

Quantized Adam with Error Feedback

CONGLIANG CHEN, The Chinese University of Hong Kong, Shenzhen, China

LI SHEN*, JD Explore Academy, China

HAOZHI HUANG, Tencent AI Lab, China

WEI LIU, Tencent, China

In this paper, we present a distributed variant of adaptive stochastic gradient method for training deep neural networks in the parameter-server model. To reduce the communication cost among the workers and server, we incorporate two types of quantization schemes, i.e., gradient quantization and weight quantization, into the proposed distributed Adam. Besides, to reduce the bias introduced by quantization operations, we propose an error-feedback technique to compensate for the quantized gradient. Theoretically, in the stochastic nonconvex setting, we show that the distributed adaptive gradient method with gradient quantization and error-feedback converges to the first-order stationary point, and that the distributed adaptive gradient method with weight quantization and error-feedback converges to the point related to the quantized level under both the single-worker and multi-worker modes. At last, we apply the proposed distributed adaptive gradient methods to train deep neural networks. Experimental results demonstrate the efficacy of our methods.

Additional Key Words and Phrases: Adam, Quantized Communication, Error Feedback

ACM Reference Format:

Congliang Chen, Li Shen, Haozhi Huang, and Wei Liu. 2021. Quantized Adam with Error Feedback. 1, 1 (June 2021), 24 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Recently, deep neural networks [8, 21] achieve high performances in many applications, such as computer vision [10, 20], natural language processing [6], speech recognition [1], reinforcement learning [29, 34], etc. However, a huge deep neural network contains millions of parameters, so its training procedure requires a large amount of training data [5, 40], which may not be stored in a single machine. In addition, due to some privacy issues [27, 43], all the training data cannot be sent to a single machine but can be stored in different devices. Therefore, how to accelerate the training process by using multiple machines over large-scale data or distributed data has already been a hot topic in both industrial and academic communities [18, 23, 26, 42].

An efficient approach to tackle this problem is to develop distributed training algorithms for the huge neural networks [4]. Most of the distributed algorithms can be summarized into two categories: one is the parameter-server [36] model (or called centralized model) shown in Fig. 1, and the other is the decentralized model [25] shown in Fig. 2. For the centralized model in Fig. 1, there are one parameter server and multiple workers. In an update iteration, all workers report the

*Li Shen is the corresponding author.

Authors' addresses: Congliang Chen, chcoli2007@163.com, The Chinese University of Hong Kong, Shenzhen, Shenzhen, Guangdong, China, 518000; Li Shen, JD Explore Academy, Beijing, China, mathshenli@gmail.com; Haozhi Huang, Tencent AI Lab, Shenzhen, Guangdong, China, matthzhuang@tencent.com; Wei Liu, Tencent, Shenzhen, Guangdong, China, wl2223@columbia.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2021 Association for Computing Machinery.

XXXX-XXXX/2021/6-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

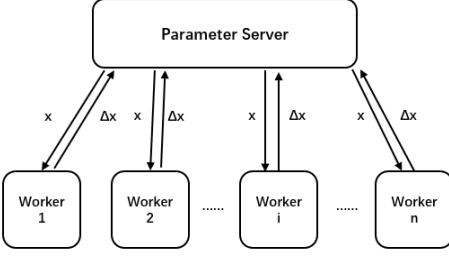


Fig. 1. The Centralized distributed model.

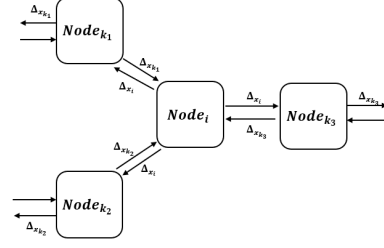


Fig. 2. The Decentralized distributed model.

update vectors to the parameter server. After gathering all the update vectors, the parameter server will update the parameters and send the parameters to all workers. While for the decentralized model in Fig. 2, there are n nodes working simultaneously. In each update iteration, each worker computes its update vector respectively and communicates with its neighbors, and then updates its own parameters. When we use a distributed training algorithm such as distributed stochastic gradient descent [23] in either the centralized model or the decentralized model, plenty of update vectors have to be communicated among different devices. Then, a communication issue emerges for huge networks.

To accelerate the distributed training process of huge deep learning models, we propose a new distributed adaptive stochastic gradient method with gradient quantization, weight quantization, and error-feedback in the parameter server model, as shown in Fig. 1. In what follows, we elaborate on each component used in the proposed method:

Quantization. Note that both gradient quantization and weight quantization are introduced in the proposed method to reduce the communication cost among the workers and the parameter server. Specifically, weight quantization is performed on the parameter server and the quantized weights are then broadcast to all the workers. Weight quantization is introduced because of the consideration of limited storage in edge devices. Meanwhile, gradient quantization is performed on each worker and then the quantized gradients are reported to the server. Thanks to the double quantization schemes, the communication cost can be largely reduced. In addition, for some resource-limited devices, storage is another issue. Weight quantization can also be used to reduce the deep neural network model size efficiently [9, 31, 45]. Especially, in federated learning, a distributed device may be smartphones or Internet of things devices, which may encounter both the storage issue and the communication issue. Thus, the weight quantization and gradient quantization schemes can jointly solve these two issues.

Adaptive learning rate. To ease the labor of tuning learning rate, we also adopt the adaptive learning rate as [3, 7, 11, 16, 33, 47] in the proposed method. Here, the adaptive learning rate is calculated by a similar definition to those in RMSProp [11] and Adam [16], except that the noisy gradients are estimated with quantized weights. Moreover, to guarantee the convergence of the proposed method, we set the exponential moving average parameter in estimating the adaptive learning rate the same as that used in Zou et al.[48].

Error-feedback. In the proposed method, an error-feedback technique is leveraged to reduce the bias introduced by gradient quantization. The error-feedback technique is also performed on each worker. Actually, the error-feedback technique is motivated by Karimireddy et al. [14] by introducing an additional term as the compensation term for the quantized gradient. However, due to the introduced adaptive learning rate and momentum, the compensation term is slightly different from that in Karimireddy et al. [14]. To the best of our knowledge, this is the first work that simultaneously employs the adaptive learning rate and the error-feedback technique.

Besides, we establish the convergence rate of the proposed algorithm. In the stochastic nonconvex setting, we show that the distributed adaptive stochastic gradient method with gradient quantization and error-feedback converges to the first-order stationary point, and that the distributed adaptive stochastic gradient method with weight quantization and error-feedback converges to the point related to the quantized level under both the single-worker and multi-worker modes. At last, we apply the proposed distributed adaptive method to train deep learning models, such as LeNet [22] on the MNIST dataset [22] and ResNet-101 [10] on the CIFAR100 dataset [19], respectively. The experimental results demonstrate the effectiveness of weight quantization, gradient quantization, and the error-feedback technique working in concert with distributed adaptive stochastic gradient method. Here, we summarize our contributions in three-fold:

- We propose a distributed variant of the adaptive stochastic gradient method to train deep learning models. The proposed approach exploits gradient quantization, weight quantization, and the error-feedback technique to accelerate the training process.
- We establish the convergence rates of the proposed distributed adaptive stochastic gradient algorithms with weight quantization, gradient quantization, and error-feedback in the non-convex stochastic setting under the single-worker and multi-worker environments, which are far different from the stochastic gradient setting because adaptive learning rate is introduced into the algorithms.
- We apply the proposed algorithms to train deep learning models including LeNet and ResNet-101. The experiments demonstrate the efficacy of the proposed algorithms.

2 RELATED WORKS

In this section, we enumerate several works that are most related to this work. We split the related works into two categories: distributed quantized algorithms and adaptive learning rate.

2.1 Distributed Quantized Algorithms

The quantization functions can be divided into two categories: unbiased quantization functions and biased quantization functions. For unbiased quantization functions, Wen et al. [39] showed that with an unbiased ternary quantization function, the distributed stochastic gradient descent algorithm can almost surely converge to a minimum point. Jiang et al. [13] showed with an unbiased quantization function, the centralized distributed stochastic gradient descent algorithm can converge with convergence rate $O(1/\sqrt{T})$. Besides, Hou et al. [12] showed that in the stochastic convex setting, with gradient quantization solely, the algorithm they proposed will converge to the optimal solution, while with weight quantization the algorithm will converge to the point near the optimal solution which is related to the weight quantization level. However, they can only deal with the unbiased quantization function, which limits the use of both algorithms and theorems.

For biased quantization functions, the main issue is to eliminate the biased error during optimization. A common technique to tackle this issue is error-feedback, where each worker stores the error of the quantization and adds the error term to the next communication before quantization. Based on the decentralized model in Fig. 2, Tang et al. [37] and Koloskova et al. [17] showed that distributed stochastic gradient descent with quantized communication and error-feedback can converge to a stationary point in the nonconvex setting with convergence rate $O(1/\sqrt{T})$. Based on the centralized model in Fig. 1, Zhou et al. [45] and Wu et al. [41] showed that few bits or integer networks can be trained empirically. Zheng et al. [44] showed the convergence of the algorithm with a block quantization function in the nonconvex setting.

Among the above-mentioned algorithms, Hou et al. [12] is the most related work to our proposed algorithm. However, their proposed algorithms do not adopt unbiased quantization on gradients.

Moreover, they do not incorporate momentum acceleration terms into their algorithm to accelerate its piratical performance. In addition, the convergence analysis in Hou et al. [12] is merely restricted to the stochastic convex setting, which makes their algorithm heuristic when it is applied to train deep learning models. By contrast, the convergence rates of our algorithms are established in the more difficult nonconvex setting. In this work, we first extend the error-feedback technique to adaptive stochastic gradient method (Adam) and then establish its convergence in the nonconvex setting, and we compare the most related works in Table 1.

Method	gradient quantization	weight quantization	convexity	communication	convergence rate
Wen et al. [39]	unbiased	no	nonconvex	centralize	$O(1/\sqrt{T})$
Zheng et al. [44]	biased	no	nonconvex	centralize	$O(1/\sqrt{T})$
Tang et al. [37]	biased	no	nonconvex	decentralize	$O(1/\sqrt{T})$
Koloskova et al. [17]	biased	no	nonconvex	decentralize	$O(1/\sqrt{T})$
Hou et al. [12]	unbiased	yes	convex	centralize	$O(1/\sqrt{T})$
Ours	biased	yes	nonconvex	centralize	$O(1/\sqrt{T})$

Table 1. Comparison among different methods.

2.2 Adaptive Learning Rate

Adaptive learning rate, as a popular optimization technique for training deep learning models, has attracted much attention. Numerous papers have studied the convergences of adaptive stochastic gradient methods, such as AdaGrad [7, 28], RMSprop [11], Adam [16], and AMSGrad [33]. Besides the counterexample of divergence when using the Adam algorithm in the convex case in [33], various works have proposed different conditions to make Adam-type methods converge to first-order stationary points. For example, [24, 38, 47] establish global convergence of AdaGrad in the nonconvex setting; Reddi et al. [33] check the difference between learning rates of two adjacent iterations and proposes a new variant called AMSGrad; Chen et al. [3] establish the convergence of AMSGrad in the nonconvex setting; Basu et al. [2] show that Adam converges when a full-batch gradient is used; Zhou et al. [46] check the independence between gradient square and learning rate to ensure the convergence for the counterexamples in [33], and Zou et al. [48] check the parameter setting to give a sufficient condition to guarantee the convergences of both Adam and RMSProp. Also, Reddi et al. [32] introduce distributed stochastic adaptive gradient methods in the centralized model and Nazari et al. [30] introduce a decentralized adaptive gradient method. In this paper, we propose a distributed variant of Adam method by incorporating quantization and error-feedback techniques. We show that the proposed method converges to a saddle point with quantized update vectors, and will be close to a saddle point when we quantize the weights of a certain network.

3 MAIN RESULTS

Throughout this paper, we consider the following stochastic nonconvex optimization:

$$\min_{x \in \mathbb{R}^n} f(x) = \mathbb{E}_{\xi \sim \mathbb{P}} [f(x, \xi)], \quad (1)$$

where ξ is a random variable with an unknown distribution \mathbb{P} , and $f: \mathbb{R}^n \rightarrow (-\infty, +\infty)$ is a lower bounded nonconvex smooth function, i.e., $f^* = \min_{x \in \mathbb{R}^d} f(x) > -\infty$.

Due to the absence of a probability distribution \mathbb{P} of the random variable ξ , the access of exact gradient $\nabla f(x)$ may be impossible, which leads to constructing an unbiased noisy estimation $\nabla f(x_t, \xi_t)$ for full gradient $\nabla f(x_t)$ at point x_t with the given sampled sequence $\{\xi_t\}$. For convenience, we denote g_t as the abbreviation, i.e., $g_t := \nabla f(x_t, \xi_t)$. Moreover, throughout this paper, we assume that objective function f is a gradient Lipschitz function. These two requirements on gradient estimation are summarized into the following assumption:

ASSUMPTION 1. Gradient ∇f is L -Lipschitz continuous, i.e., $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$. Moreover, the noisy gradient estimation g_t is upper bounded and unbiased, i.e., $E[g_t] = \nabla f(x_t)$ and $\|g_t\| \leq G$.

Below, we introduce the gradient quantization operator $Q_g(\cdot)$ and weight quantization operator $Q_x(\cdot)$ that satisfy the following assumptions, respectively.

