

433 A Proof of Auxiliary Lemmas

434 **Lemma 1.** *For the sequence defined in (17), we have*

$$Z_{t+1} - Z_t = \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}}. \quad (6)$$

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

$$\begin{aligned} \bar{X}_{t+1} &= \frac{1}{N} \sum_{i=1}^N x_{t+1,i} \\ &= \frac{1}{N} \sum_{i=1}^N \left(x_{t+0.5,i} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) \end{aligned} \quad (7)$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\sum_{j=1}^N W_{ij} x_{t,j} - \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) \quad (8)$$

$$\stackrel{(i)}{=} \left(\frac{1}{N} \sum_{j=1}^N x_{t,j} \right) - \frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \quad (9)$$

$$= \bar{X}_t - \frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}}, \quad (10)$$

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

437 Then, we have

$$\begin{aligned} Z_{t+1} - Z_t &= \bar{X}_{t+1} - \bar{X}_t + \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) - \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) \\ &= \frac{1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) - \frac{\beta_1}{1 - \beta_1} (\bar{X}_{t+1} - \bar{X}_t) \end{aligned} \quad (11)$$

$$= \frac{1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t,i}}{\sqrt{u_{t,i}}} \right) - \frac{\beta_1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right) \quad (12)$$

$$= \frac{1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{\beta_1 m_{t-1,i} + (1 - \beta_1) g_{t,i}}{\sqrt{u_{t,i}}} \right) - \frac{\beta_1}{1 - \beta_1} \left(-\frac{1}{N} \sum_{i=1}^N \alpha \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right) \quad (13)$$

$$= \alpha \frac{\beta_1}{1 - \beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}}, \quad (14)$$

438 which is the desired result. \square

439 **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*
 440 *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*
 441 *have*

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| \geq \sum_{i=1}^n |b_i(r') - \bar{b}(r')| \quad (15)$$

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

$$\sum_{i=1}^n |b_i(r) - \bar{b}(r)| = \sum_{i=1}^n |a_i - \bar{a}|. \quad (16)$$

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

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

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

447 Next, we will prove that $h(r)$ is non-increasing in each piece. Define $l(r)$ to be the largest index with
 448 $a(l(r)) < r$, and $s(r)$ to be the largest index with $a_{s(r)} < \bar{b}(r)$. Note that we have $b_i(r) = r, \forall i \leq l(r)$
 449 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

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

450 and

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

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

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

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

$$h'(r) \leq 0,$$

453 which means $h(r)$ is non-increasing in each piece. Combining with the fact that $h(r)$ is continuous,
 454 (15) is proven. When $r \leq a(i)$, we have $b(i) = \max(a_i, r) = r, \forall r \in [n]$ and $\bar{b}(r) = \frac{1}{n} \sum_{i=1}^n a_i = \bar{a}$
 455 which proves (16). \square

456 B Proof of Theorem 2

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

$$Z_t = \bar{X}_t + \frac{\beta_1}{1 - \beta_1} (\bar{X}_t - \bar{X}_{t-1}) \quad (17)$$

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

459 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
460 have

$$\mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right] = \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \quad (18)$$

461 By assuming smoothness (A1) we have

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

462 Using Lemma 1 into the above inequality and take expectation over G_t given $G_{1:t-1}$, we have

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

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

$$\alpha \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\rangle \right] \quad (20)$$

$$\begin{aligned} &\leq \mathbb{E}[f(Z_t)] - \mathbb{E}[f(Z_{t+1})] + \frac{L}{2} \mathbb{E} [\|Z_{t+1} - Z_t\|^2] \\ &\quad + \alpha \frac{\beta_1}{1 - \beta_1} \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right]. \end{aligned} \quad (21)$$

464 In addition, we have

$$\begin{aligned} &\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\rangle \\ &= \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\rangle + \left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\rangle \end{aligned} \quad (22)$$

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

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

$$\begin{aligned}
& -\frac{1}{2} \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{1}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \\
& \geq \frac{1}{2} \left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 - \frac{3}{2} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2, \quad (23)
\end{aligned}$$

