
Convergent Adaptive Gradient Methods in Decentralized Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

Adaptive gradient methods including Adam, AdaGrad, and their variants have been very successful for training deep learning models, such as neural networks, in the past few years. Meanwhile, given the need for distributed training procedure, the realm of distributed optimization algorithms is at the center of attention. With the growth of computing power and the need for using machine learning models on mobile devices, the communication cost of distributed training algorithms needs careful consideration. In response to this, more and more attention is shifted from the traditional parameter server training paradigm to the decentralized one, which usually requires lower communication costs. In this paper, we try to rigorously incorporate adaptive gradient methods into decentralized training procedures, coming up with convergent decentralized adaptive gradient methods. Specifically, we propose a general algorithmic framework that can convert existing adaptive gradient methods to their decentralized counterparts. In addition, we rigorously analyze the convergence behavior of the proposed algorithmic framework and show that if a given adaptive gradient method converges, under some specific conditions, then its decentralized counterpart is also convergent.

1 Introduction

Distributed training of machine learning models is drawing growing attention in the past few years due to its practical benefits and necessities. Given the evolution of computing capabilities of CPUs and GPUs, computation time in distributed setting is gradually dominated by the communication time in many circumstances [Chilimbi et al., 2014, McMahan et al., 2016]. As a result, a large amount of recent works has been focussing on reducing communication cost for distributed learning [Alistarh et al., 2017, Lin et al., 2017, Wangni et al., 2018, Stich et al., 2018, Wang et al., 2018, Tang et al., 2019]. In the traditional parameter (central) server setting, where a parameter server is employed to manage communication in the whole network, many effective communication reductions have been proposed based on gradient compression [Aji and Heafield, 2017] and quantization [Chen et al., 2010, Ge et al., 2013, Jegou et al., 2010]. Despite these communication reduction techniques, its cost still, usually, scales linearly with the number of workers. Due to this limitation and with the sheer size of decentralized devices, the *decentralized training paradigm* [Duchi et al., 2011b], where the parameter server is removed and each node only communicates with its neighbors, is drawing attention. It has been shown in Lian et al. [2017] that decentralized training algorithms can outperform parameter server-based algorithms when the training bottleneck is the communication cost. The decentralized training paradigm is also naturally preferred when a parameter server is not available.

In light of recent advances in nonconvex optimization, an effective way to accelerate training is by using adaptive gradient methods like AdaGrad [Duchi et al., 2011a], Adam [Kingma and Ba, 2014] or AMSGrad [Reddi et al., 2019]. Their practical benefits are proven by their popularity in training neural networks, featured by faster convergence and ease of parameter tuning compared with SGD. Despite a large amount of literature in distributed optimization, there have been few

works considering bringing adaptive gradient methods into distributed training, largely due to the lack of understanding of adaptive gradient methods convergence behaviors. Notably, Reddi et al. [2020] develop the first decentralized ADAM method for distributed optimization problems with a direct application to federated learning. An inner loop is employed to compute mini-batch gradients on each worker nodes and a global adaptive step is done to update the global parameter at each central-server iteration. Yet, in the settings of our paper, nodes can only communicate with their neighbors while a server/worker communication is needed in [Reddi et al., 2020]. Designing adaptive methods in such settings is highly non-trivial due to the already complicated update rules and the interaction between the effect of using adaptive learning rates and the decentralized communication protocols.

This paper is an attempt at bridging the gap between both realms in nonconvex optimization. Our contributions are summarized as follows:

- In this paper, we investigate the possibility of using any adaptive gradient methods in the decentralized training paradigm. We develop a general technique that can convert an adaptive gradient method from a centralized method to a decentralized method.
- By using our proposed technique, we present a new decentralized optimization algorithm, called decentralized AMSGrad, as the decentralized counterpart of AMSGrad.
- We provide a theoretical verification interface for analyzing the behavior of decentralized adaptive gradient methods obtained as a result of our technique. Built upon our proposed analysis framework for that type of decentralized algorithms, we can characterize the convergence rate of decentralized AMSGrad, which is the first convergent decentralized adaptive gradient method.

A *novel technique* in our framework is a mechanism to enforce a consensus on adaptive learning rates at different nodes. We show the importance of consensus on adaptive learning rates by proving a divergent problem instance for a recently proposed decentralized adaptive gradient method DADAM [Nazari et al., 2019], a decentralized version of ADAM, which lacks consensus mechanisms on adaptive learning rates.

After presenting related work and important concepts of decentralized adaptive methods in Section 2, we develop our general framework for converting any adaptive gradient algorithm in its decentralized counterpart along with their rigorous finite-time convergence analysis in Section 3. Section 4 and 5 conclude our work with illustrative examples of our framework.

Notations: $x_{t,i}$ denotes variable x at node i and iteration t . $\|\cdot\|_{abs}$ denotes the entry-wise L_1 norm of a matrix, i.e. $\|A\|_{abs} = \sum_{i,j} A_{i,j}$. We introduce important notations used throughout the paper: for any $t > 0$, $G_t := [g_{t,N}]$ where $[g_{t,N}]$ denotes the vector $[g_{t,1}, g_{t,2}, \dots, g_{t,N}]$, $M_t := [m_{t,N}]$, $X_t := [x_{t,N}]$, $\bar{\nabla}f(X_t) := \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i})$, $U_t := [u_{t,N}]$, $\tilde{U}_t := [\tilde{u}_{t,N}]$, $V_t := [v_{t,N}]$, $\hat{V}_t := [\hat{v}_{t,N}]$, $\bar{X}_t := \frac{1}{N} \sum_{i=1}^N x_{t,i}$, $\bar{U}_t := \frac{1}{N} \sum_{i=1}^N u_{t,i}$ and $\tilde{\bar{U}}_t := \frac{1}{N} \sum_{i=1}^N \tilde{u}_{t,i}$.

2 Decentralized Adaptive Training and Divergence of DADAM

2.1 Related Work

Decentralized optimization: Traditional decentralized optimization methods include well-known algorithms such as ADMM [Boyd et al., 2011], dual averaging [Duchi et al., 2011b], distributed subgradient descent [Nedic and Ozdaglar, 2009]. More recent algorithms include Extra [Shi et al., 2015], Next [Di Lorenzo and Scutari, 2016] and Prox-PDA [Hong et al., 2017]. While these algorithms were commonly used in applications other than deep learning, recent algorithmic advances in the machine learning community have shown that decentralized optimization can be useful for training deep models such as neural networks too. Lian et al. [2017] show that a stochastic version of decentralized subgradient descent can outperform parameter server-based algorithms when the communication cost is high. Tang et al. [2018] propose the D^2 algorithm improving the convergence rate over stochastic subgradient descent. Assran et al. [2018] propose the Stochastic Gradient Push that is more robust to network failures for training neural networks. The study of decentralized training in the machine learning community is only at its initial stage. No one has seriously considered designing adaptive gradient methods in the setting of decentralized training until a recent work [Nazari

et al., 2019] where a decentralized version of AMSGrad [Reddi et al., 2019] is proposed is proven to satisfy some non-standard regret.

Adaptive gradient methods: Adaptive gradient methods have been popular in recent years due to their superior performance in training neural networks. Most used adaptive methods include AdaGrad [Duchi et al., 2011a] or Adam [Kingma and Ba, 2014] and their variants. Key features of such methods lie in the use of momentum and adaptive learning rates (which means the learning rate is changing during optimization and are anisotropic, i.e. depend on the dimension). The method of reference, Adam, has been analyzed in [Reddi et al., 2019] where the authors point out an error in previous convergence analyses. Since then, a range of works have been focussing on analyzing the convergence behavior of the various existing adaptive gradient methods. Ward et al. [2018], Li and Orabona [2018] derive convergence guarantees for a variant of AdaGrad without coordinate-wise learning rates. Chen et al. [2018] analyze the convergence behavior of a broad class of algorithms including AMSGrad [Reddi et al., 2019] and AdaGrad. Zou and Shen [2018] provide a unified convergence analysis for AdaGrad with momentum. A few recent adaptive gradient methods can be found in [Agarwal et al., 2018, Luo et al., 2019, Zaheer et al., 2018].

2.2 Decentralized Optimization

In distributed optimization (with N nodes), we aim at solving the following problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N f_i(x), \quad (1)$$

where x is the vector of parameters and f_i is only accessible by the i th node. Through the prism of neural network training procedure, f_i can be viewed as the average loss of the data samples located at node i . Throughout the paper, we make the following assumptions for analyzing the convergence behavior of the different algorithms.

A1. For all $i \in [N]$, f_i is differentiable and the gradients is L -Lipschitz, i.e. $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|, \forall x, y$.

A2. We assume at iteration t , node i can access a stochastic gradient $g_{t,i}$. In addition, the stochastic gradients have bounded L_∞ norm and the gradients of f_i are also bounded, i.e. $\|g_{t,i}\| \leq G_\infty$, $\|\nabla f_i(x)\|_\infty \leq G_\infty$.

A3. The gradient estimators are unbiased and each coordinate have bounded variance, i.e. $\mathbb{E}[g_{t,i}] = \nabla f_i(x_{t,i})$ and $\mathbb{E}[(g_{t,i} - \nabla f_i(x_{t,i}))_j^2] \leq \sigma^2, \forall t, i, j$.

The assumptions A1 and A3 are standard in distributed optimization. A2 is slightly stronger than the traditional assumption that the estimator has bounded variance, but is commonly used for the analysis of adaptive gradient methods [Chen et al., 2018, Ward et al., 2018]. One thing that should be noted is that the bounded gradient estimator assumption in A2 implies the bounded variance assumption in A3. We denote the variance bound and the estimator bound differently to avoid confusion when we use them for different purposes. In decentralized optimization, the nodes are connected as a graph and each node only communicates to its neighbors. In such cases, one usually constructs a $N \times N$ matrix W for information sharing when designing new algorithms. We denote λ_i to be its i th largest eigenvalue and define $\lambda \triangleq \max(|\lambda_2|, |\lambda_N|)$. As can be expected, W cannot be arbitrary, the key properties required for W are listed in A4.

A4. The matrix W satisfies: (I) $\sum_{j=1}^N W_{i,j} = 1, \sum_{i=1}^N W_{i,j} = 1, W_{i,j} \geq 0$, (II) $\lambda_1 = 1, |\lambda_2| < 1, |\lambda_N| < 1$ and (III) $W_{i,j} = 0$ if node i and node j are not neighbors.

