

---

# Sparsified Distributed Adaptive Learning with Error Feedback

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 To be completed...

## 2 1 Introduction

3 Deep neural network has achieved the state-of-the-art learning performance on numerous AI appli-  
4 cations, e.g., computer vision [16, 19, 36], Natural Language Processing [18, 42, 43], Reinforcement  
5 Learning [28, 34] and recommendation systems [10, 38]. With the increasing size of both data and  
6 deep networks, standard single machine training confronts with at least two major challenges:

- 7 • Due to the limited computing power of a single machine, it would take a long time to  
8 process the massive number of data samples—training would be slow.
- 9 • In many practical scenarios, data are typically stored in multiple servers, possibly at differ-  
10 ent locations, due to the storage constraints (massive user behavior data, Internet images,  
11 etc.) or privacy reasons [7]. Transmitting data might be costly.

12 *Distributed learning* framework [12] has been a common training strategy to tackle the above two  
13 issues. For example, in centralized distributed stochastic gradient descent (SGD) protocol, data are  
14 located at  $N$  local nodes, at which the gradients of the model are computed in parallel. In each  
15 iteration, a central server aggregates the local gradients, updates the global model, and transmits  
16 back the updated model to the local nodes for subsequent gradient computation. As we can see, this  
17 setting naturally solves aforementioned issues: 1) We use  $N$  computing nodes to train the model, so  
18 the time per training epoch can be largely reduced; 2) There is no need to transmit the local data to  
19 central server. Besides, distributed training also provides stronger error tolerance since the training  
20 process could continue even one local machine breaks down. As a result of these advantages, there  
21 has been a surge of study and applications on distributed systems [6, 30, 13, 17, 20, 26, 25].

22 Among many optimization strategies, SGD is still the most popular prototype in distributed training  
23 for its simplicity and effectiveness [9, 1, 27]. Yet, when the deep learning model is very large, the  
24 communication between local nodes and central server could be expensive. Burdensome gradient  
25 transmission would slow down the whole training system, or even be impossible because of the lim-  
26 ited bandwidth in some applications. Thus, reducing the communication cost in distributed SGD has  
27 become an active topic, and an important ingredient of large-scale distributed systems (e.g. [32]).  
28 Solutions based on quantization, sparsification and other compression techniques of the local gradi-  
29 ents are proposed, e.g., [3, 39, 37, 35, 2, 5, 11, 41, 21]. As one would expect, in most approaches,  
30 there exists a trade-off between compression and model accuracy. In particular, larger bias of the  
31 compressed gradients usually brings more significant performance downgrade. Interestingly, [23]  
32 shows that the technique of *error feedback* is able to remedy the issue of such biased compressors,  
33 achieving same convergence rate and learning performance as full-gradient SGD.

34 On the other hand, in recent years, adaptive optimization algorithms (e.g. AdaGrad [14], Adam [24]  
35 and AMSGrad [31]) have become popular because of their superior empirical performance. These

methods use different implicit learning rates for different coordinates that keep changing adaptively throughout the training process, based on the learning trajectory. In many learning problems, adaptive methods have been shown to converge faster than SGD, sometimes with better generalization as well. However, the body of literature that combines adaptive methods with distributed training is still very limited. In this paper, we propose a distributed optimization algorithm with AMSGrad as the backbone, along with TopK sparsification to reduce the communication cost.

## 1.1 Our contributions

## 2 Related Work

### 2.1 Communication-efficient distributed SGD

**Quantization.** As we mentioned before, SGD is the most commonly adopted optimization method in distributed training of deep neural nets. To reduce the extensive communication in large-scale distributed systems, extensive works have considered various compression techniques applied to the gradient transaction procedure. The first strategy is quantization. [?] condenses 32-bit floating numbers into 8-bits when representing the gradients. [32, 5, 23?] use the extreme 1-bit information (sign) of the gradients, combined with tricks like momentum, majority vote and memory. Other quantization-based methods include QSGD [3, 40?] and LPC-SVRG [?], leveraging stochastic quantization. The saving in communication of quantization methods is moderate: for example, 8-bit quantization reduces the cost to 25% (compared with 32-bit full-precision). Even in the extreme 1-bit case, the largest compression ratio is around  $1/32 \approx 3.1\%$ .

**Sparsification.** Gradient sparsification is another popular solution which may provide higher compression rate. Instead of commuting the full gradient, each local worker only passes a few coordinates to the central server. Thus, we can more freely choose higher compression ratio (e.g., 1%, 0.1%), still achieving impressive performance in many applications [?]. Stochastic sparsification methods, including Random- $k$  and variance-based sparsification [37], select coordinates based on some sampling probability yielding unbiased gradient compressors. Deterministic methods are simpler, e.g., Top- $k$  [35, 33] (selecting  $k$  elements with largest magnitude), Deep Gradient Compression [?], but usually lead to biased gradient estimation. In [21], the central server identifies heavy-hitters from the count-sketch of the local gradients, which can be regarded as a noisy variant of Top- $k$  strategy. More applications and analysis of compressed distributed SGD can be found in [22?, 4? ?], among others.

**Error Feedback.** Biased gradient estimation, which is a consequence of many aforementioned methods (e.g., signSGD, Top- $k$ ), undermines the model training, both theoretically and empirically, with slower convergence and worse generalization. The technique of *error feedback* is able to “correct for the bias” and fix the convergence issue. In this procedure, the difference between the true stochastic gradient and the compressed one is accumulated locally, which is then added back to the local gradients in later iterations. [35, 23] prove the  $\mathcal{O}(\frac{1}{T})$  and  $\mathcal{O}(\frac{1}{\sqrt{T}})$  convergence rate of EF-SGD in strongly convex and non-convex setting respectively, matching the rates of vanilla SGD [? 15].

### 2.2 Adaptive optimization

When a large number of compute engines is available, being able to train global machine learning models while mutualizing the available and *decentralized* source of computation has been a growing focus for the community.

Decentralized optimization methods include methods such as ADMM [6], Distributed Subgradient Descent [30], Dual Averaging [13], Prox-PDA [20], GNSD [26], and Choco-SGD [25].

A recent work [8], which focuses on adaptive gradient methods, namely the Adam [24] and the AMSGrad [31] optimization methods, develops a decentralized variant of gradient based and adaptive methods in the context of gossip protocols. To date, very few contributions provided attempt to efficiently run adaptive gradient method in such a distributed setting. Apart from [8], (author?) [29] proposes a decentralized version of AMSGrad [31] which provably satisfies some non-standard

84 regret. Though, no sparsified variants of them have been proposed for practical purposes nor been  
85 studied in the literature.

### 86 3 Method

87 Most modern machine learning tasks can be casted as a large finite-sum optimization problem writ-  
88 ten as:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta) \quad (1)$$

89 where  $n$  denotes the number of workers,  $f_i$  represents the average loss for worker  $i$  and  $\theta$  the global  
90 model parameter taking value in  $\Theta$ , a subset of  $\mathbb{R}^d$ .

91 Some related work:

92 [23] develops variant of signSGD (as a biased compression schemes) for distributed optimization.  
93 Contributions are mainly on this error feedback variant. In [33], the authors provide theoretical  
94 results on the convergence of sparse Gradient SGD for distributed optimization (we want that for  
95 AMS here). [35] develops a variant of distributed SGD with sparse gradients too. Contributions  
96 include a memory term used while compressing the gradient (using top k for instance). Speeding up  
97 the convergence in  $\frac{1}{T^3}$ .

98 Consider standard synchronous distributed optimization setting. AMSGrad is used as the prototype,  
99 and the local workers is only in charge of gradient computation.

