# Distributed Adaptive Learning with Gradient Compression

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

This paper presents new algorithms – SPAMS and dist-SPAMS – for tackling single-machine and distributed optimization. Unlike prior works which rely on full gradient communication between the workers and the parameter-server, we design a distributed adaptive optimization method with gradient compression coupled with an error-feedback technique to alleviate the bias introduced by the compression. While the former permits to transmit fewer bits of gradient vectors to the server, we show that using the latter, which correct for the bias, our methods reach a stationary point in $\mathcal{O}(1/\sqrt{T})$ iterations, matching that of state-of-the-art single-machine and distributed methods, without any error-feedback. We illustrate our theoretical results by showing the effectiveness of our method both under the single-machine and distributed settings on various benchmark datasets.

## 1 Introduction

Deep neural network has achieved the state-of-the-art learning performance on numerous AI applications, e.g., computer vision [23, 26, 47], Natural Language Processing [25, 54, 58], Reinforcement Learning [37, 45] and recommendation systems [16, 49]. With the increasing size of both data and deep networks, standard single machine training confronts with at least two major challenges:

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

*Distributed learning* framework [18] has been a common training strategy to tackle the above two issues. For example, in centralized distributed stochastic gradient descent (SGD) protocol, data are located at $n$ local nodes, at which the gradients of the model are computed in parallel. In each iteration, a central server aggregates the local gradients, updates the global model, and transmits back the updated model to the local nodes for subsequent gradient computation. As we can see, this setting naturally solves aforementioned issues: 1) We use $n$ computing nodes to train the model, so the time per training epoch can be largely reduced; 2) There is no need to transmit the local data to central server. Besides, distributed training also provides stronger error tolerance since the training process could continue even one local machine breaks down. As a result of these advantages, there has been a surge of study and applications on distributed systems [10, 39, 20, 24, 27, 35, 33].

Among many optimization strategies, SGD is still the most popular prototype in distributed training for its simplicity and effectiveness [14, 1, 36]. Yet, when the deep learning model is very large, the communication between local nodes and central server could be expensive. Burdensome gradient transmission would slow down the whole training system, or even be impossible because of

the limited bandwidth in some applications. Thus, reducing the communication cost in distributed SGD has become an active topic, and an important ingredient of large-scale distributed systems (e.g. [42]). Solutions based on quantization, sparsification and other compression techniques of the local gradients are proposed, e.g., [4, 50, 48, 46, 3, 7, 17, 52, 28]. As one would expect, in most approaches, there exists a trade-off between compression and learning performance. In general, larger bias and variance of the compressed gradients usually bring more significant performance downgrade in terms of convergence [46, 2]. Interestingly, studies (e.g., [31]) show that the technique of *error feedback* is able to remedy the issue of such biased compressors, achieving same convergence rate as full-gradient SGD.

On the other hand, in recent years, adaptive optimization algorithms (e.g. AdaGrad [21], Adam [32] and AMSGrad [41]) 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 papar, we propose a distributed optimization algorithm with AMSGrad as the backbone, along with Top-$k$ sparsification to reduce the communication cost.

### 1.1 Our contributions

We develop a simple optimization leveraging the adaptivity of AMSGrad, and the computational virtue of TopK sparsification, for tackling a large finite-sum of nonconvex objective functions.

Our technique is shown to be both theoretically and empirically effective under *the classical centralized setting* and *the distributed setting*.

In this contribution,

- We derive a sparsified AMSGrad with error feedback, called SPAMS, with a single machine and provide its decentralized counter part.
- We provide a non-asymptotic convergence rate under each setting,
- We highlight the effectiveness of both methods through several numerical experiments

## 2 Related Work

### 2.1 Distributed SGD with compressed gradients

**Quantization.** As we mentioned before, SGD is the most commonly adopted optimization method in distributed training of deep neural nets. To reduce the expensive communication in large-scale distributed systems, extensive works have considered various compression techniques applied to the gradient transaction procedure. The first strategy is quantization. [19] condenses 32-bit floating numbers into 8-bits when representing the gradients. [42, 7, 31, 8] 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 [4, 51, 57] and LPC-SVRG [55], leveraging unbiased 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 and zeros out the others. Thus, we can more freely choose higher compression ratio (e.g., 1%, 0.1%), still achieving impressive performance in many applications [34]. Stochastic sparsification methods, including uniform sampling and magnitude based sampling [48], select coordinates based on some sampling probability yielding unbiased gradient compressors. Deterministic methods are simpler, e.g., Random-$k$, Top-$k$ [46, 44] (selecting $k$ elements with largest magnitude), Deep Gradient Compression [34], but usually lead to biased gradient estimation. In [28], the central server identifies heavy-hitters from the count-sketch [12] 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 [30, 43, 5, 6, 29], 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 [2, 9]. The technique of *error feedback* is able to "correct for the bias" and fix the problems. 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. [46, 31] 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 [40, 22].

## 2.2 Adaptive optimization

In each SGD update, all the gradient coordinates share a same learning rate, either constant or decreasing over iterations. Adaptive optimization methods cast different learning rate on each dimension. AdaGrad [21] divides the gradient element-wisely by $\sqrt{\sum_{t=1}^{T} g_t^2} \in \mathbb{R}^d$, where $g_t \in \mathbb{R}^d$ is the gradient vector at time $t$ and $d$ is the model dimensionality. Thus, it intrinsically assigns different learning rates to different coordinates throughout the training—elements with smaller previous gradient magnitude tend to move a larger step. AdaGrad has been shown to perform well especially under some sparsity structure. AdaDelta [56] and Adam [32] introduce momentum and moving average of second moment estimation into AdaGrad which lead to better performance. AMSGrad [41] fixes the potential convergence issue of Adam, which will serve as the prototype in this paper. We present the psudocode in Algorithm . In general, adaptive optimization methods are easier to tune in practice, and usually exhibit faster convergence than SGD. Thus, they have been widely used in training deep learning models in language and computer vision applications, e.g., [15, 53, 59]. In distributed setting, the work [38] proposes a decentralized system in online optimization. However, communication efficiency is not considered. The recent work [13] is the most relevant to our paper. Yet, their method is based on Adam, and requires every local node to store a local estimation of first and second moment, thus being less efficient. We will present more detailed comparison in Section 3.

# 3 Communication-Efficient Adaptive Optimization

## 3.1 Gradient Compressors

In this paper, we mainly consider deterministic $q$-deviate compressors defined as below.

**Assumption 1.** *We say a compressor $\mathcal{C} : \mathbb{R}^d \mapsto \mathbb{R}^d$ is q-deviate if for $\forall x \in \mathbb{R}^d$, $\exists\, 0 \leq q < 1$ such that $\|\mathcal{C}(x) - x\| \leq q\,\|x\|$.*

Note that, smaller $q$ indicates better approximation of the true gradient, and $q = 0$ implies no compression, i.e. $\mathcal{C}(x) = x$. We give two popular and highly efficient $q$-deviate compressors that will be compared in this paper.

**Definition 1** (Top-$k$)**.** *For $x \in \mathbb{R}^d$, denote $\mathcal{S}$ as the size-k set of $i \in [d]$ with largest k magnitude $|x_i|$. The **Top-**$k$ compressor is defined as $\mathcal{C}(x)_i = x_i$, if $i \in \mathcal{S}$; $\mathcal{C}(x)_i = 0$ otherwise.*

**Definition 2** (SIGN)**.** *For $x \in \mathbb{R}^d$, define the **SIGN** compressor as $\mathcal{C}(x) = sign(x) \times \frac{1}{d} \sum_{i=1}^{d} |x_i|$.*

**Remark 1.** *Here the scalar, mean magnitude, multiplied to $sign(x)$ ensures $0 \leq q < 1$ as required by Assumption 1, which can be shown by Cauchy-Schwartz inequality. In implementation, this scalar can be arbitrary since we can offset its influence by tuning the learning rate.*

Most modern machine learning tasks can be casted as a large finite-sum optimization problem written as:

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

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

Some related work:

[31] develops variant of signSGD (as a biased compression schemes) for distributed optimization. Contributions are mainly on this error feedback variant. In [44], the authors provide theoretical

results on the convergence of sparse Gradient SGD for distributed optimization (we want that for AMS here). [46] develops a variant of distributed SGD with sparse gradients too. Contributions include a memory term used while compressing the gradient (using top k for instance). Speeding up the convergence in $\frac{1}{T^3}$.

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

## 3.2 SPAMS with Error Feedback

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

---

**Algorithm 1** Distributed SPAMS with error-feedback

1: **Input**: parameter $\beta_1$, $\beta_2$, learning rate $\eta_t$.
2: Initialize: central server parameter $\theta_1 \in \Theta \subseteq \mathbb{R}^d$; $e_{1,i} = 0$ 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 $\tilde{g}_{t,i} = TopK(g_{t,i} + e_{t,i}, k)$
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 the global model $\theta_{t+1} = \theta_t - \eta_t \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}$
17: **end for**

---

# 4 Non-Asymptotic Convergence Analysis for the Single Machine and Decentralized settings

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

Check [13] starting with single machine and extending to distributed settings (several machines).

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

We begin by making the following assumptions.

**Assumption 2.** *(Smoothness) For $i \in [\![n]\!]$, $f_i$ is L-smooth: $\|\nabla f_i(\theta) - \nabla f_i(\vartheta)\| \leq L\|\theta - \vartheta\|$.*

**Assumption 3.** *(Unbiased and Bounded gradient **per worker**) For any iteration index $t > 0$ and worker index $i \in [\![n]\!]$, the stochastic gradient is unbiased and bounded from above: $\mathbb{E}[g_{t,i}] = \nabla f_i(\theta_t)$ and $\|g_{t,i}\| \leq G_i$.*

**Assumption 4.** *(Bounded variance **per worker**) For any iteration index $t > 0$ and worker index $i \in [\![n]\!]$, the variance of the noisy gradient is bounded: $\mathbb{E}[|g_{t,i} - \nabla f_i(\theta_t)|^2] < \sigma_i^2$.*

4

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

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

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

where $n$ denotes the number of workers.