2.3 Divergence of DADAM

Recently, Nazari et al. [2019] initiated a trial to bring adaptive gradient methods into decentralized optimization, the resulting algorithm is DADAM, which is shown in Algorithm 1. DADAM is essentially a decentralized version of AMSGrad and the key modification is the use of a consensus step on optimization variable x to transmit information across the network, encouraging convergence. The matrix W is a doubly stochastic matrix (which satisfies A4) for achieving average consensus of x . Introducing such mixing matrix is a standard approach for decentralizing an algorithm, such as distributed gradient descent [Nedic and Ozdaglar, 2009, Yuan et al., 2016]. It is proven in Nazari et al. [2019] that DADAM admits a non-standard regret bound in the online setting, however, whether the algorithm can converge to stationary points in standard offline settings such training neural networks is still unknown.

140 In the following, we show the DADAM may fail
 141 to converge in the offline nonconvex optimiza-
 142 tion settings.

143 **Theorem 1.** *There exist a problem satisfying*
 144 *A1 – A4 where DADAM fail to converge.*

145 *Proof.* Consider a 1 dimensional optimiza-
 146 tion problem distributed on two nodes
 147 $\min_x \frac{1}{2} \sum_{i=1}^2 f_i(x)$ where $f_i(x) = \frac{1}{2}(x - a_i)^2$
 148 and $a_1 = 0, a_2 = 1$. The network contains only
 149 two nodes and the matrix W satisfies $W_{ij} = \frac{1}{2}$
 150 for all i, j . For simplicity, we consider running
 151 DADAM with $\beta_1 = \beta_2 = \beta_3 = 0$ and $\epsilon = 0.6$.
 152 Suppose we initialize DADAM at $x_{1,i} = 0$
 153 for all $i \in [N]$ and use the following learning
 154 rate $\alpha = 0.001$. We have at $x_{1,i} = 0, \nabla f_1(x_{1,1}) = 0, \nabla f_2(x_{1,2}) = 1$, leading to $\hat{v}_{1,1} = 0.6$ and
 155 $\hat{v}_{1,2} = 1$. Thus, from step 1, we will have $\hat{v}_{1,2} \geq 1$. In addition, it is can be easily proved that,
 156 with the stepsize selection, we always have $\hat{v}_{1,1} < 1$, in fact, it will not reach 0.6. Thus, in the next
 157 iterations, the gradient of losses on node 1 and 2 will be scaled differently. This scaling is equivalent
 158 to running gradient descent on a objective where the losses of the two nodes are scaled by different
 159 factors. In such case, the algorithm will converge to a stationary point of a weighted average of the
 160 loss on node 1. Recall that the problem we tackle to illustrate Theorem 1 is a quadratic problem with
 161 only one minimizer. Then, since the weight of the losses on the two nodes are different and that
 162 the unbalanced weights on the two functions yields a different minimizer, the algorithm will not
 163 converge to the unique stationary point of the original loss (which is $x = 0.5$). \square

164 Theorem 1 claims that even though DADAM is proven to satisfy some regret bounds , see [Nazari et al.,
 165 2019], it can fail to converge to stationary points in the nonconvex offline setting, which is a common
 166 setting for training neural networks. We conjecture that this inconsistency is due to the definition of
 167 the regret in [Nazari et al., 2019]. In the next section, we will design decentralized adaptive gradient
 168 methods that are guaranteed to converge to stationary points and provide a characterization of that
 169 convergence in finite-time and independently of the initialization of our methods.

170 3 Decentralized Adaptive Gradient Methods and their Convergence

171 In this section, we discuss the difficulties of designing adaptive gradient methods in decentralized
 172 optimization and introduce an algorithmic framework that convert existing convergent adaptive gradi-
 173 ent methods to their decentralized counterparts. We also develop the first convergent decentralized
 174 adaptive gradient method, converted from AMSGrad, as an instance of this proposed framework.

175 3.1 Importance and Difficulties of Consensus on Adaptive Learning Rates

176 The divergent example in the previous section implies that we should synchronize the adaptive
 177 learning rates on different nodes. This can be easily achieved in the parameter server setting where
 178 all the nodes are sending their gradients to a central server at each iteration. The parameter server
 179 can then exploit the received gradients to maintain a sequence of synchronized adaptive learning
 180 rates when updating the parameters, see [Reddi et al., 2020]. However, in our setting of decentralized
 181 training, every node can only communication with its neighbors and such central parameter server
 182 does not exist. Under that setting, the information for updating the adaptive learning rates can be only
 183 shared locally instead of broadcasted over the whole network. This makes it impossible to obtain, in a
 184 single iteration, a synchronized adaptive learning rate update using all the information in the network.

185 *Systemic Approach:* On a systemic level, one way to alleviate this bottleneck is to design communica-
 186 tion protocols to give each node access to the same aggregated gradients over the whole network at
 187 least periodically if not at every iteration. Therefore, the nodes can update their individual adaptive
 188 learning rates based on the same information. However, such solution introduce extra communication
 189 cost since it involves broadcasting over the network.

190 *Algorithmic Approach:* Our contributions being on an algorithmic level, another way to solve the
 191 aforementioned problem is by letting the sequences of adaptive learning rates, present on different

Algorithm 1 DADAM (with N nodes)

```

1: Input:  $\alpha$ , current point  $X_t, u_{\frac{1}{2},i} = \hat{v}_{0,i} = \epsilon \mathbf{1}$ ,  

    $m_0 = 0$  and mixing matrix  $W$   

2: for  $t = 1, 2, \dots, T$  do  

3:   for all  $i \in [N]$  do in parallel  

4:      $g_{t,i} \leftarrow \nabla f_i(x_{t,i}) + \xi_{t,i}$   

5:      $m_{t,i} = \beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}$   

6:      $v_{t,i} = \beta_2 v_{t-1,i} + (1 - \beta_2) g_{t,i}^2$   

7:      $\hat{v}_{t,i} = \beta_3 \hat{v}_{t-1,i} + (1 - \beta_3) \max(\hat{v}_{t-1,i}, v_{t,i})$   

8:      $x_{t+\frac{1}{2},i} = \sum_{j=1}^N W_{ij} x_{t,j}$   

9:      $x_{t+1,i} = x_{t+\frac{1}{2},i} - \alpha \frac{m_{t,i}}{\sqrt{\hat{v}_{t,i}}}$   

10:  end for

```

Algorithm 2 Decentralized Adaptive Gradient Method (with N nodes)

```

1: Input: learning rate  $\alpha$ , initial point  $x_{1,i} = x_{init}, u_{\frac{1}{2},i} = \hat{v}_{0,i}, m_{0,i} = 0$ , mixing matrix  $W$ 
2: for  $t = 1, 2, \dots, T$  do
3:   for all  $i \in [N]$  do in parallel
4:      $g_{t,i} \leftarrow \nabla f_i(x_{t,i}) + \xi_{t,i}$ 
5:      $m_{t,i} = \beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}$ 
6:      $\hat{v}_{t,i} = r_t(g_{1,i}, \dots, g_{t-1,i})$ 
7:      $x_{t+\frac{1}{2},i} = \sum_{j=1}^N W_{ij} x_{t,j}$ 
8:      $\tilde{u}_{t,i} = \sum_{j=1}^N W_{ij} \tilde{u}_{t-\frac{1}{2},j}$ 
9:      $u_{t,i} = \max(\tilde{u}_{t,i}, \epsilon)$ 
10:     $x_{t+1,i} = x_{t+\frac{1}{2},i} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}}$ 
11:     $\tilde{u}_{t+\frac{1}{2},i} = \tilde{u}_{t,i} - \hat{v}_{t-1,i} + \hat{v}_{t,i}$ 
12:  end for

```

nodes, to *consent* gradually, through the iterations. Intuitively, if the adaptive learning rates can consent fast enough, the difference among the adaptive learning rates on different nodes will not affect the convergence of the algorithm. The benefit of such approach is that we do not need to introduce extra communication cost.

3.2 Decentralized Adaptive Gradient Unifying Framework

As mentioned before, we need to choose a method to implement consensus of adaptive learning rates. While each node can have different $\hat{v}_{t,i}$ in DADAM, one can keep track of the min/max/average of these adaptive learning rates and use this quantity to update the adaptive learning rates. Also one can predefine some convergent lower and upper bounds to gradually synchronize the adaptive learning rates on different nodes as developed for AdaBound in [Luo et al., 2019]. In this paper, we opt for the average consensus on $\hat{v}_{t,i}$. Since in adaptive gradient methods such as AdaGrad or Adam, the quantity $\hat{v}_{t,i}$ approximates the second moment of the gradient estimator, the average of the estimations of those second moments from different nodes is an estimation of second moment on the whole network. Also, this design will not introduce any extra hyperparameters that can potentially complicate the tuning process. We now present the main convergence result for our class methods:

Theorem 2. Assume A1-A4. Set $\alpha = 1/\sqrt{Td}$. When $\alpha \leq \frac{\epsilon^{0.5}}{16L}$, Algorithm 2 yields the following regret bound

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(Z_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{C_2}{T} + \frac{C_3}{T^{1.5}d^{0.5}} \\ &\quad + \left(\frac{C_4}{TN^{0.5}} + \frac{C_5}{T^{1.5}d^{0.5}N^{0.5}} \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right], \end{aligned} \quad (2)$$

where $\|\cdot\|_{abs}$ denotes the entry-wise L_1 norm of a matrix (i.e. $\|A\|_{abs} = \sum_{i,j} |A_{ij}|$). The constants $C_1 = \max(4, 4L/\epsilon)$, $C_2 = 6((\beta_1/(1-\beta_1))^2 + 1/(1-\lambda)^2)LG_\infty^2/\epsilon^{1.5}$, $C_3 = 16L^2(1-\lambda)G_\infty^2/\epsilon^2$, $C_4 = 2/(\epsilon^{1.5}(1-\lambda))(\lambda + \beta_1/(1-\beta_1))G_\infty^2$, $C_5 = 2/(\epsilon^2(1-\lambda))L(\lambda + \beta_1/(1-\beta_1))G_\infty^2 + 4/(\epsilon^2(1-\lambda))LG_\infty^2$ are independent of d, T and N . In addition, $\frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 \leq \alpha^2 \left(\frac{1}{1-\lambda} \right)^2 dG_\infty^2 \frac{1}{\epsilon}$ which quantifies the consensus error.