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

467 Substituting (23) and (22) into (20), we get

$$\begin{aligned}
\frac{1}{2} \alpha \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] & \leq \mathbb{E}[f(Z_t)] - \mathbb{E}[f(Z_{t+1})] + \frac{L}{2} \mathbb{E}[\|Z_{t+1} - Z_t\|^2] \\
& + \alpha \frac{\beta_1}{1 - \beta_1} \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \\
& - \alpha \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\rangle \right] \\
& + \frac{3}{2} \alpha \mathbb{E} \left[\left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right]. \quad (24)
\end{aligned}$$

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

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] \quad (25) \\
& \leq \frac{2}{T\alpha} (\mathbb{E}[f(Z_1)] - \mathbb{E}[f(Z_{T+1})]) + \frac{L}{T\alpha} \sum_{t=1}^T \mathbb{E}[\|Z_{t+1} - Z_t\|^2] \\
& + \frac{2}{T} \frac{\beta_1}{1 - \beta_1} \sum_{t=1}^T \underbrace{\mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right]}_{T_1} \\
& + \frac{2}{T} \sum_{t=1}^T \underbrace{\mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{\bar{U}_t}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right]}_{T_2} \\
& + \frac{3}{T} \sum_{t=1}^T \underbrace{\mathbb{E} \left[\left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 + \left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right]}_{T_3}. \quad (26)
\end{aligned}$$

469 Now we need to upper bound all the terms on RHS of the above inequality to get the convergence
470 rate. For the terms composing T_3 in (25), we can upper bound them by

$$\begin{aligned}
\left\| \frac{\nabla f(Z_t) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 & \leq \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \|\nabla f(Z_t) - \nabla f(\bar{X}_t)\|^2 \\
& \leq L \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \underbrace{\|Z_t - \bar{X}_t\|^2}_{T_4} \quad (27)
\end{aligned}$$

471 and

$$\begin{aligned} \left\| \frac{\frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 &\leq \frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \frac{1}{N} \sum_{i=1}^N \|\nabla f_i(x_{t,i}) - \nabla f(\bar{X}_t)\|^2 \\ &\leq L \underbrace{\frac{1}{\min_{j \in [d]} [\bar{U}_t^{1/2}]_j} \frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2}_{T_5}. \end{aligned} \quad (28)$$

472 using Jensen's inequality, Lipschitz continuity of f_i , and the fact that $f = \frac{1}{N} \sum_{i=1}^N f_i$. Next we need
 473 to bound T_4 and T_5 . Before we proceed into bounding T_5 , we need some preparations. Let's recall
 474 the update rule of X_t , we have

$$X_t = X_{t-1}W - \alpha \frac{M_{t-1}}{\sqrt{U_{t-1}}} = X_1 W^{t-1} - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} W^k, \quad (29)$$

475 where we define $W^0 = \mathbf{I}$. Since W is a symmetric matrix, we can decompose it as $W = Q\Lambda Q^T$
 476 where Q is a orthonormal matrix and Λ is a diagonal matrix whose diagonal elements correspond
 477 to eigenvalues of W in an descending order, i.e. $\Lambda_{ii} = \lambda_i$ with λ_i being i th largest eigenvalue of
 478 W . In addition, because W is a doubly stochastic matrix, we know $\lambda_1 = 1$ and $q_1 = \frac{1}{\sqrt{N}}$. With
 479 eigen-decomposition of W , we can rewrite T_5 as

$$\sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 = \|X_t - \bar{X}_t \mathbf{1}_N^T\|_F^2 = \|X_t Q Q^T - X_t \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^T\|_F^2 = \sum_{l=2}^N \|X_t q_l\|^2. \quad (30)$$

480 In addition, we can rewrite (29) as

$$X_t = X_1 W^{t-1} - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} W^k = X_1 - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} Q \Lambda^k Q^T, \quad (31)$$