#### 100 3.1 TopK AMSGrad with Error Feedback

101 The key difference (and interesting part) of our TopK AMSGrad compared with the following arxiv  
102 paper “Quantized Adam”<https://arxiv.org/pdf/2004.14180.pdf> is that, in our model only  
103 gradients are transmitted. In “QAdam”, each local worker keeps a local copy of moment estimator  
104  $m$  and  $v$ , and compresses and transmits  $m/v$  as a whole. Thus, that method is very much like the  
105 sparsified distributed SGD, except that  $g$  is changed into  $m/v$ . In our model, the moment estimates  
106  $m$  and  $v$  are computed only at the central server, with the compressed gradients instead of the full  
107 gradient. This would be the key (and difficulty) in convergence analysis.

---

#### Algorithm 1 SPARS-AMS for Distributed Learning

---

- 1: **Input:** parameter  $\beta_1, \beta_2$ , learning rate  $\eta_t$ .
  - 2: Initialize: central server parameter  $\theta_0 \in \Theta \subseteq \mathbb{R}^d$ ;  $e_{0,i} = 0$  the error accumulator for each worker; sparsity parameter  $k$ ;  $n$  local workers;  $m_0 = 0, v_0 = 0, \hat{v}_0 = 0$
  - 3: **for**  $t = 1$  to  $T$  **do**
  - 4:   **parallel for worker**  $i \in [n]$  **do:**
  - 5:     Receive model parameter  $\theta_t$  from central server
  - 6:     Compute stochastic gradient  $g_{t,i}$  at  $\theta_t$
  - 7:     Compute  $\tilde{g}_{t,i} = \text{TopK}(g_{t,i} + e_{t,i}, k)$
  - 8:     Update the error  $e_{t+1,i} = e_{t,i} + g_{t,i} - \tilde{g}_{t,i}$
  - 9:     Send  $\tilde{g}_{t,i}$  back to central server
  - 10:   **end parallel**
  - 11:   **Central server do:**
  - 12:    $\bar{g}_t = \frac{1}{n} \sum_{i=1}^n \tilde{g}_{t,i}$
  - 13:    $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \bar{g}_t$
  - 14:    $v_t = \beta_2 v_{t-1} + (1 - \beta_2) \bar{g}_t^2$
  - 15:    $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$
  - 16:   Update global model  $\theta_t = \theta_{t-1} - \eta_t \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}$
  - 17: **end for**
-

### 108 3.2 Convergence Analysis

109 Several mild assumptions to make: Nonconvex and smooth loss function, unbiased stochastic gradi-  
 110 ent, bounded variance of the gradient, bounded norm of the gradient, control of the distance between  
 111 the true gradient and its sparse variant.

112 Check [8] starting with single machine and extending to distributed settings (several machines).

113 Under the distributed setting, the goal is to derive an upper bound to the second order moment of  
 114 the gradient of the objective function at some iteration  $T_f \in [1, T]$ .

### 115 3.3 Mild Assumptions

116 We begin by making the following assumptions.

117 **A 1. (Smoothness)** For  $i \in \llbracket n \rrbracket$ ,  $f_i$  is  $L$ -smooth:  $\|\nabla f_i(\theta) - \nabla f_i(\vartheta)\| \leq L \|\theta - \vartheta\|$ .

118 **A 2. (Unbiased and Bounded gradient *per worker*)** For any iteration index  $t > 0$  and worker index  
 119  $i \in \llbracket n \rrbracket$ , the stochastic gradient is unbiased and bounded from above:  $\mathbb{E}[g_{t,i}] = \nabla f_i(\theta_t)$  and  
 120  $\|g_{t,i}\| \leq G_i$ .

121 **A 3. (Bounded variance *per worker*)** For any iteration index  $t > 0$  and worker index  $i \in \llbracket n \rrbracket$ , the  
 122 variance of the noisy gradient is bounded:  $\mathbb{E}[|g_{t,i} - \nabla f_i(\theta_t)|^2] < \sigma_i^2$ .

123 Denote by  $Q(\cdot)$  the quantization operator Line 7 of Algorithm 1, which takes as input a gradient  
 124 vector and returns a quantized version of it, and note  $\tilde{g} := Q(g)$ . Assume that

125 **A 4. (Bounded Quantization)** For any iteration  $t > 0$ , there exists a constant  $0 < q < 1$  such that  
 126  $\|g_{t,i} - \tilde{g}_{t,i}\| \leq q \|g_{t,i}\|$ , where  $g_{t,i}$  is the stochastic gradient computed at iteration  $t$  for worker  $i$   
 127 and  $\tilde{g}_{t,i}$  is its quantized counterpart. (high  $q$  means large quantization so loss of precision on the  
 128 true gradient)

129 Denote for all  $\theta \in \Theta$ :

$$f(\theta) := \frac{1}{n} \sum_{i=1}^n f_i(\theta), \quad (2)$$

130 where  $n$  denotes the number of workers.

### 131 3.4 Intermediary Lemmas

132 **Lemma 1.** Under Assumption 2 and Assumption 4 we have for any iteration  $t > 0$ :

$$\|m_t\|^2 \leq (q^2 + 1)G^2 \quad \text{and} \quad \hat{v}_t \leq (q^2 + 1)G^2 \quad (3)$$

133 where  $m_t$  and  $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$  are defined Line 15 of Algorithm 1 and  $G^2 = \frac{1}{n} \sum_{i=1}^n G_i^2$ .

134 **Lemma 2.** Under A1 to A4, with a decreasing sequence of stepsize  $\{\eta_t\}_{t>0}$ , we have:

$$-\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \text{Id})^{-1/2} \bar{g}_t \right\rangle \right] \leq -\frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2} \quad (4)$$

135 where  $\text{Id}$  is the identity matrix,  $\hat{V}_t$  the diagonal matrix which diagonal entries are  $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$   
 136 defined Line 15 of Algorithm 1 and  $\bar{g}_t$  is the aggregation of all **quantized** gradients from the workers.

137 **Lemma 3.** Under A1 to A4, with a decreasing sequence of stepsize  $\{\eta_t\}_{t>0}$ , we have:

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] &\leq -\frac{\eta_{t+1}(1-\beta_1)}{2} \left( \epsilon + \frac{(q^2+1)G^2}{1-\beta_2} \right)^{-\frac{1}{2}} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2} \\
&\quad - \eta_{t+1} \beta_1 \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\
&\quad + \left( \frac{L}{2} + \beta_1 L \right) \|\theta_t - \theta_{t-1}\|^2 \\
&\quad + \eta_{t+1} G^2 \mathbb{E} \left[ \sum_{j=1}^d \left[ (\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2} \right] \right]
\end{aligned} \tag{5}$$

138 where  $d$  denotes the dimension of the parameter vector

139 The main theorem in the decentralized setting reads:

140 **Theorem 1.** Under A1 to A4, with a constant stepsize  $\eta_t = \eta = \frac{L}{\sqrt{T_m}}$ , we have:

$$\frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \frac{\mathbb{E}[f(\theta_0) - f(\theta_{T_m})]}{L \Delta_1 \sqrt{T_m}} + d \frac{L \Delta_3}{\Delta_1 \sqrt{T_m}} + \frac{\Delta_2}{\eta \Delta_1 T_m} + \frac{1-\beta_1}{\Delta_1} \epsilon^{-\frac{1}{2}} \sqrt{(q^2+1)} G^2 \tag{6}$$

141 where

$$\begin{aligned}
\Delta_1 &:= \frac{(1-\beta_1)}{2} \left( \epsilon + \frac{(q^2+1)G^2}{1-\beta_2} \right)^{-\frac{1}{2}}, \quad \Delta_2 := q^2 + \sum_{k=t+1}^{\infty} \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2} \\
\Delta_3 &:= \left( \frac{L}{2} + 1 + \frac{\beta_1 L}{1-\beta_1} \right) (1-\beta_2)^{-1} \left( 1 - \frac{\beta_1^2}{\beta_2} \right)^{-1}
\end{aligned} \tag{7}$$

142 We remark from this bound in Theorem 1, that the more quantization we apply to our gradient  
143 vectors ( $q \uparrow$ ), the larger the upper bound of the stationary condition is, *i.e.*, the slower the algorithm  
144 is. This is intuitive as using compressed quantities will definitely impact the algorithm speed. We  
145 will observe in the numerical section below that a trade-off on the level of quantization  $q$  can be  
146 found to achieve similar speed of convergence with less computation resources used throughout the  
147 training.