Theorem 2 shows that if $\mathbb{E}[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs}] = o(T)$ and \bar{U}_t is upper bounded, then Algorithm 2 is guaranteed to converge to stationary points of the regret function. Intuitively, this means that if the adaptive learning rates on different nodes do not change too fast, the algorithm can converge. This is true as stated in [Chen et al., 2018] where it is shown that if such condition is violated, the algorithm can diverge. Furthermore, Theorem 2 conveys the benefits of using more nodes. As N becomes larger, the term σ^2/N will be small. This is also strengthened by the fact that with the growth of N , the training process tends to be more stable. We now present, in Algorithm 3, a notable special case of our algorithmic framework, namely Decentralized AMSGrad, which is a decentralized variant of AMSGrad.

223 Compared with DADAM, the above algorithm
 224 leverages a dynamic average consensus mech-
 225 anism to keep track of average of $\{\hat{v}_{t,i}\}_{i=1}^N$,
 226 stored as $\tilde{u}_{t,i}$ on i th node, and uses $u_{t,i} =$
 227 $\max(\tilde{u}_{t,i}, \epsilon)$ for updating the adaptive learn-
 228 ing rate for i th node. As the number of iter-
 229 ation grows, even though $\hat{v}_{t,i}$ on different
 230 nodes can converge to different constants, all
 231 the $u_{t,i}$ will be converge to the same number
 232 $\lim_{t \rightarrow \infty} 1/N \sum_{i=1}^N \hat{v}_{t,i}$ if the limit exists. The
 233 use of this average consensus mechanism en-
 234 ables the consensus of adaptive learning rates
 235 on different nodes, which consequentially guar-
 236 antees convergence to stationary points. The
 237 consensus of adaptive learning rates is the
 238 key difference between decentralized AMSGrad
 239 and DADAM and is the reason why decentral-
 240 ized AMSGrad is a convergent algorithm while
 241 DADAM is not. One may noticed that decentral-
 242 ized AMSGrad does not deduce to AMSGrad
 243 because $u_{t,i}$ in line 10 is calculated based on $v_{t-1,i}$ instead of $v_{t,i}$. This encourages parallel execution
 244 of gradient computation and communication. Specifically, line 4-7 in Algorithm 3 and Algorithm
 245 2 can be executed in parallel with line 8-9 to overlap communication and computation time. If
 246 $u_{t,i}$ depends on $v_{t,i}$ which in turn depends on $g_{t,i}$, the gradient computation must finish before the
 247 consensus step of adaptive learning rate line in 9. This can slow down per-iteration running time
 248 of the algorithm. To avoid such delayed adaptive learning, adding $\tilde{u}_{t-\frac{1}{2},i} = \tilde{u}_{t,i} - \hat{v}_{t-1,i} + \hat{v}_{t,i}$
 249 before line 9 and get rid of line 12 in Algorithm 2 is an option. Similar convergence guarantees will
 250 hold since one can easily modify our proof of Theorem 2 for such update rule. As stated above,
 251 Algorithm 3 converges, with the following rate:

252 **Theorem 3.** Assume A1-A4. Set $\alpha = 1/\sqrt{Td}$. When $\alpha \leq \frac{\epsilon^{0.5}}{16L}$, Algorithm 3 yields the following
 253 regret bound

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{U_t^{1/4}} \right\|^2 \right] \leq C'_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(Z_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{C'_2}{T} + \frac{d}{T} \sqrt{N} C'_4 + \frac{\sqrt{d}}{T^{1.5}} \sqrt{N} C'_5,$$

254 where $C'_1 = C_1$, $C'_2 = C_2$, $C'_3 = C_3$, $C'_4 = C_4 G_\infty^2$ and $C'_5 = C_5 G_\infty^2$. C_1, C_2, C_3, C_4, C_5 are
 255 constants independent of d, T and N defined in Theorem 2. In addition, the consensus of variables at
 256 different nodes is given by $\frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 \leq \frac{1}{T} \left(\frac{1}{1-\lambda} \right)^2 G_\infty^2 \frac{1}{\epsilon}$.

257 Theorem 3 shows that Algorithm 3 converges with a rate of $\mathcal{O}(\sqrt{d}/\sqrt{T})$ when T is large, which is the
 258 best known convergence rate under the given assumptions. Note that in some related works, SGD
 259 admits a convergence rate of $\mathcal{O}(1/\sqrt{T})$ without any dependence on the dimensions. Such improved
 260 convergence rate is under the assumption that the gradient estimator have a bounded L_2 norm, which
 261 can thus hide a dependency of \sqrt{d} in the final convergence rate.

262 3.3 Convergence Analysis

263 **Proof of Theorem 2.** The detailed proof of this section is reported in the supplementary material.
 264 We now present a proof sketch for our main convergence result of Algorithm 2.

265 *Step 1: Reparameterization.* Similarly to [Yan et al., 2018, Chen et al., 2018] with SGD (with
 266 momentum) and centralized adaptive gradient methods, define the following auxiliary sequence:

$$Z_t = \bar{X}_t + \frac{\beta_1}{1 - \beta_1} (\bar{X}_t - \bar{X}_{t-1}), \quad (3)$$

267 with $\bar{X}_0 \triangleq \bar{X}_1$. Such an auxiliary sequence can help us deal with the bias brought by the momentum
 268 and simplifies the convergence analysis. An intermediary result needed to conduct our proof reads:

Algorithm 3 Decentralized AMSGrad (with N nodes)

```

1: Input: learning rate  $\alpha$ , initial point  $x_{1,i} =$ 
    $x_{init}, u_{\frac{1}{2},i} = \hat{v}_{0,i} = \epsilon \mathbf{1}$  (with  $\epsilon \geq 0$ ),  $m_{0,i} =$ 
    $0$ , mixing matrix  $W$ 
2: for  $t = 1, 2, \dots, T$  do
3:   for all  $i \in [N]$  do in parallel
4:      $g_{t,i} \leftarrow \nabla f_i(x_{t,i}) + \xi_{t,i}$ 
5:      $m_{t,i} = \beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}$ 
6:      $v_{t,i} = \beta_2 v_{t-1,i} + (1 - \beta_2) g_{t,i}^2$ 
7:      $\hat{v}_{t,i} = \max(\hat{v}_{t-1,i}, v_{t,i})$ 
8:      $x_{t+\frac{1}{2},i} = \sum_{j=1}^N W_{ij} x_{t,j}$ 
9:      $\tilde{u}_{t,i} = \sum_{j=1}^N W_{ij} \tilde{u}_{t-\frac{1}{2},j}$ 
10:     $u_{t,i} = \max(\tilde{u}_{t,i}, \epsilon)$ 
11:     $x_{t+1,i} = x_{t+\frac{1}{2},i} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}}$ 
12:     $\tilde{u}_{t+\frac{1}{2},i} = \tilde{u}_{t,i} - \hat{v}_{t-1,i} + \hat{v}_{t,i}$ 
13:  end for
```

269 **Lemma 1.** For the sequence defined in (3), we have

$$Z_{t+1} - Z_t = \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \quad (4)$$

270 Lemma 1 does not display any momentum term in $\frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}}$. This simplification is convenient
271 since it is directly related to the current gradients instead of the exponential average of past gradients.

272 **Step 2: Smoothness.** Using smoothness assumption A1 involves the following scalar product term:
273 $\kappa_t := \langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) / \sqrt{\bar{U}_t} \rangle$ which can be lower bounded by:

$$\kappa_t \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2. \quad (5)$$

274 The above inequality substituted in the smoothness condition $f(Z_{t+1}) \leq f(Z_t) + \langle \nabla f(Z_t), Z_{t+1} -$
275 $Z_t \rangle + \frac{L}{2} \|Z_{t+1} - Z_t\|^2$ yields:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \leq \frac{2}{T\alpha} (\mathbb{E}[\Delta_f]) + \frac{L}{T\alpha} \sum_{t=1}^T \mathbb{E} [\|Z_{t+1} - Z_t\|^2] + \frac{2}{T} \frac{\beta_1}{1 - \beta_1} T_1 + \frac{2}{T} T_2 + \frac{3}{T} T_3, \quad (6)$$

276 where $\Delta_f := \mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]$ T_1, T_2 and T_3 are three terms, defined in the supplementary
277 material, and which can be tightly bounded from above. We first bound T_3 using the following
278 quantities of interest:

$$\sum_{t=1}^T \|Z_t - \bar{X}_t\|^2 \leq T \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \alpha^2 d \frac{G_\infty^2}{\epsilon} \quad \text{and} \quad \sum_{t=1}^T \frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 \leq T \alpha^2 \left(\frac{1}{1 - \lambda} \right)^2 d G_\infty^2 \frac{1}{\epsilon}. \quad (7)$$

279 where $\lambda = \max(|\lambda_2|, |\lambda_N|)$ and recall that λ_i is i th largest eigenvalue of W .

280 Then, concerning the term T_2 , few derivations, not detailed here for simplicity, yields:

$$T_2 \leq \frac{G_\infty^2}{N} \mathbb{E} \left[\sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| -\sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \right] \quad (8)$$

281 where q_l is the eigenvector corresponding to l th largest eigenvalue of W and $\|\cdot\|_{abs}$ is the entry-wise
282 L_1 norm of matrices. We can also show that

$$\sum_{t=1}^T \left\| -\sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \leq \sqrt{N} \sum_{o=0}^{T-1} \frac{\lambda}{1 - \lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \quad (9)$$

283 resulting in an upper bound for T_2 proportional to $\sum_{o=0}^{T-1} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs}$. Similarly:

$$T_1 \leq G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{\sqrt{N}} \mathbb{E} \left[\frac{1}{1 - \lambda} \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \quad (10)$$

284 **Step 3: Bounding the drift term variance.** An important term that needs upper bounding in our proof
285 is the variance of the gradients multiplied (element-wise) by the adaptive learning rate:

$$\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \leq \mathbb{E}[\|\Gamma_u^f\|^2] + \frac{d}{N} \frac{\sigma^2}{\epsilon} \quad (11)$$

286 where $\Gamma_u^f := 1/N \sum_{i=1}^N \nabla f_i(x_{t,i}) / \sqrt{u_{t,i}}$. Two consecutive and simple bounding of the above yields:

$$\sum_{t=1}^T \mathbb{E}[\|\Gamma_u^f\|^2] \leq 2 \sum_{t=1}^T \mathbb{E}[\|\Gamma_{\bar{U}}^f\|^2] + 2 \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right\| \right] \quad (12)$$

287 and

$$\sum_{t=1}^T \mathbb{E}[\|\Gamma_{\bar{U}}^f\|^2] \leq 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right]. \quad (13)$$

288 Then, by plugging the LHS of (13) in (6), and further bounding as operated for T_2, T_3 (see supple-
289 ment), we obtain the bound in Theorem 2.