**Decentralized Workers Setting:** The main theorem in the decentralized setting reads:

**Theorem 1.** *Under Assumption 2 to Assumption 4, the sequence of iterates $\{\theta_t\}_{t>0}$ output from Algorithm 1 satisfies:*

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \frac{\mathbb{E}[f(\theta_1) - f(\theta_{T+1})]}{\Delta_1 \eta_t T} + d \frac{\Delta_3}{\Delta_1 \eta_t T} + \frac{\Delta_2}{\Delta_1 T} + \frac{1 - \beta_1}{\Delta_1} \epsilon^{-\frac{1}{2}} \sqrt{(q^2 + 1)\sigma^2} \tag{3}$$

*where $\{\eta_t\}_{t>0}$ is the sequence of stepsizes and:*

$$\Delta_1 := \frac{(1 - \beta_1)}{2} (\epsilon + \frac{(q^2 + 1)\sigma^2}{1 - \beta_2})^{-\frac{1}{2}} \quad , \quad \Delta_2 := q^2 + \frac{G^2}{\epsilon 2 n^2} \overline{\beta}_1$$
$$\Delta_3 := \left( \frac{L}{2} + 1 + \frac{\beta_1 L}{1 - \beta_1} \right) (1 - \beta_2)^{-1} (1 - \frac{\beta_1^2}{\beta_2})^{-1} \tag{4}$$

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

**Corollary 1.** *Under Assumption 2 to Assumption 4, setting the stepsize as $\eta_t = L\sqrt{\frac{n}{T}}$, the sequence of iterates $\{\theta_t\}_{t>0}$ output from Algorithm 1 satisfies:*

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \mathcal{O}(\frac{1}{L\sqrt{nT}} + d\frac{L}{\sqrt{nT}} + \frac{1}{T}),$$

**Single Machine Setting:** We first provide the formulation of our method in the single machine settings in Algorithm 2. Here, the data and the computation are all performed on a single machine.

---

**Algorithm 2** SPAMS with error-feedback for a single machine

1: **Input**: parameter $\beta_1$, $\beta_2$, learning rate $\eta_t$.
2: Initialize: central server parameter $\theta_1 \in \Theta \subseteq \mathbb{R}^d$; $e_1 = 0$ 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 $\tilde{g}_t = TopK(g_t + e_t, k)$
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 the global model $\theta_{t+1} = \theta_t - \eta_t \frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}$
11: **end for**

---

The convergence rate of the vector of parameters estimated via Algorithm 2 is given below:

180 **Theorem 2.** *Under Assumption 2 to Assumption 4, with a decreasing sequence of stepsize*
181 $\{\eta_t\}_{t>0} = \frac{1}{\sqrt{T}}$, *the sequence of iterates* $\{\theta_t\}_{t>0}$ *output from Algorithm 2 satisfies:*

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \mathcal{O}(\frac{1}{\sqrt{T}} + \frac{1}{T}),$$

182 matching the convergence rate of SGD with error feedback [31].

## 5  Experiments

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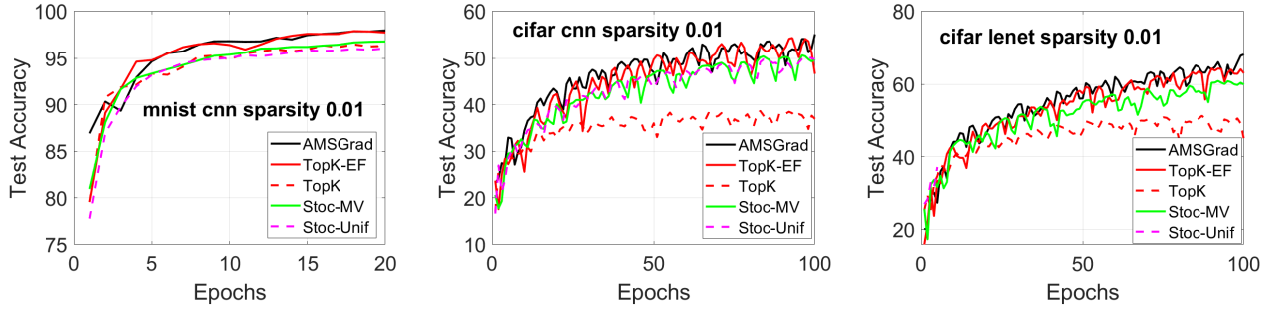


Figure 1: Test accuracy.

## 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] Ahmad Ajalloeian and Sebastian U Stich. Analysis of sgd with biased gradient estimators. *arXiv preprint arXiv:2008.00051*, 2020.

[3] Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*, 2017.

[4] 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.

[5] 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.

[6] Debraj Basu, Deepesh Data, Can Karakus, and Suhas N. Diggavi. Qsparse-local-sgd: Distributed SGD with quantization, sparsification and local computations. 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 14668–14679, 2019.

[7] 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.

[8] Jeremy Bernstein, Jiawei Zhao, Kamyar Azizzadenesheli, and Anima Anandkumar. signsgd with majority vote is communication efficient and fault tolerant. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[9] Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression for distributed learning. *CoRR*, abs/2002.12410, 2020.

[10] 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.

[11] 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.

[12] Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Automata, Languages and Programming, 29th International Colloquium, ICALP 2002, Malaga, Spain, July 8-13, 2002, Proceedings*, volume 2380 of *Lecture Notes in Computer Science*, pages 693–703. Springer, 2002.

[13] Congliang Chen, Li Shen, Haozhi Huang, Qi Wu, and Wei Liu. Quantized adam with error feedback. *arXiv preprint arXiv:2004.14180*, 2020.

[14] 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.

[15] Dami Choi, Christopher J. Shallue, Zachary Nado, Jaehoon Lee, Chris J. Maddison, and George E. Dahl. On empirical comparisons of optimizers for deep learning. *CoRR*, abs/1910.05446, 2019.

[16] 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.

[17] 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.

[18] 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.

[19] Tim Dettmers. 8-bit approximations for parallelism in deep learning. In Yoshua Bengio and Yann LeCun, editors, *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.

[20] 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.

[21] 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.

[22] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.

[23] 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.

[24] 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.

[25] 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.

[26] 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.

[27] 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.

[28] 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.

[29] Jiawei Jiang, Fangcheng Fu, Tong Yang, and Bin Cui. Sketchml: Accelerating distributed machine learning with data sketches. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, pages 1269–1284. ACM, 2018.

[30] 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.

[31] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U Stich, and Martin Jaggi. Error feed-back fixes signsgd and other gradient compression schemes. *arXiv preprint arXiv:1901.09847*, 2019.

[32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[33] Anastasia Koloskova, Sebastian U Stich, and Martin Jaggi. Decentralized stochastic optimiza-tion and gossip algorithms with compressed communication. In *International Conference on Machine Learning*, pages 3478–3487, 2019.

[34] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and Bill Dally. Deep gradient compression: Re-ducing the communication bandwidth for distributed training. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[35] 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.

[36] 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.

[37] 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.

[38] 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.

[39] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent opti-mization. *IEEE Transactions on Automatic Control*, 54(1):48, 2009.

[40] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochas-tic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.

[41] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.

[42] 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, Sin-gapore, September 14-18, 2014*, pages 1058–1062. ISCA, 2014.

[43] Zebang Shen, Aryan Mokhtari, Tengfei Zhou, Peilin Zhao, and Hui Qian. Towards more ef-ficient stochastic decentralized learning: Faster convergence and sparse communication. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stock-holmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 4631–4640. PMLR, 2018.

[44] 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.

[45] 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.

[46] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. In *Advances in Neural Information Processing Systems*, pages 4447–4458, 2018.

[47] 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.

[48] 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.

[49] 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.

[50] 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.

[51] Jiaxiang 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.

[52] 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.

[53] Yang You, Jing Li, Sashank J. Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training BERT in 76 minutes. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[54] 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.

[55] Yue Yu, Jiaxiang Wu, and Junzhou Huang. Exploring fast and communication-efficient algorithms in large-scale distributed networks. In *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pages 674–683. PMLR, 2019.

[56] Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.

[57] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. Zipml: Training linear models with end-to-end low precision, and a little bit of deep learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 4035–4043. PMLR, 2017.

[58] Lei Zhang, Shuai Wang, and Bing Liu. Deep learning for sentiment analysis: A survey. *Wiley Interdiscip. Rev. Data Min. Knowl. Discov.*, 8(4), 2018.

[59] Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q. Weinberger, and Yoav Artzi. Revisiting few-sample BERT fine-tuning. *CoRR*, abs/2006.05987, 2020.