148 **Belhal Try for Single Machine Setting:**

149 Define the auxiliary model

$$\begin{aligned}
\theta'_{t+1} &:= \theta_{t+1} - e_{t+1} \\
&= \theta_t - \eta a_t - e_{t+1} \\
&= \theta_t - \eta a_t - e_t - g_t + \tilde{g}_t \\
&= \theta_t - \eta a_t - e_t - \Delta_t \\
&= \theta'_t - \eta a_t - \Delta_t
\end{aligned}$$

150 where  $a_t := \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}$  and  $\Delta_t := g_t - \tilde{g}_t$ . By smoothness assumption we have

$$f(\theta'_{t+1}) \leq f(\theta'_t) - \langle \nabla f(\theta'_t), \eta a_t + \Delta_t \rangle + \frac{L}{2} \|\theta'_{t+1} - \theta'_t\|^2.$$

151 Thus,

$$\begin{aligned}
\mathbb{E}[f(\theta'_{t+1}) - f(\theta'_t)] &\leq -\mathbb{E}[\langle \nabla f(\theta'_t), \eta a_t + \Delta_t \rangle] + \frac{L}{2} \mathbb{E}[\|\eta a_t + \Delta_t\|^2] \\
&\leq \eta \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), \eta a_t + \Delta_t \rangle] - \mathbb{E}[\langle \nabla f(\theta_t), \eta a_t + \Delta_t \rangle] + \frac{L}{2} \mathbb{E}[\|\eta a_t + \Delta_t\|^2]
\end{aligned}$$

152 Using the smoothness assumption A1 we have

$$\mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), \eta a_t + \Delta_t \rangle] \leq L \mathbb{E}[\|\theta_t - \theta'_t\|] \mathbb{E}[\|\eta a_t + \Delta_t\|]$$

153 Hence,

$$\begin{aligned} \mathbb{E}[f(\theta'_{t+1}) - f(\theta'_t)] &\leq -\mathbb{E}[\langle \nabla f(\theta'_t), \eta a_t + \Delta_t \rangle] + \frac{L}{2} \mathbb{E}[\|\eta a_t + \Delta_t\|^2] \\ &\leq -\left(\eta \frac{1}{\sqrt{G^2 + \epsilon}} + q\right) \mathbb{E}[\|\nabla f(\theta_t)\|^2] + L \mathbb{E}[\|\theta_t - \theta'_t\|] \mathbb{E}[\|\eta a_t + \Delta_t\|] + \frac{L}{2} \mathbb{E}[\|\eta a_t + \Delta_t\|^2] \\ &\leq -\left(\eta \frac{1}{\sqrt{G^2 + \epsilon}} + q\right) \mathbb{E}[\|\nabla f(\theta_t)\|^2] + L \mathbb{E}[\|e_t\| \|\eta a_t + \Delta_t\|] + \frac{L}{2} \mathbb{E}[\|\eta a_t + \Delta_t\|^2] \end{aligned}$$

154 Summing from  $t = 0$  to  $t = T_m - 1$  and divide it by  $T_m$  yields:

$$\begin{aligned} &\left(\eta \frac{1}{\sqrt{G^2 + \epsilon}} + q\right) \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\ &\leq \sum_{t=0}^{T_m-1} \frac{\mathbb{E}[f(\theta'_t) - f(\theta'_{t+1})]}{T_m} + \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|e_t\| \|\eta a_t + \Delta_t\|] + \frac{L}{2T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\eta a_t + \Delta_t\|^2] \end{aligned}$$

155 **Bounding**  $\frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|e_t\| \|\eta a_t + \Delta_t\|]$ :

156 To begin with

$$\begin{aligned} \|e_t\| &= \|e_{t-1} + g_{t-1} - \tilde{g}_{t-1}\| \\ &= \|g_{t-1} + e_{t-1} - \text{TopK}(g_{t-1} + e_{t-1}, k)\| \\ &\leq q \|g_{t-1} + e_{t-1}\| \\ &\leq q \|g_{t-1}\| + q \|e_{t-1}\| \\ &\leq \sum_{k=1}^t q^{t-k} \|g_k\| \end{aligned}$$

157 using A4.

158 **Bounding**  $\frac{L}{2T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\eta a_t + \Delta_t\|^2]$ :

## 159 4 Sequential Model

160 Single machine method

---

### Algorithm 2 SPARS-AMS : Single machine setting

---

- 1: **Input:** parameter  $\beta_1, \beta_2$ , learning rate  $\eta_t$ .
  - 2: Initialize: central server parameter  $\theta_0 \in \Theta \subseteq \mathbb{R}^d$ ;  $e_0 = 0$  the error accumulator; sparsity parameter  $k$ ;  $m_0 = 0, v_0 = 0, \hat{v}_0 = 0$
  - 3: **for**  $t = 1$  to  $T$  **do**
  - 4:   Compute stochastic gradient  $g_t = g_{t,i_t}$  at  $\theta_t$  for randomly sampled index  $i_t$
  - 5:   Compute  $\tilde{g}_t = \text{TopK}(g_t + e_t, k)$
  - 6:   Update the error  $e_{t+1} = e_t + g_t - \tilde{g}_t$
  - 7:    $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \tilde{g}_t$
  - 8:    $v_t = \beta_2 v_{t-1} + (1 - \beta_2) \tilde{g}_t^2$
  - 9:    $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$
  - 10:   Update global model  $\theta_t = \theta_{t-1} - \eta_t \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}$
  - 11: **end for**
-

161 Let  $m'_t$  and  $\hat{v}'_t$  be the first and second moment moving average of standard AMSGrad using full  
 162 gradients. Denote

$$a_t = \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}, \quad a'_t = \frac{m'_t}{\sqrt{\hat{v}'_t + \epsilon}}.$$

163 Define the sequence

$$\mathcal{E}_{t+1} = \mathcal{E}_t + a'_t - a_t,$$

164 such that the auxiliary model

$$\begin{aligned} \theta'_{t+1} &:= \theta_{t+1} - \eta \mathcal{E}_{t+1} \\ &= \theta_t - \eta a_t - \eta \mathcal{E}_{t+1} \\ &= \theta_t - \eta a_t - \eta(\mathcal{E}_t + a'_t - a_t) \\ &= \theta'_t - \eta a'_t \end{aligned}$$

165 follows the update of full-gradient AMSGrad. By smoothness assumption we have

$$f(\theta'_{t+1}) \leq f(\theta'_t) - \eta \langle \nabla f(\theta'_t), a'_t \rangle + \frac{L}{2} \|\theta'_{t+1} - \theta'_t\|^2.$$