ASSUMPTION 2. Let $Q_g(\cdot)$ be the gradient quantization operator defined by Definition 1. We assume that there exists a constant $\delta_g \geq 0$ such that the inequality holds $\|g - Q_g(g)\| \leq (1 - \delta_g)\|g\|$.

ASSUMPTION 3. Let $Q_x(\cdot)$ be the weight quantization operator defined by Definition 1. First, we assume that the noisy gradient estimation at point x_t is an unbiased estimation of $\nabla f(Q_x(x_t))$, i.e., $E[g_t] = \nabla f(Q_x(x_t))$. In addition, we assume that there exists $\delta_x \geq 0$ such that $\|x - Q_x(x)\| \leq \delta_x$.

Assumption 1 is commonly used in analyzing adaptive stochastic type methods [3, 7, 33]. Especially, for the gradient quantization and weight quantization, the Lipschitz continuity conditions are used for bounding the error term introduced by quantization. For the weight quantization, the unbiased estimation condition $E[g_t] = \nabla f(Q_x(x_t))$ is used, which has also been used in Hou et al. [12]. All the detailed proof procedures are placed in Section 4.

3.1 Single-Machine Analysis

In this subsection, we first present the quantized Generic Adam with weight quantization, gradient quantization, and the error-feedback technique working on a single machine. Then, to show the influence on convergence related to gradient quantization or weight quantization, we establish its convergence rate with either gradient quantization or weight quantization.

Algorithm 1 unifies weight quantization, gradient quantization, and the error-feedback technique into Adam, in which $Q_x(\cdot)$ denotes the weight quantization operator, $Q_g(\cdot)$ denotes the gradient quantization operator, and e_t denotes the error-feedback term. In addition, to establish the convergence rate of Algorithm 1 in the nonconvex setting, we make the following assumptions on momentum parameter β_t , exponential moving average parameter θ_t , and base learning rate α_t .

ASSUMPTION 4. Assume that momentum parameter β_t , exponential moving average parameter θ_t , and base learning rate α_t satisfy $\beta_t \in [0, \beta]$ with $0 < \beta < 1$, $\theta_t = 1 - \theta/t$, and $\alpha_t = \alpha/\sqrt{t}$, respectively. Furthermore, we denote $\gamma = \beta/\theta'$ and $C_1 = \prod_{j=1}^N \frac{\theta_j}{\theta'}$ with $N = \max\{j | \theta_j < \theta'\}$ and θ' satisfies $\beta^2 < \theta' < 1$.

The above assumption on the hyperparameters is used to establish the convergence of adaptive stochastic type gradient method like Zou et al. [48]. In this paper, we use a simplified setting for momentum parameter β_t , exponential moving average parameter θ_t , and the base learning rate α_t to simplify the convergence analysis, compared with the sufficient condition in Zou et al. [48].

3.1.1 Gradient Quantization. Let $Q_x(x) = x$. The quantized generic Adam reduces to be generic Adam with gradient quantization and error-feedback. Below, we present the convergence rate of Algorithm 1 in the single-machine mode.

THEOREM 3.1. Let $\{x_t\}$ be the point generated by Algorithm 1 with $Q_x(x) = x$. In addition, let x_τ^T represent random variable x_τ with τ taking from $\{1, 2, \dots, T\}$ with the same probability. If Assumptions

Algorithm 1 Quantized Generic Adam

Parameters: Choose parameters $\{\alpha_t\}$, $\{\beta_t\}$, $\{\theta_t\}$, and $x_1 \in \mathbb{R}^d$, quantization functions $Q_x(\cdot)$, and $Q_g(\cdot)$. Set initial values $m_0 = 0$, $v_0 = 0$, and $e_1 = 0$.
1: **for** $t = 1, 2, \dots, T$ **do**
2: Sample a stochastic gradient of $f(Q_x(x_t))$ as g_t ;
3: $v_t = \theta_t v_{t-1} + (1 - \theta_t) g_t$;
4: $m_t = \beta_t m_{t-1} + (1 - \beta_t) g_t$;
5: $x_{t+1} = x_t - Q_g\left(\alpha_t \frac{m_t}{\sqrt{v_t + \epsilon}} + e_t\right)$;
6: $e_{t+1} = \alpha_t \frac{m_t}{\sqrt{v_t + \epsilon}} + e_t - Q_g\left(\alpha_t \frac{m_t}{\sqrt{v_t + \epsilon}} + e_t\right)$;
7: **end for**

1, 2, 4 further hold, the convergence result of Algorithm 1 holds as follows:

$$E \left[\left\| \nabla f \left(\left(x_\tau^T \right) \right) \right\|^2 \right] \leq \frac{C + C' \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}}, \quad (2)$$

where $C = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C' = \frac{2\sqrt{G^2+\epsilon d}C_3}{(1-\beta)\alpha}$, and $C_3 = \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right)$.

This theoretical result shows that with gradient quantization and error-feedback the proposed algorithm can converge to the first-order stationary point in the nonconvex setting. In addition, the convergence rate is of the same order as the original Adam in the nonconvex setting [48]. Besides, paying attention to the constant C_3 , it can be seen that the constant factor appearing in the convergence rate is related to the quantized level.

COROLLARY 3.1.1. *For given precision ξ , when we choose hyperparameters $\alpha_t = 1/\sqrt{T}$ and $\theta_t = 1 - 1/T$, to achieve $\mathbb{E}[\|\nabla f(x_\tau^T)\|^2] \leq \xi$, we have $T = O\left(\frac{1}{\xi^2}\right)$.*

REMARK. *This result shows Algorithm 1 with gradient quantization and error feedback technique can convergence in the same order as some popular method such as stochastic gradient descent and vanilla Adam.*

3.1.2 Weight Quantization. In this subsection, we set $Q_g(g) = g$ in Algorithm 1. The proposed quantized generic Adam reduces to generic Adam with the weight quantization. In this situation, to establish the convergence rate of Algorithm 1 are given below.

THEOREM 3.2. *Let $\{x_t\}$ be the point generated by Algorithm 1 with $Q_g(g) = g$. In addition, let x_τ^T represent random variable x_τ with τ taking from $\{1, 2, \dots, T\}$ with the same probability. If Assumptions 1, 3, 4 further hold, the convergence result of Algorithm 1 holds as follows:*

$$E \left[\left\| \nabla f \left(Q_x \left(x_\tau^T \right) \right) \right\|^2 \right] \leq \frac{C_5 + C_6 \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}} + C_7,$$

where $C_5 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C_6 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{LG^2\alpha^2}{\epsilon} + C_2\theta \right)$, $C_7 = \frac{8\delta_x\sqrt{G^2+\epsilon d}LG}{(1-\beta)\sqrt{\epsilon C_1}(1-\sqrt{\gamma})}$, and C_2 is defined in Theorem 3.1.

This theoretical result shows that with weight quantization the algorithm will converge to the point related to the quantized level, and when we don't use quantization the proposed algorithm will converge to the first-order stationary point by setting $\delta_x = 0$ directly. In addition, weight quantization on stochastic type gradient methods has already been considered in Khaled et al.[15], in which the authors also showed that weight quantized SGD converges to a point near the global optimum. However, the analysis of weight quantized SGD in Khaled et al. [15] is merely restricted to the strongly convex setting.

COROLLARY 3.2.1. *For given precision ξ , when we choose hyperparameters $\alpha_t = 1/\sqrt{T}$ and $\theta_t = 1 - 1/T$, to achieve $\mathbb{E}[\|\nabla f(Q_x(x_\tau^T))\|^2] \leq C'_7 + \xi$, we have $T = O\left(\frac{1}{\xi^2}\right)$, where $C'_7 = \frac{4\delta_x\sqrt{G^2+\epsilon d}LG}{(1-\beta)\sqrt{\epsilon C_1}(1-\sqrt{\gamma})}$.*

REMARK. *This result shows Algorithm 1 with weight quantization can convergence to the point near the stationary point due to quantization, but the speed to near stationary is in the same order as some popular method such as stochastic gradient descent and vallina Adam.*

3.2 Multi-Worker Analysis

In this subsection, we extend Algorithm 1 to the multi-worker setting via the parameter server model. Below, we use Algorithm 2 to represent the iteration schemes of the distributed quantized generic Adam algorithm in the parameter server, and Algorithm 3 to represent the iteration schemes in all workers, respectively. Here, we assume that all the workers work independently.

Note that communicated information \hat{x}_t and δ_t^i between the server and works is all quantized in order to improve the communication efficiency. The weight quantization procedure is performed on the server, while the gradient quantization and error-feedback procedures are performed on the workers. Below, we establish the convergence rates of distributed Adam with weight quantization, gradient quantization, and error-feedback in Algorithms 2-3 in the parameter server model.

THEOREM 3.3. *Let $\{x_t\}$ be the point generated by Algorithms 2-3. In addition, let \hat{x}_τ^T be the random variable \hat{x}_τ with τ taking from $\{1, 2, \dots, T\}$ with the same probability. If Assumptions 1-4 hold and the iterates $\|x_t\| \leq D$ are upper bounded, the convergence result of Algorithms 2-3 holds as follows:*

$$E \left[\|\nabla f(\hat{x}_\tau^T)\|^2 \right] \leq \frac{C_8 + C_9 \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}} + C_{10},$$

where $C_8 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C_9 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right)$, $C_{10} = \frac{4\sqrt{G^2+\epsilon d}\delta_x LG}{\sqrt{C_1}(1-\sqrt{\gamma})\sqrt{\epsilon}(1-\beta)}$, and C_2 is defined in Theorem 3.1.

In Algorithms 2-3, both the gradient quantization and weight quantization schemes are applied. We also show that the proposed algorithms converge to a point near the saddle point of problem (1) up to a constant. It is noted that the constant is affected by both the gradient quantized level δ_g and the weight quantized level δ_x . In addition, the limit point of the generate iterates will be influenced merely by the weight quantized level. Once gradient quantization and weight quantization reduce to identity mappings, Algorithms 2-3 reduce to the distributed Adam in the parameter server model and Theorem 3 provides their convergence rates.

COROLLARY 3.3.1. *For given precision ξ , when we choose hyperparameters $\alpha_t = 1/\sqrt{T}$ and $\theta_t = 1 - 1/T$, to achieve $E[\|\nabla f(\hat{x}_\tau^T)\|^2] \leq C'_{10} + \xi$, we have $T = O\left(\frac{1}{\xi^2}\right)$, where $C'_{10} = \frac{4\sqrt{G^2+\epsilon d}\delta_x LG}{\sqrt{C_1}(1-\sqrt{\gamma})\sqrt{\epsilon}(1-\beta)}$.*

To close this section, we give several comments on the proposed Algorithms 1-3. Different from distributed Adam where each worker transmits gradient to the parameter server and the parameter server calculates learning rate and update vector, we calculate the learning rates and update vector in local. Therefore, the error feedback technique can be applied to the adaptive algorithm. However,

Algorithm 2 The Parameter Server

Parameters: Choose $x_1 \in R^d$, and quantization function $Q_x(\cdot)$.

- 1: **for** $t = 1, 2, \dots, T$ **do**
- 2: Broadcasting $Q_x(x_t)$;
- 3: Gathering all updates from workers $\hat{\delta}_t = \frac{1}{N} \sum_{i=1}^N \delta_t^{(i)}$;
- 4: $x_{t+1} = x_t + \hat{\delta}_t$;
- 5: **end for**
- 6: **Output** $Q_x(x_t)$.

Algorithm 3 The i -th Worker

Parameters: Choose parameters $\{\alpha_t\}, \{\beta\}, \{\theta_t\}$, and quantization function $Q_g(\cdot)$. Set initial values $m_0^{(i)} = 0, v_0^{(i)} = 0, e_1^{(i)} = 0$, and $\hat{x}_0 = 0$.

- 1: **for** $t = 1, 2, \dots, T$ **do**
- 2: Receiving \hat{x}_t from the server;
- 3: Sample a stochastic gradient of $f(\hat{x}_t)$ as $g_t^{(i)}$;
- 4: $v_t^{(i)} = \theta_t v_{t-1}^{(i)} + (1 - \theta_t) [g_t^{(i)}]^2$;
- 5: $m_t^{(i)} = \beta m_{t-1}^{(i)} + (1 - \beta) g_t^{(i)}$;
- 6: Sending $\delta_t^{(i)} = Q_g \left(\alpha_t \frac{m_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon}} + e_t^{(i)} \right)$;
- 7: $e_{t+1}^{(i)} = \alpha_t \frac{m_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon}} + e_t^{(i)} - \delta_t^{(i)}$;
- 8: **end for**

the proof will be complicated due to N different learning rates being involved in the algorithm, and the following section will give a detailed proof of the above theorems.

4 PROOF DETAILS

In this section, we provide the detailed proof procedures of Theorem 3.1, Theorem 3.2, Theorem 3.3 and the related corollaries of the main theorems.

4.1 Proof of Theorem 3.1

Before providing the detailed proof of Theorem 3.1, we first denote several useful notations. Then, we provide several lemmas that are used to split the main proof of Theorem 3.1 for better readability.

