beta_kA)Y_k, \beta_kAW E_k \rangle_F. \end{aligned}$$

Taking the conditional expectation with respect to the filtration  $\mathcal{F}_k$ , we obtain:

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k}(\diamond) &\stackrel{(a)}{=} \|(I - (I - A)\beta_k)Y_k\|_F^2 + \beta_k^2 \mathbf{E}_{\mathcal{F}_k} \|AW E_k\|_F^2 \\ &\stackrel{(b)}{\leq} (1 - (1 - \sigma_2)\beta_k)^2 \|Y_k\|_F^2 + n^2 \beta_k^2 \Delta_k^2, \end{aligned}$$

where (a) follows from the unbiasedness of the random quantization scheme, (b) follows from

$$\|AW E_k\|_F^2 \leq \|AW\|_F^2 \|E_k\|_F^2 \stackrel{(c)}{\leq} n^2 \Delta_k^2, \quad (36)$$

and (c) follows from

$$\|AW\|_F^2 \leq \|W\|_F^2 = n - 1 \leq n. \quad (37)$$

Therefore,

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} \|Y_{k+1}\|^2 &\leq (1 + \eta) \left[ (1 - (1 - \sigma_2)\beta_k)^2 \|Y_k\|_F^2 + n^2 \beta_k^2 \Delta_k^2 \right] \\ &\quad + \left(1 + \frac{1}{\eta}\right) \alpha_k^2 L^2 \end{aligned} \quad (38)$$

Setting  $\eta = (1 - \sigma_2)\beta_k$ , we obtain

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} \|Y_{k+1}\|^2 &\leq (1 + (1 - \sigma_2)\beta_k) \left[ (1 - (1 - \sigma_2)\beta_k)^2 \|Y_k\|_F^2 \right. \\ &\quad \left. + 2\beta_k^2 n^2 \Delta_k^2 \right] + \left(1 + \frac{1}{(1 - \sigma_2)\beta_k}\right) \alpha_k^2 L^2. \end{aligned} \quad (39)$$

Next, consider the following sequence of inequalities:

$$\|WG_k\|_F^2 \leq \|G_k\|_F^2 \leq \sum_{i=1}^n L_i^2 \leq L^2. \quad (40)$$

After some algebraic manipulations, we get

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} \|Y_{k+1}\|^2 &\leq (1 - (1 - \sigma_2)^2 \beta_k) (1 - (1 - \sigma_2) \beta_k)^2 \|Y_k\|_F^2 \\ &\quad + (1 + (1 - \sigma_2) \beta_k) n^2 \beta_k^2 \Delta_k^2 \\ &\quad + \left( \frac{(1 - \sigma_2) \beta_k + 1}{(1 - \sigma_2)} \right) L^2 \frac{\alpha_k^2}{\beta_k}. \end{aligned} \quad (41)$$

Finally, using the fact that  $\{\beta_k\}$  is a nonincreasing sequence, we obtain the following inequality

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} \|Y_{k+1}\|^2 &\leq (1 - (1 - \sigma_2) \beta_k) \|Y_k\|_F^2 \\ &\quad + (1 - (1 - \sigma_2) \beta_0) n^2 \beta_k^2 \Delta_k^2 \\ &\quad + \left( \frac{(1 - \sigma_2) \beta_0 + 1}{(1 - \sigma_2)} \right) L^2 \frac{\alpha_k^2}{\beta_k}. \end{aligned} \quad (42)$$

■

For the second technical result in this section, we will obtain an upper bound on the conditional expectation of  $r_{k+1}$  with respect to the filtration  $\mathcal{F}_k$ . To do so, we assume strong convexity of the global objective function,  $f$ .