481 where the last equality is because $x_{1,i} = x_{1,j}$, $\forall i, j$ and thus $X_1 W = X_1$. Then we have when $l > 1$,

$$X_t q_l = (X_1 - \alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} Q \Lambda^k Q^T) q_l = -\alpha \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} q_l \lambda_l^k, \quad (32)$$

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

483 Combining (30) and (32), we have

$$T_5 = \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 = \sum_{l=2}^N \|X_t q_l\|^2 = \sum_{l=2}^N \alpha^2 \left\| \sum_{k=0}^{t-2} \frac{M_{t-k-1}}{\sqrt{U_{t-k-1}}} \lambda_l^k q_l \right\|^2 \leq \alpha^2 \left(\frac{1}{1-\lambda} \right)^2 N d G_\infty^2 \frac{1}{\epsilon}, \quad (33)$$

484 where the last inequality follows from the fact that $g_{t,i} \leq G_\infty$, $\|q_l\| = 1$, and $|\lambda_l| \leq \lambda < 1$. Now let
 485 us turn to T_4 , it can be rewritten as

$$\begin{aligned} \|Z_t - \bar{X}_t\|^2 &= \left\| \frac{\beta_1}{1-\beta_1} (\bar{X}_t - \bar{X}_{t-1}) \right\|^2 = \left(\frac{\beta_1}{1-\beta_1} \right)^2 \alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{m_{t-1,i}}{\sqrt{u_{t-1,i}}} \right\|^2 \\ &\leq \left(\frac{\beta_1}{1-\beta_1} \right)^2 \alpha^2 d \frac{G_\infty^2}{\epsilon}. \end{aligned} \quad (34)$$

486 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
 487 $\mathcal{O}(\alpha^2)$. Next we will bound T_2 and T_1 . Define $G_1 \triangleq \max_{t \in [T]} \max_{i \in [N]} \|\nabla f_i(x_{t,i})\|_\infty$, $G_2 \triangleq$
 488 $\max_{t \in [T]} \|\nabla f(Z_t)\|_\infty$, $G_3 \triangleq \max_{t \in [T]} \max_{i \in [N]} \|g_{t,i}\|_\infty$ and $G_\infty = \max(G_1, G_2, G_3)$. Then we

489 have

$$\begin{aligned}
T_2 &= \sum_{t=1}^T \mathbb{E} \left[\left\langle \nabla f(Z_t), \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{\bar{U}_t}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\rangle \right] \\
&\leq \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{\sqrt{[\bar{U}_t]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{1}{\sqrt{[\bar{U}_t]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| \frac{\sqrt{[\bar{U}_t]_j} + \sqrt{[u_{t,i}]_j}}{\sqrt{[\bar{U}_t]_j} + \sqrt{[u_{t,i}]_j}} \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{[\bar{U}_t]_j - [u_{t,i}]_j}{[\bar{U}_t]_j \sqrt{[u_{t,i}]_j} + \sqrt{[\bar{U}_t]_j} [u_{t,i}]_j} \right| \right] \\
&\leq \underbrace{\mathbb{E} \left[\sum_{t=1}^T G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \left| \frac{[\bar{U}_t]_j - [u_{t,i}]_j}{2\epsilon^{1.5}} \right| \right]}_{T_6}, \tag{35}
\end{aligned}$$

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

491 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 ,
492 then we have

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

493 where the second inequality is due to Lemma 2, introduced Section A, and the fact that $U_t =$
494 $\max(\tilde{U}_t, \epsilon)$ element-wisely. Recall from update rule of U_t , by defining $\hat{V}_{-1} \triangleq \hat{V}_0$ and $U_0 \triangleq U_{1/2}$,
495 we have $\forall t \geq 0$

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

496 and thus

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

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

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

498 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
499 $q_1 = \frac{1}{\sqrt{N}} \mathbf{1}_N^T$) and $q_i^T q_j = 0$ when $i \neq j$. Note by definition of $\|\cdot\|_{abs}$, we have $\forall A, B, \|A+B\|_{abs} \leq$