290 4 Numerical Experiments

291 In this section, we conduct experiments to test the performance of Decentralized AMSGrad, see
 292 Algorithm 3, on both homogeneous data and heterogeneous data distribution (i.e. the data generating
 293 distribution on different nodes are assumed to be different). We compare it with DADAM and the
 294 decentralized stochastic gradient descent (DGD) developed in [Lian et al., 2017]. The task consists
 295 of training a Convolutional Neural Network (CNN) with 3 convolution layers followed by a fully
 296 connected layer on MNIST [LeCun, 1998]. We set $\epsilon = 10^{-6}$ for both Decentralized AMSGrad and
 297 DADAM, the learning rate is chosen from the grid $[10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}]$ based on
 298 validation accuracy for all algorithms. In all the following experiments, the graph contains 5 nodes
 299 and the nodes form a ring, each node can only talk with its two adjacent neighbors. We set $W_{ij} = 1/3$
 300 if there nodes i and j are neighbors and $W_{ij} = 0$ otherwise for the mixing matrix. More details and
 301 experiments can be found in Appendix A.4.

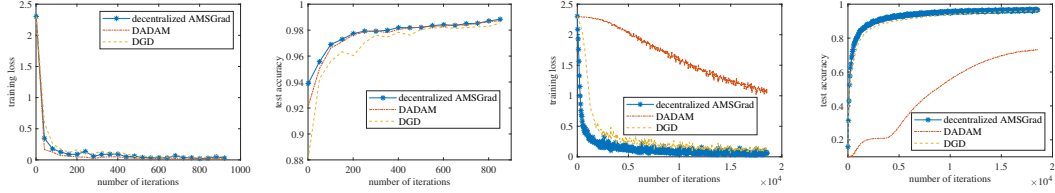


Figure 1: Performance comparison on homogeneous and heterogeneous data

302 Figure 1 shows the performance of different algorithms on homogeneous data. The whole dataset
 303 is shuffled evenly split to different nodes. We can see that decentralized AMSGrad and DADAM
 304 performs quite similarly and DGD is slower compared with them in terms of both training loss and
 305 test accuracy. Though we have prove in previous sections that DADAM is not a convergent algorithm,
 306 its performance is still quite good on homogeneous data. The reason is that the adaptive learning
 307 rates tend to be similar on different nodes when we have homogeneous data distribution. However,
 308 this is usually not true when we have heterogeneous data distribution. This motivates us to compare
 309 the performance of the algorithms on a different data distribution.

310 In Figure 1, we compare the performance of different algorithms on heterogeneous data. In this
 311 case, each node only contains training data with two labels out of ten. We can see that all algorithm
 312 converges significantly slower compared with the case with homogeneous data. Especially, the
 313 performance of DADAM deteriorates significantly. decentralized AMSGrad achieves the best training
 314 and testing performance in this experiment.

315 5 Conclusion

316 This paper studies the problem of designing adaptive gradient methods for decentralized training. We
 317 propose a unifying algorithmic framework that can convert existing adaptive gradient methods to
 318 decentralized settings. With rigorous convergence analysis, we show that if the original algorithm
 319 satisfies converges under some minor conditions, the converted algorithm obtained using our proposed
 320 framework is guaranteed to converge to stationary points of the regret function. By applying
 321 our framework to AMSGrad, we propose the first convergent adaptive gradient methods, namely
 322 Decentralized AMSGrad. Experiments show that the proposed algorithm achieves better performance
 323 than the baselines.

6 Broader Impact of Our Work

References

- Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang. The case for full-matrix adaptive regularization. *arXiv preprint arXiv:1806.02958*, 2018.
- Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*, 2017.
- 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.
- Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Michael Rabbat. Stochastic gradient push for distributed deep learning. *arXiv preprint arXiv:1811.10792*, 2018.
- 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.
- Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of adam-type algorithms for non-convex optimization. *arXiv preprint arXiv:1808.02941*, 2018.
- Yongjian Chen, Tao Guan, and Cheng Wang. Approximate nearest neighbor search by residual vector quantization. *Sensors*, 10(12):11259–11273, 2010.
- Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project adam: Building an efficient and scalable deep learning training system. In *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, pages 571–582, 2014.
- Paolo Di Lorenzo and Gesualdo Scutari. Next: In-network nonconvex optimization. *IEEE Transactions on Signal and Information Processing over Networks*, 2(2):120–136, 2016.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011a.
- 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, 2011b.
- Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. Optimized product quantization for approximate nearest neighbor search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2946–2953, 2013.
- 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 *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1529–1538. JMLR.org, 2017.
- Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2010.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Xiaoyu Li and Francesco Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. *arXiv preprint arXiv:1805.08114*, 2018.
- Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pages 5330–5340, 2017.

369 Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J Dally. Deep gradient compression:
370 Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*,
371 2017.

372 Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic
373 bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019.

374 H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient
375 learning of deep networks from decentralized data. *arXiv preprint arXiv:1602.05629*, 2016.

376 Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. Dadam: A consensus-based
377 distributed adaptive gradient method for online optimization. *arXiv preprint arXiv:1901.09109*,
378 2019.

379 Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization.
380 *IEEE Transactions on Automatic Control*, 54(1):48, 2009.

381 Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
382 Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint*
383 *arXiv:2003.00295*, 2020.

384 Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv*
385 *preprint arXiv:1904.09237*, 2019.

386 Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentralized
387 consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.

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

390 Hanlin Tang, Xiangru Lian, Ming Yan, Ce Zhang, and Ji Liu. D^2 : Decentralized training over
391 decentralized data. *arXiv preprint arXiv:1803.07068*, 2018.

392 Hanlin Tang, Xiangru Lian, Tong Zhang, and Ji Liu. Doublesqueeze: Parallel stochastic gradient
393 descent with double-pass error-compensated compression. *arXiv preprint arXiv:1905.05957*, 2019.

394 Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen
395 Wright. Atomo: Communication-efficient learning via atomic sparsification. In *Advances in*
396 *Neural Information Processing Systems*, pages 9850–9861, 2018.

397 Jianqiao Wangni, Jiale Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-
398 efficient distributed optimization. In *Advances in Neural Information Processing Systems*, pages
399 1299–1309, 2018.

400 Rachel Ward, Xiaoxia Wu, and Leon Bottou. Adagrad stepsizes: Sharp convergence over nonconvex
401 landscapes, from any initialization. *arXiv preprint arXiv:1806.01811*, 2018.

402 Yan Yan, Tianbao Yang, Zhe Li, Qihang Lin, and Yi Yang. A unified analysis of stochastic momentum
403 methods for deep learning. *arXiv preprint arXiv:1808.10396*, 2018.

404 Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent. *SIAM*
405 *Journal on Optimization*, 26(3):1835–1854, 2016.

406 Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive methods
407 for nonconvex optimization. In *Advances in Neural Information Processing Systems*, pages
408 9793–9803, 2018.

409 Fangyu Zou and Li Shen. On the convergence of weighted adagrad with momentum for training deep
410 neural networks. *arXiv preprint arXiv:1808.03408*, 2018.

411 A Appendix

412 A.1 Proof of Theorem 2

413 To prove convergence of the algorithm, we first define an auxiliary sequence

$$Z_t = \bar{X}_t + \frac{\beta_1}{1 - \beta_1} (\bar{X}_t - \bar{X}_{t-1}) \quad (14)$$

414 with $\bar{X}_0 \triangleq \bar{X}_1$.

415 Then we have the following Lemma to characterize the difference of iterations of sequence Z_t .

416 **Lemma.** *For the sequence defined in (14), we have*

$$Z_{t+1} - Z_t = \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \quad (15)$$

417 **Proof:** See Appendix A.3. □

418 Since $\mathbb{E}[g_{t,i}] = \nabla f(x_{t,i})$ and $u_{t,i}$ is a function of $G_{1:t-1}$ (which denotes G_1, G_2, \dots, G_{t-1}), we
419 have

$$\mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \quad (16)$$

420 By assuming smoothness (A1) we have

$$f(Z_{t+1}) \leq f(Z_t) + \langle \nabla f(Z_t), Z_{t+1} - Z_t \rangle + \frac{L}{2} \|Z_{t+1} - Z_t\|^2 \quad (17)$$

421 Substitute (61) into the above inequality and take expectation over G_t given $G_{1:t-1}$, we have

$$\begin{aligned} \mathbb{E}_{G_t|G_{1:t-1}} [f(Z_{t+1})] &\leq f(Z_t) - \alpha \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\rangle + \frac{L}{2} \mathbb{E}_{G_t|G_{1:t-1}} [\|Z_{t+1} - Z_t\|^2] \\ &\quad + \alpha \frac{\beta_1}{1 - \beta_1} \mathbb{E}_{G_t|G_{1:t-1}} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \end{aligned} \quad (18)$$

422 Then take expectation over $G_{1:t-1}$ and rearrange, we have

$$\begin{aligned} \alpha \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\rangle \right] &\leq \mathbb{E}[f(Z_t)] - \mathbb{E}[f(Z_{t+1})] + \frac{L}{2} \mathbb{E} [\|Z_{t+1} - Z_t\|^2] \\ &\quad + \alpha \frac{\beta_1}{1 - \beta_1} \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \end{aligned} \quad (19)$$

423 In addition, we have

$$\begin{aligned} &\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\rangle \\ &= \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\rangle + \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\rangle \end{aligned} \quad (20)$$

424 and the first term on RHS of the equality can be lower bounded as

$$\begin{aligned}
& \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\rangle \\
&= \frac{1}{2} \left\| \frac{\nabla f(Z_t)}{\bar{U}_t^{1/4}} \right\|^2 + \frac{1}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i})}{\bar{U}_t^{1/4}} \right\|^2 - \frac{1}{2} \left\| \frac{\nabla f(Z_t) - \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i})}{\bar{U}_t^{1/4}} \right\|^2 \\
&\geq \frac{1}{4} \left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 + \frac{1}{4} \left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{1}{2} \left\| \frac{\nabla f(Z_t) - \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i})}{\bar{U}_t^{1/4}} \right\|^2 \\
&\quad - \frac{1}{2} \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{1}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \\
&\geq \frac{1}{2} \left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \quad (21)
\end{aligned}$$

425 where the inequalities are all due to Cauchy-Schwartz.