NOTATION 1. Denote $\sigma_t^2 = \mathbb{E}_t[g_t^2]$ where $\mathbb{E}_t[\cdot]$ is the conditional expectation on the random variables $\{x_t, v_{t-1}, m_{t-1}\}$. Denote $\hat{v}_t = \theta_t v_{t-1} + (1 - \theta_t) \sigma_t^2$, $\hat{\eta}_t = \alpha_t / \sqrt{\hat{v}_t}$, $\Delta_t = -\alpha_t m_t / \sqrt{v_t + \epsilon}$, $M_t = \mathbb{E}[\langle \nabla f(x_t), \Delta_t \rangle + L(2 - \delta_g) \|\Delta_t\|^2 + L(2 - \delta_g) \|e_t\| \|\Delta_t\|]$ and $\|x\|_{\hat{\eta}_t}^2 = \sum_{i=1}^d \hat{\eta}_t^{(i)} x^{(i)2}$. In addition, let $\tilde{x}_t = x_t - e_t$. Then, it holds that $\tilde{x}_{t+1} = x_t - Q_g(-\Delta_t + e_t) - e_{t+1} = x_t + \Delta_t - e_t = \tilde{x}_t + \Delta_t$.

LEMMA 4.1. Given two positive sequences $\{a_i\}_{i=1}^n$ and $\{b_i\}_{i=1}^n$, it holds that

$$\left(\sum_{i=1}^n a_i b_i \right)^2 \leq \left(\sum_{i=1}^n a_i^2 \right) \left(\sum_{i=1}^n b_i^2 \right).$$

PROOF.

$$\begin{aligned} \left(\sum_{i=1}^n a_i b_i \right)^2 &= \sum_{i=1}^n \sum_{j=1}^n a_i b_i a_j b_j \leq \sum_{i=1}^n \sum_{j=1}^n \frac{1}{2} (a_i^2 b_j^2 + a_j^2 b_i^2) = \frac{1}{2} \sum_{i=1}^n a_i^2 \left(\sum_{j=1}^n b_j^2 \right) + \frac{1}{2} \sum_{i=1}^n b_i^2 \left(\sum_{j=1}^n a_j^2 \right) \\ &= \left(\sum_{i=1}^n a_i^2 \right) \left(\sum_{i=1}^n b_i^2 \right), \end{aligned}$$

where the second inequality is the arithmetic inequality with positive numbers $a_i b_j$ and $a_j b_i$. \square

LEMMA 4.2. By using Notation 1 and the iteration scheme of Theorem 3.1, for all $t \geq 1$ the following inequality holds:

$$m_t^2 \leq \frac{1}{C_1 (1 - \gamma) (1 - \theta_t)} v_t.$$

PROOF. By using the definition of m_t in Theorem 3.1, it directly holds that $m_t = \sum_{i=1}^t \beta^{t-i} (1 - \beta) g_i$. Let $\Theta(t, i) = \prod_{j=i+1}^t \theta_j$ for $i < t$, and $\Theta(i, i) = 1$. According to the definition of v_t , it holds that $v_t = \sum_{i=1}^t \left(\prod_{j=i+1}^t \theta_j \right) (1 - \theta_i) g_i^2 = \sum_{i=1}^t \Theta(t, i) (1 - \theta_i) g_i^2$.

With Lemma 4.1, it holds that

$$\begin{aligned} m_t^2 &= \left(\sum_{i=1}^t \frac{\beta^{t-i} (1 - \beta)}{\sqrt{\Theta(t, i) (1 - \theta_i)}} \sqrt{\Theta(t, i) (1 - \theta_i)} g_i \right)^2 \leq \sum_{i=1}^t \frac{\beta^{2(t-i)} (1 - \beta)^2}{\Theta(t, i) (1 - \theta_i)} \sum_{i=1}^t \Theta(t, i) (1 - \theta_i) g_i^2 \\ &= \sum_{i=1}^t \frac{\beta^{2(t-i)} (1 - \beta)^2}{\Theta(t, i) (1 - \theta_i)} v_t \leq \sum_{i=1}^t \frac{\beta^{2(t-i)}}{\left(\prod_{j=i+1}^N \theta_j / \theta' \right) (\theta')^{t-i} (1 - \theta_i)} v_t \\ &\leq \sum_{i=1}^t \frac{\beta^{2(t-i)}}{\left(\prod_{j=1}^N \theta_j / \theta' \right) (\theta')^{t-i} (1 - \theta_i)} v_t \leq \frac{1}{C_1 (1 - \theta_t)} \sum_{i=1}^t \left(\frac{\beta^2}{\theta'} \right)^{t-i} v_t \leq \frac{1}{C_1 (1 - \gamma) (1 - \theta_t)} v_t. \end{aligned}$$

Then, we obtain the targeted result. \square

LEMMA 4.3. Let τ be randomly chosen from $\{1, 2, \dots, T\}$ with equal probabilities $p_\tau = \frac{1}{T}$. We have the following estimate:

$$\mathbb{E}[\|\nabla f(x_\tau)\|^2] \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right].$$

PROOF. Note that $\|\hat{v}_t\|_1 = \theta_t \|v_{t-1}\|_1 + (1 - \theta_t) \|\sigma_t\|^2$ and $\|g_t\| \leq G$. It is straightforward to prove $\|v_t\|_1 \leq G^2$. Hence, we have $\|\hat{v}_t + \epsilon\|_1 \leq G^2 + \epsilon d$.

Utilizing this inequality, we have

$$\begin{aligned} \|\nabla f(x_t)\|^2 &= \frac{\|\nabla f(x_t)\|^2}{\sqrt{\|\hat{v}_t + \epsilon\|_1}} \sqrt{\|\hat{v}_t + \epsilon\|_1} = \sqrt{\|\hat{v}_t + \epsilon\|_1} \sum_{k=1}^d \frac{|\nabla_k f(x_t)|^2}{\sqrt{\sum_{l=1}^d \hat{v}_{t,l} + \epsilon}} \\ &\leq \sqrt{\|\hat{v}_t + \epsilon\|_1} \alpha_t^{-1} \sum_{k=1}^d \frac{\alpha_t}{\sqrt{\hat{v}_{t,k} + \epsilon}} |\nabla_k f(x_t)|^2 = \sqrt{\|\hat{v}_t + \epsilon\|_1} \alpha_t^{-1} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \\ &\leq \sqrt{G^2 + \epsilon d} \alpha_t^{-1} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha_T} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2. \end{aligned}$$

Then, by using the definition of x_τ , we obtain

$$\mathbb{E}[\|\nabla f(x_\tau)\|^2] = \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla f(x_t)\|^2] \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right].$$

Thus, the desired result is obtained. \square

LEMMA 4.4. By using Notation 1, the following inequality holds:

$$\sum_{t=1}^T \|\Delta_t\|^2 \leq \frac{G^2}{\epsilon} \sum_{t=1}^T \frac{\alpha^2}{t}.$$

PROOF. By using the definition of m_t , it holds $\|m_t\|^2 \leq G^2$.

Then, $\|\Delta_t\|^2 = \left\| \frac{\alpha_t m_t}{\sqrt{\hat{v}_t + \epsilon}} \right\|^2 \leq \frac{G^2}{\epsilon} \alpha_t^2$ by using the definition of Δ_t .

Therefore, $\sum_{t=1}^T \|\Delta_t\|^2 \leq \frac{G^2}{\epsilon} \sum_{t=1}^T \frac{\alpha^2}{t}$. \square

LEMMA 4.5. By the iteration scheme of Algorithm 1, it holds that

$$\sum_{t=1}^T \|e_t\| \|\Delta_t\| \leq \sum_{t=1}^T \frac{1 - \delta_g}{\delta_g} \|\Delta_t\|^2.$$

PROOF. By the definition of noisy term e_t and Δ_t , it holds

$$\begin{aligned} \|e_t\| &= \|\Delta_t + e_{t-1} - Q_g(\Delta_t + e_{t-1})\| \leq (1 - \delta_g) \|\Delta_t + e_{t-1}\| \leq (1 - \delta_g) \|\Delta_t\| + (1 - \delta_g) \|e_{t-1}\| \\ &\leq \sum_{i=1}^t (1 - \delta_g)^{t-i+1} \|\Delta_i\|. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \sum_{t=1}^T \|e_t\| \|\Delta_t\| &\leq \sum_{t=1}^T \sum_{i=1}^t (1 - \delta_g)^{t-i+1} \|\Delta_i\| \|\Delta_t\| \leq \frac{1}{2} \sum_{t=1}^T \sum_{i=1}^t (1 - \delta_g)^{t-i+1} (\|\Delta_i\|^2 + \|\Delta_t\|^2) \\ &\leq \frac{1 - \delta_g}{2\delta_g} \sum_{t=1}^T \|\Delta_t\|^2 + \frac{1 - \delta_g}{2\delta_g} \sum_{i=1}^T \|\Delta_i\|^2 = \frac{1 - \delta_g}{\delta_g} \sum_{t=1}^T \|\Delta_t\|^2. \end{aligned}$$

Hence, we obtain the desired result. \square

LEMMA 4.6. By the definition of M_k , it holds that

$$\sum_{t=1}^T M_t \leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} - \frac{1-\beta}{2} \mathbb{E} \left[\left\| \sum_{t=1}^T \nabla f(x_t) \right\|_{\hat{\eta}_t}^2 \right],$$

where

$$\begin{aligned} C_2 = & \frac{5\alpha G^3(1-\beta)}{2\epsilon\sqrt{\theta}} \left(\frac{\beta}{(1-\beta)\sqrt{\theta_1 C_1}(1-\gamma)} + 1 \right)^2 + \frac{5\alpha G^3}{2\epsilon\sqrt{\theta}} + \frac{5\beta^2 \alpha d \sqrt{\epsilon}}{2\sqrt{\theta}(1-\beta)\theta_1 C_1(1-\gamma)} \\ & + \frac{5\alpha\sqrt{G^2+\epsilon}G^2\beta^2}{2(1-\beta)\sqrt{\theta}\theta_1 C_1(1-\gamma)\epsilon} + \frac{5\alpha\sqrt{G^2+\epsilon}\beta^2 d}{2(1-\beta)\sqrt{\theta}\theta_1 C_1(1-\gamma)}. \end{aligned}$$

PROOF. To split M_t , first we introduce the following two equalities. Using the definitions of v_t and \hat{v}_t , we obtain

$$\begin{aligned} \frac{(1-\beta)\alpha_t g_t}{\sqrt{v_t+\epsilon}} &= \frac{(1-\beta)\alpha_t g_t}{\sqrt{\hat{v}_t+\epsilon}} + (1-\beta)\alpha_t g_t \left(\frac{1}{\sqrt{v_t+\epsilon}} - \frac{1}{\sqrt{\hat{v}_t+\epsilon}} \right) \\ &= (1-\beta)\hat{\eta}_t g_t + (1-\beta)\alpha_t g_t \frac{(1-\theta_t)(\sigma_t^2 - g_t^2)}{\sqrt{v_t+\epsilon}\sqrt{\hat{v}_t+\epsilon}(\sqrt{v_t+\epsilon} + \sqrt{\hat{v}_t+\epsilon})} \\ &= (1-\beta)\hat{\eta}_t g_t + \hat{\eta}_t \sigma_t \frac{(1-\theta_t)g_t}{\sqrt{v_t+\epsilon}} \frac{(1-\beta)\sigma_t}{\sqrt{v_t+\epsilon} + \sqrt{\hat{v}_t+\epsilon}} - \hat{\eta}_t g_t \frac{(1-\theta_t)g_t}{\sqrt{v_t+\epsilon}} \frac{(1-\beta)g_t}{\sqrt{v_t+\epsilon} + \sqrt{\hat{v}_t+\epsilon}}. \end{aligned}$$

In addition, it is not hard to check that the following equality holds:

$$\begin{aligned} \beta\alpha_t m_{t-1} \left(\frac{1}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} - \frac{1}{\sqrt{v_t + \epsilon}} \right) &= \beta\alpha_t m_{t-1} \frac{(1-\theta_t)(g_t^2 + \epsilon)}{\sqrt{v_t + \epsilon}\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}(\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon})} \\ &= \hat{\eta}_t g_t \frac{(1-\theta_t)g_t}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} + \hat{\eta}_t \epsilon \frac{(1-\theta_t)}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \\ &\quad + \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1-\theta_t)g_t}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}g_t}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \\ &\quad + \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1-\theta_t)\sqrt{\epsilon}}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\epsilon}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}}. \end{aligned}$$

For convenience, we denote

$$\begin{aligned} A_t^1 &= \hat{\eta}_t g_t \frac{(1-\theta_t)g_t}{\sqrt{v_t + \epsilon}} \left(\frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} + \frac{(1-\beta)g_t}{\sqrt{v_t + \epsilon} + \sqrt{\hat{v}_t + \epsilon}} \right), \\ A_t^2 &= -\hat{\eta}_t \sigma_t \frac{(1-\theta_t)g_t}{\sqrt{v_t + \epsilon}} \frac{(1-\beta)\sigma_t}{\sqrt{v_t + \epsilon} + \sqrt{\hat{v}_t + \epsilon}}, \quad A_t^3 = \hat{\eta}_t \epsilon \frac{(1-\theta_t)}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}}, \\ A_t^4 &= \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1-\theta_t)g_t}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}g_t}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}}, \\ A_t^5 &= \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1-\theta_t)\sqrt{\epsilon}}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\epsilon}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t}\sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}}. \end{aligned}$$