*Lemma 2:* Assume that  $f$  is a  $\mu$ -strongly convex function with  $L$ -Lipschitz continuous gradient. Let  $\{\alpha_k\}$  be a non-increasing, non-negative sequence, and

$$r_{k+1} = \|\bar{x}_{k+1} - x^*\|^2. \quad (43)$$

Then, the following inequality holds:

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k} [r_{k+1}] &\leq \left(1 - \frac{\mu}{2} \alpha_k\right) r_k + \alpha_k^2 L^2 + \beta_k^2 \Delta_k^2 \\ &\quad + 2\alpha_k ((f(x^*) - f(x_k^l)) + \alpha_k (L + \frac{8L^2}{\mu})) \|Y_k\|_F^2 \end{aligned} \quad (44)$$

for all  $l \in \mathcal{V}$ .

*Proof:* We begin with the following identity:

$$\begin{aligned} \|\bar{x}_{k+1} - x^*\|^2 &= \|(1 - \beta_k) \bar{x}_k + \beta_k \bar{q}_k - \alpha_k \bar{g}_k - x^*\|^2 \\ &= \underbrace{\|\bar{x}_k - x^* - \alpha_k \bar{g}_k\|^2}_{\stackrel{\text{def}}{=} A} + \underbrace{\beta_k^2 \|\bar{e}_k\|^2}_{\stackrel{\text{def}}{=} B} \\ &\quad + \underbrace{2\beta_k (\bar{x}_k - x^* - \alpha_k \bar{g}_k)^T \bar{e}_k}_{\stackrel{\text{def}}{=} C}. \end{aligned} \quad (45)$$

We proceed to analyse each of the terms  $A$ ,  $B$ , and  $C$  above.

First, consider

$$A \stackrel{\text{def}}{=} \|\bar{x}_k - x^* - \alpha_k \bar{g}_k\|^2 \quad (46)$$

Then, the following identity holds:

$$\begin{aligned} A &= \|\bar{x}_k - x^*\|^2 + \alpha_k^2 \|\bar{g}_k\|^2 \\ &\quad - \underbrace{2\alpha_k \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_k^i)^T (\bar{x}_k - x^*)}_{\stackrel{\text{def}}{=} A_1}. \end{aligned} \quad (47)$$

Consider the following identity

$$\begin{aligned} A_1 &= -2\alpha_k \frac{1}{n} \sum_{i=1}^n \underbrace{\left( \nabla f_i(x_k^i) - \nabla f_i(\bar{x}_k) \right)^T (\bar{x}_k - x^*)}_{\stackrel{\text{def}}{=} A_{11}} \\ &\quad - \underbrace{2\alpha_k \nabla f(\bar{x}_k)^T (\bar{x}_k - x^*)}_{\stackrel{\text{def}}{=} A_{12}}. \end{aligned} \quad (48)$$

Then the following sequence of inequalities hold:

$$A_{11} \stackrel{(a)}{\leq} 2\alpha_k \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k^i) - \nabla f_i(\bar{x}_k)\| \|\bar{x}_k - x^*\| \quad (49)$$

$$\stackrel{(b)}{\leq} 2\alpha_k \frac{1}{n} \sum_{i=1}^n L \|x_k^i - \bar{x}_k\| \|\bar{x}_k - x^*\| \quad (50)$$

$$\stackrel{(c)}{\leq} \alpha_k \frac{L}{n} \sum_{i=1}^n \left( \eta \|x_k^i - \bar{x}_k\|^2 + \frac{1}{\eta} \|\bar{x}_k - x^*\|^2 \right) \quad (51)$$

$$\stackrel{(d)}{=} \alpha_k \eta L \|X_k - \mathbf{1} \bar{x}_k\|_F^2 + \alpha_k L \frac{1}{\eta} \|\bar{x}_k - x^*\|^2, \quad (52)$$

where (a) follows from the Cauchy-Schwarz inequality, (b) follows from the Lipschitz continuity of  $f_i$ ,  $i = 1, \dots, n$ , (c) follows from Proposition 2, and (d) follows from Eq. (30). Setting  $\eta = 4L/\mu$ , we obtain:

$$A_{11} \leq \alpha_k \frac{4L^2}{\mu} \|Y_k\|_F^2 + \alpha_k \frac{4}{\mu} \|\bar{x}_k - x^*\|^2. \quad (53)$$