## A  Some Important Notations

For the following proofs, denote

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\tilde{g}_t \quad \text{and} \quad m'_t = \beta_1 m'_{t-1} + (1 - \beta_1)g_t$$

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

## B  Single Machine Setting

### B.1  $\beta_1 = 0$

*Proof.* Denote the following auxiliary sequences,

$$\theta'_t := \theta_t - \eta \frac{e_t}{\sqrt{\hat{v}_{t-1} + \epsilon}},$$

such that

$$\begin{aligned}
\theta'_{t+1} &= \theta_{t+1} - \eta \frac{e_{t+1}}{\sqrt{\hat{v}_t + \epsilon}} \\
&= \theta_t - \eta \frac{\tilde{g}_t + e_{t+1}}{\sqrt{\hat{v}_t + \epsilon}} \\
&= \theta_t - \eta \frac{e_t}{\sqrt{\hat{v}_t + \epsilon}} - \eta \frac{g_t}{\sqrt{\hat{v}_t + \epsilon}} \\
&= \theta'_t - \eta \frac{g_t}{\sqrt{\hat{v}_t + \epsilon}}.
\end{aligned}$$

where (a) uses the fact that $\tilde{g}_t + e_{t+1} = g_t + e_t$. By Assumption 2 we have

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

Taking expectation regarding the randomness at step $t$,

$$\begin{aligned}
\mathbb{E}[f(\theta'_{t+1})] - f(\theta'_t) &\le -\eta\mathbb{E}[\langle \nabla f(\theta'_t), a'_t\rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] \\
&= -\eta\mathbb{E}[\langle \nabla f(\theta_t), a'_t\rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] + \eta\mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), a'_t\rangle]. \quad (5)
\end{aligned}$$

**The first term in (19).** We have

$$\begin{aligned}
M_t &:= -\mathbb{E}[\langle \nabla f(\theta_t), a'_t\rangle] = -\mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_t + \epsilon}}\rangle] \\
&= \underbrace{-\mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_{t-1} + \epsilon}}\rangle]}_{I} + \underbrace{\mathbb{E}[\langle \nabla f(\theta_t), (\frac{1}{\sqrt{\hat{v}_{t-1} + \epsilon}} - \frac{1}{\sqrt{\hat{v}_t + \epsilon}})m'_t\rangle]}_{II}.
\end{aligned}$$

To bound I, note that

$$\begin{aligned}
I &= -\mathbb{E}[\langle \nabla f(\theta_t), \frac{g_t}{\sqrt{\hat{v}_{t-1} + \epsilon}}\rangle] \\
&= -\mathbb{E}\mathbb{E}[\langle \nabla f(\theta_t), \frac{g_t}{\sqrt{\hat{v}_{t-1} + \epsilon}}\rangle | \mathcal{F}_{t-1}] \\
&\le -\frac{1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}\mathbb{E}[\|\nabla f(\theta_t)\|^2], \quad (6)
\end{aligned}$$

13

where the last inequality follows from Lemma 6. Regarding the second term in (19), we have

$$II \leq G^2 \mathbb{E}[\sum_{i=1}^{d} |\frac{1}{\sqrt{\hat{v}_{t-1,i} + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{t,i} + \epsilon}}|].$$

Summing over $t = 1, ..., T$, we obtain

$$\sum_{t=1}^{T} M_t \leq G^2 \mathbb{E}[\sum_{t=1}^{T} \sum_{i=1}^{d} |\frac{1}{\sqrt{\hat{v}_{t-1,i} + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{t,i} + \epsilon}}|] - \frac{1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\leq G^2 \sum_{i=1}^{d} (\frac{1}{\sqrt{\hat{v}_{0,i} + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{T,i} - \epsilon}}) - \frac{1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\leq \frac{G^2 d}{\sqrt{\epsilon}} - \frac{1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

where the second inequality holds by cancelling terms as $\hat{v}_t$ is a non-decreasing sequence.

**Bounding the last two terms in in (19).** For the second term in (19) we have

$$\mathbb{E}[\|a_t'\|^2] = \mathbb{E}[\|\frac{m_t'}{\sqrt{\hat{v}_t + \epsilon}}\|^2] \leq \frac{1}{\epsilon} \mathbb{E}[\|g_t\|^2]$$

$$\leq \frac{1}{\epsilon} (\sigma^2 + \mathbb{E}[\|\nabla f(\theta_t)\|^2]),$$

by Assumption 4. For the last term,

$$\mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), a_t' \rangle] \tag{7}$$

$$= \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), \frac{g_t'}{\sqrt{\hat{v}_t + \epsilon}} \rangle] \tag{8}$$

$$\leq \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), \frac{\nabla f(\theta_t)}{\sqrt{\hat{v}_{t-1} + \epsilon}} \rangle] + \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), \frac{g_t'}{\sqrt{\hat{v}_t + \epsilon}} - \frac{g_t'}{\sqrt{\hat{v}_{t-1} + \epsilon}} \rangle] \tag{9}$$

$$\overset{(a)}{\leq} \frac{\eta^2 L^2 \rho}{2} \mathbb{E}[\|\frac{e_t}{\sqrt{\hat{v}_{t-1} + \epsilon}}\|^2 + \frac{1}{2\rho\epsilon} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \eta L G \mathbb{E}[\|\frac{e_t}{\sqrt{\hat{v}_{t-1} + \epsilon}}\| \|\frac{1}{\sqrt{\hat{v}_t + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{t-1} + \epsilon}}\|] \tag{10}$$

$$\overset{(b)}{\leq} \frac{\eta^2 L^2 \rho}{2} \frac{4q^2}{(1-q^2)^2 \epsilon} + \frac{1}{2\rho\epsilon} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{2\eta L G^2 q}{(1-q^2)\sqrt{\epsilon}} \mathbb{E}[\|\frac{1}{\sqrt{\hat{v}_t + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{t-1} + \epsilon}}\|_1], \tag{11}$$

where (a) uses Young's inequality, Assumption 2 and Assumption 3, and (b) is due to the property that $l_2$ norm is smaller than $l_1$ norm and Lemma 4. Choosing $\rho = \frac{2\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{\epsilon}$, we obtain

$$\mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), a_t' \rangle]$$

$$\leq \frac{4T\eta^2 q^2 L^2 \sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{(1-q^2)^2 \epsilon^2} G^2 + \frac{1}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$+ \frac{2\eta L G^2 q}{(1-q^2)\sqrt{\epsilon}} \mathbb{E}[\|\frac{1}{\sqrt{\hat{v}_t + \epsilon}} - \frac{1}{\sqrt{\hat{v}_{t-1} + \epsilon}}\|_1].$$

398 Summing over $t = 1, .., T$ gives

$$\sum_{t=1}^{T} \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_t'), a_t' \rangle]$$

$$\leq \frac{4T\eta^2 q^2 L^2 \sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{(1-q^2)^2 \epsilon^2} G^2 + \frac{1}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$+ \frac{2\eta L G^2 q}{(1-q^2)\sqrt{\epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\frac{1}{\sqrt{\hat{v}_t} + \epsilon} - \frac{1}{\sqrt{\hat{v}_{t-1}} + \epsilon}\|_1]$$

$$\leq \frac{4T\eta^2 q^2 L^2 \sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{(1-q^2)^2 \epsilon^2} G^2 + \frac{1}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{2\eta L q G^2 d}{(1-q^2)\epsilon}.$$

399 Putting it all together we have

$$\mathbb{E}[f(\theta_{T+1}') - f(\theta_1')]$$

$$\leq \frac{\eta G^2 d}{\sqrt{\epsilon}} - \frac{\eta}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{\eta^2 L}{2\epsilon}(T\sigma^2 + \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2])$$

$$+ \frac{4T\eta^3 q^2 L^2 \sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{(1-q^2)^2 \epsilon^2} G^2 + \frac{\eta}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{2\eta^2 L q G^2 d}{(1-q^2)\epsilon}$$

$$\leq \eta \left[ \frac{\eta L}{2\epsilon} - \frac{3}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}} \right] \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{T\eta^2 L\sigma^2}{2\epsilon} + \frac{4T\eta^3 q^2 L \sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}{(1-q^2)^2 \epsilon^2} G^2$$

$$+ \frac{\eta G^2 d}{\sqrt{\epsilon}} + \frac{2\eta^2 L q G^2 d}{(1-q^2)\epsilon}.$$

400 Setting $\eta \leq \frac{3\epsilon}{2L\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}}$ and re-arranging terms, we arrive at

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq C_1 \frac{\mathbb{E}[f(\theta_1') - f(\theta_{T+1}')]}{T\eta} + \frac{\eta C_1 L\sigma^2}{2\epsilon} + \frac{2\eta^2 C_1^2 q^2 L^2 G^2}{(1-q^2)^2 \epsilon^2} + \frac{C_1 G^2 d}{T\epsilon} + \frac{2\eta C_1 L q G^2 d}{T(1-q^2)\epsilon}$$

$$\leq C_1 \frac{\mathbb{E}[f(\theta_1) - f(\theta^*)]}{T\eta} + \frac{\eta C_1 L\sigma^2}{2\epsilon} + \frac{2\eta^2 C_1^2 q^2 L^2 G^2}{(1-q^2)^2 \epsilon^2} + \frac{C_1 G^2 d}{T\epsilon} + \frac{2\eta C_1 L q G^2 d}{T(1-q^2)\epsilon},$$

401 where $C_1 = 2\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2 + \epsilon}$. The last inequality is because $\theta_1' = \theta_1$, and $\theta^* = \arg\min_\theta f(\theta)$.
402 The proof is complete.

403 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

404 With the learning rate $\sqrt{1/T}$ we get

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \frac{\mathbb{E}[f(\theta_{T+1}') - f(\theta_1')]}{C_1\sqrt{T}} + \frac{L}{\sqrt{T}2\epsilon C_1}\frac{4q^2}{(1-q^2)^2}\sigma^2 + G^2\frac{d}{TC_1\sqrt{\epsilon}}$$

405 T is finite so ok, with a lot of workers we overpass the curse of dimensionality ($d/\sqrt{n}$ term). The
406 $\frac{\sqrt{n}L}{\sqrt{T}2\epsilon C_1}$ term with the variance is problematic.

## B.2 Intermediary Lemmas

**Lemma 1.** *Under Assumption1 to Assumption4 we have:*

$$\mathbb{E}\|m_t'\|^2 \le C\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2],$$

$$\mathbb{E}[\|m_t\|^2] \le (3q^2 + \frac{4q^2(6q^2+3)}{(1-q^2)^2} + 1)C\sigma^2 + (6q^2+3)C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2],$$

*where $C_1 = (1-\beta_1^2)(1 + \frac{1}{4(1-\beta_1^2)})$ and $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}$.*

*Proof.* We have by Young's inequality

$$\mathbb{E}[\|m_t'\|^2] = \mathbb{E}[\|\beta_1 m_{t-1}' + (1-\beta_1)g_t\|^2]$$
$$\le (1+\frac{\rho}{2})\beta_1^2 \mathbb{E}[\|m_{t-1}'\|^2] + (1+\frac{1}{2\rho})(1-\beta_1)^2 \mathbb{E}[\|g_t\|^2].$$

Since $\mathbb{E}[\|g_t\|^2] \le \sigma^2 + \mathbb{E}[\|\nabla f(\theta_t)\|^2]$, by choosing $\rho = 2(1-\beta_1^2)$, we derive

$$\mathbb{E}[\|m_t'\|^2] \le \beta_1^2(2-\beta_1^2)\mathbb{E}[\|m_{t-1}'\|^2] + (1-\beta_1)^2(1 + \frac{1}{4(1-\beta_1^2)})\mathbb{E}[\|g_t\|^2] \tag{12}$$

$$\le \frac{(1-\beta_1)^2}{1-\beta_1^2(2-\beta_1^2)}(1 + \frac{1}{4(1-\beta_1^2)})\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2] \tag{13}$$

$$:= C\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2], \tag{14}$$

due to $\beta_1 < 1$, $m_0' = 0$ and the bounded variance assumption. Here $C_1 = (1-\beta_1^2)(1 + \frac{1}{4(1-\beta_1^2)})$ and $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}$.