Then, we obtain

$$\begin{aligned}
 \Delta_t - \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}\Delta_{t-1} &= -\frac{\alpha_t m_t}{\sqrt{v_t + \epsilon}} + \frac{\beta\alpha_t m_{t-1}}{\sqrt{\theta_t}(v_{t-1} + \epsilon)} \\
 &= -\frac{(1-\beta)\alpha_t g_t}{\sqrt{v_t + \epsilon}} + \beta\alpha_t m_{t-1} \left(\frac{1}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} - \frac{1}{\sqrt{v_t + \epsilon}} \right) \\
 &= -(1-\beta)\hat{\eta}_t g_t + A_t^1 + A_t^2 + A_t^3 + A_t^4 + A_t^5.
 \end{aligned}$$

By using the above inequalities, we derive the upper bound for the term $\mathbb{E}\langle \nabla f(x_t), \Delta_t \rangle$:

$$\begin{aligned}
 \mathbb{E}\langle \nabla f(x_t), \Delta_t \rangle &= \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}\mathbb{E}\langle \nabla f(x_t), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(x_t), \Delta_t - \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}\Delta_{t-1} \rangle \\
 &= \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(x_t) - \nabla f(x_{t-1}), \Delta_{t-1} \rangle) \\
 &\quad + \mathbb{E}\langle \nabla f(x_t), -(1-\beta)\hat{\eta}_t g_t \rangle + \mathbb{E}\langle \nabla f(x_t), A_t^1 \rangle + \mathbb{E}\langle \nabla f(x_t), A_t^2 \rangle \\
 &\quad + \mathbb{E}\langle \nabla f(x_t), A_t^3 \rangle + \mathbb{E}\langle \nabla f(x_t), A_t^4 \rangle + \mathbb{E}\langle \nabla f(x_t), A_t^5 \rangle.
 \end{aligned} \tag{3}$$

For the first term in Eq (3), we have

$$\begin{aligned}
 &\frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(x_t) - \nabla f(x_{t-1}), \Delta_{t-1} \rangle) \\
 &\leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + \mathbb{E}[L\|x_t - x_{t-1}\|\|\Delta_{t-1}\|]) \\
 &= \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + L\mathbb{E}[\|\Delta_{t-1}\|\|Q_g(\Delta_{t-1} + e_{t-1})\|]) \\
 &\leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + L\mathbb{E}[\|\Delta_{t-1}\|(\|Q_g(\Delta_{t-1} + e_{t-1}) - (\Delta_{t-1} + e_{t-1})\| + \|\Delta_{t-1} + e_{t-1}\|)]) \\
 &\leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + L\mathbb{E}[\|\Delta_{t-1}\|((1-\delta_g)\|\Delta_{t-1} + e_{t-1}\| + \|\Delta_{t-1} + e_{t-1}\|)]) \\
 &\leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + L\mathbb{E}[\|\Delta_{t-1}\|((2-\delta_g)\|\Delta_{t-1}\| + (2-\delta_g)\|e_{t-1}\|)]) \\
 &\leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}(\mathbb{E}\langle \nabla f(x_{t-1}), \Delta_{t-1} \rangle + L(2-\delta_g)\mathbb{E}[\|\Delta_{t-1}\|^2] + L(2-\delta_g)\mathbb{E}[\|e_{t-1}\|\|\Delta_{t-1}\|]) \\
 &= \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}M_{t-1}.
 \end{aligned}$$

For the second term in Eq. (3), we have

$$\begin{aligned}
 -(1-\beta)\mathbb{E}\langle \nabla f(x_t), \hat{\eta}_t g_t \rangle &= -(1-\beta)\mathbb{E}[\mathbb{E}_t\langle \nabla f(x_t), \hat{\eta}_t g_t \rangle] = -(1-\beta)\mathbb{E}\langle \nabla f(x_t), \hat{\eta}_t \mathbb{E}_t[g_t] \rangle \\
 &= -(1-\beta)\mathbb{E}\|\nabla f(x_t)\|_{\hat{\eta}_t}^2,
 \end{aligned}$$

where the second equality is because given $\{x_t, v_{t-1}, m_{t-1}\}$, $\nabla f(x_t)$ and $\hat{\eta}_t$ are independent to g_t . Besides, $\mathbb{E}_t[\nabla f(x_t)] = \nabla f(x_t)$ and $\mathbb{E}_t[\hat{\eta}_t] = \hat{\eta}_t$.

For the third term in Eq. (3), we have

$$\begin{aligned}
 \mathbb{E}\langle \nabla f(x_t), A_t^1 \rangle &\leq \mathbb{E} \left\langle \frac{\sqrt{\hat{\eta}_t} |\nabla f(x_t)| |g_t|}{\sigma_t}, \frac{\sqrt{\hat{\eta}_t} \sigma_t (1 - \theta_t) |g_t| \left| \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} + \frac{(1-\beta)g_t}{\sqrt{v_t + \epsilon}} \right|}{\sqrt{v_t + \epsilon}} \right\rangle \\
 &\leq \mathbb{E} \left\langle \frac{\sqrt{\hat{\eta}_t} |\nabla f(x_t)| |g_t|}{\sigma_t}, \sqrt{\frac{\alpha G}{\sqrt{\theta}}} \frac{(1 - \theta_t) |g_t| \left(\left| \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \right| + \left| \frac{(1-\beta)g_t}{\sqrt{v_t + \epsilon}} \right| \right)}{\sqrt{v_t + \epsilon}} \right\rangle \\
 &\leq \mathbb{E} \left\langle \frac{\sqrt{\hat{\eta}_t} |\nabla f(x_t)| |g_t|}{\sigma_t}, \sqrt{\frac{\alpha G}{\sqrt{\theta}}} \left(\frac{\beta}{\sqrt{\theta_1 C_1 (1 - \gamma)}} + (1 - \beta) \right) \frac{\sqrt{1 - \theta_t} |g_t|}{\sqrt{v_t + \epsilon}} \right\rangle \\
 &\leq \frac{1 - \beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{5\alpha G^3 (1 - \beta) (1 - \theta_t)}{2\epsilon \sqrt{\theta}} \left(\frac{\beta}{(1 - \beta) \sqrt{\theta_1 C_1 (1 - \gamma)}} + 1 \right)^2,
 \end{aligned}$$

where the equalities hold according to the following inequities, respectively,

$$\begin{aligned}
 \sqrt{\hat{\eta}_t} \sigma_t &= \sqrt{\frac{\alpha_t \sigma_t^2}{\sqrt{\hat{\delta}_t} + \epsilon}} \leq \sqrt{\frac{\alpha G}{\sqrt{\theta}}}, \\
 \left| \frac{(1 - \beta) g_t}{\sqrt{v_t + \epsilon}} \right| &\leq \left| \frac{(1 - \beta) g_t}{\sqrt{(1 - \theta_t) g_t^2}} \right| = \frac{1 - \beta}{\sqrt{1 - \theta_t}}, \\
 \left| \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \right| &\leq \left| \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1}}} \right| \leq \frac{\beta}{\sqrt{\theta_1 (1 - \theta_t) C_1 (1 - \gamma)}}.
 \end{aligned}$$

For the fourth term in Eq. (3), we have

$$\begin{aligned}
 \mathbb{E}\langle \nabla f(x_t), A_t^2 \rangle &\leq \mathbb{E} \left\langle |\nabla f(x_t)|, \left| \hat{\eta}_t \sigma_t \frac{(1 - \theta_t) g_t}{\sqrt{v_t + \epsilon}} \frac{(1 - \beta) \sigma_t}{\sqrt{v_t + \epsilon} + \sqrt{\hat{\delta}_t} + \epsilon} \right| \right\rangle \\
 &\leq \mathbb{E} \left\langle \sqrt{\hat{\eta}_t} |\nabla f(x_t)|, \sqrt{\hat{\eta}_t} \sigma_t \left| \frac{(1 - \theta_t) g_t}{\sqrt{v_t + \epsilon}} \right| \left| \frac{(1 - \beta) \sigma_t}{\sqrt{v_t + \epsilon} + \sqrt{\hat{\delta}_t} + \epsilon} \right| \right\rangle \\
 &\leq \frac{1 - \beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{2\alpha G}{5\theta} \mathbb{E} \left\| \frac{\sqrt{1 - \theta_t} g_t}{\sqrt{v_t + \epsilon}} \right\|^2 \leq \frac{1 - \beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{5\alpha G^3 (1 - \theta_t)}{2\sqrt{\theta} \epsilon},
 \end{aligned}$$

where the equality holds according to $\left| \frac{(1-\beta)\sigma_t}{\sqrt{v_t+\epsilon}+\sqrt{\hat{\delta}_t}+\epsilon} \right| \leq \left| \frac{(1-\beta)\sigma_t}{\sqrt{v_t+\epsilon}} \right| \leq \frac{1-\beta}{\sqrt{1-\theta_t}}$.

For the fifth term in Eq. (3), we have

$$\begin{aligned}
 \mathbb{E}\langle \nabla f(x_t), A_t^3 \rangle &\leq \mathbb{E} \left\langle |\nabla f(x_t)|, \left| \hat{\eta}_t \epsilon \frac{(1 - \theta_t)}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \right| \right\rangle \\
 &\leq \mathbb{E} \left\langle \sqrt{\hat{\eta}_t} |\nabla f(x_t)|, \sqrt{\hat{\eta}_t} \epsilon \frac{(1 - \theta_t) \beta |m_{t-1}|}{\sqrt{\epsilon} \sqrt{\theta_t v_{t-1}}} \right\rangle \leq \frac{1 - \beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{5(1 - \theta_t) \beta^2 \alpha d \sqrt{\epsilon}}{2\sqrt{\theta} (1 - \beta) \theta_1 C_1 (1 - \gamma)},
 \end{aligned}$$

where the inequality holds according to $\sqrt{\hat{\eta}_t} \epsilon \leq \sqrt{\frac{\alpha_t \epsilon^2}{\sqrt{\hat{\delta}_t} + \epsilon}} \leq \sqrt{\frac{\alpha \epsilon^{3/2}}{\sqrt{\theta}}}$.

For the sixth term in Eq. (3), we have

$$\begin{aligned} & \mathbb{E} \langle \nabla f(x_t), A_t^4 \rangle \\ & \leq \mathbb{E} \left\langle |\nabla f(x_t)|, \left| \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1-\theta_t) g_t}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} g_t}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \right| \right\rangle \\ & \leq \frac{1-\beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{5\alpha\sqrt{G^2 + \epsilon} G^2 \beta^2 (1-\theta_t)}{2(1-\beta) \sqrt{\theta} \theta_1 C_1 (1-\gamma) \epsilon}, \end{aligned}$$

where the equalities hold according to $\sqrt{\hat{\eta}_t} \sqrt{\sigma_t^2 + \epsilon} = \sqrt{\frac{\alpha_t(\sigma_t^2 + \epsilon)}{\sqrt{\hat{v}_t + \epsilon}}} \leq \sqrt{\frac{\alpha\sqrt{G^2 + \epsilon}}{\sqrt{\theta}}}$, and

$$\begin{aligned} & \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} g_t}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \\ & \leq \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} g_t}{\sqrt{v_t + \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon}} \leq \frac{\beta}{\sqrt{\theta_1 (1-\theta_t) C_1 (1-\gamma)}}. \end{aligned}$$

For the seventh term in Eq. (3), we have

$$\begin{aligned} & \mathbb{E} \langle \nabla f(x_t), A_t^5 \rangle \\ & \leq \mathbb{E} \left\langle |\nabla f(x_t)|, \left| \hat{\eta}_t \sqrt{\sigma_t^2 + \epsilon} \frac{(1+\theta_t) \sqrt{\epsilon}}{\sqrt{v_t + \epsilon}} \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\epsilon}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \right| \right\rangle \\ & \leq \mathbb{E} \left\langle \sqrt{\hat{\eta}_t} |\nabla f(x_t)|, \sqrt{\frac{\alpha\sqrt{G^2 + \epsilon}}{\sqrt{\theta}}} \frac{\beta \sqrt{1-\theta_t} \sqrt{\epsilon}}{\sqrt{\theta_1 C_1 (1-\gamma) \sqrt{v_t + \epsilon}}} \right\rangle \leq \frac{1-\beta}{10} \mathbb{E} \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 + \frac{5\alpha\sqrt{G^2 + \epsilon} \beta^2 (1-\theta_t) d}{2(1-\beta) \sqrt{\theta} \theta_1 C_1 (1-\gamma)}, \end{aligned}$$

where the inequality holds according to

$$\begin{aligned} & \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\epsilon}}{\sqrt{v_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon} + \sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \\ & \leq \frac{\beta m_{t-1}}{\sqrt{\theta_t v_{t-1} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\epsilon}}{\sqrt{(1-\theta_t) \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^2 + \epsilon}}{\sqrt{\hat{v}_t + \epsilon}} \leq \frac{\beta}{\sqrt{\theta_1 (1-\theta_t) C_1 (1-\gamma)}}. \end{aligned}$$

Therefore, by combining the above upper estimations for the seven terms in Eq. (3), we obtain

$$\mathbb{E} \langle \nabla f(x_t), \Delta_t \rangle \leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} M_{t-1} + C_2 (1-\theta_t) - \frac{1-\beta}{2} \mathbb{E} [\|\nabla f(x_t)\|_{\hat{\eta}_t}^2],$$

where $C_2 = \frac{5\alpha G^3 (1-\beta)}{2\epsilon \sqrt{\theta}} \left(\frac{\beta}{(1-\beta) \sqrt{\theta_1 C_1 (1-\gamma)}} + 1 \right)^2 + \frac{5\alpha G^3}{2\epsilon \sqrt{\theta}} + \frac{5\beta^2 \alpha d \sqrt{\epsilon}}{2\sqrt{\theta} (1-\beta) \theta_1 C_1 (1-\gamma)} + \frac{5\alpha \sqrt{G^2 + \epsilon} G^2 \beta^2}{2(1-\beta) \sqrt{\theta} \theta_1 C_1 (1-\gamma) \epsilon} + \frac{5\alpha \sqrt{G^2 + \epsilon} \beta^2 d}{2(1-\beta) \sqrt{\theta} \theta_1 C_1 (1-\gamma)}$.