426 Substituting (21) and (20) into (19), we get

$$\begin{aligned}
\frac{1}{2} \alpha \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq \mathbb{E}[f(Z_t)] - \mathbb{E}[f(Z_{t+1})] + \frac{L}{2} \mathbb{E}[\|Z_{t+1} - Z_t\|^2] \\
&\quad + \alpha \frac{\beta_1}{1 - \beta_1} \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \\
&\quad - \alpha \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\rangle \right] \\
&\quad + \frac{3}{2} \alpha \mathbb{E} \left[\left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \quad (22)
\end{aligned}$$

427 Then sum over the above inequality from $t = 1$ to T and divide both sides by $T\alpha/2$, we have

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq \frac{2}{T\alpha} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) + \frac{L}{T\alpha} \sum_{t=1}^T \mathbb{E}[\|Z_{t+1} - Z_t\|^2] \\
&\quad + \frac{2}{T} \frac{\beta_1}{1 - \beta_1} \underbrace{\sum_{t=1}^T \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right]}_{T_1} \\
&\quad + \frac{2}{T} \sum_{t=1}^T \underbrace{\mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{\bar{U}_t}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right]}_{T_2} \\
&\quad + \frac{3}{T} \sum_{t=1}^T \underbrace{\mathbb{E} \left[\left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right]}_{T_3} \quad (23)
\end{aligned}$$

428 Now we need to upper bound all the terms on RHS of the above inequality to get the convergence
429 rate.

430 For terms in T_3 in (23), we can upper bound them by

$$\left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \leq \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \|\nabla f(Z_t) - \nabla f(\bar{X}_t)\|^2 \leq L \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \underbrace{\|Z_t - \bar{X}_t\|^2}_{T_4} \quad (24)$$

431 and

$$\begin{aligned} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 &\leq \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \frac{1}{N} \sum_{i=1}^N \|\nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)\|^2 \\ &\leq L \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \frac{1}{N} \underbrace{\sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2}_{T_5} \end{aligned} \quad (25)$$

432 using Jensen's inequality, Lipschitz continuity of f_i , and the fact that $f = \frac{1}{N} \sum_{i=1}^N f_i$.

433 What we need to do next is to bound T_4 and T_5 and we will bound T_5 first.

434 Before we proceed into bounding T_5 , we need some preparations. Let's recall the update rule of X_t ,
435 we have

$$X_t = X_{t-1}W - \alpha \frac{M_{t-1}}{\sqrt{U_{t-1}}} = X_1 W^{t-1} - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} W^k \quad (26)$$

436 where we define $W^0 = \mathbf{I}$.

437 Since W is a symmetric matrix, we can decompose it as $W = Q\Lambda Q^T$ where Q is a orthonormal
438 matrix and Λ is a diagonal matrix whose diagonal elements correspond to eigenvalues of W in an
439 descending order, i.e. $\Lambda_{ii} = \lambda_i$ with λ_i being i th largest eigenvalue of W . In addition, because W is
440 a doubly stochastic matrix, we know $\lambda_1 = 1$ and $q_1 = \frac{1}{\sqrt{N}}$

441 With eigen-decomposition of W , we can rewrite T_5 as

$$\sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 = \|X_t - \bar{X}_t \mathbf{1}_N^T\|_F^2 = \|X_t Q Q^T - X_t \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T\|_F^2 = \sum_{l=2}^N \|X_t q_l\|^2 \quad (27)$$

442 In addition, we can rewrite (26) as

$$X_t = X_1 W^{t-1} - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} W^k = X_1 - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} Q \Lambda^k Q^T \quad (28)$$

443 where the last equality is because $x_{1,i} = x_{1,j}$, $\forall i, j$ and thus $X_1 W = X_1$.

444 Then we have when $l > 1$,

$$X_t q_l = (X_1 - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} Q \Lambda^k Q^T) q_l = -\alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} q_l \lambda_l^k \quad (29)$$

445 because Q is orthonormal and $X_1 q_l = x_{1,1} \mathbf{1}_N^T q_l = x_{1,1} \sqrt{N} q_1^T q_l = 0, \forall l \neq 1$.

446 Combining (27) and (29), we have

$$T_5 = \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 = \sum_{l=2}^N \|X_t q_l\|^2 = \sum_{l=2}^N \alpha^2 \left\| \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} \lambda_l^k q_l \right\|^2 \leq \alpha^2 \left(\frac{1}{1-\lambda} \right)^2 N d G_\infty^2 \frac{1}{\epsilon} \quad (30)$$

447 where the last inequality follows from the fact that $g_{t,i} \leq G_\infty$, $\|q_l\| = 1$, and $|\lambda_l| \leq \lambda < 1$.

448 Now let us turn to T_4 , it can be rewritten as

$$\|Z_t - \bar{X}_t\|^2 = \left\| \frac{\beta_1}{1 - \beta_1} (\bar{X}_t - \bar{X}_{t-1}) \right\|^2 = \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right\|^2 \leq \left(\frac{\beta_1}{1 - \beta_1} \right)^2 \alpha^2 d \frac{G_\infty^2}{\epsilon} \quad (31)$$

449 Now we know both T_4 and T_5 are in the order of $\mathcal{O}(\alpha^2)$ and thus T_3 is in the order of $\mathcal{O}(\alpha^2)$.

450 Next we will bound T_2 and T_1 . Define $G_1 \triangleq \max_{t \in [T]} \max_{i \in [N]} \|\nabla f_i(x_{t,i})\|_\infty$, $G_2 \triangleq$
 451 $\max_{t \in [T]} \|\nabla f(Z_t)\|_\infty$, $G_3 \triangleq \max_{t \in [T]} \max_{i \in [N]} \|g_{t,i}\|_\infty$ and $G_\infty = \max(G_1, G_2, G_3)$

452 Then we have

$$\begin{aligned} T_2 &= \sum_{t=1}^T \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{\bar{U}_t}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \\ &\leq \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{\sqrt{[\bar{U}_t]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{\sqrt{[\bar{U}_t]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| \frac{\sqrt{[\bar{U}_t]_j} + \sqrt{[u_{t,i}]_j}}{\sqrt{[\bar{U}_t]_j} + \sqrt{[u_{t,i}]_j}} \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{[\bar{U}_t]_j - [u_{t,i}]_j}{[\bar{U}_t]_j \sqrt{[u_{t,i}]_j} + \sqrt{[\bar{U}_t]_j} [u_{t,i}]_j} \right| \right] \\ &\leq \mathbb{E} \left[\underbrace{\sum_{t=1}^T G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{[\bar{U}_t]_j - [u_{t,i}]_j}{2\epsilon^{1.5}} \right|}_{T_6} \right] \end{aligned} \quad (32)$$

453 where the last inequality is due to $[u_{t,i}]_j \geq \epsilon$, $\forall t, i, j$.

454 To simplify notations, let's define $\|A\|_{abs} = \sum_{i,j} |A_{ij}|$ to be the entry-wise L_1 norm of a matrix A ,
 455 then we have

$$\begin{aligned} T_6 &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \|\bar{U}_t \mathbf{1}^T - U_t\|_{abs} \\ &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \|\tilde{U}_t \mathbf{1}^T - \tilde{U}_t\|_{abs} \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \|\tilde{U}_t \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T - \tilde{U}_t Q Q^T\|_{abs} \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \tilde{U}_t \sum_{l=2}^N q_l q_l^T \right\|_{abs} \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \end{aligned}$$

456 where the second inequality is due to Lemma 2 and the fact that $U_t = \max(\tilde{U}_t, \epsilon)$ element-wisely.

457 **Lemma 2.** Given a set of numbers a_1, \dots, a_n and denote their mean to be $\bar{a} = \frac{1}{n} \sum_{i=1}^n a_i$. In
 458 addition, define $b_i(r) \triangleq \max(a_i, r)$ and $\bar{b}(r) = \frac{1}{n} \sum_{i=1}^n b_i(r)$. For any r and r' with $r' \geq r$ we
 459 have

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| \geq \sum_{i=1}^n |b_i(r') - \bar{b}(r')| \quad (33)$$

460 and when $r \leq \min_{i \in [n]} a_i$, we have

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| = \sum_{i=1}^n |a_i - \bar{a}| \quad (34)$$

461 **Proof:** See Appendix A.3. □

462 Recall from update rule of U_t , by defining $\hat{V}_{-1} \triangleq \hat{V}_0$ and $U_0 \triangleq U_{1/2}$, we have $\forall t \geq 0$

$$\tilde{U}_{t+1} = (\tilde{U}_t - \hat{V}_{t-1} + \hat{V}_t)W \quad (35)$$

463 and thus

$$\tilde{U}_t = \tilde{U}_0 W^t + \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k})W^k = \tilde{U}_0 + \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k})Q\Lambda^k Q^T \quad (36)$$

464 Then we further have when $l \neq 1$,

$$\tilde{U}_{tql} = (\tilde{U}_0 + \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k})Q\Lambda^k Q^T)q_l = \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k})q_l \lambda_l^k \quad (37)$$

465 where the last equality is due to the definition $\tilde{U}_0 \triangleq U_{1/2} = \epsilon \mathbf{1}_d \mathbf{1}_N^T = \sqrt{N} \epsilon \mathbf{1}_d \mathbf{1}_N^T$ (recall that
466 $q_1 = \frac{1}{\sqrt{N}} \mathbf{1}_N^T$) and $q_i^T q_j = 0$ when $i \neq j$.