500 $\|A\|_{abs} + \|B\|_{abs}$, then we have

$$\begin{aligned}
T_6 &\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \\
&= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \left\| - \sum_{k=1}^t (-\hat{V}_{t-1-k} + \hat{V}_{t-k}) \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_{abs} \\
&\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \|(-\hat{V}_{t-1-k} + \hat{V}_{t-k}) \sum_{l=2}^N q_l \lambda_l^k q_l^T\|_{abs} \\
&= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T (-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j \right\|_1 \\
&\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_1 \|(-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j\|_1 \\
&\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \sqrt{N} \left\| \sum_{l=2}^N q_l \lambda_l^k q_l^T \right\|_2 \|(-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j\|_1 \\
&\leq \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \sum_{j=1}^d \|(-\hat{V}_{t-1-k} + \hat{V}_{t-k})^T e_j\|_1 \sqrt{N} \lambda^k \\
&= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{k=1}^t \|(-\hat{V}_{t-1-k} + \hat{V}_{t-k})\|_{abs} \sqrt{N} \lambda^k \\
&= \frac{G_\infty^2}{N} \sum_{t=1}^T \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{t-1} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \sqrt{N} \lambda^{t-o} \\
&= \frac{G_\infty^2}{N} \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{T-1} \sum_{t=o+1}^T \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \sqrt{N} \lambda^{t-o} \\
&\leq \frac{G_\infty^2}{\sqrt{N}} \frac{1}{2\epsilon^{1.5}} \sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs}, \tag{36}
\end{aligned}$$

501 where $\lambda = \max(|\lambda_2|, |\lambda_N|)$. Combining (35) and (36), we have

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

502 Now we need to bound T_1 , we have

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

$$\begin{aligned}
&\stackrel{(a)}{\leq} \sum_{t=1}^T \mathbb{E} \left[G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{2\epsilon^{1.5}} |([\tilde{u}_{t-1,i}]_j - [\tilde{u}_{t,i}]_j)| \right] \\
&= G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{N} \mathbb{E} \left[\sum_{t=1}^T \|\tilde{U}_{t-1} - \tilde{U}_t\|_{abs} \right]
\end{aligned} \tag{37}$$

503 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. In
504 addition, by update rule of \tilde{U}_t , we have

$$\begin{aligned}
&\sum_{t=1}^T \|\tilde{U}_{t-1} - \tilde{U}_t\|_{abs} \\
&= \sum_{t=1}^T \|\tilde{U}_{t-1} - (\tilde{U}_{t-1} - \hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\
&= \sum_{t=1}^T \|\tilde{U}_{t-1}(QQ^T - Q\Lambda Q^T) + (-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\
&= \sum_{t=1}^T \|\tilde{U}_{t-1}(\sum_{l=2}^N q_l(1 - \lambda_l)q_l^T) + (-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\
&\leq \sum_{t=1}^T \left\| \sum_{k=1}^{t-1} (-\hat{V}_{t-2-k} + \hat{V}_{t-1-k}) \sum_{l=2}^N q_l \lambda_l^k (1 - \lambda_l) q_l^T \right\|_{abs} + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})W\|_{abs} \\
&\leq \sum_{t=1}^T \left(\sum_{k=1}^{t-1} \|-\hat{V}_{t-2-k} + \hat{V}_{t-1-k}\|_{abs} \sqrt{N} \lambda^k \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \\
&= \sum_{t=1}^T \left(\sum_{o=1}^{t-1} \|-\hat{V}_{o-2} + \hat{V}_{o-1}\|_{abs} \sqrt{N} \lambda^{t-o} \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \\
&= \sum_{o=1}^{T-1} \sum_{t=o+1}^T \left(\|-\hat{V}_{o-2} + \hat{V}_{o-1}\|_{abs} \sqrt{N} \lambda^{t-o} \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \\
&\leq \sum_{o=1}^{T-1} \frac{\lambda}{1-\lambda} \left(\|-\hat{V}_{o-2} + \hat{V}_{o-1}\|_{abs} \sqrt{N} \right) + \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \\
&\leq \frac{1}{1-\lambda} \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \sqrt{N}.
\end{aligned} \tag{38}$$