We proceed to bounding  $A_{12}$  as follows:

$$A_{12} = 2\alpha_k \nabla f(\bar{x}_k)^T (x^* - \bar{x}_k) \quad (54)$$

$$\stackrel{(a)}{\leq} 2\alpha_k \left( f(x^*) - f(\bar{x}_k) - \frac{\mu}{2} \|x^* - \bar{x}_k\|^2 \right) \quad (55)$$

$$\stackrel{(b)}{=} 2\alpha_k (f(x^*) - f(x_k^l) + f(x_k^l) - f(\bar{x}_k)) - \alpha_k \mu \|x^* - \bar{x}_k\|^2 \quad (56)$$

$$\stackrel{(c)}{\leq} 2\alpha_k (f(x^*) - f(x_k^l)) + 2\alpha_k \left( \nabla f(\bar{x}_k)^T (x_k^l - \bar{x}_k) + \frac{L}{2} \|x_k^l - \bar{x}_k\|^2 \right) - \alpha_k \mu \|x^* - \bar{x}_k\|^2 \quad (57)$$

$$\stackrel{(d)}{=} 2\alpha_k (f(x^*) - f(x_k^l)) + 2\alpha_k \left( (\nabla f(\bar{x}_k) - \nabla f(x^*))^T (x_k^l - \bar{x}_k) \right) + \alpha_k L \|x_k^l - \bar{x}_k\|^2 - \alpha_k \mu \|x^* - \bar{x}_k\|^2 \quad (58)$$

$$\stackrel{(e)}{\leq} 2\alpha_k (f(x^*) - f(x_k^l)) + 2\alpha_k L \|\bar{x}_k - x^*\| \|x_k^l - \bar{x}_k\| + \alpha_k L \|x_k^l - \bar{x}_k\|^2 - \alpha_k \mu \|x^* - \bar{x}_k\|^2 \quad (59)$$

$$\stackrel{(f)}{\leq} 2\alpha_k (f(x^*) - f(x_k^l)) + \alpha_k L \left( \eta \|\bar{x}_k - x^*\|^2 + \frac{1}{\eta} \|x_k^l - \bar{x}_k\|^2 \right) + \alpha_k L \|x_k^l - \bar{x}_k\|^2 - \alpha_k \mu \|x^* - \bar{x}_k\|^2, \quad (60)$$

where (a) follows from  $f$  being  $\mu$ -strongly convex, (b) follows from adding and subtracting  $f(x_k^l)$ , (c) follows the  $L$ -Lipschitz continuity of the gradient, and (d) follows from

the first order optimality condition  $\nabla f(x^*) = 0$ , (e) follows from the Cauchy-Schwarz inequality, and (f) follows from Proposition 2. Setting  $\eta = 4L/\mu$ , we obtain:

$$A_{12} \leq 2\alpha_k(f(x^*) - f(x_k^l)) - \alpha_k \frac{3\mu}{4} \|x^* - \bar{x}_k\|^2 + \alpha_k \left(L + \frac{4L^2}{\mu}\right) \|x_k^l - \bar{x}_k\|^2. \quad (61)$$

Therefore,

$$A_1 \leq 2\alpha_k((f(x^*) - f(x_k^l)) - \alpha_k \frac{\mu}{2} \|\bar{x}_k - x^*\|^2 + \alpha_k(L + \frac{4L^2}{\mu}) \|x_k^l - \bar{x}_k\|^2 + \alpha_k \frac{4L^2}{\mu} \|Y_k\|_F^2). \quad (62)$$

We next turn our attention to  $B$ . From the convexity of  $\|\cdot\|^2$ , we obtain the following upper bound:

$$B = \beta_k^2 \left\| \frac{1}{n} \sum_{i=1}^n e_k^i \right\|^2 \leq \beta_k^2 \frac{1}{n} \sum_{i=1}^n \|e_k^i\|^2. \quad (63)$$