166 Thus,

$$\begin{aligned} \mathbb{E}[f(\theta'_{t+1}) - f(\theta'_t)] &\leq -\eta \mathbb{E}[\langle \nabla f(\theta'_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] \\ &= -\eta \mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] + \eta \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), a'_t \rangle] \\ &\leq -\eta \mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] + \eta \mathbb{E}[\frac{\eta^2 \rho}{2} \|\mathcal{E}_t\|^2 + \frac{1}{2\rho} \|a'_t\|^2] \\ &\leq -\eta \frac{\mathbb{E}\|\nabla f(\theta_t)\|^2}{\sqrt{G^2 + \epsilon}} + \frac{\eta}{2\rho} \frac{\mathbb{E}\|\nabla f(\theta_t)\|^2}{\epsilon} + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] + \frac{\eta^3 \rho}{2} \mathbb{E}\|\mathcal{E}_t\|^2, \end{aligned}$$

167 when  $\beta_1 = 0$  for example. We may discard this assumption and use more complicated bound on the  
 168 first two terms. The third term can be bounded by constant yielding  $O(1/\sqrt{T})$  rate eventually when  
 169 taking decreasing learning rate. The key is to get a good bound on the cumulative error sequence,  
 170  $\mathcal{E}_t$ . We have the following:

$$\begin{aligned} \mathbb{E}\|\mathcal{E}_{t+1}\|^2 &= \mathbb{E}\|\mathcal{E}_t + a'_t - a_t + \text{TopK}(\mathcal{E}_t + a'_t) - \text{TopK}(\mathcal{E}_t + a'_t)\|^2 \\ &\leq 2\mathbb{E}\|\mathcal{E}_t + a'_t - \text{TopK}(\mathcal{E}_t + a'_t)\|^2 + 2\mathbb{E}\|a_t - \text{TopK}(\mathcal{E}_t + a'_t)\|^2 \\ &\stackrel{(a)}{\leq} 2q\mathbb{E}\|\mathcal{E}_t + a'_t\| + 2\mathbb{E}\|a_t - \text{TopK}(\mathcal{E}_t + a'_t)\|^2 \\ &\leq 2q[(1+r)\mathbb{E}\|\mathcal{E}_t\|^2 + (1+\frac{1}{r})\mathbb{E}\|a'_t\|^2] + 2\mathbb{E}\|a_t - \text{TopK}(\mathcal{E}_t + a'_t)\|^2. \end{aligned}$$

171 where (a) uses A3. Current try: If we can bound the last term in the same form as the first two terms,  
 172 then we can use recursion to get the desired result. We can have

$$\mathbb{E}\|a_t - \text{TopK}(\mathcal{E}_t + a'_t)\|^2 = \mathbb{E}\|\frac{\tilde{m}_t}{\sqrt{\hat{v}_t + \epsilon}} - \|^2$$

## 173 5 Experiments

174 Our proposed TopK-EF with AMSGrad matches that of full AMSGrad, in distributed learning.  
175 Number of local workers is 20. Error feedback fixes the convergence issue of using solely the  
176 TopK gradient.

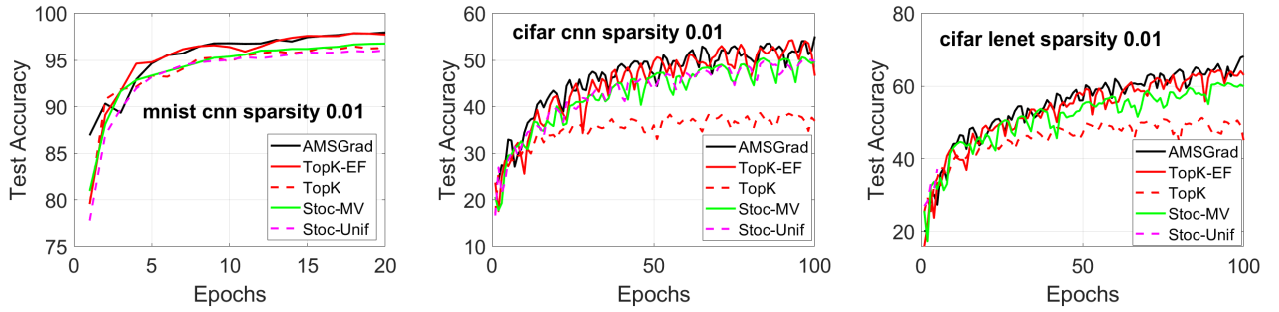


Figure 1: Test accuracy.

## 177 6 Conclusion



## References

- [1] Naman Agarwal, Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and Brendan McMahan. cpsgd: Communication-efficient and differentially-private distributed SGD. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 7575–7586, 2018.
- [2] Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*, 2017.
- [3] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- [4] Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric Renggli. The convergence of sparsified gradient methods. *arXiv preprint arXiv:1809.10505*, 2018.
- [5] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pages 560–569. PMLR, 2018.
- [6] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- [7] Ken Chang, Niranjan Balachandar, Carson K. Lam, Darvin Yi, James M. Brown, Andrew Beers, Bruce R. Rosen, Daniel L. Rubin, and Jayashree Kalpathy-Cramer. Distributed deep learning networks among institutions for medical imaging. *J. Am. Medical Informatics Assoc.*, 25(8):945–954, 2018.
- [8] Congliang Chen, Li Shen, Haozhi Huang, Qi Wu, and Wei Liu. Quantized adam with error feedback. *arXiv preprint arXiv:2004.14180*, 2020.
- [9] Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project adam: Building an efficient and scalable deep learning training system. In *Symposium on Operating Systems Design and Implementation*, pages 571–582, 2014.
- [10] Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19, 2016*, pages 191–198. ACM, 2016.
- [11] Christopher De Sa, Matthew Feldman, Christopher Ré, and Kunle Olukotun. Understanding and optimizing asynchronous low-precision stochastic gradient descent. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*, pages 561–574, 2017.
- [12] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc’Aurelio Ranzato, Andrew W. Senior, Paul A. Tucker, Ke Yang, and Andrew Y. Ng. Large scale distributed deep networks. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*, pages 1232–1240, 2012.
- [13] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- [14] John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. In *COLT 2010 - The 23rd Conference on Learning Theory, Haifa, Israel, June 27-29, 2010*, pages 257–269, 2010.

- [15] Saeed Ghadimi and Guanghai Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [16] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2672–2680, 2014.
- [17] Priya Goyal, Piotr Dollár, Ross B. Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: training imagenet in 1 hour. *CoRR*, abs/1706.02677, 2017.
- [18] Alex Graves, Abdel-rahman Mohamed, and Geoffrey E. Hinton. Speech recognition with deep recurrent neural networks. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2013, Vancouver, BC, Canada, May 26-31, 2013*, pages 6645–6649. IEEE, 2013.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016.
- [20] Mingyi Hong, Davood Hajinezhad, and Ming-Min Zhao. Prox-pda: The proximal primal-dual algorithm for fast distributed nonconvex optimization and learning over networks. In *International Conference on Machine Learning*, pages 1529–1538, 2017.
- [21] Nikita Ivkin, Daniel Rothchild, Enayat Ullah, Vladimir Braverman, Ion Stoica, and Raman Arora. Communication-efficient distributed SGD with sketching. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 13144–13154, 2019.
- [22] Peng Jiang and Gagan Agrawal. 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*, pages 2530–2541, 2018.
- [23] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. *arXiv preprint arXiv:1901.09847*, 2019.
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [25] Anastasia Koloskova, Sebastian U Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. In *International Conference on Machine Learning*, pages 3478–3487, 2019.
- [26] Songtao Lu, Xinwei Zhang, Haoran Sun, and Mingyi Hong. Gnsd: A gradient-tracking based nonconvex stochastic algorithm for decentralized optimization. In *2019 IEEE Data Science Workshop (DSW)*, pages 315–321, 2019.
- [27] Hiroaki Mikami, Hisahiro Suganuma, Yoshiki Tanaka, Yuichi Kageyama, et al. Massively distributed sgd: Imagenet/resnet-50 training in a flash. *arXiv preprint arXiv:1811.05233*, 2018.
- [28] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.
- [29] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. Dadam: A consensus-based distributed adaptive gradient method for online optimization. *arXiv preprint arXiv:1901.09109*, 2019.