505 Combining (37) and (38), we have

$$T_1 \leq G_\infty^2 \frac{1}{2\epsilon^{1.5}} \frac{1}{N} \mathbb{E} \left[\frac{1}{1-\lambda} \sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \sqrt{N} \right]. \tag{39}$$

506 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

$$\begin{aligned}
&\|Z_{t+1} - Z_t\|^2 \\
&= \left\| \alpha \frac{\beta_1}{1-\beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) - \alpha \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left\| \frac{\beta_1}{1-\beta_1} \frac{1}{N} \sum_{i=1}^N m_{t-1,i} \odot \left(\frac{1}{\sqrt{u_{t-1,i}}} - \frac{1}{\sqrt{u_{t,i}}} \right) \right\|^2 + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2
\end{aligned}$$

$$\begin{aligned}
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{\sqrt{\epsilon}} \left| \frac{1}{\sqrt{[u_{t-1,i}]_j}} - \frac{1}{\sqrt{[u_{t,i}]_j}} \right| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{\sqrt{\epsilon}} \left| \frac{[u_{t,i}]_j - [u_{t-1,i}]_j}{2\epsilon^{1.5}} \right| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d \frac{1}{2\epsilon^2} |[u_{t,i}]_j - [u_{t-1,i}]_j| + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \\
&= 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \frac{1}{2\epsilon^2} \|\tilde{U}_t - \tilde{U}_{t-1}\|_{abs} + 2\alpha^2 \left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2, \tag{40}
\end{aligned}$$

507 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
508 $\max(\cdot, \epsilon)$ is 1-Lipschitz.

509 Then, we have

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E}[\|Z_{t+1} - Z_t\|^2] \\
&\leq 2\alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 G_\infty^2 \frac{1}{N} \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{t=1}^T \|\tilde{U}_t - \tilde{U}_{t-1}\|_{abs} \right] + 2\alpha^2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&\leq \alpha^2 \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{G_\infty^2}{\sqrt{N}} \frac{1}{\epsilon^2} \frac{1}{1-\lambda} \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] + 2\alpha^2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \tag{41}
\end{aligned}$$

510 where the last inequality is due to (38).

511 We now bound the last term on RHS of the above inequality. A trivial bound can be

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

512 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
513 assumption that $[v_{t,i}]_j \geq \epsilon, \forall i$). However, the above bound is independent of N , to get a better bound,
514 we need a more involved analysis to show its dependency on N . To do this, we first notice that

$$\begin{aligned}
&\mathbb{E}_{G_t|G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&= \mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left\langle \frac{\nabla f_i(x_{t,i}) + \xi_{t,i}}{\sqrt{u_{t,i}}}, \frac{\nabla f_j(x_{t,j}) + \xi_{t,j}}{\sqrt{u_{t,j}}} \right\rangle \right] \\
&\stackrel{(a)}{=} \mathbb{E}_{G_t|G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] + \mathbb{E}_{G_t|G_{1:t-1}} \left[\frac{1}{N^2} \sum_{i=1}^N \left\| \frac{\xi_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&\stackrel{(b)}{=} \left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{1}{N^2} \sum_{i=1}^N \sum_{l=1}^d \frac{\mathbb{E}_{G_t|G_{1:t-1}}[[\xi_{t,i}]_l^2]}{[u_{t,i}]_l} \\
&\stackrel{(c)}{\leq} \left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{d}{N} \frac{\sigma^2}{\epsilon} \tag{42}
\end{aligned}$$

515 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 $\xi_j, \forall j \neq i$,
516 (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
517 $[u_{t,i}]_l \geq \epsilon$ by definition.