On the other hand, let $N_t = L(2-\delta_g) \mathbb{E}[\|\Delta_t\|^2] + L(2-\delta_g) \mathbb{E}[\|e_t\| \|\Delta_t\|] + C_2(1-\theta_t)$.

Recalling the definition of M_t . For M_1 , we have

$$\begin{aligned} M_1 &= \mathbb{E} \left[- \left\langle \nabla f(x_1), \frac{\alpha_1 m_1}{\sqrt{\theta_1 + \epsilon}} \right\rangle + L(2-\delta_g) \|\Delta_1\|^2 + L(2-\delta_g) \|e_1\| \|\Delta_1\| \right] \\ &= \mathbb{E} \left[- \left\langle \nabla f(x_1), (1-\beta_1) \hat{\eta}_1 g_1 + \hat{\eta}_1 \sigma_1 \frac{(1-\theta_1) g_1}{\sqrt{v_1 + \epsilon}} \frac{(1-\beta) \sigma_1}{\sqrt{v_1 + \epsilon} + \sqrt{\hat{v}_1 + \epsilon}} - \hat{\eta}_1 g_1 \frac{(1-\theta_1) g_1}{\sqrt{v_1 + \epsilon}} \frac{(1-\beta) g_1}{\sqrt{v_1 + \epsilon} + \sqrt{\hat{v}_1 + \epsilon}} \right\rangle \right] \\ &\quad + \mathbb{E} [L(2-\delta_g) \|\Delta_1\|^2] + L(2-\delta_g) \mathbb{E} [\|e_1\| \|\Delta_1\|] \leq N_1, \end{aligned}$$

where the last inequality holds by using A_t^1 and A_t^2 .

Let $M_0 = 0$. Then we have $M_t \leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}M_{t-1} + N_t - \frac{1-\beta}{2}\mathbb{E}[\|\nabla f(x_t)\|_{\hat{\eta}_t}^2] \leq \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_{t-1}}M_{t-1} + N_t$. By induction we have

$$\begin{aligned} M_t &\leq \sum_{i=1}^t \frac{\alpha_t \beta^{t-i}}{\alpha_i \sqrt{\Theta(t, i)}} N_i - \frac{1-\beta}{2} \mathbb{E}[\|\nabla f(x_t)\|_{\hat{\eta}_t}^2] \leq \frac{1}{\sqrt{C_1}} \sum_{i=1}^t \left(\frac{\beta}{\sqrt{\theta^t}}\right)^{t-i} N_i - \frac{1-\beta}{2} \mathbb{E}[\|\nabla f(x_t)\|_{\hat{\eta}_t}^2] \\ &\leq \frac{1}{\sqrt{C_1}} \sum_{i=1}^t \sqrt{\gamma}^{t-i} N_i - \frac{1-\beta}{2} \mathbb{E}[\|\nabla f(x_t)\|_{\hat{\eta}_t}^2]. \end{aligned}$$

Therefore,

$$\begin{aligned} \sum_{t=1}^T M_t &\leq \frac{1}{\sqrt{C_1}} \sum_{t=1}^T \sum_{i=1}^t \sqrt{\gamma}^{t-i} N_i - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \\ &\leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \sum_{t=1}^T N_t - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \\ &= \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \sum_{t=1}^T [L(2-\delta_g) \mathbb{E}\|\Delta_t\|^2 + L(2-\delta_g) \mathbb{E}\|e_t\| \|\Delta_t\| + \frac{C_2\theta}{t}] - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \\ &\leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right]. \end{aligned}$$

Hence, the targeted result holds. \square

PROOF OF THEOREM 3.1. Based on Notation (1) and the above lemmas, then we can prove Theorem 1. First, according to the gradient Lipschitz condition of f , it holds

$$\begin{aligned} f(\tilde{x}_{t+1}) &\leq f(\tilde{x}_t) + \langle \nabla f(\tilde{x}_t), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 \\ &= f(\tilde{x}_t) + \langle \nabla f(x_t), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 + \langle \nabla f(\tilde{x}_t) - \nabla f(x_t), \Delta_t \rangle \\ &\leq f(\tilde{x}_t) + \langle \nabla f(x_t), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 + \|\nabla f(\tilde{x}_t) - \nabla f(x_t)\| \|\Delta_t\| \\ &\leq f(\tilde{x}_t) + \langle \nabla f(x_t), \Delta_t \rangle + L(2-\delta_g) \|\Delta_t\|^2 + L(2-\delta_g) \|e_t\| \|\Delta_t\|. \end{aligned}$$

Recall that $M_t = \mathbb{E}[\langle \nabla f(x_t), \Delta_t \rangle + L(2-\delta_g) \|\Delta_t\|^2 + L(2-\delta_g) \|e_t\| \|\Delta_t\|]$. Then we have

$$\begin{aligned} f^* &\leq \mathbb{E}[f(\tilde{x}_{T+1})] \leq f(x_1) + \sum_{t=1}^T M_t \\ &\leq f(x_1) + \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right]. \end{aligned}$$

Using the above lemmas and arranging the corresponding terms, we have

$$\begin{aligned} \left(\mathbb{E}[\|\nabla f(x_t^T)\|^2] \right) &\leq \frac{\sqrt{G^2+\epsilon d}}{\alpha\sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \leq \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha\sqrt{T}} \left(f(x_1) - f^* + C_3 \sum_{t=1}^T \frac{1}{t} \right) \\ &\leq \frac{C + C' \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}}, \end{aligned}$$

where $C_3 = \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right)$, $C = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C' = \frac{2\sqrt{G^2+\epsilon d}C_3}{(1-\beta)\alpha}$, respectively. Hence, the proof is completed. \square

PROOF OF COROLLARY 3.1.1. For a fixed iteration T , let $\alpha_t = \frac{\alpha}{\sqrt{t}}$ and $\theta = 1 - \frac{\theta}{T}$. When $\alpha_t = \frac{\alpha}{\sqrt{t}}$, Lemma 4.4 will update to $\sum_{t=1}^T \|\Delta_t\|^2 \leq \frac{G^2 \alpha^2}{\epsilon}$. By the same proof in Lemma 4.6, we have

$$\begin{aligned} \sum_{t=1}^T M_t &\leq \frac{1}{\sqrt{C_1} (1 - \sqrt{\gamma})} \sum_{t=1}^T [L (2 - \delta_g) \mathbb{E} \|\Delta_t\|^2 + L (2 - \delta_g) \mathbb{E} \|e_t\| \|\Delta_t\| + \frac{C_2 \theta}{T}] - \frac{1 - \beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \\ &\leq \frac{L (2 - \delta_g) G^2 \alpha^2}{\sqrt{C_1} (1 - \sqrt{\gamma}) \epsilon \delta_g} + \frac{C_2 \theta}{\sqrt{C_1} (1 - \sqrt{\gamma})} - \frac{1 - \beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right]. \end{aligned}$$

Based on the proof of Theorem 3.1, we have

$$\begin{aligned} (\mathbb{E} [\|\nabla f(x_t^T)\|^2]) &\leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(x_t)\|_{\hat{\eta}_t}^2 \right] \\ &\leq \frac{2\sqrt{G^2 + \epsilon d}}{(1 - \beta) \alpha \sqrt{T}} \left(f(x_1) - f^* + \frac{L (2 - \delta_g) G^2 \alpha^2}{\sqrt{C_1} (1 - \sqrt{\gamma}) \epsilon \delta_g} + \frac{C_2 \theta}{\sqrt{C_1} (1 - \sqrt{\gamma})} \right). \end{aligned}$$

By bounding $\frac{2\sqrt{G^2 + \epsilon d}}{(1 - \beta) \alpha \sqrt{T}} \left(f(x_1) - f^* + \frac{L (2 - \delta_g) G^2 \alpha^2}{\sqrt{C_1} (1 - \sqrt{\gamma}) \epsilon \delta_g} + \frac{C_2 \theta}{\sqrt{C_1} (1 - \sqrt{\gamma})} \right)$ by ξ , we obtain $T = O\left(\frac{1}{\xi^2}\right)$. \square

4.2 Proof of Theorem 3.2

To prove Theorem 3.2, we first define some notations and provide several useful lemmas.

NOTATION 2. Let $\mathbb{E}_t[\cdot]$ be the conditional expectation conditioned on $\{x_t, v_{t-1}, m_{t-1}\}$. Denote $\sigma_t^2 = \mathbb{E}_t[g_t^2]$, $\hat{v}_t = \theta_t v_{t-1} + (1 - \theta_t) \sigma_t^2$, $\hat{\eta}_t = \alpha_t / \sqrt{\hat{v}_t + \epsilon}$, $\Delta_t = -\alpha_t m_t / \sqrt{v_t + \epsilon}$, and $M_t = \mathbb{E}[\langle \nabla f(Q_x(x_t)), \Delta_t \rangle + L \|\Delta_t\|^2 + 2L \|\Delta_t\|]$.

LEMMA 4.7. By using Notation 2, the following inequality holds:

$$\sum_{t=1}^T \|\Delta_t\| \leq \frac{2G\alpha\sqrt{T}}{\sqrt{\epsilon}}.$$

PROOF.

$$\sum_{t=1}^T \|\Delta_t\| = \sum_{t=1}^T \alpha_t \left\| \frac{m_t}{\sqrt{v_t + \epsilon}} \right\| \leq \sum_{t=1}^T \frac{G\alpha}{\sqrt{t\epsilon}} \leq \frac{2G\alpha\sqrt{T}}{\sqrt{\epsilon}}.$$

\square

LEMMA 4.8. Based on Notation 2, we have the following upper estimation for M_t ,

$$\sum_{t=1}^T M_t \leq \frac{1}{\sqrt{C_1} (1 - \sqrt{\gamma})} \left(\frac{LG^2 \alpha^2}{\epsilon} + C_2 \theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4L\delta_x G\alpha\sqrt{T}}{\sqrt{\epsilon C_1} (1 - \sqrt{\gamma})} - \frac{1 - \beta}{2} \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(Q_x(x_t))\|_{\hat{\eta}_t}^2 \right],$$

where

$$\begin{aligned} C_2 &= \frac{5\alpha G^3 (1 - \beta)}{2\epsilon \sqrt{\theta}} \left(\frac{\beta}{(1 - \beta) \sqrt{\theta_1 C_1} (1 - \gamma)} + 1 \right)^2 + \frac{5\alpha G^3}{2\epsilon \sqrt{\theta}} + \frac{5\beta^2 \alpha d \sqrt{\epsilon}}{2\sqrt{\theta} (1 - \beta) \theta_1 C_1 (1 - \gamma)} \\ &\quad + \frac{5\alpha \sqrt{G^2 + \epsilon G^2 \beta^2}}{2(1 - \beta) \sqrt{\theta} \theta_1 C_1 (1 - \gamma) \epsilon} + \frac{5\alpha \sqrt{G^2 + \epsilon \beta^2 d}}{2(1 - \beta) \sqrt{\theta} \theta_1 C_1 (1 - \gamma)}. \end{aligned}$$

PROOF. Denoting the same notations $A_t^1, A_t^2, A_t^3, A_t^4, A_t^5$ in Lemma 4.6, we have

$$\begin{aligned}\mathbb{E}\langle \nabla f(Q_x(x_t)), \Delta_t \rangle &= \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} \mathbb{E}\langle \nabla f(Q_x(x_t)), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)), \Delta_t - \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_t} \Delta_{t-1} \rangle \\ &= \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)) - \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle) \\ &\quad + \mathbb{E}\langle \nabla f(Q_x(x_t)), - (1 - \beta) \hat{\eta}_t g_t \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)), A_t^1 \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)), A_t^2 \rangle \\ &\quad + \mathbb{E}\langle \nabla f(Q_x(x_t)), A_t^3 \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)), A_t^4 \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)), A_t^5 \rangle.\end{aligned}\quad (4)$$

For the first term in Eq. (4), we have

$$\begin{aligned}&\frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + \mathbb{E}\langle \nabla f(Q_x(x_t)) - \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle) \\ &\leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + \mathbb{E}[L \|Q_x(x_t) - Q_x(x_{t-1})\| \|\Delta_{t-1}\|]) \\ &\leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + \mathbb{E}[L (\|Q_x(x_t) - x_t\| + \|x_t - x_{t-1}\| + \|Q_x(x_{t-1}) - x_{t-1}\|) \|\Delta_{t-1}\|]) \\ &\leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + \mathbb{E}[L (\delta_x + \|\Delta_{t-1}\| + \delta_x) \|\Delta_{t-1}\|]) \\ &\leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} (\mathbb{E}\langle \nabla f(Q_x(x_{t-1})), \Delta_{t-1} \rangle + L \mathbb{E}[\|\Delta_{t-1}\|^2] + 2L\delta_x \mathbb{E}[\|\Delta_{t-1}\|]) = \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} M_{t-1}.\end{aligned}$$