- [30] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48, 2009.
- [31] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [32] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014*, pages 1058–1062. ISCA, 2014.
- [33] Shaohuai Shi, Kaiyong Zhao, Qiang Wang, Zhenheng Tang, and Xiaowen Chu. A convergence analysis of distributed sgd with communication-efficient gradient sparsification. In *IJCAI*, pages 3411–3417, 2019.
- [34] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy P. Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nat.*, 550(7676):354–359, 2017.
- [35] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. In *Advances in Neural Information Processing Systems*, pages 4447–4458, 2018.
- [36] Athanasios Voulodimos, Nikolaos Doulamis, Anastasios D. Doulamis, and Eftychios Protopapadakis. Deep learning for computer vision: A brief review. *Comput. Intell. Neurosci.*, 2018:7068349:1–7068349:13, 2018.
- [37] Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems*, pages 1299–1309, 2018.
- [38] Jian Wei, Jianhua He, Kai Chen, Yi Zhou, and Zuoyin Tang. Collaborative filtering and deep learning based recommendation system for cold start items. *Expert Systems with Applications*, 69:29–39, 2017.
- [39] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. *arXiv preprint arXiv:1705.07878*, 2017.
- [40] Jiayang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang. Error compensated quantized SGD and its applications to large-scale distributed optimization. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5321–5329. PMLR, 2018.
- [41] Guandao Yang, Tianyi Zhang, Polina Kirichenko, Junwen Bai, Andrew Gordon Wilson, and Chris De Sa. Swalp: Stochastic weight averaging in low precision training. In *International Conference on Machine Learning*, pages 7015–7024. PMLR, 2019.
- [42] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing [review article]. *IEEE Comput. Intell. Mag.*, 13(3):55–75, 2018.
- [43] Lei Zhang, Shuai Wang, and Bing Liu. Deep learning for sentiment analysis: A survey. *Wiley Interdiscip. Rev. Data Min. Knowl. Discov.*, 8(4), 2018.

## 314 A Appendix

## 315 B Proofs

### 316 B.1 Proof of Lemmas

317 **Lemma.** Under Assumption 2 and Assumption 4 we have for any iteration  $t > 0$ :

$$\|m_t\|^2 \leq (q^2 + 1)G^2 \quad \text{and} \quad \hat{v}_t \leq (q^2 + 1)G^2 \quad (8)$$

318 where  $m_t$  and  $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$  are defined Line 15 of Algorithm 1 and  $G^2 = \frac{1}{n} \sum_{i=1}^N G_i^2$ .

319 *Proof.* We start by writing

$$\|\bar{g}_t\|^2 = \left\| \frac{1}{n} \sum_{i=1}^N \tilde{g}_{t,i} \right\|^2 \leq \frac{1}{n} \sum_{i=1}^N \|\tilde{g}_{t,i}\|^2 \quad (9)$$

320 Though, using Assumption 2 and Assumption 4 we have:

$$\|\tilde{g}_{t,i}\|^2 = \|g_{t,i} + \tilde{g}_{t,i} - g_{t,i}\|^2 \leq \|g_{t,i}\|^2 + \|\tilde{g}_{t,i} - g_{t,i}\|^2 \leq (q^2 + 1)G_i^2 \quad (10)$$

321 Hence

$$\|\bar{g}_t\|^2 \leq (q^2 + 1)G^2 \quad (11)$$

322 where  $G^2 = \frac{1}{n} \sum_{i=1}^N G_i^2$ . Then, by construction in Algorithm 1:

$$\|m_t\|^2 \leq \beta_1^2 \|m_{t-1}\|^2 + (1 - \beta_1)^2 \|\bar{g}_t\|^2 \leq \beta_1^2 \|m_{t-1}\|^2 + (1 - \beta_1)^2 (q^2 + 1)G^2 \quad (12)$$

323 Since we have by initialization that  $\|m_0\|^2 \leq G^2$ , then we prove by induction that  $\|m_t\|^2 \leq (q^2 + 1)G^2$ .

325 Similarly

$$\hat{v}_t = \max(v_t, \hat{v}_{t-1}) = \max(\hat{v}_{t-1}, \beta_2 v_{t-1} + (1 - \beta_2) \bar{g}_t^2) \leq \max(\hat{v}_{t-1}, \beta_2 v_{t-1} + (1 - \beta_2)(q^2 + 1)G^2) \quad (13)$$

326  $\square$

327 **Lemma.** Under A1 to A4, with a decreasing sequence of stepsize  $\{\eta_t\}_{t>0}$ , we have:

$$-\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \right\rangle \right] \leq -\frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2} \quad (14)$$

328 where  $\mathbf{I}_d$  is the identity matrix,  $\hat{V}_t$  the diagonal matrix which diagonal entries are  $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$   
329 defined Line 15 of Algorithm 1 and  $\bar{g}_t$  is the aggregation of all **quantized** gradients from the workers.

330 *Proof.* We first decompose  $\bar{g}_t$  as the sum of the unbiased stochastic gradients and its quantized  
331 versions as computed Line 7 of Algorithm 1:

$$\bar{g}_t = \frac{1}{n} \sum_{i=1}^N \tilde{g}_{t,i} = \frac{1}{n} \sum_{i=1}^N [g_{t,i} + \tilde{g}_{t,i} - g_{t,i}] \quad (15)$$

332 Hence,

$$\begin{aligned} T_1 &:= -\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \right\rangle \right] \\ &= \underbrace{-\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \frac{1}{n} \sum_{i=1}^N g_{t,i} \right\rangle \right]}_{t_1} - \underbrace{\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \frac{1}{n} \sum_{i=1}^N \tilde{g}_{t,i} - g_{t,i} \right\rangle \right]}_{t_2} \end{aligned} \quad (16)$$

333 **Bounding  $t_1$ :** Using the Tower rule, we have:

$$\begin{aligned}
t_1 &:= -\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \frac{1}{n} \sum_{i=1}^N g_{t,i} \right\rangle \right] \\
&= -\eta_{t+1} \mathbb{E} \left[ \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \frac{1}{n} \sum_{i=1}^N g_{t,i} \right\rangle \mid \mathcal{F}_t \right] \right] \\
&= -\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^N g_{t,i} \mid \mathcal{F}_t \right] \right\rangle \right]
\end{aligned} \tag{17}$$

334 Using Assumption 2 and Lemma 1, we have that

$$\begin{aligned}
t_1 &:= -\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \frac{1}{n} \sum_{i=1}^N g_{t,i} \right\rangle \right] \\
&\leq -\eta_{t+1} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2]
\end{aligned} \tag{18}$$

335 **Bounding  $t_2$ :**

336 We first recall Young's inequality with a constant  $\delta \in (0, 1)$  as follows:

$$\langle X \mid Y \rangle \leq \frac{1}{\delta} \|X\|^2 + \delta \|Y\|^2. \tag{19}$$