Finally, taking the expectation of  $\|\bar{x}_{k+1} - x^*\|^2$  with respect to  $\mathcal{F}_k$ , the term  $C$  vanishes due to the fact that the random quantizer is unbiased Eq. (8). Moreover, due to Eq. (10) the term  $B$  is upper bounded by

$$\mathbf{E}_{\mathcal{F}_k}[B] \leq \beta_k^2 \Delta_k^2. \quad (64)$$

Thus,

$$\mathbf{E}_{\mathcal{F}_k} \|\bar{x}_{k+1} - x^*\|^2 \leq \|\bar{x}_k - x^*\|^2 + \alpha_k^2 \|\bar{g}_k\|^2 + \beta_k^2 \Delta_k^2 + 2\alpha_k((f(x^*) - f(x_k^l)) - \alpha_k \frac{\mu}{2} \|\bar{x}_k - x^*\|^2 + \alpha_k(L + \frac{4L^2}{\mu}) \|x_k^l - \bar{x}_k\|^2 + \alpha_k \frac{4L^2}{\mu} \|Y_k\|_F^2). \quad (65)$$

Finally, since

$$\left\| \frac{1}{n} \sum_{i=1}^n g_k^i \right\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|g_k^i\|^2 \leq L^2, \quad (66)$$

we have:

$$\mathbf{E}_{\mathcal{F}_k}[r_{k+1}] \leq \left(1 - \frac{\mu}{2}\alpha_k\right)r_k + \alpha_k^2 L^2 + \beta_k^2 \Delta_k^2 + 2\alpha_k((f(x^*) - f(x_k^l)) + \alpha_k(L + \frac{4L^2}{\mu}) \|x_k^l - \bar{x}_k\|^2 + \alpha_k \frac{4L^2}{\mu} \|Y_k\|_F^2). \quad (67)$$

■

## V. PROOF OF THEOREM 1 AND CONVERGENCE RATE IMPROVEMENT

We begin by specifying the step-size sequences

$$\alpha_k = \frac{4/\mu}{k+1} \quad (68)$$

and

$$\beta_k = \frac{4/(1-\sigma_2)}{(k+1)^{3/4}}. \quad (69)$$

*Lemma 3:* Let  $\{\eta_k\}$  be a non-increasing sequence defined as:

$$\eta_k \stackrel{\text{def}}{=} \eta \frac{\alpha_k}{\beta_k}. \quad (70)$$

Define the following Lyapunov function:

$$\mathcal{V}_{k+1} \stackrel{\text{def}}{=} r_{k+1} + \eta_{k+1} \|Y_{k+1}\|_F^2. \quad (71)$$

The following inequality holds

$$\mathbf{E}[\mathcal{V}_{k+1}] \leq \left(\frac{k-1}{k+1}\right) \mathbf{E}[\mathcal{V}_k] - \left(\frac{8/\mu}{k+1}\right) (f(x_k^l) - f(x^*)) + \Gamma_k, \quad (72)$$

where

$$\Gamma_k = \frac{16/\mu^2}{(k+1)^2} + \frac{40L^2(L + 8L^2/\mu)/\mu^3}{(k+1)^{3/2}} + \left[ \frac{4}{(1-\sigma_2)(k+1)^{3/2}} + \frac{320(L + 8L^2/\mu)n^2}{(1-\sigma_2)^2(k+1)^{7/4}} \right] \times \left(\frac{Cd}{2^b-1}\right) a^2 \left(\sum_{t=0}^{k-1} \alpha_k\right)^2. \quad (73)$$

*Proof:* The proof is in Appendix B. ■

Equipped with Lemmas 1, 2, and 3, we are now ready to prove Theorem 1.