For $m_t$ which consists of the compressed stochastic gradients, first note that

$$\mathbb{E}[\|\tilde{g}_t\|^2] = \mathbb{E}[\|\mathcal{C}(g_t + e_t) - (g_t + e_t) + g_t + e_t - \nabla f(\theta_t) + \nabla f(\theta_t)\|^2]$$
$$\le \sigma^2 + 3\mathbb{E}[q^2\|g_t + e_t - \nabla f(\theta_t) + \nabla f(\theta_t)\|^2 + \|e_t\|^2 + \|\nabla f(\theta_t)\|^2]$$
$$\le (3q^2 + 1)\sigma^2 + (6q^2 + 3)\mathbb{E}[\|e_t\|^2 + \|\nabla f(\theta_t)\|^2]$$
$$\le (3q^2 + \frac{4q^2(6q^2+3)}{(1-q^2)^2} + 1)\sigma^2 + (6q^2+3)\mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

where the first inequality is because of Assumption 1 and that the stochastic error $(g_t - \nabla f(\theta_t))$ is mean-zero and independent of other terms. The bound on $\|e_t\|^2$ in the last inequality is due to Lemma 3 of [31]. Then by similar induction we can obtain

$$\mathbb{E}[\|m_t\|^2] \le (3q^2 + \frac{4q^2(6q^2+3)}{(1-q^2)^2} + 1)C\sigma^2 + (6q^2+3)C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2].$$

**Lemma 2.** *Suppose $\gamma = \beta_1/\beta_2 < 1$. Then, for $\forall t$,*

$$\|a_t\|^2 := \|\frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}\|^2 \le \frac{(1-\beta_1)d}{(1-\beta_2)(1-\gamma)}.$$

419 *Proof.* We have

$$
\begin{aligned}
\|\frac{m_t}{\sqrt{\hat{v}_t + \epsilon}}\|^2 &= \sum_{i=1}^{d} \frac{m_{t,i}^2}{\hat{v}_{t,i} + \epsilon} \\
&\leq \frac{(1-\beta_1)^2}{1-\beta_2} \sum_{i=1}^{d} \frac{(\sum_{\tau=1}^{t} \beta_1^{t-\tau} \tilde{g}_{\tau,i})^2}{\sum_{\tau=1}^{t} \beta_2^{t-\tau} \tilde{g}_{\tau,i}^2} \\
&\overset{(a)}{\leq} \frac{(1-\beta_1)^2}{1-\beta_2} \sum_{i=1}^{d} \frac{(\sum_{\tau=1}^{t} \beta_1^{t-\tau})(\sum_{\tau=1}^{t} \beta_1^{t-\tau} \tilde{g}_{\tau,i}^2)}{\sum_{\tau=1}^{t} \beta_2^{t-\tau} \tilde{g}_{\tau,i}^2} \\
&\leq \frac{1-\beta_1}{1-\beta_2} \sum_{i=1}^{d} \frac{\sum_{\tau=1}^{t} \beta_1^{t-\tau} \tilde{g}_{\tau,i}^2}{\sum_{\tau=1}^{t} \beta_2^{t-\tau} \tilde{g}_{\tau,i}^2} \\
&\leq \frac{(1-\beta_1)d}{1-\beta_2} \sum_{\tau=1}^{t} \gamma^{\tau} \\
&\leq \frac{(1-\beta_1)d}{(1-\beta_2)(1-\gamma)},
\end{aligned}
$$

420 where (a) is a consequence of Cauchy-Schwartz inequality. $\qquad\square$

421 **Lemma 3.** *Define*

$$
H_t := \mathbb{E}[\sum_{i=1}^{d} |\frac{1}{\sqrt{\hat{v}_{t-1} + \epsilon}} - \frac{1}{\sqrt{\hat{v}_t + \epsilon}}|]
$$

$$
S_t := \sum_{\tau=1}^{t} (\beta_1^2 (2 - \beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2])
$$

422 *then the following inequalities hold:*

$$
\sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^{\tau} S_{t-\tau} \leq \frac{1}{(1-\beta_1)(1-\beta_1^2(2-\beta_1^2))} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]
$$

$$
\sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^{\tau} H_{t-\tau} \leq \frac{d}{(1-\beta)\sqrt{\epsilon}}.
$$

423 *Proof.* By arranging terms, it holds that

$$
\begin{aligned}
\sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^{\tau} S_{t-\tau} &\leq \sum_{t=2}^{T} (\sum_{\tau=0}^{T-t} \beta_1^{T-t-\tau}) S_t \\
&\leq \frac{1}{1-\beta_1} \sum_{t=2}^{T} \sum_{\tau=1}^{t} (\beta_1^2 (2 - \beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2]) \\
&\leq \frac{1}{1-\beta_1} \sum_{t=1}^{T} (\sum_{\tau=0}^{T-t-1} (\beta_1^2 (2 - \beta_1^2))^{T-t-\tau}) \mathbb{E}[\|\nabla f(\theta_t)\|^2] \\
&\leq \frac{1}{(1-\beta_1)(1-\beta_1^2(2-\beta_1^2))} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2].
\end{aligned}
$$

17

Using similar strategy, we can write

$$\sum_{t=2}^{T}\sum_{\tau=0}^{t-2}\beta_1^\tau H_{t-\tau} \le \sum_{t=2}^{T}(\sum_{\tau=0}^{T-t}\beta_1^{T-t-\tau})H_t$$

$$\le \frac{1}{1-\beta}\sum_{t=2}^{T}\mathbb{E}[\sum_{i=1}^{d}|\frac{1}{\sqrt{\hat{v}_{t-1}+\epsilon}} - \frac{1}{\sqrt{\hat{v}_t+\epsilon}}|$$

$$\le \frac{d}{(1-\beta)\sqrt{\epsilon}},$$

where the last inequality is derived by cancelling terms due to the fact that $\{\hat{v}_t\}_{t>0}$ is a non-decreasing sequence, hence $\hat{v}_t \le \hat{v}_{t-1}$. This completes the proof of the lemma. $\square$

**Lemma 4.** *For the error sequence $e_t$ in* SPAMS*, under Assumption 4, we have for $\forall t$,*

$$\mathbb{E}[\|e_{t+1}\|^2] \le \frac{4q^2}{(1-q^2)^2}\sigma^2 + \frac{2q^2}{1-q^2}\sum_{\tau=1}^{t}(\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2].$$

*Proof.* We start by using Assumption 1 and Young's inequality to get

$$\|e_{t+1}\|^2 = \|g_t + e_t - \mathcal{C}(g_t + e_t)\|^2$$

$$\le q^2\|g_t + e_t\|^2$$

$$\le q^2(1+\rho)\|e_t\|^2 + q^2(1+\frac{1}{\rho})\|g_t\|^2$$

$$\le \frac{1+q^2}{2}\|e_t\|^2 + \frac{2q^2}{1-q^2}\|g_t\|^2,$$

by choosing $\rho = \frac{1-q^2}{2q^2}$. Now by recursion and the initialization $e_1 = 0$, we have

$$\mathbb{E}[\|e_{t+1}\|^2] \le \frac{2q^2}{1-q^2}\sum_{\tau=1}^{t}(\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\|g_\tau\|^2]$$

$$\le \frac{4q^2}{(1-q^2)^2}\sigma^2 + \frac{2q^2}{1-q^2}\sum_{\tau=1}^{t}(\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2],$$

which proves the lemma. Meanwhile, we also have the absolute bound $\|e_t\|^2 \le \frac{4q^2}{(1-q^2)^2}G^2$. $\square$

**Lemma 5.** *For the moving average error sequence $\mathcal{E}_t$, it holds that*

$$\sum_{t=1}^{T}\mathbb{E}[\|\mathcal{E}_t\|^2] \le \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{4q^2}{(1-q^2)^2\epsilon}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2].$$

*Proof.* Denote $K_t := \sum_{\tau=1}^{t}(\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2]$ and $K_0 = 0$. We have

$$\mathbb{E}[\|\mathcal{E}_t\|^2] = \mathbb{E}[\|\frac{(1-\beta_1)\sum_{\tau=1}^{t}\beta_1^{t-\tau}e_\tau}{\sqrt{\hat{v}_t+\epsilon}}\|^2]$$

$$\le \frac{(1-\beta_1)^2}{\epsilon}\sum_{i=1}^{d}\mathbb{E}[(\sum_{\tau=1}^{t}\beta_1^{t-\tau}e_{\tau,i})^2]$$

$$\overset{(a)}{\le} \frac{(1-\beta_1)^2}{\epsilon}\sum_{i=1}^{d}\mathbb{E}[(\sum_{\tau=1}^{t}\beta_1^{t-\tau})(\sum_{\tau=1}^{t}\beta_1^{t-\tau}e_{\tau,i}^2)]$$

$$\le \frac{1-\beta_1}{\epsilon}\sum_{\tau=1}^{t}\beta_1^{t-\tau}\mathbb{E}[\|e_\tau\|^2]$$

$$\overset{(b)}{\le} \frac{4q^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2(1-\beta_1)}{(1-q^2)\epsilon}\sum_{\tau=1}^{t}\beta_1^{t-\tau}K_\tau,$$

18

where (a) is due to Cauchy-Schwartz and (b) is a result of Lemma 4. Summing over $t = 1, ..., T$ and using the similar technique as in Lemma 3 leads to

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathcal{E}_t\|^2] = \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2(1-\beta_1)}{(1-q^2)\epsilon}\sum_{t=1}^{T}\sum_{\tau=1}^{t}\beta_1^{t-\tau}K_\tau$$

$$\leq \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2}{(1-q^2)\epsilon}\sum_{t=1}^{T}\sum_{\tau=1}^{t}(\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2]$$

$$\leq \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{4q^2}{(1-q^2)^2\epsilon}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

which gives the desired result.