467 Note by definition of $\|\cdot\|_{abs}$, we have $\forall A, B, \|A + B\|_{abs} \leq \|A\|_{abs} + \|B\|_{abs}$, then we have

$$\begin{aligned} T_6 &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \sum_{l=2}^N \tilde{U}_{tql} q_l^T \right\|_{abs} \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k}) \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_{abs} \\ &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \left\| (-\hat{V}_{t-1-k} + \hat{V}_{t-k}) \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_{abs} \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T (-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j \right\|_1 \\ &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_1 \left\| (-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j \right\|_1 \\ &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \sqrt{N} \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_2 \left\| (-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j \right\|_1 \\ &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \left\| (-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j \right\|_1 \sqrt{N} \lambda^k \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \left\| (-\hat{V}_{t-1-k} + \hat{V}_{t-k}) \right\|_{abs} \sqrt{N} \lambda^k \\ &= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{t-1} \left\| (-\hat{V}_{o-1} + \hat{V}_o) \right\|_{abs} \sqrt{N} \lambda^{t-o} \\ &= \frac{G_\infty^2}{N} \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{T-1} \sum_{t=o+1}^T \left\| (-\hat{V}_{o-1} + \hat{V}_o) \right\|_{abs} \sqrt{N} \lambda^{t-o} \\ &\leq \frac{G_\infty^2}{\sqrt{N}} \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \left\| (-\hat{V}_{o-1} + \hat{V}_o) \right\|_{abs} \end{aligned} \quad (38)$$

468 where $\lambda = \max(|\lambda_2|, |\lambda_N|)$.

469 Combining (32) and (38), we have

$$T_2 \leq \frac{G_\infty^2}{\sqrt{N}} \frac{1}{2\epsilon^{1.5}} \frac{\lambda}{1-\lambda} \mathbb{E} \left[\sum_{o=0}^{T-1} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] \quad (39)$$

470 Now we need to bound T_1 , we have

$$\begin{aligned} T_1 &= \sum_{t=1}^T \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \\ &\leq \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{\sqrt{[u_{t-1,i}]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \left(\frac{1}{\sqrt{[u_{t-1,i}]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right) \frac{\sqrt{[u_{t,i}]_j} + \sqrt{[u_{t-1,i}]_j}}{\sqrt{[u_{t,i}]_j} + \sqrt{[u_{t-1,i}]_j}} \right| \right] \\ &\leq \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{2\epsilon^{1.5}} ([u_{t-1,i}]_j - [u_{t,i}]_j) \right| \right] \\ &\stackrel{(a)}{\leq} \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{2\epsilon^{1.5}} |([\tilde{u}_{t-1,i}]_j - [\tilde{u}_{t,i}]_j)| \right] \\ &= G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{N} \mathbb{E} \left[\sum_{t=1}^T \|\tilde{U}_{t-1} - \tilde{U}_t\|_{abs} \right] \end{aligned} \quad (40)$$

471 where (a) is due to $[\tilde{u}_{t-1,i}]_j = \max([u_{t-1,i}]_j, \epsilon)$ and the function $\max(\cdot, \epsilon)$ is 1-Lipschitz.

472 In addition, by update rule of U_t , we have

$$\begin{aligned} &\sum_{t=1}^T \|\tilde{U}_{t-1} - \tilde{U}_t\|_{abs} \\ &= \sum_{t=1}^T \|\tilde{U}_{t-1} - (\tilde{U}_{t-1} - \hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\ &= \sum_{t=1}^T \|\tilde{U}_{t-1}(I - W) + (-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\ &= \sum_{t=1}^T \|\tilde{U}_{t-1}(QQ^T - Q\Lambda Q^T) + (-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\ &= \sum_{t=1}^T \|\tilde{U}_{t-1}(\sum_{l=2}^N q_l(1 - \lambda_l)q_l^T) + (-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\ &\leq \sum_{t=1}^T \left\| \sum_{k=1}^{t-1} (-\hat{V}_{t-2-k} + \hat{V}_{t-1-k}) \sum_{l=2}^N q_l \lambda_l^k (1 - \lambda_l) q_l^T \right\|_{abs} + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\ &\leq \sum_{t=1}^T \left(\sum_{k=1}^{t-1} \|-\hat{V}_{t-2-k} + \hat{V}_{t-1-k}\|_{abs} \sqrt{N} \lambda^k \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \\ &= \sum_{t=1}^T \left(\sum_{o=1}^{t-1} \|-\hat{V}_{o-2} + \hat{V}_{o-1}\|_{abs} \sqrt{N} \lambda^{t-o} \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \end{aligned}$$

$$\begin{aligned}
&= \sum_{o=1}^{T-1} \sum_{t=o+1}^T \left(\| -\hat{V}_{o-2} + \hat{V}_{o-1} \|_{abs} \sqrt{N} \lambda^{t-o} \right) + \sum_{t=1}^T \| (-\hat{V}_{t-2} + \hat{V}_{t-1}) \|_{abs} \\
&\leq \sum_{o=1}^{T-1} \frac{\lambda}{1-\lambda} \left(\| -\hat{V}_{o-2} + \hat{V}_{o-1} \|_{abs} \sqrt{N} \right) + \sum_{t=1}^T \| (-\hat{V}_{t-2} + \hat{V}_{t-1}) \|_{abs} \\
&\leq \frac{1}{1-\lambda} \sum_{t=1}^T \| (-\hat{V}_{t-2} + \hat{V}_{t-1}) \|_{abs} \sqrt{N}
\end{aligned} \tag{41}$$

473 Combining (40) and (41), we have

$$T_1 \leq G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{N} \mathbb{E} \left[\frac{1}{1-\lambda} \sum_{t=1}^T \| (-\hat{V}_{t-2} + \hat{V}_{t-1}) \|_{abs} \sqrt{N} \right] \tag{42}$$

474 What remains is to bound $\sum_{t=1}^T \mathbb{E} [\|Z_{t+1} - Z_t\|^2]$. By update rule of Z_t , we have

$$\begin{aligned}
&\|Z_{t+1} - Z_t\|^2 \\
&= \left\| \alpha \frac{\beta_1}{1-\beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left\| \frac{\beta_1}{1-\beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\|^2 + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{\sqrt{\epsilon}} \left| \frac{1}{\sqrt{[u_{t-1,i}]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{\sqrt{\epsilon}} \left| \frac{[u_{t,i}]_j - [u_{t-1,i}]_j}{2\epsilon^{1.5}} \right| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{2\epsilon^2} |\tilde{u}_{t,i,j} - \tilde{u}_{t-1,i,j}| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&= 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \frac{1}{2\epsilon^2} \|\tilde{U}_t - \tilde{U}_{t-1}\|_{abs} + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2
\end{aligned} \tag{43}$$

475 where the last inequality is again due to the definition that $[\tilde{u}_{t,i}]_j = \max([u_{t,i}]_j, \epsilon)$ and the fact that
476 $\max(\cdot, \epsilon)$ is 1-Lipschitz.

477 Then, we have

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} [\|Z_{t+1} - Z_t\|^2] \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{t=1}^T \|\tilde{U}_t - \tilde{U}_{t-1}\|_{abs} \right] + 2\alpha^2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&\leq \alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{G_\infty^2}{\sqrt{N}} \frac{1}{\epsilon^2} \frac{1}{1-\lambda} \mathbb{E} \left[\sum_{t=1}^T \| (-\hat{V}_{t-2} + \hat{V}_{t-1}) \|_{abs} \right] + 2\alpha^2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right]
\end{aligned} \tag{44}$$

478 where the last inequality is due to (41).

479 Now let's bound the last term on RHS of the above inequality. A trivial bound can be

$$\sum_{t=1}^T \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \leq \sum_{t=1}^T d G_\infty^2 \frac{1}{\epsilon}$$

480 due to $\|g_{t,i}\| \leq G_\infty$ and $[u_{t,i}]_j \geq \epsilon, \forall j$ (this is easy to verify from update rule of $u_{t,i}$ and the
 481 assumption that $[v_{t,i}]_j \geq \epsilon, \forall i$). However, the above bound is independent of N , to get a better bound,
 482 we need a more involved analysis to show its dependency on N . To do this, we first notice that

$$\begin{aligned}
 & \mathbb{E}_{G_t|G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
 &= \mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left\langle \frac{\nabla f_i(x_{t,i}) + \xi_{t,i}}{\sqrt{u_{t,i}}}, \frac{\nabla f_j(x_{t,j}) + \xi_{t,j}}{\sqrt{u_{t,j}}} \right\rangle \right] \\
 &\stackrel{(a)}{=} \mathbb{E}_{G_t|G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] + \mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N^2} \sum_{i=1}^N \left\| \frac{\xi_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
 &\stackrel{(b)}{=} \left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{1}{N^2} \sum_{i=1}^N \sum_{l=1}^d \frac{\mathbb{E}_{G_t|G_{1:t-1}} [[\xi_{t,i}]_l^2]}{[u_{t,i}]_l} \\
 &\stackrel{(c)}{\leq} \left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{d}{N} \frac{\sigma^2}{\epsilon}
 \end{aligned} \tag{45}$$

483 where (a) is due to $\mathbb{E}_{G_t|G_{1:t-1}} [\xi_{t,i}] = 0$ and $\xi_{t,i}$ is independent of $x_{t,j}, \forall j, u_{t,j}, \forall j$, and $\xi_j, \forall j \neq i$,
 484 (b) comes from the fact that $x_{t,i}, u_{t,i}$ are fixed given $G_{1:t}$, (c) is due to $\mathbb{E}_{G_t|G_{1:t-1}} [[\xi_{t,i}]_l^2] \leq \sigma^2$ and
 485 $[u_{t,i}]_l \geq \epsilon$ by definition.

