

**Supplementary materials for
“CADA: Communication-Adaptive Distributed Adam”**

In this supplementary document, we first present some basic inequalities that will be used frequently in this document, and then present the missing derivations of some claims, as well as the proofs of all the lemmas and theorems in the paper, which is followed by details on our experiments. The content of this supplementary document is summarized as follows.

Table of Contents

6	Supporting Lemmas	13
7	Missing Derivations in Section 2.2	16
7.1	Derivations of (6)	16
7.2	Derivations of (9)	17
7.3	Derivations of (11)	17
8	Proof of Lemma 1	18
9	Proof of Lemma 2	19
10	Proof of Lemma 3	21
11	Proof of Theorem 1	23
12	Proof of Theorem 2	25
13	Additional Numerical Results	26
13.1	Logistic regression.	26
13.2	Training neural networks.	27

6 Supporting Lemmas

Define the σ -algebra $\Theta^k = \{\theta^l, 1 \leq l \leq k\}$. For convenience, we also initialize parameters as $\theta^{-D}, \theta^{-D+1}, \dots, \theta^{-1} = \theta^0$. Some basic facts used in the proof are reviewed as follows.

Fact 1. Assume that $X_1, X_2, \dots, X_n \in \mathbb{R}^p$ are independent random variables, and $EX_1 = \dots = EX_n = 0$. Then

$$\mathbb{E} \left[\left\| \sum_{i=1}^n X_i \right\|^2 \right] = \sum_{i=1}^n \mathbb{E} [\|X_i\|^2]. \tag{18}$$

Fact 2. (Young’s inequality) For any $\theta_1, \theta_2 \in \mathbb{R}^p, \varepsilon > 0$,

$$\langle \theta_1, \theta_2 \rangle \leq \frac{\|\theta_1\|^2}{2\varepsilon} + \frac{\varepsilon\|\theta_2\|^2}{2}. \tag{19}$$

As a consequence, we have

$$\|\theta_1 + \theta_2\|^2 \leq \left(1 + \frac{1}{\varepsilon}\right)\|\theta_1\|^2 + (1 + \varepsilon)\|\theta_2\|^2. \tag{20}$$

Fact 3. (Cauchy-Schwarz inequality) For any $\theta_1, \theta_2, \dots, \theta_n \in \mathbb{R}^p$, we have

$$\left\| \sum_{i=1}^n \theta_i \right\|^2 \leq n \sum_{i=1}^n \|\theta_i\|^2. \tag{21}$$

Lemma 4. For $k - \tau_{\max} \leq l \leq k - D$, if $\{\theta^k\}$ are the iterates generated by CADA, we have

$$\begin{aligned} & \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\rangle \right] \\ & \leq \frac{L\epsilon^{-\frac{1}{2}}}{12\alpha_k} \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] + 6DL\alpha_k\epsilon^{-\frac{1}{2}}\sigma_m^2 \end{aligned} \quad (22)$$

and similarly, we have

$$\begin{aligned} & \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \mathcal{L}_m(\theta^l) - \nabla \ell(\theta^l; \theta^{k-\tau_m^k}) \right) \right\rangle \right] \\ & \leq \frac{L\epsilon^{-\frac{1}{2}}}{12\alpha_k} \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] + 3DL\alpha_k\epsilon^{-\frac{1}{2}}\sigma_m^2. \end{aligned} \quad (23)$$

Proof: We first show the following holds.

$$\begin{aligned} & \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^l), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\rangle \right] \\ & \stackrel{(a)}{=} \mathbb{E} \left[\mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^l), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\rangle \middle| \Theta^l \right] \right] \\ & \stackrel{(b)}{=} \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^l), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \mathbb{E} \left[\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \middle| \Theta^l \right] \right\rangle \right] \\ & = \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^l), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \mathcal{L}_m(\theta^l) - \nabla \mathcal{L}_m(\theta^l) \right) \right\rangle \right] = 0 \end{aligned} \quad (24)$$

where (a) follows from the law of total probability, and (b) holds because \hat{V}^{k-D} is deterministic conditioned on Θ^l when $k - D \leq l$.

We first prove (22) by decomposing it as

$$\begin{aligned} & \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\rangle \right] \\ & \stackrel{(c)}{=} \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k) - \nabla \mathcal{L}(\theta^l), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\rangle \right] \\ & \stackrel{(d)}{\leq} L\mathbb{E} \left[\left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{4}} \left\| \theta^k - \theta^l \right\| \left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{4}} \left(\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right) \right\| \right\| \right] \\ & \stackrel{(e)}{\leq} \frac{L\epsilon^{-\frac{1}{2}}}{12D\alpha_k} \underbrace{\mathbb{E} [\|\theta^k - \theta^l\|^2]}_{I_1} + \frac{6DL\alpha_k\epsilon^{-\frac{1}{2}}}{2} \underbrace{\mathbb{E} [\|\nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k})\|^2]}_{I_2} \end{aligned} \quad (25)$$

where (c) holds due to (24), (d) uses Assumption 1, and (e) applies the Young's inequality.