337 Using Young's inequality (19) with parameter equal to 1:

$$\begin{aligned}
t_2 &\leq \frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + \frac{\eta_{t+1}}{2n^2} \mathbb{E} [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \sum_{i=1}^N \{\tilde{g}_{t,i} - g_{t,i}\}^2] \\
&\stackrel{(a)}{\leq} \frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + \frac{\eta_{t+1}}{2n^2} \mathbb{E} [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2}]^2 \sum_{i=1}^N \{\tilde{g}_{t,i} - g_{t,i}\}^2 \\
&\stackrel{(b)}{\leq} \frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + \frac{\eta_{t+1}}{2n^2} \mathbb{E} [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2}]^2 \mathbb{E} \left[ \sum_{i=1}^N \{\tilde{g}_{t,i} - g_{t,i}\}^2 \right] \\
&\stackrel{(c)}{\leq} \frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + \frac{\eta_{t+1}}{\epsilon 2n^2} \mathbb{E} \left[ \sum_{i=1}^N \tilde{g}_{t,i}^2 \right] \\
&\stackrel{(d)}{\leq} \frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2}
\end{aligned} \tag{20}$$

338 where (a) uses the Cauchy-Schwartz inequality, (b) is due to the non-negativeness of both  $\hat{V}_{t+1}$   
339 and  $\|\sum_{i=1}^N \{g_{t,i} + \tilde{g}_{t,i} - g_{t,i}\}\|^2$  and (c) uses the Triangle inequality. We use Assumption 3 and  
340 Assumption 4 in (d).

341 Finally, combining (18) and (20) yields

$$-\eta_{t+1} \mathbb{E} \left[ \left\langle \nabla f(\theta_t) \mid (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \tilde{g}_t \right\rangle \right] \leq -\frac{\eta_{t+1}}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E} [\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2} \tag{21}$$

342  $\square$

343 **Lemma.** Under A1 to A4, with a decreasing sequence of stepsize  $\{\eta_t\}_{t>0}$ , we have:

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] &\leq -\frac{\eta_{t+1}(1 - \beta_1)}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2 \frac{G^2 \eta_{t+1}}{\epsilon 2n^2} \\
&\quad - \eta_{t+1} \beta_1 \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\
&\quad + \left( \frac{L}{2} + \beta_1 L \right) \|\theta_t - \theta_{t-1}\|^2 \\
&\quad + \eta_{t+1} G^2 \mathbb{E} \left[ \sum_{j=1}^d \left[ (\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2} \right] \right]
\end{aligned} \tag{22}$$

344 where  $d$  denotes the dimension of the parameter vector

345 *Proof.* Denote the following auxiliary variables at iteration  $t + 1$

$$z_{t+1} = \theta_{t+1} + \frac{\beta_1}{1 - \beta_1} (\theta_{t+1} - \theta_t) \tag{23}$$

346 By assumption Assumption 1, we can write the smoothness condition on the overall objective (2),  
347 between iteration  $t$  and  $t + 1$ :

$$f(\theta_{t+1}) \leq f(\theta_t) + \langle \nabla f(\theta_t) | \theta_{t+1} - \theta_t \rangle + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \tag{24}$$

348 Denote by  $\hat{V}_t$  the diagonal matrix which diagonal entries are  $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$  defined Line 15 of  
349 Algorithm 1. Hence, we obtain,

$$f(\theta_{t+1}) \leq f(\theta_t) - \eta_{t+1} \langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \tag{25}$$

350 where  $\mathbf{I}_d$  denotes the identity matrix.

351 We now take the expectation of those various terms conditioned on the filtration  $\mathcal{F}_t$  of the total  
352 randomness up to iteration  $t$ .

$$\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] \leq -\eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] + \frac{L}{2} \mathbb{E}[\|\theta_{t+1} - \theta_t\|^2] \tag{26}$$

353 We now focus on the computation of the inner product obtained in the equation above. We have

$$\begin{aligned}
&\eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] \\
&= \eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} + (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} - (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] \\
&= \eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] + \eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} - (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2}] m_{t+1} \rangle] \\
&= \eta_{t+1} \beta_1 \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] + \eta_{t+1} (1 - \beta_1) \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \rangle] \\
&\quad + \eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} - (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2}] m_{t+1} \rangle]
\end{aligned} \tag{28}$$

354 where  $\bar{g}_t$  is the aggregated gradients from all workers.

355 Plugging the above in (26) yields:

$$\begin{aligned}
& \mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] \\
& \leq \underbrace{-\beta_1 \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle]}_{A_t} \eta_{t+1} \\
& \quad \underbrace{-\mathbb{E}[\langle \nabla f(\theta_t) | [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} - (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2}] m_{t+1} \rangle]}_{B_t} \eta_{t+1} \\
& \quad \underbrace{-(1 - \beta_1) \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \rangle]}_{C_t} \eta_{t+1} + \frac{L}{2} \mathbb{E}[\|\theta_{t+1} - \theta_t\|^2]
\end{aligned} \tag{29}$$

356 To begin with, by the tower rule, we have that

$$A_t = -\beta_1 \mathbb{E}[\mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle | \mathcal{F}_t]] \tag{30}$$

$$= -\beta_1 \langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle - \beta_1 \langle \nabla f(\theta_t) - \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle \tag{31}$$

$$\tag{32}$$

where we recognize the first term as the term in (27), at iteration  $t - 1$  and hence apply the same decomposition as in (28). Coupling with the smoothness of  $f$ , which gives that

$$-\beta_1 \langle \nabla f(\theta_t) - \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle \leq \frac{\beta_1 L}{\eta_{t-1}} \|\theta_t - \theta_{t-1}\|^2$$

357 we obtain,

$$\begin{aligned}
A_t &= -\beta_1 \mathbb{E}[\mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle | \mathcal{F}_t]] \\
&\leq \eta_{t+1} \beta_1 (A_{t-1} + B_{t-1} + C_{t-1}) + \eta_{t+1} \frac{\beta_1 L}{\eta_{t-1}} \|\theta_t - \theta_{t-1}\|^2
\end{aligned} \tag{33}$$

358 Then,

$$\begin{aligned}
B_t &= -\mathbb{E}[\langle \nabla f(\theta_t) | [(\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} - (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2}] m_{t+1} \rangle] \\
&= \mathbb{E}[\sum_{j=1}^d \nabla^j f(\theta_t) m_{t+1}^j [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]] \\
&\stackrel{(a)}{\leq} \mathbb{E}[\|\nabla f(\theta_t)\| \|m_{t+1}\| \sum_{j=1}^d [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]] \\
&\stackrel{(b)}{\leq} G^2 \mathbb{E}[\sum_{j=1}^d [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]]
\end{aligned} \tag{34}$$

359 where  $\nabla^j f(\theta_t)$  denotes the  $j$ -th component of the gradient vector  $\nabla f(\theta_t)$ , (a) uses of the Cauchy-  
360 Schwartz inequality and (b) boils down from the norm of the gradient vector boundedness assump-  
361 tion 2, denoting  $G := \frac{1}{n} \sum_{i=1}^n G_i$ .