486 Then we have

$$\begin{aligned}
 \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] &= \mathbb{E}_{G_{1:t-1}} \left[\mathbb{E}_{G_t|G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \right] \\
 &\leq \mathbb{E}_{G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{d}{N} \frac{\sigma^2}{\epsilon} \right] \\
 &= \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] + \frac{d}{N} \frac{\sigma^2}{\epsilon}
 \end{aligned} \tag{46}$$

487 In traditional analysis of SGD-like distributed algorithms, the term corresponding to
 488 $\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right]$ will be merged with the first order descent when the stepsize is cho-
 489 sen to be small enough. However, in our case, the term cannot be merged because it is different from
 490 the first order descent in our algorithm. A brute-force upper bound is possible but this will lead to a
 491 worse convergence rate in terms of N . Thus, we need a more detailed analysis for the term in the
 492 following.

$$\begin{aligned}
 \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] &= \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} + \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right] \\
 &\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right] \\
 &\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left\| \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right]
 \end{aligned}$$

$$\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] + 2\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{U_t}} \right\|_1 \right] \quad (47)$$

Summing over T , we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] \\ & \leq 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{U_t}} \right\|_1 \right] \end{aligned} \quad (48)$$

For the last term on RHS of (48), we can bound it similarly as what we did for T_2 from (32) to (38), which yields

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{U_t}} \right\|_1 \right] \\ & \leq \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \frac{1}{2\epsilon^{1.5}} \|u_{t,i} - U_t\|_1 \right] \\ & = \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} G_\infty^2 \frac{1}{2\epsilon^2} \|\bar{U}_t \mathbf{1}^T - U_t\|_{abs} \right] \\ & \leq \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} G_\infty^2 \frac{1}{2\epsilon^2} \left\| -\sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \right] \\ & \leq \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] \end{aligned} \quad (49)$$

Further, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] \\ & \leq 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t)}{\sqrt{U_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] \\ & = 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\sqrt{U_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] \end{aligned} \quad (50)$$

and the last term on RHS of the above inequality can be bounded following similar procedures from (25) to (30), as what we did for T_3 . Completing the procedures yields

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{U_t}} \right\|^2 \right] \\ & \leq \sum_{t=1}^T \mathbb{E} \left[L \frac{1}{\epsilon} \frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 \right] \\ & \leq \sum_{t=1}^T \mathbb{E} \left[L \frac{1}{\epsilon} \frac{1}{N} \alpha^2 \left(\frac{1}{1-\lambda} \right) N d G_\infty^2 \frac{1}{\epsilon} \right] \\ & = T L \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) d G_\infty^2 \end{aligned} \quad (51)$$

499 Finally, combining (46) to (51), we get

$$\begin{aligned}
& \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
& \leq 4 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 4TL \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) dG_\infty^2 \\
& \quad + 2 \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] + T \frac{d}{N} \frac{\sigma^2}{\epsilon} \\
& \leq 4 \frac{1}{\sqrt{\epsilon}} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] + 4TL \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) dG_\infty^2 \\
& \quad + 2 \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] + T \frac{d}{N} \frac{\sigma^2}{\epsilon}. \tag{52}
\end{aligned}$$

500 where the last inequality is due to each element of \bar{U}_t is lower bounded by ϵ by definition.

501 Combining all above, we can have

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \\
& \leq \frac{2}{T\alpha} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) \\
& \quad + \frac{L}{T} \alpha \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{G_\infty^2}{\sqrt{N}} \frac{1}{\epsilon^2} \frac{1}{1-\lambda} \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
& \quad + \frac{8L}{T} \alpha \frac{1}{\sqrt{\epsilon}} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] + 8L^2 \alpha \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) dG_\infty^2 \\
& \quad + \frac{4L}{T} \alpha \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] + 2L\alpha \frac{d}{N} \frac{\sigma^2}{\epsilon} \\
& \quad + \frac{2}{T} \frac{\beta_1}{1-\beta_1} G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{\sqrt{N}} \mathbb{E} \left[\frac{1}{1-\lambda} \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
& \quad + \frac{2}{T} \frac{G_\infty^2}{\sqrt{N}} \frac{1}{2\epsilon^{1.5}} \frac{\lambda}{1-\lambda} \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
& \quad + \frac{3}{T} \left(\sum_{t=1}^T L \left(\frac{1}{1-\lambda} \right)^2 \alpha^2 dG_\infty^2 \frac{1}{\epsilon^{1.5}} + \sum_{t=1}^T L \left(\frac{\beta_1}{1-\beta_1} \right)^2 \alpha^2 d \frac{G_\infty^2}{\epsilon^{1.5}} \right) \\
& = \frac{2}{T\alpha} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) + 2L\alpha \frac{d}{N} \frac{\sigma^2}{\epsilon} + 8L\alpha \frac{1}{\sqrt{\epsilon}} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \\
& \quad + 3\alpha^2 d \left(\left(\frac{\beta_1}{1-\beta_1} \right)^2 + \left(\frac{1}{1-\lambda} \right)^2 \right) L \frac{G_\infty^2}{\epsilon^{1.5}} + 8\alpha^3 L^2 \left(\frac{1}{1-\lambda} \right) d \frac{G_\infty^2}{\epsilon^2} \\
& \quad + \frac{1}{T\epsilon^{1.5}} \frac{G_\infty^2}{\sqrt{N}} \frac{1}{1-\lambda} \left(L\alpha \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{1}{\epsilon^{0.5}} + \lambda + \frac{\beta_1}{1-\beta_1} + 2L\alpha \frac{1}{\epsilon^{0.5}} \lambda \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right]. \tag{53}
\end{aligned}$$

502 Set $\alpha = \frac{1}{\sqrt{dT}}$ and when $\alpha \leq \frac{\epsilon^{0.5}}{16L}$, we further have

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \\
& \leq \frac{4}{T\alpha} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) + 4L\alpha \frac{d}{N} \frac{\sigma^2}{\epsilon} \\
& \quad + 6\alpha^2 d \left(\left(\frac{\beta_1}{1-\beta_1} \right)^2 + \left(\frac{1}{1-\lambda} \right)^2 \right) L \frac{G_\infty^2}{\epsilon^{1.5}} + 16\alpha^3 L^2 \left(\frac{1}{1-\lambda} \right) d \frac{G_\infty^2}{\epsilon^2} \\
& \quad + \frac{2}{T\epsilon^{1.5}} \frac{G_\infty^2}{\sqrt{N}} \frac{1}{1-\lambda} \left(L\alpha \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{1}{\epsilon^{0.5}} + \lambda + \frac{\beta_1}{1-\beta_1} + 2L\alpha \frac{1}{\epsilon^{0.5}} \lambda \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
& = \frac{4\sqrt{d}}{\sqrt{T}} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) + 4L \frac{\sqrt{d}}{\sqrt{T}} \frac{1}{N} \frac{\sigma^2}{\epsilon} \\
& \quad + 6 \frac{1}{T} \left(\left(\frac{\beta_1}{1-\beta_1} \right)^2 + \left(\frac{1}{1-\lambda} \right)^2 \right) L \frac{G_\infty^2}{\epsilon^{1.5}} + 16 \frac{1}{T^{1.5}d^{0.5}} L^2 \left(\frac{1}{1-\lambda} \right) \frac{G_\infty^2}{\epsilon^2} \\
& \quad + \frac{2}{T\epsilon^{1.5}} \frac{G_\infty^2}{\sqrt{N}} \frac{1}{1-\lambda} \left(\frac{L}{\sqrt{Td}} \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{1}{\epsilon^{0.5}} + \lambda + \frac{\beta_1}{1-\beta_1} + 2 \frac{L}{\sqrt{Td}} \frac{1}{\epsilon^{0.5}} \lambda \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
& \leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(Z_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5}d^{0.5}} C_3 \\
& \quad + \left(\frac{1}{TN^{0.5}} C_4 + \frac{1}{T^{1.5}d^{0.5}N^{0.5}} C_5 \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \tag{54}
\end{aligned}$$

503 where the first inequality is obtained by moving the term $8L\alpha \frac{1}{\sqrt{\epsilon}} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right]$ on the

504 RHS of (53) to the LHS to cancel it using the assumption $8L\alpha \frac{1}{\sqrt{\epsilon}} \leq \frac{1}{2}$ followed by multiplying both

505 sides by 2, and the constants introduced in the last step are defined as following

$$\begin{aligned}
C_1 &= \max(4, 4L/\epsilon) \\
C_2 &= 6 \left(\left(\frac{\beta_1}{1-\beta_1} \right)^2 + \left(\frac{1}{1-\lambda} \right)^2 \right) L \frac{G_\infty^2}{\epsilon^{1.5}} \\
C_3 &= 16L^2 \left(\frac{1}{1-\lambda} \right) \frac{G_\infty^2}{\epsilon^2} \\
C_4 &= \frac{2}{\epsilon^{1.5}} \frac{1}{1-\lambda} \left(\lambda + \frac{\beta_1}{1-\beta_1} \right) G_\infty^2 \\
C_5 &= \frac{2}{\epsilon^2} \frac{1}{1-\lambda} L \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 + \frac{4}{\epsilon^2} \frac{\lambda}{1-\lambda} LG_\infty^2. \tag{55}
\end{aligned}$$

506 Substituting into $Z_1 = \bar{X}_1$ completes the proof \square

507 A.2 Proof of Theorem 3

508 By Theorem 2, we know under the assumptions of the theorem, we have

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5}d^{0.5}} C_3 \\
& \quad + \left(\frac{1}{TN^{0.5}} C_4 + \frac{1}{T^{1.5}d^{0.5}N^{0.5}} C_5 \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \tag{56}
\end{aligned}$$

where $\|\cdot\|_{abs}$ denotes the entry-wise L_1 norm of a matrix (i.e. $\|A\|_{abs} = \sum_{i,j} |A_{ij}|$) and C_1, C_2, C_3, C_4, C_5 are defined in Theorem 2.

Since Algorithm 3 is a special case of 2, building on result of Theorem 2, we just need to characterize the growth speed of $\mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right]$ to prove convergence of Algorithm 3. By the update rule of Algorithm 3, we know \hat{V}_t is non decreasing and thus

$$\begin{aligned}
& \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
&= \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^d |-\hat{v}_{t-2,i,j} + \hat{v}_{t-1,i,j}| \right] \\
&= \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^d (-\hat{v}_{t-2,i,j} + \hat{v}_{t-1,i,j}) \right] \\
&= \mathbb{E} \left[\sum_{i=1}^N \sum_{j=1}^d (-\hat{v}_{-1,i,j} + \hat{v}_{T-1,i,j}) \right] \\
&= \mathbb{E} \left[\sum_{i=1}^N \sum_{j=1}^d (-\hat{v}_{0,i,j} + \hat{v}_{T-1,i,j}) \right]
\end{aligned} \tag{57}$$

where the last equality is because we defined $\hat{V}_{-1} \triangleq \hat{V}_0$ previously.

Further, because $\|g_{t,i}\|_\infty \leq G_\infty, \forall t, i$ and $v_{t,i}$ is a exponential moving average of $g_{k,i}^2, k = 1, 2, \dots, t$, we know $|\hat{v}_{t,i,j}| \leq G_\infty^2, \forall t, i, j$. In addition, by update rule of \hat{V}_t , we also know each element of \hat{V}_t also cannot be greater than G_∞^2 , i.e. $|\hat{v}_{t,i,j}| \leq G_\infty^2, \forall t, i, j$.