The rest six parts remain the same as Lemma 4.6. Then we have

$$\mathbb{E}\langle \nabla f(Q_x(x_t)), \Delta_t \rangle \leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} M_{t-1} + C_2 (1 - \theta_t) - \frac{1 - \beta}{2} \mathbb{E}[\|\nabla f(Q_x(x_t))\|_{\hat{\eta}_t}^2],$$

$$\text{where } C_2 = \frac{5\alpha G^3(1-\beta)}{2\epsilon\sqrt{\theta}} \left(\frac{\beta}{(1-\beta)\sqrt{\theta_1}C_1(1-\gamma)} + 1 \right)^2 + \frac{5\alpha G^3}{2\epsilon\sqrt{\theta}} + \frac{5\beta^2\alpha d\sqrt{\epsilon}}{2\sqrt{\theta}(1-\beta)\theta_1C_1(1-\gamma)} + \frac{5\alpha\sqrt{G^2+\epsilon}G^2\beta^2}{2(1-\beta)\sqrt{\theta}\theta_1C_1(1-\gamma)\epsilon} + \frac{5\alpha\sqrt{G^2+\epsilon}\beta^2d}{2(1-\beta)\sqrt{\theta}\theta_1C_1(1-\gamma)}.$$

Define $N_t = L \mathbb{E}[\|\Delta_{t-1}\|^2] + 2L\delta_x \mathbb{E}[\|\Delta_{t-1}\|] + C_2 (1 - \theta_t)$. By applying the same induction in Lemma 4.6, we have

$$\begin{aligned}\sum_{t=1}^T M_t &\leq \frac{1}{\sqrt{C_1} (1 - \sqrt{\gamma})} \sum_{t=1}^T N_t - \frac{1 - \beta}{2} \mathbb{E}[\sum_{t=1}^T \|\nabla f(Q_x(x_t))\|_{\hat{\eta}_t}^2] \\ &= \frac{1}{\sqrt{C_1} (1 - \sqrt{\gamma})} \left(\frac{LG^2\alpha^2}{\epsilon} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4L\delta_x G\alpha\sqrt{T}}{\sqrt{\epsilon}C_1 (1 - \sqrt{\gamma})} - \frac{1 - \beta}{2} \mathbb{E}[\sum_{t=1}^T \|\nabla f(Q_x(x_t))\|_{\hat{\eta}_t}^2].\end{aligned}$$

Hence, we obtain the desired result. \square

PROOF OF THEOREM 3.2. By the Lipschitz continuity of the gradient, we have

$$\begin{aligned}f(x_{t+1}) &\leq f(x_t) + \langle \nabla f(x_t), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 \\ &= f(x_t) + \langle \nabla f(Q_x(x_t)), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 + \langle \nabla f(x_t) - \nabla f(Q_x(x_t)), \Delta_t \rangle \\ &\leq f(x_t) + \langle \nabla f(Q_x(x_t)), \Delta_t \rangle + \frac{L}{2} \|\Delta_t\|^2 + L\delta_x \|\Delta_t\|.\end{aligned}$$

Let $M_t = \mathbb{E}[\langle \nabla f(Q_t(x_t)), \Delta_t \rangle + L\|\Delta_t\|^2 + 2L\delta_x\|\Delta_t\|]$. Then we have

$$\begin{aligned}f^* &\leq \mathbb{E}[f(x_{T+1})] \leq f(x_1) + \sum_{t=1}^T M_t \\ &\leq f(x_1) + \frac{1}{\sqrt{C_1} (1 - \sqrt{\gamma})} \left(\frac{LG^2\alpha^2}{\epsilon} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4L\delta_x G\alpha\sqrt{T}}{\sqrt{\epsilon}C_1 (1 - \sqrt{\gamma})} - \frac{1 - \beta}{2} \mathbb{E}[\sum_{t=1}^T \|\nabla f(Q_x(x_t))\|_{\hat{\eta}_t}^2].\end{aligned}$$

Arranging the corresponding terms suitably, we obtain

$$\begin{aligned} \mathbb{E} \left[\left\| \nabla f(Q_x(x_t^T)) \right\|^2 \right] &\leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(Q(x_t)) \right\|_{\eta_t}^2 \right] \\ &\leq \frac{2\sqrt{G^2 + \epsilon d}}{(1-\beta)\alpha\sqrt{T}} \left(f(x_1) - f^* + \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{LG^2\alpha^2}{\epsilon} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4L\delta_x G\alpha\sqrt{T}}{\sqrt{\epsilon C_1}(1-\sqrt{\gamma})} \right) \\ &= \frac{C_5 + C_6 \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}} + C_7, \end{aligned}$$

where C_5, C_6 , and C_7 are defined as $C_5 = \frac{2\sqrt{G^2 + \epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C_6 = \frac{2\sqrt{G^2 + \epsilon d}}{(1-\beta)\alpha\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{LG^2\alpha^2}{\epsilon} + C_2\theta \right)$, and $C_7 = \frac{8\delta_x\sqrt{G^2 + \epsilon d}LG}{(1-\beta)\sqrt{\epsilon C_1}(1-\sqrt{\gamma})}$, respectively.

Hence, the proof is completed. \square

PROOF OF COROLLARY 3.2.1. For a fixed iteration T , let $\alpha_t = \frac{\alpha}{\sqrt{t}}$ and $\theta_t = 1 - \frac{\theta}{T}$. It is not hard to check that Lemma 4.7 updates to $\sum_{t=1}^T \|\Delta_t\| \leq \frac{\alpha G \sqrt{T}}{\sqrt{\epsilon}}$, and Lemma 4.8 updates to

$$\sum_{t=1}^T M_t \leq \frac{LG^2\alpha^2}{\sqrt{C_1}(1-\sqrt{\gamma})\epsilon} + \frac{C_2\theta}{\sqrt{C_1}(1-\sqrt{\gamma})} + \frac{2L\delta_x G\alpha\sqrt{T}}{\sqrt{\epsilon C_1}(1-\sqrt{\gamma})} - \frac{1-\beta}{2} \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(Q_x(x_t)) \right\|_{\eta_t}^2 \right].$$

Then, we have

$$\mathbb{E} \left[\left\| \nabla f(Q_x(x_t^T)) \right\|^2 \right] \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E} \left[\sum_{t=1}^T \left\| \nabla f(Q(x_t)) \right\|_{\eta_t}^2 \right] \leq \frac{C'_5}{\sqrt{T}} + C'_7,$$

where $C'_5 = \frac{2\sqrt{G^2 + \epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^* + \frac{LG^2\alpha^2}{\sqrt{C_1}(1-\sqrt{\gamma})\epsilon} + \frac{C_2\theta}{\sqrt{C_1}(1-\sqrt{\gamma})})$, $C'_7 = \frac{4\delta_x\sqrt{G^2 + \epsilon d}LG}{(1-\beta)\sqrt{\epsilon C_1}(1-\sqrt{\gamma})}$.

Hence, by bounding $\frac{C'_5}{\sqrt{T}} + C'_7 \leq C'_7 + \xi$, we obtain the desired result. \square

4.3 Proof of Theorem 3.3

To prove Theorem 3, we first denote a few notations and provide several useful lemmas.

NOTATION 3. Let $\mathbb{E}_t[\cdot]$ be the conditional expectation conditioned on $\{x_t, v_{t-1}, m_{t-1}\}$. Denote $\sigma_t^{(i)^2} = \mathbb{E}_t[g_t^{(i)^2}]$, $\hat{v}_t^{(i)} = \theta_t v_{t-1}^{(i)} + (1-\theta_t)\sigma_t^{(i)^2}$, $\hat{\eta}_t^{(i)} = \alpha_t / \sqrt{\hat{v}_t^{(i)}} + \epsilon$, $\Delta_t^{(i)} = -\alpha_t m_t^{(i)} / \sqrt{v_t^{(i)}} + \epsilon$, $\hat{\Delta}_t = \frac{1}{N} \sum_{i=1}^N \Delta_t^{(i)}$, and $M_t = \frac{1}{N} \sum_{i=1}^N \mathbb{E} \langle \nabla f(\hat{x}_t), \Delta_t^{(i)} \rangle + \frac{1}{N^2} L(2-\delta_g) \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} \|\Delta_t^{(i)}\| \|\Delta_t^{(j)}\| + \frac{1}{N^2} L(2-\delta_g) \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} \|\Delta_t^{(i)}\| \|\Delta_t^{(j)}\| + \frac{1}{N} 2L \sum_{i=1}^N \mathbb{E} \|\Delta_t^{(i)}\|$. In addition, let $\tilde{x}_t = x_t - \frac{1}{N} \sum_{i=1}^N e_t^{(i)}$. Then

$$\tilde{x}_{t+1} = x_t - \frac{1}{N} \sum_{i=1}^N Q_g(-\Delta_t^{(i)} + e_t^{(i)}) - \frac{1}{N} \sum_{i=1}^N e_{t+1}^{(i)} = x_t + \frac{1}{N} \sum_{i=1}^N \Delta_t^{(i)} - \frac{1}{N} \sum_{i=1}^N e_t = \tilde{x}_t + \frac{1}{N} \sum_{i=1}^N \Delta_t^{(i)}.$$

LEMMA 4.9. With Notation 3, we derive an upper estimation for $\hat{\Delta}_t$ as: $\sum_{t=1}^T \|\hat{\Delta}_t\|^2 \leq \frac{G^2}{\epsilon} \sum_{t=1}^T \frac{\alpha^2}{t}$.

PROOF. By using the definition of $\hat{\Delta}_t$, it is not hard to check that the following equations hold:

$$\sum_{t=1}^T \|\hat{\Delta}_t\|^2 = \sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \Delta_t^{(i)} \right\|^2 \leq \sum_{t=1}^T \frac{1}{N} \sum_{i=1}^N \|\Delta_t^{(i)}\|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \|\Delta_t^{(i)}\|^2 \leq \frac{G^2}{\epsilon} \sum_{t=1}^T \frac{\alpha^2}{t}.$$

\square

LEMMA 4.10. Let e_i be the noisy term in Algorithm2 and $\hat{\Delta}_t$ be the term defined in Notation 3. Then it holds that

$$\mathbb{E} \left[\frac{1}{N^2} \sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \|e_t^{(i)}\| \|\Delta_t^{(j)}\| \right] \leq \frac{1-\delta_g}{\delta_g N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\Delta_t^{(i)}\|^2 \right].$$

PROOF. Using the definition of the noisy term e_t , the following holds:

$$\|e_t^{(j)}\| = \|\Delta_t^{(j)} + e_{t-1}^{(j)} - Q_g(\Delta_t^{(j)} + e_{t-1}^{(j)})\| \leq (1-\delta_g) (\|\Delta_t^{(j)}\| + \|e_{t-1}^{(j)}\|) \leq \sum_{k=1}^t (1-\delta_g)^{t-k+1} \|\Delta_k^{(j)}\|.$$

Then, we have

$$\begin{aligned} \mathbb{E} \left[\frac{1}{N^2} \sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \|e_t^{(i)}\| \|\Delta_t^{(j)}\| \right] &\leq \mathbb{E} \left[\frac{1}{N^2} \sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^t (1-\delta_g)^{t-k+1} \|\Delta_k^{(i)}\| \|\Delta_t^{(j)}\| \right] \\ &\leq \mathbb{E} \left[\frac{1}{N^2} \sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^t \frac{(1-\delta_g)^{t-k+1}}{2} (\|\Delta_k^{(i)}\|^2 + \|\Delta_t^{(j)}\|^2) \right] \\ &= \mathbb{E} \left[\frac{1}{N} \sum_{t=1}^T \sum_{i=1}^N \sum_{k=1}^t \frac{(1-\delta_g)^{t-k+1}}{2} (\|\Delta_k^{(i)}\|^2) + \frac{1}{N} \sum_{t=1}^T \sum_{j=1}^N \sum_{k=1}^t \frac{(1-\delta_g)^{t-k+1}}{2} (\|\Delta_t^{(j)}\|^2) \right] \\ &\leq \mathbb{E} \left[\frac{1-\delta_g}{2\delta_g N} \sum_{k=1}^T \sum_{i=1}^N \|\Delta_k^{(i)}\|^2 + \frac{1-\delta_g}{2\delta_g N} \sum_{t=1}^T \sum_{j=1}^N \|\Delta_t^{(j)}\|^2 \right] = \mathbb{E} \left[\frac{1-\delta_g}{\delta_g N} \sum_{t=1}^T \sum_{i=1}^N \|\Delta_t^{(i)}\|^2 \right]. \end{aligned}$$

□

LEMMA 4.11. Let τ be randomly chosen from $\{1, 2, \dots, T\}$ with equal probabilities $p_\tau = \frac{1}{T}$. We have the following estimate:

$$\mathbb{E}[\|\nabla f(\hat{x}_\tau)\|^2] \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T} N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right].$$

PROOF. By Lemma 4.3, $\mathbb{E}[\|\nabla f(\hat{x}_\tau)\|^2] \leq \frac{\sqrt{G^2 + \epsilon d}}{\alpha \sqrt{T}} \mathbb{E}[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2]$ holds for any i . Hence, the proof is finished. □