362 Plugging the above into (29) yields

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] &\leq \eta_{t+1}(A_t + B_t + C_t) + \frac{L}{2}\mathbb{E}[\|\theta_{t+1} - \theta_t\|^2] \\
&\leq -\eta_{t+1}\beta_1\mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\
&\quad + \eta_{t+1}G^2\mathbb{E}\left[\sum_{j=1}^d [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]\right] \\
&\quad + \left(\frac{L}{2} + \eta_{t+1}\frac{\beta_1 L}{\eta_{t-1}}\right)\|\theta_t - \theta_{t-1}\|^2 \\
&\quad - \eta_{t+1}(1 - \beta_1)\mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \rangle]
\end{aligned} \tag{35}$$

363 We bound the last term on the RHS,  $-\eta_{t+1}\mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} \bar{g}_t \rangle]$  with Lemma 2

364 Under the assumption that we use a decreasing stepsize such that  $\eta_{t+1} \leq \eta_t$ , and given that according  
365 to Line 15 we have that  $\hat{v}_{t+1} \geq \hat{v}_t$  by construction, we obtain

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] &\leq -\frac{\eta_{t+1}(1 - \beta_1)}{2}\left(\epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2}\right)^{-\frac{1}{2}}\mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2\frac{G^2\eta_{t+1}}{\epsilon 2n^2} \\
&\quad - \eta_{t+1}\beta_1\mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\
&\quad + \left(\frac{L}{2} + \beta_1 L\right)\|\theta_t - \theta_{t-1}\|^2 \\
&\quad + \eta_{t+1}G^2\mathbb{E}\left[\sum_{j=1}^d [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]\right]
\end{aligned} \tag{36}$$

366 Finally, using Lemma 2, we obtain the desired result.  $\square$

## 367 B.2 Proof of Theorem 1

368 **Theorem.** Under A1 to A4, with a constant stepsize  $\eta_t = \eta = \frac{L}{\sqrt{T_m}}$ , we have:

$$\frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \frac{\mathbb{E}[f(\theta_0) - f(\theta_{T_m})]}{L\Delta_1\sqrt{T_m}} + d\frac{L\Delta_3}{\Delta_1\sqrt{T_m}} + \frac{\Delta_2}{\eta\Delta_1 T_m} + \frac{1 - \beta_1}{\Delta_1}\epsilon^{-\frac{1}{2}}\sqrt{(q^2 + 1)}G^2 \tag{37}$$

369 where

$$\begin{aligned}
\Delta_1 &:= \frac{(1 - \beta_1)}{2}\left(\epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2}\right)^{-\frac{1}{2}}, \quad \Delta_2 := q^2 + \sum_{k=t+1}^{\infty} \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2} \\
\Delta_3 &:= \left(\frac{L}{2} + 1 + \frac{\beta_1 L}{1 - \beta_1}\right)(1 - \beta_2)^{-1}\left(1 - \frac{\beta_1^2}{\beta_2}\right)^{-1}
\end{aligned} \tag{38}$$

370 *Proof.* By Lemma 3 we have

$$\begin{aligned}
\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] &\leq -\frac{\eta_{t+1}(1 - \beta_1)}{2}\left(\epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2}\right)^{-\frac{1}{2}}\mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2\frac{G^2\eta_{t+1}}{\epsilon 2n^2} \\
&\quad - \eta_{t+1}\beta_1\mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\
&\quad + \left(\frac{L}{2} + \beta_1 L\right)\|\theta_t - \theta_{t-1}\|^2 \\
&\quad + \eta_{t+1}G^2\mathbb{E}\left[\sum_{j=1}^d [(\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2}]\right]
\end{aligned} \tag{39}$$



371 Let us consider the following sequence, defined for all  $t > 0$ :

$$R_t := f(\theta_t) - \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \quad (40)$$

372 We compute the following expectation:

$$\begin{aligned} \mathbb{E}[R_{t+1}] - \mathbb{E}[R_t] &= \mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] - \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] \\ &\quad + \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \end{aligned} \quad (41)$$

373 Using the Assumption 1, we note that:

$$\mathbb{E}[f(\theta_{t+1}) - f(\theta_t)] \leq -\eta_{t+1} \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \quad (42)$$

374 which yields

$$\begin{aligned} \mathbb{E}[R_{t+1}] - \mathbb{E}[R_t] &= -(\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \mathbb{E}[\langle \nabla f(\theta_t) | (\hat{V}_{t+1} + \epsilon \mathbf{I}_d)^{-1/2} m_{t+1} \rangle] \\ &\quad + \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] \\ &\quad + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \\ &\leq (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \mathbb{E}[A_t + B_t + C_t] \\ &\quad - \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \mathbb{E}[A_{t-1} + B_{t-1} + C_{t-1}] \\ &\quad + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \end{aligned} \quad (43)$$

375 where  $A_t, B_t, C_t$  are defined in (29).

376 We use (33) and (34) to bound  $A_t$  and  $B_t$ , and Lemma 2 to bound  $C_t$  where we precise that the  
377 learning rate  $\eta_{t+1}$  becomes  $\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}$ . Hence

$$\begin{aligned} \mathbb{E}[R_{t+1}] - \mathbb{E}[R_t] &\leq \left( (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \beta_1 - \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \right) \mathbb{E}[A_{t-1} + B_{t-1} + C_{t-1}] \\ &\quad + (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) G^2 \mathbb{E} \left[ \sum_{j=1}^d \left[ (\hat{v}_{t+1}^j + \epsilon)^{-1/2} - (\hat{v}_t^j + \epsilon)^{-1/2} \right]^2 \right] \\ &\quad + \left( \frac{L}{2} + (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \frac{\beta_1 L}{\eta_{t-1}} \right) \|\theta_{t+1} - \theta_t\|^2 \\ &\quad - (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \frac{(1 - \beta_1)}{2} \left( \epsilon + \frac{(q^2 + 1)G^2}{1 - \beta_2} \right)^{-\frac{1}{2}} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\ &\quad + q^2 \eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2} \end{aligned} \quad (44)$$

378 where the last term in the LHS is due to Lemma 3.

379 By assumption, we have that for all  $t > 0$ ,  $\eta_{t+1} \leq \eta_t$ . Also, set the tuning parameters such that

$$\eta_t + \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \leq \frac{\eta_t}{1 - \beta_1} \quad (45)$$

380 so that

$$\begin{aligned} & (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \beta_1 - \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} = 0 \\ \iff & (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \beta_1 = \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \end{aligned} \quad (46)$$

381 Note that  $-(\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) \frac{(1-\beta_1)}{2} (\epsilon + \frac{(q^2+1)G^2}{1-\beta_2})^{-\frac{1}{2}} \leq -\eta_{t+1} \frac{(1-\beta_1)}{2} (\epsilon + \frac{(q^2+1)G^2}{1-\beta_2})^{-\frac{1}{2}}$   
 382 since  $\sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2} \geq 0$ .

383 The above coupled with (44) yields

$$\begin{aligned} \mathbb{E}[R_{t+1}] - \mathbb{E}[R_t] & \leq -\eta_{t+1} \frac{(1-\beta_1)}{2} (\epsilon + \frac{(q^2+1)G^2}{1-\beta_2})^{-\frac{1}{2}} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + q^2 \eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2} \\ & \quad - (\eta_{t+1} + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2}) G^2 \mathbb{E}[\sum_{j=1}^d [(\hat{v}_t^j + \epsilon)^{-1/2} - (\hat{v}_{t+1}^j + \epsilon)^{-1/2}]] \\ & \quad + \left( \frac{L}{2} + 1 + \frac{\beta_1 L}{1 - \beta_1} \right) \mathbb{E}[\|\theta_{t+1} - \theta_t\|^2] \end{aligned} \quad (47)$$

384 We now sum from  $t = 0$  to  $t = T_m - 1$  the inequality in (47), and divide it by  $T_m$ :

$$\begin{aligned} & \eta \frac{(1-\beta_1)}{2} (\epsilon + \frac{(q^2+1)G^2}{1-\beta_2})^{-\frac{1}{2}} \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\ & \leq \frac{\mathbb{E}[R_0] - \mathbb{E}[R_{T_m}]}{T_m} + \frac{q^2 \eta + \sum_{k=t+1}^{\infty} \eta_k \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2}}{T_m} \\ & \quad + \left( \frac{L}{2} + 1 + \frac{\beta_1 L}{1 - \beta_1} \right) \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\theta_{t+1} - \theta_t\|^2] \end{aligned} \quad (48)$$

385 where we have used the fact that  $(\hat{v}_t^j + \epsilon)^{-1/2} - (\hat{v}_{t+1}^j + \epsilon)^{-1/2} \geq 0$  for all dimension  $j \in [d]$  by  
 386 construction of  $\hat{v}_{t+1}^j$ .