*Proof: (Proof of Theorem 1)*

From Lemma 3, multiplying both sides of Eq. (72) by  $(k+1)^2$  yields:

$$(k+1)^2 \mathbf{E}[\mathcal{V}_{k+1}] \leq (k^2-1) \mathbf{E}[\mathcal{V}_k] - (8/\mu)(k+1)(f(x_k^l) - f(x^*)) + (k+1)^2 \Gamma_k, \quad (74)$$

which, due to the non-negativity of  $\mathcal{V}_k$ , can be further relaxed to:

$$(k+1)^2 \mathbf{E}[\mathcal{V}_{k+1}] \leq k^2 \mathbf{E}[\mathcal{V}_k] - (8/\mu)(k+1)(f(x_k^l) - f(x^*)) + (k+1)^2 \Gamma_k. \quad (75)$$

Summing both sides of Eq. (76) from  $k=1$  to  $T > 1$ , we get:

$$(T+1)^2 \mathbf{E}[\mathcal{V}_{T+1}] \leq \mathbf{E}[\mathcal{V}_1] - (8/\mu) \sum_{k=1}^T (k+1)(f(x_k^l) - f(x^*)) + \sum_{k=1}^T (k+1)^2 \Gamma_k. \quad (76)$$

We will now analyze  $\sum_{k=1}^T (k+1)^2 \Gamma_k$ . First, notice that:

$$\sum_{t=0}^{T-1} \alpha_t = \frac{4}{\mu} \sum_{t=0}^{T-1} \frac{1}{t+1} \stackrel{(a)}{\leq} \frac{4}{\mu} \int_0^{T-1} \frac{1}{t+1} dt = \frac{4}{\mu} \ln T, \quad (77)$$

where (a) follows from the Cauchy-Maclaurin *integral test*. Then, we consider

$$\sum_{k=1}^T (k+1)^{1/2} \stackrel{(b)}{\leq} \int_1^T (t+1)^{1/2} dt \stackrel{(c)}{\leq} \frac{2}{3} (T+1)^{3/2}, \quad (78)$$

where (b) follows from the integral test, and (c) is obtained by dropping the term corresponding to the lower-end of the integration interval. Therefore, we have:

$$\begin{aligned} \sum_{k=1}^T (k+1)^{1/2} \left( \sum_{t=0}^{k-1} \alpha_t \right)^2 &\stackrel{(d)}{\leq} \sum_{k=1}^T (k+1)^{1/2} \left( \sum_{t=0}^{T-1} \alpha_t \right)^2 \quad (79) \\ &\stackrel{(e)}{\leq} \frac{32}{3\mu^2} (\ln T)^2 (T+1)^{3/2}. \quad (80) \end{aligned}$$

where (d) follows from  $k \leq T$ , and (e) follows from the integral test.

Similarly, we obtain

$$\sum_{k=1}^T (k+1)^{1/4} \left( \sum_{t=0}^{k-1} \alpha_t \right)^2 \leq \frac{64}{5\mu^2} (\ln T)^2 (T+1)^{5/4}, \quad (81)$$

Using the inequalities above, we get:

$$\begin{aligned} \sum_{k=1}^T (k+1)^2 \Gamma_k &\leq \frac{16}{\mu} T + \left[ \frac{(128/3\mu^2)}{(1-\sigma_2)} (\ln T)^2 (T+1)^{3/2} \right. \\ &+ \left. \frac{(256n^2(L+8L^2/\mu)/\mu)}{(1-\sigma_2)^2} (\ln T)^2 (T+1)^{5/4} \right] \left( \frac{Cd}{2^b-1} \right)^2 \\ &+ \frac{(8L(L+8L^2/\mu))}{3\mu^3} (T+1)^{3/2}. \quad (82) \end{aligned}$$

Using Eq. (82) in Eq. (76) and dividing both sides of the resulting inequality by  $\sum_{k=1}^T (k+1) \approx (T+1)^2$ , we get:

$$\begin{aligned} \mathbf{E}[\mathcal{V}_{T+1}] &\leq \frac{\mathbf{E}[\mathcal{V}_1]}{(T+1)^2} - \frac{\frac{8}{\mu} \sum_{k=1}^T (k+1) (f(x_k^l) - f(x^*))}{\sum_{k=1}^T (k+1)} \\ &+ \frac{16}{\mu} \frac{T}{(T+1)^2} + \left[ \frac{(128/3\mu^2)}{(1-\sigma_2)} \frac{(\ln T)^2}{(T+1)^{1/2}} \right. \\ &+ \left. \frac{(256n^2(L+8L^2/\mu)/\mu)}{(1-\sigma_2)^2} \frac{(\ln T)^2}{(T+1)^{3/4}} \right] \left( \frac{Cd}{2^b-1} \right)^2 \\ &+ \frac{(8L(L+8L^2/\mu))}{3\mu^3} \frac{1}{(T+1)^{1/2}}, \end{aligned}$$

which by rearranging both sides, dropping the positive term  $\mathbf{E}[\mathcal{V}_{T+1}]$ , and using Jensen's inequality, we obtain

$$\begin{aligned} \mathbf{E}[f(z_i^k)] - f^* &\leq \frac{\sum_{k=1}^T (k+1) (f(x_k^l) - f(x^*))}{\sum_{k=1}^T (k+1)} \\ &\leq \frac{\mu \mathbf{E}[\mathcal{V}_1]}{8(T+1)^2} + \frac{2}{T+1} + \frac{16}{3\mu(1-\sigma_2)} \left( \frac{Cd}{2^b-1} \right)^2 \frac{(\ln T)^2}{(T+1)^{1/2}} \\ &+ \frac{4n^2(L+8L^2)}{(1-\sigma_2)^2} \left( \frac{Cd}{2^b-1} \right)^2 \frac{(\ln T)^2}{(T+1)^{3/4}} \\ &+ \frac{(8L(L+8L^2/\mu))}{3\mu^3} \frac{1}{(T+1)^{1/2}}, \quad (83) \end{aligned}$$

where  $z_i^k$  is defined in Algorithm (1). Clearly, the right-hand side of Eq. (83) is dominated by a term that decays with

$$\mathcal{O}\left(\frac{(\ln T)^2}{(1-\sigma_2)(T+1)^{1/2}}\right).$$

## VI. A NUMERICAL EXAMPLE

In this section, we provide a numerical experiment to illustrate for the performance of Eq. (12) under the random quantization scheme given in Section II-C. Our simulation is similar to the ones in [7], however, we do not require any projection to a predefined compact set. As a result, the quantization intervals at each step in our algorithm are different and increasing in  $k$ .

We apply Eq. (12) for solving the distributed variant of the popular linear regression problem. The goal of this problem is to find a linear relationship between a set of variables and some real valued outcome. That is, given a training set  $S = \{(w_i, b_i) \in \mathbb{R}^d \times \mathbb{R}\}$  for  $i = 1, \dots, n$ . We want to find the solution of the following optimization problem

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^n (w_i^T x - b_i)^2, \quad (84)$$

In this case, we know that the optimal centralized solution can be efficiently found by solving a least squares problem. The challenge is to how to solve it efficiently in a decentralized fashion, under strict quantization constraints.

Suppose that  $w_i \in \mathbb{R}^5$ . We consider simulated training data sets where  $(w_i, b_i)$ ,  $i \in \{1, \dots, n\}$ , are generated randomly and independently according to the following uniform distributions:

$$w_i(j) \sim \mathcal{U}[0, 0.65], \quad j \in \{1, \dots, d\} \quad (85)$$

and

$$b_i \sim \mathcal{U}[0, 0.45]. \quad (86)$$

We study the performance of our distributed gradient method on the undirected connected graph of 40 nodes shown in Fig. 1, i.e.,  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $n = |\mathcal{V}| = 40$ . Our graph is generated randomly, similar to the ones in [7]. Finally, the adjacency matrix  $A$  is chosen as a Lazy Metropolis matrix corresponding to  $\mathcal{G}$ , i.e.,

$$A = [a_{ij}] = \begin{cases} \frac{1}{2(\max\{|\mathcal{N}_i|, |\mathcal{N}_j|\})}, & \text{if } (i, j) \in \mathcal{E} \\ 0, & \text{if } (i, j) \notin \mathcal{E} \text{ and } i \neq j \\ 1 - \sum_{j \in \mathcal{N}_i} a_{ij}, & \text{if } i = j \end{cases}$$