518 Then we have

$$\begin{aligned}
\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] &= \mathbb{E}_{G_{1:t-1}} \left[\mathbb{E}_{G_t | G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \right] \\
&\leq \mathbb{E}_{G_{1:t-1}} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 + \frac{d}{N} \frac{\sigma^2}{\epsilon} \right] \\
&= \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] + \frac{d}{N} \frac{\sigma^2}{\epsilon} \tag{43}
\end{aligned}$$

519 In traditional analysis of SGD-like distributed algorithms, the term corresponding to
520 $\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right]$ will be merged with the first order descent when the stepsize is cho-
521 sen to be small enough. However, in our case, the term cannot be merged because it is different from
522 the first order descent in our algorithm. A brute-force upper bound is possible but this will lead to a
523 worse convergence rate in terms of N . Thus, we need a more detailed analysis for the term in the
524 following.

$$\begin{aligned}
&\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] \tag{44} \\
&= \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} + \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right] \\
&\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right] \\
&\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \left\| \nabla f_i(x_{t,i}) \odot \left(\frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right) \right\|^2 \right] \\
&\leq 2\mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right\|_1 \right] \tag{45}
\end{aligned}$$

525 Summing over T , we have

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&\leq 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right\|_1 \right] \tag{46}
\end{aligned}$$

526 For the last term on RHS of (46), we can bound it similarly as what we did for T_2 from (35) to (36),
527 which yields

$$\begin{aligned}
\sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \left\| \frac{1}{\sqrt{u_{t,i}}} - \frac{1}{\sqrt{\bar{U}_t}} \right\|_1 \right] &\leq \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N G_\infty^2 \frac{1}{\sqrt{\epsilon}} \frac{1}{2\epsilon^{1.5}} \|u_{t,i} - \bar{U}_t\|_1 \right] \\
&= \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} G_\infty^2 \frac{1}{2\epsilon^2} \|\bar{U}_t \mathbf{1}^T - U_t\|_{abs} \right]
\end{aligned}$$

$$\begin{aligned}
&\leq \sum_{t=1}^T \mathbb{E} \left[\frac{1}{N} G_\infty^2 \frac{1}{2\epsilon^2} \left\| - \sum_{l=2}^N \tilde{U}_t q_l q_l^T \right\|_{abs} \right] \\
&\leq \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right]
\end{aligned} \tag{47}$$

528 Further, we have

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] \\
&\leq 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t)}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] \\
&= 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 2 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right]
\end{aligned} \tag{48}$$

529 and the last term on RHS of the above inequality can be bounded following similar procedures from
530 (28) to (33), as what we did for T_3 . Completing the procedures yields

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{\nabla f_i(\bar{X}_t) - \nabla f_i(x_{t,i})}{\sqrt{\bar{U}_t}} \right\|^2 \right] \\
&\leq \sum_{t=1}^T \mathbb{E} \left[L \frac{1}{\epsilon} \frac{1}{N} \sum_{i=1}^N \|x_{t,i} - \bar{X}_t\|^2 \right] \\
&\leq \sum_{t=1}^T \mathbb{E} \left[L \frac{1}{\epsilon} \frac{1}{N} \alpha^2 \left(\frac{1}{1-\lambda} \right) N d G_\infty^2 \frac{1}{\epsilon} \right] \\
&= T L \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) d G_\infty^2
\end{aligned} \tag{49}$$

531 Finally, combining (43) to (49), we get

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \left[\left\| \frac{1}{N} \sum_{i=1}^N \frac{g_{t,i}}{\sqrt{u_{t,i}}} \right\|^2 \right] \\
&\leq 4 \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\sqrt{\bar{U}_t}} \right\|^2 \right] + 4 T L \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) d G_\infty^2 \\
&\quad + 2 \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] + T \frac{d}{N} \frac{\sigma^2}{\epsilon} \\
&\leq 4 \frac{1}{\sqrt{\epsilon}} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] + 4 T L \frac{1}{\epsilon^2} \alpha^2 \left(\frac{1}{1-\lambda} \right) d G_\infty^2 \\
&\quad + 2 \frac{1}{\sqrt{N}} G_\infty^2 \frac{1}{2\epsilon^2} \mathbb{E} \left[\sum_{o=0}^{T-1} \frac{\lambda}{1-\lambda} \|(-\hat{V}_{o-1} + \hat{V}_o)\|_{abs} \right] + T \frac{d}{N} \frac{\sigma^2}{\epsilon}.
\end{aligned} \tag{50}$$