387 We now bound the two remaining terms:

388 **Bounding**  $-\mathbb{E}[R_{T_m}]$ :

389 By definition (40) of  $R_t$  we have, using Lemma 1:

$$\begin{aligned} -\mathbb{E}[R_{T_m}] & \leq \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \mathbb{E}[\langle \nabla f(\theta_{t-1}) | (\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t \rangle] - f(\theta_{T_m}) \\ & \leq \left\| \sum_{k=t}^{\infty} \eta_k \beta_1^{k-t+1} \right\| \|\nabla f(\theta_{t-1})\| \|(\hat{V}_t + \epsilon \mathbf{I}_d)^{-1/2} m_t\| \\ & \leq \eta_{t+1} (1 - \beta_1) \epsilon^{-\frac{1}{2}} \sqrt{(q^2 + 1)G^2} - f(\theta_{T_m}) \end{aligned} \quad (49)$$

390 **Bounding**  $\sum_{t=0}^{T_m-1} \mathbb{E}[\|\theta_{t+1} - \theta_t\|^2]$ :

391 By definition in Algorithm 1:

$$\|\theta_{t+1} - \theta_t\|^2 = \eta_{t+1}^2 \left[ (\hat{V}_{t+1} + \epsilon I_d)^{-\frac{1}{2}} m_{t+1} \right]^2 = \eta_{t+1}^2 \sum_{j=1}^d \frac{|m_{t+1}^j|^2}{\hat{v}_{t+1}^j + \epsilon} \quad (50)$$

392 For any dimension  $j \in [d]$ ,

$$\begin{aligned} |m_{t+1}^j|^2 &= |\beta_1 m_t^j + (1 - \beta_1) \bar{g}_t^j|^2 \\ &\leq \beta_1 (\beta_1^2 |m_{t-1}^j|^2 + (1 - \beta_1)^2 |\bar{g}_{t-1}^j|^2) + |\bar{g}_t^j|^2 \\ &\leq \sum_{k=0}^t \beta_1^{2(t-k)} |\bar{g}_k^j|^2 \\ &\leq \sum_{k=0}^t \frac{\beta_1^{2(t-k)}}{\beta_2^{t-k}} \beta_2^{t-k} |\bar{g}_k^j|^2 \end{aligned} \quad (51)$$

393 Using Cauchy-Schwartz inequality we obtain

$$\begin{aligned} |m_{t+1}^j|^2 &\leq \sum_{k=0}^t \frac{\beta_1^{2(t-k)}}{\beta_2^{t-k}} \beta_2^{t-k} |\bar{g}_k^j|^2 \leq \sum_{k=0}^t \left( \frac{\beta_1^2}{\beta_2} \right)^{t-k} \sum_{k=0}^t \beta_2^{t-k} |\bar{g}_k^j|^2 \\ &\leq \frac{1}{1 - \frac{\beta_1^2}{\beta_2}} \sum_{k=0}^t \beta_2^{t-k} |\bar{g}_k^j|^2 \end{aligned} \quad (52)$$

394 On the other hand we have

$$\hat{v}_{t+1}^j \geq \beta_2 \hat{v}_t^j + (1 - \beta_2) (\bar{g}_t^j)^2 \quad (53)$$

395 and since it is also true for iteration  $t = 1$ , we have by induction replacing  $v_t^j$  in the above that

$$\hat{v}_{t+1}^j \geq (1 - \beta_2) \sum_{k=0}^t \beta_2^{t-k} |\bar{g}_k^j|^2 \iff \frac{\sum_{k=0}^t \beta_2^{t-k} |\bar{g}_k^j|^2}{\hat{v}_{t+1}^j} \leq (1 - \beta_2)^{-1} \quad (54)$$

396 Hence, we can derive from (50) that

$$\begin{aligned} \|\theta_{t+1} - \theta_t\|^2 &= \eta_{t+1}^2 \sum_{j=1}^d \frac{|m_{t+1}^j|^2}{\hat{v}_{t+1}^j + \epsilon} \leq \eta_{t+1}^2 \sum_{j=1}^d \frac{|m_{t+1}^j|^2}{\hat{v}_{t+1}^j} \\ &\stackrel{(a)}{\leq} \eta_{t+1}^2 \sum_{j=1}^d \frac{1}{1 - \frac{\beta_1^2}{\beta_2}} \frac{\sum_{k=0}^t \beta_2^{t-k} |\bar{g}_k^j|^2}{\hat{v}_{t+1}^j} \\ &\stackrel{(b)}{\leq} \eta_{t+1}^2 d (1 - \beta_2)^{-1} \left(1 - \frac{\beta_1^2}{\beta_2}\right)^{-1} \end{aligned} \quad (55)$$

397 where (a) uses (52) and (b) uses (54).

398 Plugging the two bounds in (48), we obtain the following bound:

$$\begin{aligned} \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] &\leq \frac{\mathbb{E}[f(\theta_0) - f(\theta_{T_m})]}{\eta \Delta_1 T_m} + \frac{q^2 \eta + \sum_{k=t+1}^{\infty} \eta \beta_1^{k-t+2} \frac{G^2}{\epsilon 2 n^2}}{\eta \Delta_1 T_m} \\ &\quad + \frac{1 - \beta_1}{\Delta_1} \epsilon^{-\frac{1}{2}} \sqrt{(q^2 + 1)} G^2 \\ &\quad + \left( \frac{L}{2} + 1 + \frac{\beta_1 L}{1 - \beta_1} \right) \frac{1}{\eta \Delta_1} \eta^2 d (1 - \beta_2)^{-1} \left(1 - \frac{\beta_1^2}{\beta_2}\right)^{-1} \end{aligned} \quad (56)$$

399 where  $\Delta_1 := \frac{(1-\beta_1)}{2}(\epsilon + \frac{(q^2+1)G^2}{1-\beta_2})^{-\frac{1}{2}}$

400 With a constant stepsize  $\eta = \frac{L}{\sqrt{T_m}}$  we get the final convergence bound as follows:

$$\begin{aligned} \frac{1}{T_m} \sum_{t=0}^{T_m-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] &\leq \frac{\mathbb{E}[f(\theta_0) - f(\theta_{T_m})]}{L\Delta_1\sqrt{T_m}} + d \frac{L\Delta_3}{\Delta_1\sqrt{T_m}} \\ &\quad + \frac{\Delta_2}{\eta\Delta_1 T_m} + \frac{1-\beta_1}{\Delta_1} \epsilon^{-\frac{1}{2}} \sqrt{(q^2+1)} G^2 \end{aligned} \tag{57}$$

401 where  $\Delta_2 := q^2 + \sum_{k=t+1}^{\infty} \beta_1^{k-t+2} \frac{G^2}{\epsilon 2n^2}$  and  $\Delta_3 := \left(\frac{L}{2} + 1 + \frac{\beta_1 L}{1-\beta_1}\right) (1-\beta_2)^{-1} (1 - \frac{\beta_1^2}{\beta_2})^{-1}$ .

402 □