In this example, the state variables  $x_k^i$  were quantized using 16 bits, without the need for periodic transmissions real numbers among agents. In Fig. 2, we clearly see the convergence of the quantized vs. the unquantized algorithm, which asymptotically achieve the same value of the objective function, but there is a slight degradation due to the quantization error. It is important to note that despite the fact that the instantaneous quantization error increases with time, our algorithm successfully compensates for that by driving the estimates swiftly to the desired optimal solution. Hence, the excellent asymptotic convergence.

## VII. CONCLUSIONS

We introduced a novel two-time scale algorithm for distributed gradient descent under random quantization. Unlike other works, our algorithm does not require periodic transmission of real valued messages that specify the quantization

■

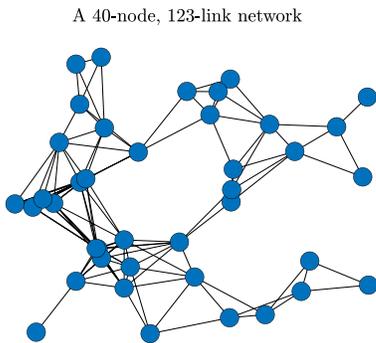


Fig. 1. Network with 40 nodes and 123 communication links used in our distributed linear regression application.

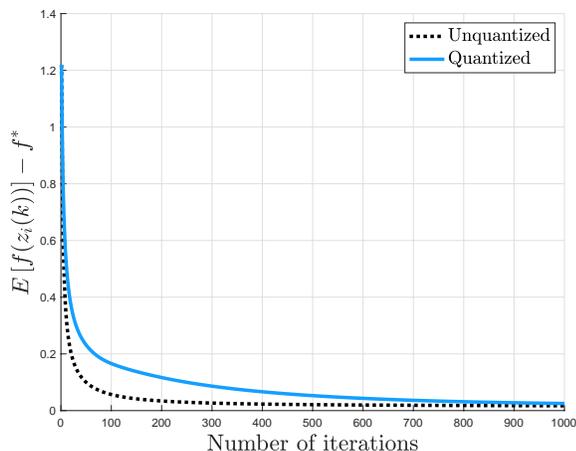


Fig. 2. Performance of our algorithm with (solid line) and without (dashed line) quantization. In this example, the state variables was quantized using 16 bits/dimension.

bins for encoding and decoding, which in principle violates the finite bandwidth constraint at each link. We showed that the convergence rate of our algorithm is  $\mathcal{O}((\ln(k))^2/\sqrt{k})$ , which improves on the current state of the art without the use of adaptive quantization.

## APPENDIX I USEFUL INEQUALITIES

*Proposition 2:* Let  $a, b \in \mathbb{R}^d$ . For  $\eta > 0$ , the following inequality holds:

$$\|a + b\|^2 \leq (1 + \eta)\|a\|^2 + \left(1 + \frac{1}{\eta}\right)\|b\|^2. \quad (87)$$

or, equivalently,

$$\langle a, b \rangle \leq \frac{1}{2} \left( \eta \|a\|^2 + \frac{1}{\eta} \|b\|^2 \right). \quad (88)$$

## APPENDIX II PROOF OF LEMMA 3

Using Lemmas 2 and 3, after some algebraic manipulations, we obtain

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k}[\mathcal{V}_{k+1}] &\leq \left(1 - \frac{\mu}{2}\alpha_k\right)\mathcal{V}_k + 2\alpha_k(f^* - f(x_k^l)) \\ &+ \left[\alpha_k^2 L^2 + (\beta_k^2 + 2\eta(1 - (1 - \sigma_2)\beta_0)\alpha_k\beta_k n^2) \left(\frac{Cd}{2^b - 1}\right)^2 \left(\sum_{t=0}^{k-1} \alpha_t\right)^2\right. \\ &\quad \left. + \eta \left(\frac{(1 - \sigma_2)\beta_0 + 1}{1 - \sigma_2}\right) L^2 \frac{\alpha_k^3}{\beta_k^2}\right] \\ &+ \|Y_k\|_F^2 \left(L + \frac{8L^2}{\mu} + \frac{\mu}{2}\eta_k - \eta(1 - \sigma_2)\right)\alpha_k. \quad (89) \end{aligned}$$