$\square$

**Lemma 6.** *It holds that* $\forall t \in [T], \forall i \in [d], \hat{v}_{t,i} \leq \frac{4(1+q^2)^3}{(1-q^2)^2}G^2.$

*Proof.* For any $t$, by Lemma 4 and Assumption 3 we have

$$\|\tilde{g}_t\|^2 = \|\mathcal{C}(g_t + e_t)\|^2$$

$$\leq \|\mathcal{C}(g_t + e_t) - (g_t + e_t) + (g_t + e_t)\|^2$$

$$\leq 2(q^2 + 1)\|g_t + e_t\|^2$$

$$\leq 4(q^2 + 1)(G^2 + \frac{4q^2}{(1-q^2)^2}G^2)$$

$$= \frac{4(1+q^2)^3}{(1-q^2)^2}G^2.$$

It's then easy to show by the updating rule of $\hat{v}_t$,

$$\hat{v}_{t,i} = (1 - \beta_2)\sum_{\tau=1}^{t}\tilde{g}_{t,i}^2 \leq \frac{4(1+q^2)^3}{(1-q^2)^2}G^2.$$

$\square$

## B.3 Proof of Theorem 3

**Theorem 3.** *Denote* $C' = \frac{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}{1-\beta_1}$, $C = \frac{(1-\beta_1)^2}{1-\beta_1^2(2-\beta_1)^2}(1 + \frac{1}{4(1-\beta_1^2)})$, *and* $\gamma = \beta_1/\beta_2 <$ 1. *Under Assumption 1 to Assumption 4, with* $\eta_t = \eta \leq \min\{\frac{1-\beta_1}{C}, \frac{(1-q^2)^2}{2q^2}\}\frac{(1-\beta_1)\epsilon}{4L\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}$,

SPAMS *satisfies*

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq C'\Big(\frac{\mathbb{E}[f(\theta_1) - f(\theta^*)]}{T\eta} + \frac{2dG^2}{T(1-\beta_1)\sqrt{\epsilon}} + \frac{\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon}$$

$$+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2\eta Lq^2\sigma^2}{(1-q^2)^2\epsilon}\Big).$$

*Proof.* Let $m_t'$ be the first moment moving average of standard AMSGrad using full gradients, *i.e.,* the gradient with respect to the index data point $i_t$ computed Line 4 of Algorithm 2 before applying any compression operator.

Denote

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\tilde{g}_t \quad \text{and} \quad m_t' = \beta_1 m_{t-1}' + (1 - \beta_1)g_t$$

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

19

449    By construction we have $m'_t = (1 - \beta_1) \sum_{i=1}^k \beta_1^{t-i} g_t$.

450    Denote the following auxiliary sequences,

$$\mathcal{E}_{t+1} := \frac{(1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t+1-\tau} e_\tau}{\sqrt{\hat{v}_t} + \epsilon}$$

$$\theta'_{t+1} := \theta_{t+1} - \eta \mathcal{E}_{t+1}.$$

451    Then,

$$
\begin{aligned}
\theta'_{t+1} &= \theta_{t+1} - \eta \mathcal{E}_{t+1} \\
&= \theta_t - \eta \frac{(1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} \tilde{g}_\tau + (1 - \beta_1) \sum_{\tau=1}^{t+1} \beta_1^{t+1-\tau} e_\tau}{\sqrt{\hat{v}_t} + \epsilon} \\
&= \theta_t - \eta \frac{(1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} (\tilde{g}_\tau + e_{\tau+1}) + (1 - \beta) \beta_1^t e_1}{\sqrt{\hat{v}_t} + \epsilon} \\
&= \theta_t - \eta \frac{(1 - \beta_1) \sum_{\tau=1}^t \beta_1^{t-\tau} e_\tau}{\sqrt{\hat{v}_t} + \epsilon} - \eta \frac{m'_t}{\sqrt{\hat{v}_t} + \epsilon} \\
&\overset{(a)}{=} \theta'_t - \eta \frac{m'_t}{\sqrt{\hat{v}_t} + \epsilon} := \theta'_t - \eta a'_t,
\end{aligned}
$$

452    where (a) uses the fact that $\tilde{g}_t + e_{t+1} = g_t + e_t$, $e_1 = 0$ at initialization. By Assumption 2 we have

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

453    Thus,

$$
\begin{aligned}
\mathbb{E}[f(\theta'_{t+1}) - f(\theta'_t)] &\le -\eta \mathbb{E}[\langle \nabla f(\theta'_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] \\
&= -\eta \mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] + \eta \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), a'_t \rangle] \\
&\le -\eta \mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \frac{\eta^2 L}{2} \mathbb{E}[\|a'_t\|^2] + \eta^2 L \mathbb{E}[\|\mathcal{E}_t\| \|a'_t\|] \\
&\le -\eta \mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \eta^2 L \mathbb{E}[\|a'_t\|^2] + \frac{\eta^2 L}{2} \mathbb{E}[\|\mathcal{E}_t\|^2]. \qquad (15)
\end{aligned}
$$

454    **Bounding the first term in (21).** We have

$$
\begin{aligned}
M_t := -\mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] &= -\mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_t} + \epsilon} \rangle] \\
&= \underbrace{-\mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle]}_{I} + \underbrace{\mathbb{E}[\langle \nabla f(\theta_t), (\frac{1}{\sqrt{\hat{v}_{t-1}} + \epsilon} - \frac{1}{\sqrt{\hat{v}_t} + \epsilon}) m'_t \rangle]}_{II}.
\end{aligned}
$$

455    To bound I, note that

$$
\begin{aligned}
I &= -\mathbb{E}[\langle \nabla f(\theta_t), \frac{(1 - \beta_1) g_t}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle] - \mathbb{E}[\langle \nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle] \\
&= -\mathbb{E}\mathbb{E}[\langle \nabla f(\theta_t), \frac{(1 - \beta_1) g_t}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle | \mathcal{F}_{t-1}] - \mathbb{E}[\langle \nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle] \\
&= -(1 - \beta_1) \mathbb{E}[\frac{\|\nabla f(\theta_t)\|^2}{\sqrt{\hat{v}_{t-1}} + \epsilon}] - \mathbb{E}[\langle \nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle] \\
&\le -\frac{1 - \beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2} G^2} + \epsilon} \mathbb{E}[\|\nabla f(\theta_t)\|^2] - \beta_1 \mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle], \qquad (16)
\end{aligned}
$$

where the last inequality follows from Lemma 6. Regarding the second term in (16), we have

$$
-\mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}+\epsilon}}\rangle]
$$

$$
= -\mathbb{E}[\langle \nabla f(\theta_{t-1}), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}+\epsilon}}\rangle] - \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_{t-1}), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}+\epsilon}}\rangle]
$$

$$
= M_{t-1} + \eta L \mathbb{E}[\|\frac{m_{t-1}}{\sqrt{\hat{v}_{t-1}+\epsilon}}\|\|\frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}+\epsilon}}\|]
$$

$$
\leq M_{t-1} + \frac{\eta L}{\epsilon}\mathbb{E}[\|m'_{t-1}\|^2] + \eta L \mathbb{E}[\|a_{t-1}\|^2] \tag{17}
$$

$$
\leq M_{t-1} + \frac{\eta L}{\epsilon}(C\sigma^2 + C_1 \sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2]) + \frac{\eta L(1-\beta_1)d}{(1-\beta_2)(1-\gamma)}, \tag{18}
$$

where Lemma 1 and Lemma 2 are used, with $C_1 = (1-\beta_1^2)(1+\frac{1}{4(1-\beta_1^2)})$ and $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}$. Putting parts together we obtain

$$
I \leq \beta_1 M_{t-1} + \frac{\eta\beta_1 LC\sigma^2}{\epsilon} + \frac{\eta\beta_1 LC_1}{\epsilon}\sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2])
$$

$$
+ \frac{\eta L\beta_1(1-\beta_1)d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}\mathbb{E}[\|\nabla f(\theta_t)\|^2].
$$

For II, it holds that

$$
II \leq G^2 \mathbb{E}[\sum_{i=1}^{d}|\frac{1}{\sqrt{\hat{v}_{t-1}+\epsilon}} - \frac{1}{\sqrt{\hat{v}_t+\epsilon}}|].
$$

Denoting $H_t := \mathbb{E}[\sum_{i=1}^{d}|\frac{1}{\sqrt{\hat{v}_{t-1}+\epsilon}} - \frac{1}{\sqrt{\hat{v}_t+\epsilon}}|]$, $S_t := \sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2])$. We arrive at

$$
M_t \leq \beta_1 M_{t-1} + \frac{\eta\beta_1 LC\sigma^2}{\epsilon} + \frac{\eta\beta_1 LC_1}{\epsilon}S_t + G^2 H_t
$$

$$
+ \frac{\eta L\beta_1(1-\beta_1)d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}\mathbb{E}[\|\nabla f(\theta_t)\|^2]
$$

$$
\leq \beta_1 M_{t-1} + \frac{\eta\beta_1 LC\sigma^2}{\epsilon} + \frac{\eta\beta_1 LC_1}{\epsilon}S_t + G^2 H_t + \frac{\eta L\beta_1(1-\beta_1)d}{(1-\beta_2)(1-\gamma)}.
$$

By induction, we have

$$
M_t \leq \beta_1^{t-1}M_1 + G^2\sum_{\tau=0}^{t-2}\beta_1^\tau H_{t-\tau} + \frac{\eta\beta_1 LC_1}{\epsilon}\sum_{\tau=0}^{t-2}\beta_1^\tau S_{t-\tau} + \frac{\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon}
$$

$$
+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}\mathbb{E}[\|\nabla f(\theta_t)\|^2],
$$

21

since $\beta_1 < 1$. Summing over $t = 1, ..., T$, we obtain

$$\sum_{t=1}^{T} M_t \leq \sum_{t=1}^{T} \beta_1^{t-1} M_1 + G^2 \sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^\tau H_{t-\tau} + \frac{\eta \beta_1 L C_1}{\epsilon} \sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^\tau S_{t-\tau}$$

$$+ \frac{T\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\overset{(a)}{\leq} \frac{2dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)}$$

$$+ \left[ \frac{\eta LC}{(1-\beta_1)\epsilon} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \right] \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\leq \frac{2dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{3(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

when $\eta$ is chosen to be $\eta \leq \frac{(1-\beta_1)^2\epsilon}{4LC\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}$. Here, (a) is due to $M_1 = \mathbb{E}[\langle \nabla f(\theta_1), a_0' \rangle] \leq$

$\beta_1 dG^2/\sqrt{\epsilon}$ and Lemma 3. It remains to bound the last two terms in (21).