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

533 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) d G_\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 d G_\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].
\end{aligned} \tag{51}$$

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

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

$$\begin{aligned}
& + \frac{2}{T\epsilon^{1.5}} \frac{G_\infty^2}{\sqrt{N}} \frac{1}{1-\lambda} \left(\frac{L}{\sqrt{Td}} \left(\frac{\beta_1}{1-\beta_1} \right)^2 \frac{1}{\epsilon^{0.5}} + \lambda + \frac{\beta_1}{1-\beta_1} + 2 \frac{L}{\sqrt{Td}} \frac{1}{\epsilon^{0.5}} \lambda \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right] \\
\leq & C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(Z_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5} d^{0.5}} C_3 \\
& + \left(\frac{1}{TN^{0.5}} C_4 + \frac{1}{T^{1.5} d^{0.5} N^{0.5}} C_5 \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right],
\end{aligned}$$

535 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
536 RHS of (51) to the LHS to cancel it using the assumption $8L\alpha \frac{1}{\sqrt{\epsilon}} \leq \frac{1}{2}$ followed by multiplying both
537 sides by 2, and the constants introduced in the last step are defined as following

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

538 Substituting into $Z_1 = \bar{X}_1$ completes the proof. □

539 C Proof of Theorem 3

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

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5} d^{0.5}} C_3 \\ &\quad + \left(\frac{1}{TN^{0.5}} C_4 + \frac{1}{T^{1.5} d^{0.5} N^{0.5}} C_5 \right) \mathbb{E} \left[\sum_{t=1}^T \|(-\hat{V}_{t-2} + \hat{V}_{t-1})\|_{abs} \right], \end{aligned} \quad (52)$$

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

543 Since Algorithm 3 is a special case of 2, building on result of Theorem 2, we just need to characterize
544 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
545 update rule of Algorithm 3, we know \hat{V}_t is non decreasing and thus

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

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

547 Further, because $\|g_{t,i}\|_\infty \leq G_\infty, \forall t, i$ and $v_{t,i}$ is a exponential moving average of $g_{k,i}^2, k =$
548 $1, 2, \dots, t$, we know $|[v_{t,i}]_j| \leq G_\infty^2, \forall t, i, j$. In addition, by update rule of \hat{V}_t , we also know
549 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
550 $[\hat{v}_{0,i}]_j \geq 0$, we have

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

551 Substituting the above into (52), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \frac{\nabla f(\bar{X}_t)}{\bar{U}_t^{1/4}} \right\|^2 \right] &\leq C_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C_2 + \frac{1}{T^{1.5} d^{0.5}} C_3 \\ &\quad + \frac{d}{T} C_4 \sqrt{N} G_\infty^2 + \frac{\sqrt{d}}{T^{1.5}} C_5 \sqrt{N} G_\infty^2 \\ &= C'_1 \frac{\sqrt{d}}{\sqrt{T}} \left(\mathbb{E}[f(\bar{X}_1)] - \min_z f(z) + \frac{\sigma^2}{N} \right) + \frac{1}{T} C'_2 + \frac{1}{T^{1.5} d^{0.5}} C'_3 \\ &\quad + \frac{d}{T} \sqrt{N} C'_4 + \frac{\sqrt{d}}{T^{1.5}} \sqrt{N} C'_5 \end{aligned} \quad (53)$$

552 where we have

$$C'_1 = C_1 \quad C'_2 = C_2 \quad C'_3 = C_3 \quad C'_4 = C_4 G_\infty^2 \quad C'_5 = C_5 G_\infty^2. \quad (54)$$

553 and concluding our proof. \square

554 D Additional Experiments and Details

555 In this section, we compare the learning curves of different algorithms with different stepsizes on
 556 heterogeneous data distribution. We use 5 nodes and the heterogeneous data distribution is created
 557 by assigning each node with data of only two labels and there are no overlapping labels between
 558 different nodes. For all algorithms, we compare stepsizes in the set $[1e-1, 1e-2, 1e-3, 1e-4, 1e-5,$
 559 $1e-6]$.