Given the fact that $\hat{v}_{0,i,j} \geq 0$, we have

$$\mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] = \mathbb{E} \left[\sum_{i=1}^N \sum_{j=1}^d (-\hat{v}_{0,i,j} + \hat{v}_{T-1,i,j}) \right] \leq \mathbb{E} \left[\sum_{i=1}^N \sum_{j=1}^d G_\infty^2 \right] = NdG_\infty^2$$

Substituting the above into (56), we have

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5} d^{0.5}} C_3 \\
&\quad + \frac{d}{T} C_4 \sqrt{N} G_\infty^2 + \frac{\sqrt{d}}{T^{1.5}} C_5 \sqrt{N} G_\infty^2 \\
&= C'_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C'_2 + \frac{1}{T^{1.5} d^{0.5}} C'_3 \\
&\quad + \frac{d}{T} \sqrt{N} C'_4 + \frac{\sqrt{d}}{T^{1.5}} \sqrt{N} C'_5
\end{aligned} \tag{58}$$

where we have

$$\begin{aligned}
C'_1 &= C_1 \\
C'_2 &= C_2 \\
C'_3 &= C_3 \\
C'_4 &= C_4 G_\infty^2 \\
C'_5 &= C_5 G_\infty^2
\end{aligned} \tag{59}$$

The proof is complete. \square

522 A.3 Proof of Lemmas

523 **Lemma 1.** For the sequence defined in (14), we have

$$Z_{t+1} - Z_t = \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \quad (60)$$

524 **Proof:** By update rule of Algorithm 2, we first have

$$\begin{aligned} \bar{X}_{t+1} &= \frac{1}{N} \sum_{i=1}^N x_{t+1,i} \\ &= \frac{1}{N} \sum_{i=1}^N \left(x_{t+0.5,i} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left(\sum_{j=1}^N W_{ij} x_{t,j} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) \\ &\stackrel{(i)}{=} \left(\frac{1}{N} \sum_{j=1}^N x_{t,j} \right) - \frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \\ &= \bar{X}_t - \frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \end{aligned} \quad (61)$$

525 where (i) is due to an interchange of summation and $\sum_{i=1}^N W_{ij} = 1$.

526 Then, we have

$$\begin{aligned} Z_{t+1} - Z_t &= \bar{X}_{t+1} - \bar{X}_t + \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) - \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) \\ &= \frac{1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) - \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) \\ &= \frac{1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) - \frac{\beta_1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right) \\ &= \frac{1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{\beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}}{\sqrt{u_{t,i}}} \right) - \frac{\beta_1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right) \\ &= \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \end{aligned} \quad (62)$$

527 which is the desired result. \square

528 **Lemma 2.** Given a set of numbers a_1, \dots, a_n and denote their mean to be $\bar{a} = \frac{1}{n} \sum_{i=1}^n a_i$. In
529 addition, define $b_i(r) \triangleq \max(a_i, r)$ and $\bar{b}(r) = \frac{1}{n} \sum_{i=1}^n b_i(r)$. For any r and r' with $r' \geq r$ we
530 have

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| \geq \sum_{i=1}^n |b_i(r') - \bar{b}(r')| \quad (63)$$

531 and when $r \leq \min_{i \in [n]} a_i$, we have

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| = \sum_{i=1}^n |a_i - \bar{a}| \quad (64)$$

532 **Proof:** Without loss of generality, let's assume $a_i \leq a_j$ when $i < j$, i.e. a_i is a non-decreasing
533 sequence. Define

$$h(r) = \sum_{i=1}^n |b_i(r) - \bar{b}(r)| = \sum_{i=1}^n \left| \max(a_i, r) - \frac{1}{n} \sum_{j=1}^n \max(a_j, r) \right|, \quad (65)$$

we need to prove that h is a non-increasing function of r . First, it is easy to see that h is a continuous function of r with non-differentiable points $r = a_i, i \in [n]$, thus h is a piece-wise linear function.

Next, we will prove that $h(r)$ is non-increasing in each piece. Define $l(r)$ to be the largest index with $a(l(r)) < r$, and $s(r)$ to be the largest index with $a_{s(r)} < \bar{b}(r)$. Note that we have $b_i(r) = r, \forall i \leq l(r)$ and $b_i(r) - \bar{b}(r) \leq 0, \forall i \leq s(r)$ because a_i is a non-decreasing sequence. Therefore, we have

$$h(r) = \sum_{i=1}^{l(r)} (\bar{b}(r) - r) + \sum_{i=l(r)+1}^{s(r)} (\bar{b}(r) - a_i) + \sum_{i=s(r)+1}^n (a_i - \bar{b}(r)). \quad (66)$$

and

$$\bar{b}(r) = \frac{1}{n} \left(l(r)r + \sum_{i=l(r)+1}^n a_i \right) \quad (67)$$

Taking derivative of the above form, we know the derivative of $h(r)$ at differentiable points is

$$\begin{aligned} h'(r) &= l(r) \left(\frac{l(r)}{n} - 1 \right) + (s(r) - l(r)) \frac{l(r)}{n} - (n - s(r)) \frac{l(r)}{n} \\ &= \frac{l(r)}{n} ((l(r) - n) + (s(r) - l(r)) - (n - s(r))) \end{aligned} \quad (68)$$

Since we have $s(r) \leq n$ we know $(l(r) - n) + (s(r) - l(r)) - (n - s(r)) \leq 0$ and thus

$$h'(r) \leq 0 \quad (69)$$

which means $h(r)$ is non-increasing in each piece. Combining with the fact that $h(r)$ is continuous, (64) is proven.

When $r \leq a(i)$, we have $b(i) = \max(a_i, r) = r, \forall r \in [n]$ and $\bar{b}(r) = \frac{1}{n} \sum_{i=1}^n a_i = \bar{a}$ which proves (65). \square

A.4 Additional experiments and details

In this section, we compare the learning curves of different algorithms with different stepsizes on heterogeneous data distribution. We use 5 nodes and the heterogeneous data distribution is created by assigning each node with data of only two labels and there are no overlapping labels between different nodes. For all algorithms, we compare stepsizes in the set [1e-1, 1e-2, 1e-3, 1e-4, 1e-5, 1e-6].

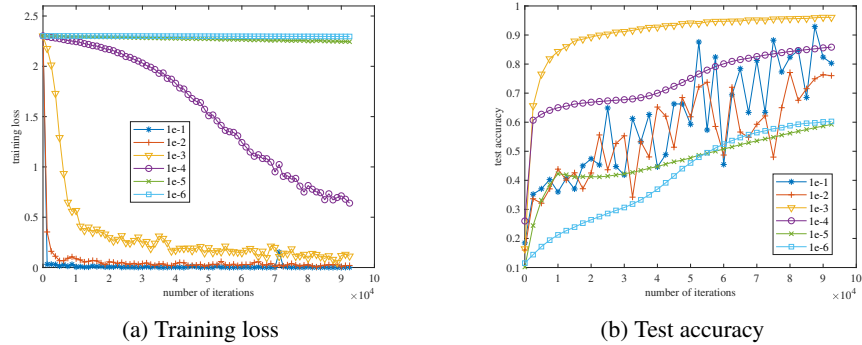


Figure 2: Performance comparison of different stepsizes for DGD

Figure 2 shows the training loss and test accuracy of DGD, it can be seen that the stepsize 1e-3 works best for DGD in terms of test accuracy and 1e-1 works best in terms of training loss. The difference is caused by the inconsistency among the value of parameters on different nodes when the stepsize is large. The training loss is calculated as the average of the loss value of different local models evaluated on their local training batch. Thus, though the training loss is small evaluated at a particular

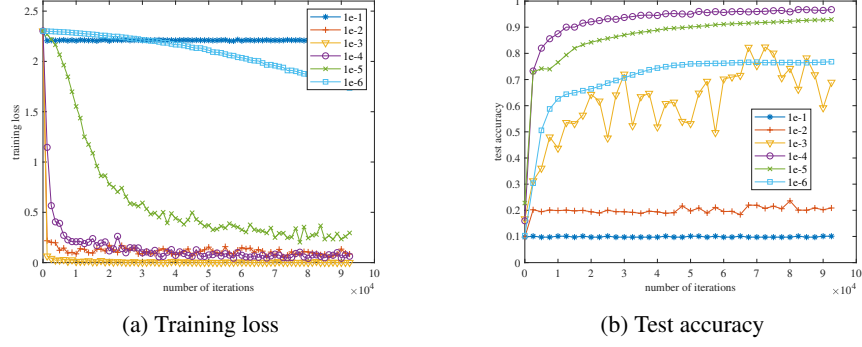


Figure 3: Performance comparison of different stepsizes for decentralized AMSGrad

node, the test accuracy will be low when evaluating data with labels not seen by the node (recall that each node contains data with different labels).

Figure 3 shows the performance of decentralized AMSGrad with different stepsizes, we can see its best performance is better than DGD and the performance is stabler (the test performance is less sensitive to stepsize choice).

Figure 4 shows the performance of DADM, as it can be expected, the performance of DADAM is not as good as DGD and decentralized AMSGrad since it is not a convergent algorithm and the heterogeneity in data amplified the non-convergence issue of DADAM.

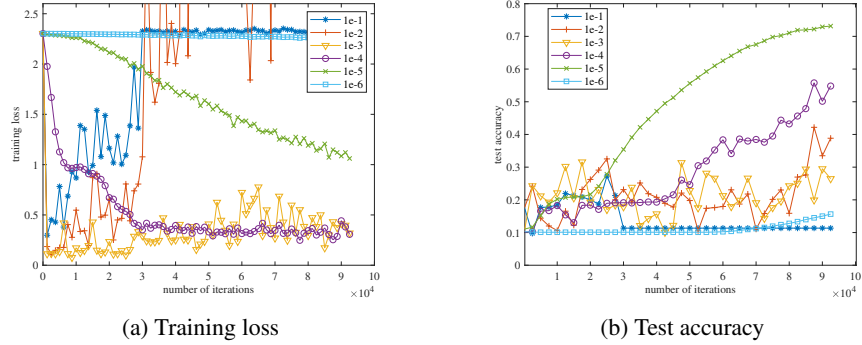


Figure 4: Performance comparison of different stepsizes for DADAM

From the experiments above, we can see the advantages of decentralized AMSGrad in terms of both performance and ease of parameter tuning, and the importance of ensuring the theoretical convergence of algorithms.