**Bounding the last two terms in in (21).** We have

$$\mathbb{E}[\|a_t'\|^2] = \mathbb{E}[\|\frac{m_t'}{\sqrt{\hat{v}_t} + \epsilon}\|^2] \leq \frac{1}{\epsilon}\mathbb{E}[\|m_t'\|^2].$$

By Lemma 1, it follows that

$$\mathbb{E}[\|a_t'\|^2] \leq \frac{1}{\epsilon}(C\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2]).$$

Summing over $t = 1, ..., T$, we obtain

$$\sum_{t=1}^{T} \|a_t'\|^2 \leq \frac{TC\sigma^2}{\epsilon} + \frac{C}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]),$$

where the last inequality can be derived similar to Lemma 3.

For the last term in (21), we have by Lemma 5

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathcal{E}_t\|^2] \leq \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{4q^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2].$$

22

471 **Completing the proof.** Summing (21) over $t = 1, ..., T$ and integrating things together, we have

$$\mathbb{E}[f(\theta'_{T+1}) - f(\theta'_1)]$$

$$\leq \eta \sum_{t=1}^{T} M_t + \frac{T\eta^2 CL\sigma^2}{\epsilon} + \frac{C\eta^2 L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2])$$

$$+ \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon} + \frac{2\eta^2 Lq^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\leq \frac{2\eta dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta^2\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta^2 L\beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{3\eta(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$+ \frac{T\eta^2 CL\sigma^2}{\epsilon} + \Big[\frac{C\eta^2 L}{\epsilon} + \frac{2\eta^2 Lq^2}{(1-q^2)^2\epsilon}\Big] \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]) + \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon}$$

$$\leq - \frac{\eta(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{2\eta dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta^2\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon}$$

$$+ \frac{T\eta^2 L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon},$$

472 when $\eta \leq \frac{(1-q^2)^2(1-\beta_1)\epsilon}{8Lq^2\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}$, where the last line is because $C\eta L \leq \frac{(1-\beta_1)\epsilon}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}$ also holds.

473 Re-arranging terms, we get that when $\eta \leq \min\{\frac{1-\beta_1}{C}, \frac{(1-q^2)^2}{2q^2}\}\frac{(1-\beta_1)\epsilon}{4L\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}$,

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq C'\Big(\frac{\mathbb{E}[f(\theta'_1) - f(\theta'_{T+1})]}{T\eta} + \frac{2dG^2}{T(1-\beta_1)\sqrt{\epsilon}} + \frac{\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon}$$

$$+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2\eta Lq^2\sigma^2}{(1-q^2)^2\epsilon}\Big)$$

$$\leq C'\Big(\frac{\mathbb{E}[f(\theta_1) - f(\theta^*)]}{T\eta} + \frac{2dG^2}{T(1-\beta_1)\sqrt{\epsilon}} + \frac{\eta\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon}$$

$$+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2\eta Lq^2\sigma^2}{(1-q^2)^2\epsilon}\Big).$$

474 where $C' = \frac{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}{1-\beta_1}$, and $C = \frac{(1-\beta_1)^2}{1-\beta_1^2(2-\beta_1)^2}(1 + \frac{1}{4(1-\beta_1^2)})$. The last inequality is because
475 $\theta'_1 = \theta_1$, and $\theta^* = \arg\min_\theta f(\theta)$. The proof is complete.

476 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

477 **Corollary 2.** *Under the setting in Theorem 3, if the learning rate is chosen to be $\eta \leq$*
478 $\min\{\min\{\frac{1-\beta_1}{C}, \frac{(1-q^2)^2}{2q^2}\}\frac{(1-\beta_1)\epsilon}{4L\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2+\epsilon}}, \frac{1}{\sqrt{T}}\}$, *then the convergence rate of* SPAMS *admits*

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq \mathcal{O}(\frac{1}{\sqrt{T}} + \frac{1}{T}).$$

## C Distributed setting Xiaoyun

### C.1 $\beta_1 = 0$

*Proof.* Using the smoothness assumption we have

$$\mathbb{E}[f(\theta'_{t+1})] - f(\theta'_t) \leq -\eta\mathbb{E}[\langle \nabla f(\theta'_t), a'_t \rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2]$$
$$= -\eta\mathbb{E}[\langle \nabla f(\theta_t), a'_t \rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] + \eta\mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta'_t), a'_t \rangle]. \quad (19)$$

$\square$

### C.2 $\beta_1 \neq 0$

**Assumption 5.** *The true gradient deviation is bounded by* $\frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_t) - \nabla f(\theta_t)\|^2 \leq \sigma_g^2, \forall t.$

**Lemma 7.** *For the distributed* SPAMS *with* $n$ *local workers, we have*

$$\mathbb{E}\|\bar{m}'_t\|^2 \leq \frac{C\sigma^2}{n} + C_1\sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2],$$

$$\mathbb{E}[\|\bar{m}_t\|^2] \leq \frac{C\sigma^2}{n} + (3q^2 + \frac{4q^2(6q^2+3)}{(1-q^2)^2})C\sigma^2 + (6q^2+3)C_1\sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2],$$

*where* $C_1 = (1-\beta_1^2)(1 + \frac{1}{4(1-\beta_1^2)})$ *and* $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}.$

*Proof.* First we investigate the variance of average gradients. It holds that

$$\mathbb{E}[\|\bar{g}_t\|^2] = \mathbb{E}\left[\|\frac{1}{n}\sum_{i=1}^{n}g_{t,i}\|^2\right]$$
$$= \frac{1}{n^2}\mathbb{E}\left[\|\sum_{i=1}^{n}(g_{t,i} - \nabla f_i(\theta_t) + \nabla f_i(\theta_t))\|^2\right]$$
$$\leq \frac{\sigma^2}{n} + \left\|\frac{1}{n}\sum_{i=1}^{n}\nabla f_i(\theta_t)\right\|^2 = \frac{\sigma^2}{n} + \|\nabla f(\theta_t)\|^2,$$

as $g_{t,i} - \nabla f_i(\theta_t), i \in [n]$ are mean-zero and independent random variables. Analogous to Lemma 1, we have

$$\mathbb{E}[\|m'_t\|^2] \leq \frac{C\sigma^2}{n} + C_1\sum_{\tau=1}^{t}(\beta_1^2(2-\beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2], \quad (20)$$

with $C_1 = (1-\beta_1^2)(1 + \frac{1}{4(1-\beta_1^2)})$ and $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}.$

24

For $\bar{m}_t$, the first moment sequence based on averaged compressed stochastic gradients, the following bound holds

$$\mathbb{E}[\|\bar{\tilde{g}}_t\|^2] = \mathbb{E}[\|\frac{1}{n}\sum_{i=1}^{n}\mathcal{C}(g_{t,i} + e_{t,i})\|^2]$$

$$= \mathbb{E}[\|\frac{1}{n}\sum_{t=1}^{N}\big(\mathcal{C}(g_{t,i} + e_{t,i}) - (g_{t,i} + e_{t,i}) + g_{t,i} + e_{t,i} - \nabla f_i(\theta_t) + \nabla f_i(\theta_t)\big)\|^2]$$

$$\leq \frac{\sigma^2}{n} + \frac{1}{n^2}\mathbb{E}[\|\sum_{t=1}^{N}(\mathcal{C}(g_{t,i} + e_{t,i}) - (g_{t,i} + e_{t,i})) + \sum_{t=1}^{N}e_{t,i} + \sum_{t=1}^{N}\nabla f_i(\theta_t)\|^2]$$

$$\leq \frac{\sigma^2}{n} + \frac{3}{n}\sum_{i=1}^{n}\mathbb{E}[q^2\|g_{t,i} + e_{t,i}\|^2 + \|e_{t,i}\|^2] + 3\|\frac{1}{n}\sum_{i=1}^{n}\nabla f_i(\theta_t)\|^2$$

$$\leq \frac{\sigma^2}{n} + (3q^2 + \frac{4q^2(6q^2 + 3)}{(1 - q^2)^2})\sigma^2 + (6q^2 + 3)\mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

where the first inequality is because of Assumption 1 and that the stochastic error $(g_t - \nabla f(\theta_t))$ is mean-zero and independent of other terms. The bound on $\|e_t\|^2$ in the last inequality is due to Lemma 3 of [31]. Then by similar induction we can obtain

$$\mathbb{E}[\|m_t\|^2] \leq \frac{C\sigma^2}{n} + (3q^2 + \frac{4q^2(6q^2 + 3)}{(1 - q^2)^2})C\sigma^2 + (6q^2 + 3)C_1\sum_{\tau=1}^{t}(\beta_1^2(2 - \beta_1^2))^{t-\tau}\mathbb{E}[\|\nabla f(\theta_\tau)\|^2].$$