LEMMA 4.12. By the definition of M_t , we obtain its upper-estimation as follows:

$$\sum_{t=1}^T M_t \leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4\delta_x LG\alpha\sqrt{T}}{\sqrt{\epsilon}} \right) \frac{1-\beta}{2N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right].$$

PROOF. Based on the similar proof of $\Delta_t - \frac{\beta\alpha_t}{\sqrt{\theta_t}\alpha_t} \Delta_{t-1}$ in Lemma 4.6, we define

$$\begin{aligned} (A_t^1)^{(i)} &= \hat{\eta}_t^{(i)} g_t^{(i)} \frac{(1-\theta_t)g_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon}} \left(\frac{\beta m_{t-1}^{(i)}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}} + \frac{(1-\beta)g_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\hat{v}_t^{(i)} + \epsilon}} \right), \\ (A_t^2)^{(i)} &= -\hat{\eta}_t^{(i)} \sigma_t^{(i)} \frac{(1-\theta_t)g_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon}} \frac{(1-\beta)\sigma_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\hat{v}_t^{(i)} + \epsilon}}, \quad (A_t^3)^{(i)} = \hat{\eta}_t^{(i)} \epsilon \frac{(1-\theta_t)}{\sqrt{v_t^{(i)} + \epsilon}} \frac{\beta m_{t-1}^{(i)}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}}, \end{aligned}$$

$$\begin{aligned}
(A_t^4)^{(i)} &= \hat{\eta}_t^{(i)} \sqrt{\sigma_t^{(i)^2} + \epsilon} \frac{(1-\theta_t) g_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon}} \frac{\beta m_{t-1}^{(i)}}{\sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} g_t^{(i)}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^{(i)^2} + \epsilon}}{\sqrt{\hat{v}_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}}, \\
(A_t^5)^{(i)} &= \hat{\eta}_t^{(i)} \sqrt{\sigma_t^{(i)^2} + \epsilon} \frac{(1-\theta_t) \sqrt{\epsilon}}{\sqrt{v_t^{(i)} + \epsilon}} \frac{\beta m_{t-1}^{(i)}}{\sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\epsilon}}{\sqrt{v_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}} \frac{\sqrt{1-\theta_t} \sqrt{\sigma_t^{(i)^2} + \epsilon}}{\sqrt{\hat{v}_t^{(i)} + \epsilon} + \sqrt{\theta_t v_{t-1}^{(i)} + \theta_t \epsilon}}.
\end{aligned}$$

Then, via the same proof as Lemma 4.6, the following equation holds

$$\begin{aligned}
\mathbb{E}\langle \nabla f(\hat{x}_t), \hat{\Delta}_t \rangle &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}\langle \nabla f(\hat{x}_t), \Delta_t^{(i)} \rangle = \frac{1}{N} \sum_{i=1}^N \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} \mathbb{E}\langle \nabla f(\hat{x}_t), \Delta_{t-1}^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t), \Delta_t^{(i)} - \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_t} \Delta_{t-1}^{(i)} \rangle \\
&= \frac{1}{N} \sum_{i=1}^N \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} \left(\mathbb{E}\langle \nabla f(\hat{x}_{t-1}), \Delta_{t-1}^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t) - \nabla f(\hat{x}_{t-1}), \Delta_{t-1}^{(i)} \rangle \right) \\
&\quad + \mathbb{E}\langle \nabla f(\hat{x}_t), -(1-\beta) \hat{\eta}_t^{(i)} g_t^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t^{(i)}), (A_t^1)^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t^{(i)}), (A_t^2)^{(i)} \rangle \\
&\quad + \mathbb{E}\langle \nabla f(\hat{x}_t^{(i)}), (A_t^3)^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t^{(i)}), (A_t^4)^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t^{(i)}), (A_t^5)^{(i)} \rangle.
\end{aligned} \tag{5}$$

For the first term in Eq. (5), it holds that

$$\begin{aligned}
&\frac{1}{N} \sum_{i=1}^N \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} \left(\mathbb{E}\langle \nabla f(\hat{x}_{t-1}), \Delta_{t-1}^{(i)} \rangle + \mathbb{E}\langle \nabla f(\hat{x}_t) - \nabla f(\hat{x}_{t-1}), \Delta_{t-1}^{(i)} \rangle \right) \\
&\leq \frac{1}{N} \sum_{i=1}^N \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} \left(\mathbb{E}\langle \nabla f(\hat{x}_{t-1}), \Delta_{t-1}^{(i)} \rangle + L \mathbb{E} \|\hat{x}_t - \hat{x}_{t-1}\| \|\Delta_{t-1}^{(i)}\| \right) = \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} M_{t-1}.
\end{aligned}$$

The remain 6 terms in Eq. (5) are the same as Lemma 4.6. Thus, we have

$$\mathbb{E}\langle \nabla f(\hat{x}_t), \hat{\Delta}_t \rangle \leq \frac{\beta \alpha_t}{\sqrt{\theta_t} \alpha_{t-1}} M_{t-1} + C_2 (1-\theta_t) - \frac{1}{N} \sum_{i=1}^N \frac{1-\beta}{2} \mathbb{E}[\|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2].$$

Let $N_t = \frac{L(2-\delta_g)}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} \|\Delta_t^{(i)}\| \|\Delta_t^{(j)}\| + \frac{L(2-\delta_g)}{N^2} \sum_{i=1}^N \sum_{j=1}^N \mathbb{E} \|\Delta_t^{(i)}\| \|e_t^{(j)}\| + \frac{2L}{N} \sum_{i=1}^N \mathbb{E} \|\Delta_t^{(i)}\| + C_2 (1-\theta_t)$. Similarly, we can acquire

$$\begin{aligned}
\sum_{t=1}^T M_t &\leq \frac{1}{\sqrt{C_1} (1-\sqrt{\gamma})} \sum_{t=1}^T N_t - \frac{1-\beta}{2N} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right] \\
&\leq \frac{1}{\sqrt{C_1} (1-\sqrt{\gamma})} \left(\left(\frac{L(2-\delta_g) G^2 \alpha^2}{\epsilon \delta_g} + C_2 \theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4\delta_x L G \alpha \sqrt{T}}{\sqrt{\epsilon}} \right) - \frac{1-\beta}{2N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right].
\end{aligned}$$

Hence, the proof is completed. \square

PROOF OF THEOREM 3.3. By using the gradient Lipschitz continuity of f , it holds that

$$\begin{aligned}
\mathbb{E}f(\tilde{x}_{t+1}) &\leq \mathbb{E} \left(f(\tilde{x}_t) + \langle \nabla f(\tilde{x}_t), \hat{\Delta}_t \rangle + \frac{L}{2} \|\hat{\Delta}_t\|^2 \right) \\
&= \mathbb{E} \left[f(\tilde{x}_t) + \langle \nabla f(\hat{x}_t), \hat{\Delta}_t \rangle + \frac{L}{2} \|\hat{\Delta}_t\|^2 + \langle \nabla f(x_t) - \nabla f(\hat{x}_t), \hat{\Delta}_t \rangle + \langle \nabla f(\tilde{x}_t) - \nabla f(x_t), \hat{\Delta}_t \rangle \right] \\
&\leq \mathbb{E}f(\tilde{x}_t) + M_t.
\end{aligned}$$

Taking summation over both sides of the above inequality, it holds that

$$\begin{aligned} f^* &\leq \mathbb{E}f(\tilde{x}_{T+1}) \leq f(x_1) + \sum_{t=1}^T M_t \\ &\leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4\delta_x LG\alpha\sqrt{T}}{\sqrt{\epsilon}} \right) - \frac{1-\beta}{2N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right]. \end{aligned}$$

Arranging the terms, it holds that

$$\begin{aligned} \mathbb{E}[\|\nabla f(\hat{x}_T)\|^2] &\leq \frac{\sqrt{G^2+\epsilon d}}{\alpha\sqrt{T}N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f((\hat{x}_t))\|_{\hat{\eta}_t^{(i)}}^2 \right], \\ &\leq \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha\sqrt{T}} \left\{ f(x_1) - f^* + \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \sum_{t=1}^T \frac{1}{t} + \frac{4\delta_x LG\alpha\sqrt{T}}{\sqrt{\epsilon}} \right) \right\} \\ &= \frac{C_8 + C_9 \sum_{t=1}^T \frac{1}{t}}{\sqrt{T}} + C_{10}, \end{aligned}$$

where $C_8 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} (f(x_1) - f^*)$, $C_9 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right)$, $C_{10} = \frac{8\sqrt{G^2+\epsilon d}\delta_x LG}{\sqrt{C_1}(1-\sqrt{\gamma})\sqrt{\epsilon}(1-\beta)}$, respectively.

Hence, the proof is completed. \square

PROOF OF COROLLARY 3.3.1. Given iteration T , let $\alpha_t = \frac{\alpha}{\sqrt{T}}$ and $\theta_t = 1 - \frac{\theta}{T}$. Lemma 4.12 updates to

$$\sum_{t=1}^T M_t \leq \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta + \frac{2\delta_x LG\alpha\sqrt{T}}{\sqrt{\epsilon}} \right) - \frac{1-\beta}{2N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f(\hat{x}_t)\|_{\hat{\eta}_t^{(i)}}^2 \right].$$

Then, we have

$$\mathbb{E}[\|\nabla f(\hat{x}_T)\|^2] \leq \frac{\sqrt{G^2+\epsilon d}}{\alpha\sqrt{T}N} \sum_{i=1}^N \mathbb{E} \left[\sum_{t=1}^T \|\nabla f((\hat{x}_t))\|_{\hat{\eta}_t^{(i)}}^2 \right] \leq \frac{C'_8}{\sqrt{T}} + C'_{10},$$

where $C'_8 = \frac{2\sqrt{G^2+\epsilon d}}{(1-\beta)\alpha} \left(f(x_1) - f^* + \frac{1}{\sqrt{C_1}(1-\sqrt{\gamma})} \left(\frac{L(2-\delta_g)G^2\alpha^2}{\epsilon\delta_g} + C_2\theta \right) \right)$, $C'_{10} = \frac{4\sqrt{G^2+\epsilon d}\delta_x LG}{\sqrt{C_1}(1-\sqrt{\gamma})\sqrt{\epsilon}(1-\beta)}$.

Hence, by bounding $\frac{C'_8}{\sqrt{T}} + C'_{10} \leq C'_{10} + \xi$, we obtain the desired result. \square

5 EXPERIMENTS

In this section, we apply the proposed algorithms to train deep neural networks including VGG16 network [35] and ResNet-101 [10] on CIFAR10 [19] and CIFAR100 [19] datasets, respectively. Below, we first describe the implementation details and then show the experimental results.

5.1 Implementation Details

We evaluate the effectiveness of the proposed algorithms via training ResNet-101 [10] on the CIFAR100 dataset. The CIFAR100 dataset contains 100 classes. In each class, there are 600 color images with the size of 32×32 , among which 500 images serve as training images and the rest are testing images. We train ResNet-101 with 8 workers, each worker will take a batch gradient with batch size 16, so the batch size for one optimization step is 128. We set the total number of training epochs as 200. In addition, we add ℓ_2 regularization with a scalar $5e-4$ as the weight decay term. Input images are all down-scaled to 1/8 of their original sizes after 100 convolutional layers and then fed into a fully-connected layer for the 100-class classifications. The output channel numbers of 1-11 conv layers, 12-23 conv layers, 24-92 conv layers, and 93-100 conv layers are 64, 128, 256, and

512, respectively. Also, we train VGG16 network [35] on the CIFAR10 dataset. CIFAR10 contains 10 classes of 32×32 images, including 50,000 training images and 10,000 testing images. We train VGG with 8 workers and each worker takes 16 samples to estimate gradient for 78,200 iterations with ℓ_2 regularization which is the same as the previous regularization term. The VGG16 contains 13 convolutional layers, and 3 fully-connected layers with 4096, 4096 and 10 neurons per layer. The output channel numbers of 1-2 conv layers, 3-4 conv layers, 5-7 layers, 8-13 layers are 64, 128, 256, 512, respectively.

Besides, in all experiments we choose β as 0.99, θ as 0.999, and ϵ as $1e-5$. Since it is natural to choose the exponential decay strategy on learning rate, we choose to reduce α by half every 50 epochs instead of using α/\sqrt{t} as [33]. We choose the starting learning rate as 0.001. The value is chosen based on grid search on the set $\{0.01, 0.001, 0.0001\}$ based on the accuracy of the full precision setting. For gradient quantization, we use function $Q_g(g) = \|g\|_\infty \arg \min_{\hat{g} \in \mathcal{G}} \|g/\|g\|_\infty - \hat{g}\|$, where $\mathcal{G} = \{-1, \dots, -2^{-k_g}, 0, 2^{-k_g}, 2^{-k_g+1}, \dots, 1\}$. For weight quantization, we use function $Q_x(x) = 0.5 \times \arg \min_{\hat{x} \in \mathcal{X}} \|2x - \hat{x}\|$, where $\mathcal{X} = \{-1, \dots, -\frac{1}{2^{k_x}}, 0, \frac{1}{2^{k_x}}, \frac{2}{2^{k_x}}, \dots, 1\}$.

We compared our method with TernGrad [39] and Zheng et al. [44] for gradient quantization, where the learning rate in these two methods is 0.1 which generated by grid search in $\{0.1, 0.05, 0.01\}$. For weight quantization, we compare with the result which quantizes the final model directly named WQuan in tables.