Choosing  $\eta = 2(L + L^2/8)/(1 - \sigma_2)$ , since  $\alpha_0/\beta_0 = \mu/(1 - \sigma_2)$ , we may drop the last term in Eq. (89). Thus,

$$\begin{aligned} \mathbf{E}_{\mathcal{F}_k}[\mathcal{V}_{k+1}] &\leq \left(1 - \frac{\mu}{2}\alpha_k\right)\mathcal{V}_k + 2\alpha_k(f^* - f(x_k^l)) \\ &+ \left[\alpha_k^2 L^2 + (\beta_k^2 + 2\eta(1 - (1 - \sigma_2)\beta_0)\alpha_k\beta_k n^2) \left(\frac{Cd}{2^b - 1}\right)^2 \left(\sum_{t=0}^{k-1} \alpha_t\right)^2\right. \\ &\quad \left. + \eta \left(\frac{(1 - \sigma_2)\beta_0 + 1}{1 - \sigma_2}\right) L^2 \frac{\alpha_k^3}{\beta_k^2}\right] \quad (90) \end{aligned}$$

Substituting the step sequences  $\alpha_k$  and  $\beta_k$  into Eq. (90), after some algebra, and taking the expectation over all possible  $\mathcal{F}_k$  on both sides, we obtain Eq. (72).

## REFERENCES

- [1] A. Nedic, "Distributed gradient methods for convex machine learning problems in networks: Distributed optimization," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 92–101, 2020.
- [2] T. Yang, X. Yi, J. Wu, Y. Yuan, D. Wu, Z. Meng, Y. Hong, H. Wang, Z. Lin, and K. H. Johansson, "A survey of distributed optimization," *Annual Reviews in Control*, vol. 47, pp. 278–305, 2019.
- [3] J. Li, G. Chen, Z. Wu, and X. He, "Distributed subgradient method for multi-agent optimization with quantized communication," *Mathematical Methods in the Applied Sciences*, vol. 40, no. 4, pp. 1201–1213, 2017.
- [4] A. Reiszadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "An exact quantized decentralized gradient descent algorithm," *IEEE Transactions on Signal Processing*, vol. 67, no. 19, pp. 4934–4947, 2019.
- [5] H. Taheri, A. Mokhtari, H. Hassani, and R. Pedarsani, "Quantized decentralized stochastic learning over directed graphs," in *Proceedings of the 37th International Conference on Machine Learning*, vol. 119, 13–18 Jul 2020, pp. 9324–9333.
- [6] T. T. Doan, S. T. Maguluri, and J. Romberg, "Fast convergence rates of distributed subgradient methods with adaptive quantization," *IEEE Transactions on Automatic Control*, pp. 1–1, 2020.
- [7] —, "Convergence rates of distributed gradient methods under random quantization: A stochastic approximation approach," *IEEE Transactions on Automatic Control*, 2020.
- [8] A. Chattopadhyay and U. Mitra, "Dynamic sensor subset selection for centralized tracking of an iid process," *IEEE Transactions on Signal Processing*, vol. 68, pp. 3209–3224, 2020.
- [9] T. T. Doan, S. T. Maguluri, and J. Romberg, "On the convergence of distributed subgradient methods under quantization," in *2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2018, pp. 567–574.
- [10] E. Cinlar, *Probability and stochastics*. Springer Science & Business Media, 2011, vol. 261.
- [11] A. Nedić, A. Olshevsky, and M. G. Rabbat, "Network topology and communication-computation tradeoffs in decentralized optimization," *Proceedings of the IEEE*, vol. 106, no. 5, pp. 953–976, 2018.