$\square$

**Lemma 8.** *For the averaged error sequence $\bar{e}_t$ in distributed* SPAMS, *under Assumption 4, for $\forall t$,*

$$\mathbb{E}[\|\bar{e}_{t+1}\|^2] \leq \frac{4q^2}{(1 - q^2)^2}\sigma^2 + \frac{2q^2}{1 - q^2}\sum_{\tau=1}^{t}(\frac{1 + q^2}{2})^{t-\tau}\mathbb{E}[\frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_\tau)\|^2].$$

*Proof.* We have

$$\mathbb{E}[\|\bar{e}_{t+1}\|^2] = \mathbb{E}[\|\frac{1}{n}\sum_{i=1}^{n}e_{t,i}\|^2]$$

$$\leq \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}[\|e_{t,i}\|^2]$$

$$\leq \frac{4q^2}{(1 - q^2)^2}\sigma^2 + \frac{2q^2}{1 - q^2}\sum_{\tau=1}^{t}(\frac{1 + q^2}{2})^{t-\tau}\mathbb{E}[\frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_\tau)\|^2],$$

where we use Lemma 4 for each local worker. $\square$

**Lemma 9.** *For the moving average error sequence $\bar{\mathcal{E}}_t$ averaged over all local workers, we have*

$$\sum_{t=1}^{T}\mathbb{E}[\|\bar{\mathcal{E}}_t\|^2] \leq \frac{4Tq^2}{(1 - q^2)^2\epsilon}(\sigma^2 + \sigma_g^2) + \frac{4q^2}{(1 - q^2)^2\epsilon}\sum_{t=1}^{T}\mathbb{E}[\|\frac{1}{n}\sum_{i=1}^{n}\nabla f_i(\theta_t)\|^2],$$

*Proof.* The proof is similar to Lemma 5. Denote $K_t := \sum_{\tau=1}^{t} (\frac{1+q^2}{2})^{t-\tau} \mathbb{E}[\frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(\theta_\tau)\|^2]$ and $K_0 = 0$. We have

$$
\begin{aligned}
\mathbb{E}[\|\bar{\mathcal{E}}_t\|^2] &= \mathbb{E}[\|\frac{(1-\beta_1)\sum_{\tau=1}^{t}\beta_1^{t-\tau}\bar{e}_\tau}{\sqrt{\hat{v}_t} + \epsilon}\|^2] \\
&\leq \frac{(1-\beta_1)^2}{\epsilon} \sum_{i=1}^{d} \mathbb{E}[(\sum_{\tau=1}^{t} \beta_1^{t-\tau}\bar{e}_{\tau,i})^2] \\
&\overset{(a)}{\leq} \frac{(1-\beta_1)^2}{\epsilon} \sum_{i=1}^{d} \mathbb{E}[(\sum_{\tau=1}^{t} \beta_1^{t-\tau})(\sum_{\tau=1}^{t} \beta_1^{t-\tau}\bar{e}_{\tau,i}^2)] \\
&\leq \frac{1-\beta_1}{\epsilon} \sum_{\tau=1}^{t} \beta_1^{t-\tau}\mathbb{E}[\|\bar{e}_\tau\|^2] \\
&\overset{(b)}{\leq} \frac{4q^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2(1-\beta_1)}{(1-q^2)\epsilon} \sum_{\tau=1}^{t} \beta_1^{t-\tau}K_\tau,
\end{aligned}
$$

where (a) is due to Cauchy-Schwartz and (b) is a result of Lemma 8. Summing over $t = 1, ..., T$ and using the similar technique as in Lemma 3 leads to

$$
\begin{aligned}
\sum_{t=1}^{T} \mathbb{E}[\|\bar{\mathcal{E}}_t\|^2] &= \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2(1-\beta_1)}{(1-q^2)\epsilon} \sum_{t=1}^{T}\sum_{\tau=1}^{t} \beta_1^{t-\tau}K_\tau \\
&\leq \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{2q^2}{(1-q^2)\epsilon} \sum_{t=1}^{T}\sum_{\tau=1}^{t} (\frac{1+q^2}{2})^{t-\tau}\mathbb{E}[\frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_\tau)\|^2] \\
&\leq \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{4q^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_t)\|^2] \\
&= \frac{4Tq^2}{(1-q^2)^2\epsilon}\sigma^2 + \frac{4q^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\frac{1}{n}\sum_{i=1}^{n}\nabla f_i(\theta_t)\|^2 + \frac{1}{n}\sum_{i=1}^{n}\|\nabla f_i(\theta_t) - \nabla f(\theta_t)\|^2] \\
&\leq \frac{4Tq^2}{(1-q^2)^2\epsilon}(\sigma^2 + \sigma_g^2) + \frac{4q^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\frac{1}{n}\sum_{i=1}^{n}\nabla f_i(\theta_t)\|^2],
\end{aligned}
$$

where the last two lines hold because of variance decomposition and Assumption 5.

$\square$

Denote the average gradient as $\bar{g}_t = \frac{1}{n}\sum_{i=1}^{n}\tilde{g}_{t,i}$, and $\bar{g}'_t = \frac{1}{n}\sum_{i=1}^{n}g_{t,i}$ be the average of true (uncompressed) local gradients. With a little change of notation, we denote $\bar{m}_0 = \bar{m}'_0 = 0$, and

$$\bar{m}_t = \beta_1\bar{m}_{t-1} + (1-\beta_1)\bar{g}_t \quad \text{and} \quad \bar{m}'_t = \beta_1\bar{m}'_{t-1} + (1-\beta_1)\bar{g}'_t$$

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

By construction we have $m'_t = (1-\beta_1)\sum_{i=1}^{k}\beta_1^{t-i}\bar{g}_t$.

Let $\bar{e}_t = \frac{1}{n}\sum_{i=1}^{n}e_{t,i}$. Denote the following auxiliary sequences,

$$
\begin{aligned}
\bar{\mathcal{E}}_{t+1} &:= \frac{(1-\beta_1)\sum_{i=1}^{t+1}\beta_1^{t+1-i}\bar{e}_i}{\sqrt{\hat{v}_t} + \epsilon} \\
\theta'_{t+1} &:= \theta_{t+1} - \eta\mathcal{E}_{t+1}.
\end{aligned}
$$

26

Then,

$$
\begin{aligned}
\theta'_{t+1} &= \theta_{t+1} - \eta\bar{\mathcal{E}}_{t+1} \\
&= \theta_t - \eta\frac{(1-\beta_1)\sum_{i=1}^t \beta_1^{t-i}\bar{g}_i + (1-\beta_1)\sum_{i=1}^{t+1}\beta_1^{t+1-i}\bar{e}_i}{\sqrt{\hat{v}_t}+\epsilon} \\
&= \theta_t - \eta\frac{(1-\beta_1)\sum_{i=1}^t \beta_1^{t-i}(\bar{g}_i + \bar{e}_{i+1}) + (1-\beta)\beta_1^t\bar{e}_1}{\sqrt{\hat{v}_t}+\epsilon} \\
&= \theta_t - \eta\frac{(1-\beta_1)\sum_{i=1}^t \beta_1^{t-i}\bar{e}_i}{\sqrt{\hat{v}_t}+\epsilon} - \eta\frac{\bar{m}'_t}{\sqrt{\hat{v}_t}+\epsilon} \\
&\stackrel{(a)}{=} \theta'_t - \eta\frac{\bar{m}'_t}{\sqrt{\hat{v}_t}+\epsilon} := \theta'_t - \eta a'_t,
\end{aligned}
$$

512 where (a) uses the fact that $\tilde{g}_{t,i} + e_{t+1,i} = g_{t,i} + e_{t,i}$ for $\forall i \in [N]$. By Assumption 2 we have

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

513 Thus,

$$
\begin{aligned}
\mathbb{E}[f(\theta'_{t+1}) - f(\theta'_t)] &\le -\eta\mathbb{E}[\langle\nabla f(\theta'_t), a'_t\rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] \\
&= -\eta\mathbb{E}[\langle\nabla f(\theta_t), a'_t\rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] + \eta\mathbb{E}[\langle\nabla f(\theta_t) - \nabla f(\theta'_t), a'_t\rangle] \\
&\le -\eta\mathbb{E}[\langle\nabla f(\theta_t), a'_t\rangle] + \frac{\eta^2 L}{2}\mathbb{E}[\|a'_t\|^2] + \eta^2 L\mathbb{E}[\|\mathcal{E}_t\|\|a'_t\|] \\
&\le -\eta\mathbb{E}[\langle\nabla f(\theta_t), a'_t\rangle] + \eta^2 L\mathbb{E}[\|a'_t\|^2] + \frac{\eta^2 L}{2}\mathbb{E}[\|\mathcal{E}_t\|^2]. \qquad (21)
\end{aligned}
$$

514 **Bounding the first term in (21).** We have

$$
\begin{aligned}
M_t &:= -\mathbb{E}[\langle\nabla f(\theta_t), a'_t\rangle] = -\mathbb{E}[\langle\nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_t}+\epsilon}\rangle] \\
&= \underbrace{-\mathbb{E}[\langle\nabla f(\theta_t), \frac{m'_t}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle]}_{I} + \underbrace{\mathbb{E}[\langle\nabla f(\theta_t), (\frac{1}{\sqrt{\hat{v}_{t-1}}+\epsilon} - \frac{1}{\sqrt{\hat{v}_t}+\epsilon})m'_t\rangle]}_{II}.
\end{aligned}
$$