5.2 Results Illustration

In this section, we will illustrate the results of training ResNet-101 on the CIFAR100 dataset and training VGG16 on the CIFAR10 dataset, respectively. Two tables show the test accuracy after 200 epochs of training, where the first column represents training methods, the second column represents the test accuracy, the third column represents bits required for gradient communication, and the last column represents the bits to save a model. As for the same method, we can set different k_x and k_g , and we can get a different number of bits needed for gradients and weights.

5.2.1 Results of Training ResNet-101 on the CIFAR100. Figure 3 and Table 2 show the result of training ResNet-101 on the CIFAR100 dataset. When the algorithm involves gradient quantization, we compare our algorithm with Zheng et al. [44] and TernGrad [39] with a different number of communication bits, which has been shown in the left figure of Figure 3 and the first 9 lines in Table 2. It can be shown that even with gradient quantization, our algorithm can outperform TernGrad [39] and Zheng et al. [44]. The middle figure in Figure 3 shows the result of using weight quantization only. Row 10-12 in Table 2 shows the comparison results between quantizing weight during training and after training, which shows when quantizing weight during the training process can achieve higher test accuracy. The right figure in Figure 3 and the last 4 rows demonstrate the results of

Method	Test Acc	Comm	Size
QADAM	77.94 \pm 0.44%	162.9	162.9
QADAM	77.44 \pm 0.40%	15.27	162.9
QADAM	78.36 \pm 0.37%	10.18	162.9
TernGrad[39]	76.69 \pm 0.16%	162.9	162.9
TernGrad[39]	76.62 \pm 0.38%	15.27	162.9
TernGrad[39]	76.00 \pm 0.15%	10.18	162.9
[44]	76.28 \pm 0.24%	162.9	162.9
[44]	76.20 \pm 0.47%	15.27	162.9
[44]	76.06 \pm 0.38%	10.18	162.9
QADAM	78.10 \pm 0.33%	162.9	81.44
QADAM	78.08 \pm 0.45%	162.9	40.72
WQuan	77.00 \pm 0.57%	162.9	81.44
WQuan	76.68 \pm 0.42%	162.9	40.72
QADAM	78.18 \pm 0.43%	15.27	81.44
QADAM	78.32 \pm 0.47%	10.18	81.44
QADAM	78.19 \pm 0.28%	15.27	40.72
QADAM	78.04 \pm 0.39%	10.18	40.72

Table 2. Test Accuracy for CIFAR100. "Comm" stands for communication cost and unit is MB per iteration. "Size" stands for model size and unit is MB. WQuan means quantizing weight after training.

combining gradient quantization and weight quantization. It shows even though we shrink model size into 1/4 of its original size and gradient size into 1/16 of its original size, respectively, it can still give comparable results.

5.2.2 Results of Training VGG16 on the CIFAR10. Figure 4 and Table 3 show the result of training VGG16 on the CIFAR10 dataset. The left figure in Figure 4 and the first 9 rows in Table 3 show results of gradient quantization comparison among our algorithm, TernGrad[39] and Zheng et al.[44]. Our results and Zheng et al.[44] can achieve similar performance, but TernGrad [39] gets worse performance due to noise to ensure unbiasedness. When considering weight quantization, shown in the middle figure of Figure 4 and row 10-13 in Table 3, performance is similar between quantizing during training or after training. Further, as it has shown in the last 4 rows with different sizes of the model and different gradient quantization, our method can still achieve high accuracy compared to the full precision version (first row in Table 3).

Method	Test Acc	Comm	Size
QADAM	91.21 \pm 0.24%	512.3	512.3
QADAM	91.29 \pm 0.26%	48.03	512.3
QADAM	90.64 \pm 0.18%	32.02	512.3
TernGrad[39]	90.60 \pm 0.13%	512.3	512.3
TernGrad[39]	90.60 \pm 0.32%	48.03	512.3
TernGrad[39]	89.91 \pm 0.28%	32.02	512.3
[44]	91.51 \pm 0.22%	512.3	512.3
[44]	91.38 \pm 0.36%	48.03	512.3
[44]	90.93 \pm 0.10%	32.02	512.3
QADAM	91.06 \pm 0.31%	512.3	256.2
QADAM	91.21 \pm 0.41%	512.3	128.1
WQuan	91.21 \pm 0.24%	512.3	256.2
WQuan	90.18 \pm 0.39%	512.3	128.1
QADAM	90.87 \pm 0.27%	48.03	256.2
QADAM	91.20 \pm 0.22%	32.02	256.2
QADAM	91.16 \pm 0.26%	48.03	128.1
QADAM	90.48 \pm 0.36%	32.02	128.1

Table 3. Test Accuracy for CIFAR10. "Comm" stands for communication cost and unit is MB per iteration. "Size" stands for model size and unit is MB. WQuan means quantizing weight with the first model.

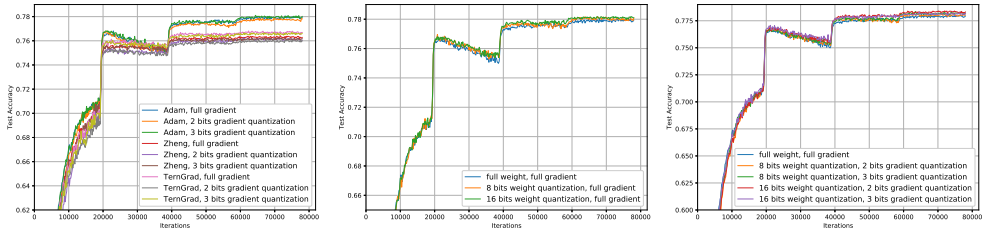


Fig. 3. Results for Training ResNet-101 on CIFAR100.

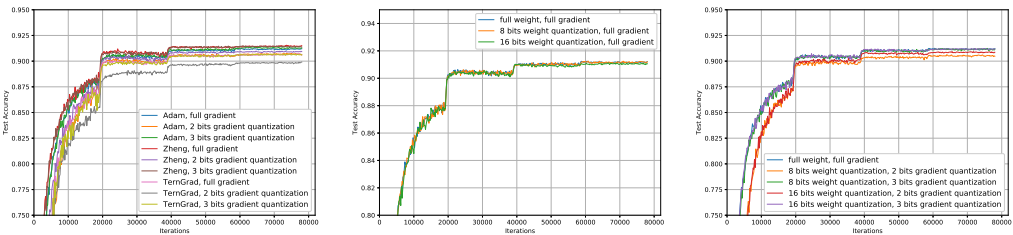


Fig. 4. Results for Training VGG16 on CIFAR10.

6 CONCLUSIONS

To accelerate the training process of deep learning models, we proposed distributed Adam with weight quantization, gradient quantization, and the error-feedback technique in the parameter-server model. Through capitalizing on the two schemes of weight quantization and gradient

quantization, the communication cost between the server and works can be significantly alleviated. In addition, the proposed error-feedback technique can suppress the bias caused by the gradient quantization step, thereby making the proposed algorithms more efficient. We further established the convergence rates of the proposed algorithms in the nonconvex stochastic setting and showed that quantized Adam with the error-feedback technique converges to the neighborhood of a stationary point under both the single-worker and multi-worker modes. Moreover, we applied the proposed algorithms to train VGG16 on the CIFAR10 dataset and ResNet-101 on the CIFAR100 dataset, respectively. The experiments demonstrate the efficacy of the proposed algorithms.

REFERENCES

- [1] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. 2016. Deep speech 2: End-to-end speech recognition in english and mandarin. In *International conference on machine learning*. 173–182.
- [2] Amitabh Basu, Soham De, Anirbit Mukherjee, and Enayat Ullah. 2018. Convergence guarantees for rmsprop and adam in non-convex optimization and their comparison to nesterov acceleration on autoencoders. *arXiv preprint arXiv:1807.06766* (2018).
- [3] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. 2018. On the convergence of a class of adam-type algorithms for non-convex optimization. *arXiv preprint arXiv:1808.02941* (2018).
- [4] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’aurelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. 2012. Large scale distributed deep networks. In *Advances in neural information processing systems*. 1223–1231.
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [7] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research* 12, Jul (2011), 2121–2159.
- [8] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT press.
- [9] Song Han, Huizi Mao, and William J Dally. 2015. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149* (2015).
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [11] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. 2012. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. (2012).
- [12] Lu Hou, Ruiliang Zhang, and James T Kwok. 2018. Analysis of quantized models. (2018).
- [13] Peng Jiang and Gagan Agrawal. 2018. A linear speedup analysis of distributed deep learning with sparse and quantized communication. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*. 2530–2541.
- [14] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. 2019. Error feedback fixes signsgd and other gradient compression schemes. In *International Conference on Machine Learning*. PMLR, 3252–3261.
- [15] Ahmed Khaled and Peter Richtárik. 2019. Gradient descent with compressed iterates. *arXiv preprint arXiv:1909.04716* (2019).
- [16] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [17] Anastasia Koloskova, Tao Lin, Sebastian U Stich, and Martin Jaggi. 2019. Decentralized Deep Learning with Arbitrary Communication Compression. In *International Conference on Learning Representations*.
- [18] Tim Kraska, Ameet S. Talwalkar, John C. Duchi, R. Griffith, M. Franklin, and Michael I. Jordan. 2013. MLbase: A Distributed Machine-learning System. In *CIDR*.
- [19] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [21] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.
- [22] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.
- [23] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, and Bor-Yiing Su. 2014. Scaling distributed machine learning with the parameter server. In *11th {USENIX}*

- Symposium on Operating Systems Design and Implementation ({OSDI} 14)*. 583–598.
- [24] Xiaoyu Li and Francesco Orabona. 2019. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 983–992.
 - [25] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. 2017. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*. 5330–5340.
 - [26] Tie-Yan Liu, Wei Chen, and Taifeng Wang. 2017. Distributed machine learning: Foundations, trends, and practices. In *Proceedings of the 26th International Conference on World Wide Web Companion*. 913–915.
 - [27] H Brendan McMahan et al. 2021. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning* 14, 1 (2021).
 - [28] H Brendan McMahan and Matthew Streeter. 2010. Adaptive bound optimization for online convex optimization. *arXiv preprint arXiv:1002.4908* (2010).
 - [29] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fiedjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533.
 - [30] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. 2019. Dadam: A consensus-based distributed adaptive gradient method for online optimization. *arXiv preprint arXiv:1901.09109* (2019).
 - [31] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. 2016. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*. Springer, 525–542.
 - [32] Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. 2020. Adaptive Federated Optimization. *arXiv preprint arXiv:2003.00295* (2020).
 - [33] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. 2019. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237* (2019).
 - [34] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature* 529, 7587 (2016), 484.
 - [35] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
 - [36] Alexander Smola and Shraavan Narayanamurthy. 2010. An architecture for parallel topic models. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 703–710.
 - [37] H Tang, X Lian, S Qiu, L Yuan, C Zhang, T Zhang, and J Liu. 2019. DeepSqueeze: Decentralized meets error-compensated compression. *arXiv preprint arXiv:1907.07346* (2019).
 - [38] Rachel Ward, Xiaoxia Wu, and Leon Bottou. 2019. AdaGrad stepsizes: Sharp convergence over nonconvex landscapes. In *International Conference on Machine Learning*. PMLR, 6677–6686.
 - [39] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. 2017. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In *Advances in neural information processing systems*. 1509–1519.
 - [40] Baoyuan Wu, Weidong Chen, Yanbo Fan, Yong Zhang, Jinlong Hou, Jie Liu, and Tong Zhang. 2019. Tencent ml-images: A large-scale multi-label image database for visual representation learning. *IEEE Access* 7 (2019), 172683–172693.
 - [41] Shuang Wu, Guoqi Li, Feng Chen, and Luping Shi. 2018. Training and inference with integers in deep neural networks. *arXiv preprint arXiv:1802.04680* (2018).
 - [42] Eric P Xing, Qirong Ho, Wei Dai, Jin Kyu Kim, Jinliang Wei, Seunghak Lee, Xun Zheng, Pengtao Xie, Abhimanu Kumar, and Yaoliang Yu. 2015. Petuum: A new platform for distributed machine learning on big data. *IEEE Transactions on Big Data* 1, 2 (2015), 49–67.
 - [43] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)* 10, 2 (2019), 1–19.
 - [44] Shuai Zheng, Ziyue Huang, and James Kwok. 2019. Communication-efficient distributed blockwise momentum SGD with error-feedback. In *Advances in Neural Information Processing Systems*. 11450–11460.
 - [45] Shuchang Zhou, Zekun Ni, Xinyu Zhou, He Wen, Yuxin Wu, and Yuheng Zou. 2016. DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients. *CoRR* abs/1606.06160 (2016). [arXiv:1606.06160](http://arxiv.org/abs/1606.06160) <http://arxiv.org/abs/1606.06160>
 - [46] Zhiming Zhou, Qingru Zhang, Guansong Lu, Hongwei Wang, Weinan Zhang, and Yong Yu. 2018. AdaShift: Decorrelation and Convergence of Adaptive Learning Rate Methods. In *International Conference on Learning Representations*.
 - [47] Fangyu Zou, Li Shen, Zequn Jie, Ju Sun, and Wei Liu. 2018. Weighted AdaGrad with unified momentum. *arXiv preprint arXiv:1808.03408* (2018).
 - [48] Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. 2019. A sufficient condition for convergences of adam and rmsprop. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 11127–11135.