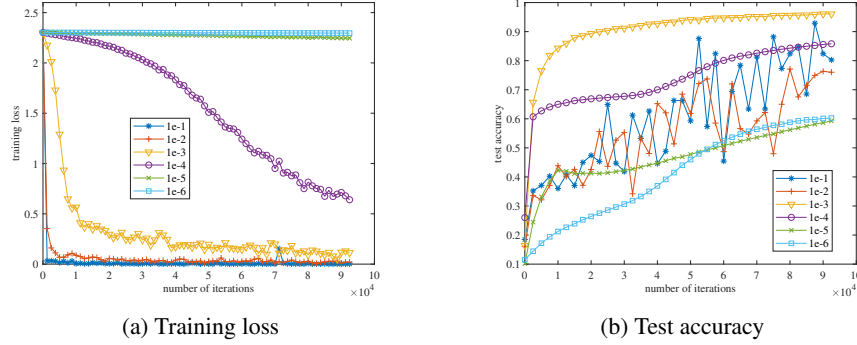


Figure 2: Performance comparison of different stepsizes for DGD

560 Figure 2 shows the training loss and test accuracy of DGD, it can be seen that the stepsize $1e-3$ works
 561 best for DGD in terms of test accuracy and $1e-1$ works best in terms of training loss. The difference
 562 is caused by the inconsistency among the value of parameters on different nodes when the stepsize
 563 is large. The training loss is calculated as the average of the loss value of different local models
 564 evaluated on their local training batch. Thus, though the training loss is small evaluated at a particular
 565 node, the test accuracy will be low when evaluating data with labels not seen by the node (recall that
 566 each node contains data with different labels).

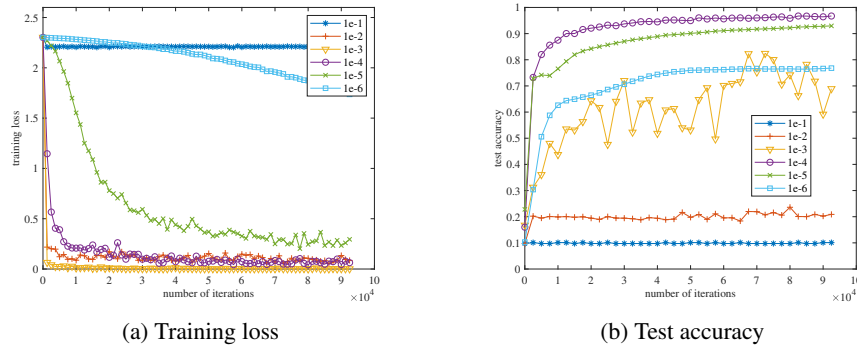


Figure 3: Performance comparison of different stepsizes for decentralized AMSGrad

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

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

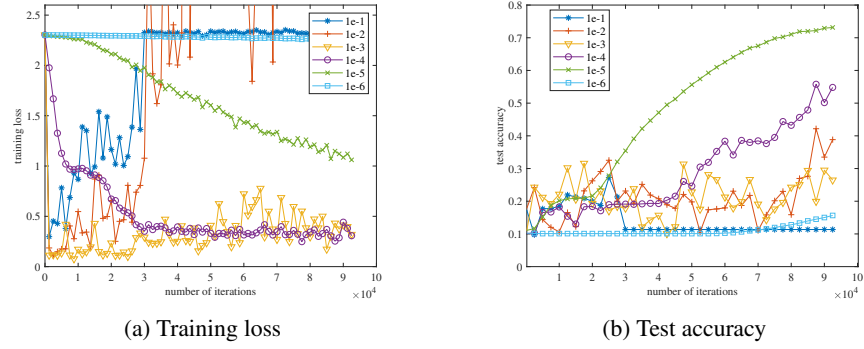


Figure 4: Performance comparison of different stepsizes for DADAM

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.