515 To bound I, note that

$$
\begin{aligned}
I &= -\mathbb{E}[\langle\nabla f(\theta_t), \frac{(1-\beta_1)g_t}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle] - \mathbb{E}[\langle\nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle] \\
&= -\mathbb{E}\mathbb{E}[\langle\nabla f(\theta_t), \frac{(1-\beta_1)g_t}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle|\mathcal{F}_{t-1}] - \mathbb{E}[\langle\nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle] \\
&= -(1-\beta_1)\mathbb{E}[\frac{\|\nabla f(\theta_t)\|^2}{\sqrt{\hat{v}_{t-1}}+\epsilon}] - \mathbb{E}[\langle\nabla f(\theta_t), \frac{\beta_1 m'_{t-1}}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle] \\
&\le -\frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2}+\epsilon}\mathbb{E}[\|\nabla f(\theta_t)\|^2] - \beta_1\mathbb{E}[\langle\nabla f(\theta_t), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}}+\epsilon}\rangle].
\end{aligned}
$$

27

Regarding the second term, we have

$$
- \mathbb{E}[\langle \nabla f(\theta_t), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle]
$$

$$
= -\mathbb{E}[\langle \nabla f(\theta_{t-1}), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle] - \mathbb{E}[\langle \nabla f(\theta_t) - \nabla f(\theta_{t-1}), \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \rangle]
$$

$$
= M_{t-1} + \eta L \mathbb{E}[\| \frac{m_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \| \| \frac{m'_{t-1}}{\sqrt{\hat{v}_{t-1}} + \epsilon} \|]
$$

$$
\leq M_{t-1} + \frac{\eta L}{\epsilon} \mathbb{E}[\|m'_{t-1}\|^2] + \eta L \mathbb{E}[\|a_{t-1}\|^2]
$$

$$
\leq M_{t-1} + \frac{\eta L}{\epsilon} (C\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2]) + \frac{\eta L(1-\beta_1)d}{(1-\beta_2)(1-\gamma)},
$$

where Lemma 1 and Lemma 2 are used, with $C_1 = (1 - \beta_1^2)(1 + \frac{1}{4(1-\beta_1^2)})$ and $C = \frac{C_1}{1-\beta_1^2(2-\beta_1^2)}$.
Putting parts together we obtain

$$
I \leq \beta_1 M_{t-1} + \frac{\eta \beta_1 L C \sigma^2}{\epsilon} + \frac{\eta \beta_1 L C_1}{\epsilon} \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2])
$$

$$
+ \frac{\eta L \beta_1 (1-\beta_1)d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \mathbb{E}[\|\nabla f(\theta_t)\|^2].
$$

For II, it holds that

$$
II \leq G^2 \mathbb{E}[\sum_{i=1}^{d} | \frac{1}{\sqrt{\hat{v}_{t-1}} + \epsilon} - \frac{1}{\sqrt{\hat{v}_t} + \epsilon} |].
$$

Denoting $H_t := \mathbb{E}[\sum_{i=1}^{d} | \frac{1}{\sqrt{\hat{v}_{t-1}}+\epsilon} - \frac{1}{\sqrt{\hat{v}_t}+\epsilon} |]$, $S_t := \sum_{\tau=1}^{t} (\beta_1^2(2-\beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2])$. We arrive at

$$
M_t \leq \beta_1 M_{t-1} + \frac{\eta \beta_1 L C \sigma^2}{\epsilon} + \frac{\eta \beta_1 L C_1}{\epsilon} S_t + G^2 H_t
$$

$$
+ \frac{\eta L \beta_1 (1-\beta_1)d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \mathbb{E}[\|\nabla f(\theta_t)\|^2]
$$

$$
\leq \beta_1 M_{t-1} + \frac{\eta \beta_1 L C \sigma^2}{\epsilon} + \frac{\eta \beta_1 L C_1}{\epsilon} S_t + G^2 H_t + \frac{\eta L \beta_1 (1-\beta_1)d}{(1-\beta_2)(1-\gamma)}.
$$

By induction, we have

$$
M_t \leq \beta_1^{t-1} M_1 + G^2 \sum_{\tau=0}^{t-2} \beta_1^\tau H_{t-\tau} + \frac{\eta \beta_1 L C_1}{\epsilon} \sum_{\tau=0}^{t-2} \beta_1^\tau S_{t-\tau} + \frac{\eta \beta_1 L C \sigma^2}{(1-\beta_1)\epsilon}
$$

$$
+ \frac{\eta L \beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \mathbb{E}[\|\nabla f(\theta_t)\|^2],
$$

since $\beta_1 < 1$. For bounding the summations, we have the following result.

28

Summing over $t = 1, ..., T$, we obtain

$$\sum_{t=1}^{T} M_t \leq \sum_{t=1}^{T} \beta_1^{t-1} M_1 + G^2 \sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^\tau H_{t-\tau} + \frac{\eta \beta_1 L C_1}{\epsilon} \sum_{t=2}^{T} \sum_{\tau=0}^{t-2} \beta_1^\tau S_{t-\tau}$$

$$+ \frac{T\eta\beta_1 L C \sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L \beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\overset{(a)}{\leq} \frac{2dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta\beta_1 L C \sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L \beta_1 d}{(1-\beta_2)(1-\gamma)}$$

$$+ \left[ \frac{\eta L C}{(1-\beta_1)\epsilon} - \frac{1-\beta_1}{\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \right] \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]$$

$$\leq \frac{2dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta\beta_1 L C \sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta L \beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{3(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2],$$

when $\eta$ is chosen to be $\eta \leq \frac{(1-\beta_1)^2 \epsilon}{4LC\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}$. Here, (a) is due to $M_1 = \mathbb{E}[\langle \nabla f(\theta_1), a_0' \rangle] \leq$

$\beta_1 dG^2/\sqrt{\epsilon}$ and Lemma 3. It remains to bound the last two terms in (21).

**Bounding the last two terms in in (21).** We have

$$\mathbb{E}[\|a_t'\|^2] = \mathbb{E}[\|\frac{m_t'}{\sqrt{\hat{v}_t} + \epsilon}\|^2] \leq \frac{1}{\epsilon} \mathbb{E}[\|m_t'\|^2].$$

By Lemma 1, it follows that

$$\mathbb{E}[\|a_t'\|^2] \leq \frac{1}{\epsilon}(C\sigma^2 + C_1 \sum_{\tau=1}^{t} (\beta_1^2(2 - \beta_1^2))^{t-\tau} \mathbb{E}[\|\nabla f(\theta_\tau)\|^2]).$$

Summing over $t = 1, ..., T$, we obtain

$$\sum_{t=1}^{T} \|a_t'\|^2 \leq \frac{TC\sigma^2}{\epsilon} + \frac{C}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]),$$

where the last inequality can be derived similar to Lemma 3.

For the last term in (21), we have by Lemma 5

$$\sum_{t=1}^{T} \mathbb{E}[\|\mathcal{E}_t\|^2] \leq \frac{4Tq^2}{(1-q^2)^2 \epsilon} \sigma^2 + \frac{4q^2}{(1-q^2)^2 \epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2].$$

29

532 **Completing the proof.** Summing (21) over $t = 1, ..., T$ and integrating things together, we have

$$
\mathbb{E}[f(\theta'_{T+1}) - f(\theta'_1)]
$$

$$
\leq \eta \sum_{t=1}^{T} M_t + \frac{T\eta^2 CL\sigma^2}{\epsilon} + \frac{C\eta^2 L}{\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2])
$$

$$
+ \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon} + \frac{2\eta^2 Lq^2}{(1-q^2)^2\epsilon} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]
$$

$$
\leq \frac{2\eta dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta^2\beta_1 LC\sigma^2}{(1-\beta_1)\epsilon} + \frac{T\eta^2 L\beta_1 d}{(1-\beta_2)(1-\gamma)} - \frac{3\eta(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]
$$

$$
+ \frac{T\eta^2 CL\sigma^2}{\epsilon} + \Big[\frac{C\eta^2 L}{\epsilon} + \frac{2\eta^2 Lq^2}{(1-q^2)^2\epsilon}\Big] \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2]) + \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon}
$$

$$
\leq -\frac{\eta(1-\beta_1)}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] + \frac{2\eta dG^2}{(1-\beta_1)\sqrt{\epsilon}} + \frac{T\eta^2 LC\sigma^2}{(1-\beta_1)\epsilon}
$$

$$
+ \frac{T\eta^2 L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2T\eta^2 Lq^2\sigma^2}{(1-q^2)^2\epsilon},
$$

533   when $\eta \leq \frac{(1-q^2)^2(1-\beta_1)\epsilon}{8Lq^2\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}$, where the last line is because $C\eta L \leq \frac{(1-\beta_1)\epsilon}{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}$ also holds.

534   Re-arranging terms, we get that when $\eta \leq \min\{\frac{1-\beta_1}{C}, \frac{(1-q^2)^2}{2q^2}\}\frac{(1-\beta_1)\epsilon}{4L\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}$,

$$
\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}[\|\nabla f(\theta_t)\|^2] \leq C'\Big(\frac{\mathbb{E}[f(\theta'_1) - f(\theta'_{T+1})]}{T\eta} + \frac{2dG^2}{T(1-\beta_1)\sqrt{\epsilon}} + \frac{\eta LC\sigma^2}{(1-\beta_1)\epsilon}
$$

$$
+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2\eta Lq^2\sigma^2}{(1-q^2)^2\epsilon}\Big)
$$

$$
\leq C'\Big(\frac{\mathbb{E}[f(\theta_1) - f(\theta^*)]}{T\eta} + \frac{2dG^2}{T(1-\beta_1)\sqrt{\epsilon}} + \frac{\eta LC\sigma^2}{(1-\beta_1)\epsilon}
$$

$$
+ \frac{\eta L\beta_1 d}{(1-\beta_2)(1-\gamma)} + \frac{2\eta Lq^2\sigma^2}{(1-q^2)^2\epsilon}\Big).
$$

535   where $C' = \frac{4\sqrt{\frac{4(1+q^2)^3}{(1-q^2)^2}G^2 + \epsilon}}{1-\beta_1}$, and $C = \frac{(1-\beta_1)^2}{1-\beta_1^2(2-\beta_1)^2}(1 + \frac{1}{4(1-\beta_1^2)})$. The last inequality is because
536   $\theta'_1 = \theta_1$, and $\theta^* = \arg\min_\theta f(\theta)$. The proof is complete.

537   $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$