Applying the Cauchy-Schwarz inequality to I_1 , we have

$$\begin{aligned} I_1 & = \mathbb{E} \left[\left\| \sum_{d=1}^{k-l} (\theta^{k+1-d} - \theta^{k-d}) \right\|^2 \right] \\ & \leq (k-l) \sum_{d=1}^{k-l} \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \leq D \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2]. \end{aligned} \quad (26)$$

Applying Assumption 2 to I_2 , we have

$$\begin{aligned} I_2 & = \mathbb{E} \left[\left\| \nabla \ell(\theta^l; \xi_m^k) - \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right\|^2 \right] \\ & = \mathbb{E} \left[\left\| \nabla \ell(\theta^l; \xi_m^k) \right\|^2 \right] + \mathbb{E} \left[\left\| \nabla \ell(\theta^l; \xi_m^{k-\tau_m^k}) \right\|^2 \right] \leq 2\sigma_m^2 \end{aligned} \quad (27)$$

where the last inequality uses Assumption 2. Plugging (26) and (27) into (25), it leads to (22).

Likewise, following the steps to (25), it can be verified that (23) also holds true.

Lemma 5. Under Assumption 2, the parameters $\{h^k, \hat{v}^k\}$ of CADA in Algorithm 1 satisfy

$$\|h^k\| \leq \sigma, \quad \forall k; \quad \hat{v}_i^k \leq \sigma^2, \quad \forall k, i \quad (28)$$

where $\sigma := \frac{1}{M} \sum_{m \in \mathcal{M}} \sigma_m$.

Proof: Using Assumption 2, it follows that

$$\|\nabla^k\| = \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right\| \leq \frac{1}{M} \sum_{m \in \mathcal{M}} \left\| \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right\| \leq \frac{1}{M} \sum_{m \in \mathcal{M}} \sigma_m = \sigma. \quad (29)$$

Therefore, from the update (2a), we have

$$\|h^{k+1}\| \leq \beta_1 \|h^k\| + (1 - \beta_1) \|\nabla^k\| \leq \beta_1 \|h^k\| + (1 - \beta_1) \sigma.$$

Since $\|h^1\| \leq \sigma$, it follows by induction that $\|h^{k+1}\| \leq \sigma, \forall k$.

Using Assumption 2, it follows that

$$\begin{aligned} (\nabla_i^k)^2 &= \left(\frac{1}{M} \sum_{m \in \mathcal{M}} \nabla_i \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right)^2 \\ &\leq \frac{1}{M} \sum_{m \in \mathcal{M}} \left(\nabla_i \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right)^2 \\ &\leq \frac{1}{M} \sum_{m \in \mathcal{M}} \left\| \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right\|^2 = \frac{1}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 \leq \sigma^2. \end{aligned} \quad (30)$$

Similarly, from the update (2b), we have

$$\hat{v}_i^{k+1} \leq \max\{\hat{v}_i^k, \beta_2 \hat{v}_i^k + (1 - \beta_2) (\nabla_i^k)^2\} \leq \max\{\hat{v}_i^k, \beta_2 \hat{v}_i^k + (1 - \beta_2) \sigma^2\}.$$

Since $\hat{v}_i^1 = \hat{v}_i^1 \leq \sigma^2$, it follows by induction that $\hat{v}_i^{k+1} \leq \sigma^2$.

Lemma 6. Under Assumption 2, the iterates $\{\theta^k\}$ of CADA in Algorithm 1 satisfy

$$\|\theta^{k+1} - \theta^k\|^2 \leq \alpha^2 p (1 - \beta_2)^{-1} (1 - \beta_3)^{-1} \quad (31)$$

where p is the dimension of θ , $\beta_1 < \sqrt{\beta_2} < 1$, and $\beta_3 := \beta_1^2 / \beta_2$.

Proof: Choosing $\beta_1 < 1$ and defining $\beta_3 := \beta_1^2 / \beta_2$, it can be verified that

$$\begin{aligned} |h_i^{k+1}| &= |\beta_1 h_i^k + (1 - \beta_1) \nabla_i^k| \leq \beta_1 |h_i^k| + |\nabla_i^k| \\ &\leq \beta_1 (\beta_1 |h_i^{k-1}| + |\nabla_i^{k-1}|) + |\nabla_i^k| \\ &\leq \sum_{l=0}^k \beta_1^{k-l} |\nabla_i^l| = \sum_{l=0}^k \sqrt{\beta_3}^{k-l} \sqrt{\beta_2}^{k-l} |\nabla_i^l| \\ &\stackrel{(a)}{\leq} \left(\sum_{l=0}^k \beta_3^{k-l} \right)^{\frac{1}{2}} \left(\sum_{l=0}^k \beta_2^{k-l} (\nabla_i^l)^2 \right)^{\frac{1}{2}} \\ &\leq (1 - \beta_3)^{-\frac{1}{2}} \left(\sum_{l=0}^k \beta_2^{k-l} (\nabla_i^l)^2 \right)^{\frac{1}{2}} \end{aligned} \quad (32)$$

where (a) follows from the Cauchy-Schwartz inequality.

For \hat{v}_i^k , first we have that $\hat{v}_i^1 \geq (1 - \beta_2) (\nabla_i^1)^2$. Then since

$$\hat{v}_i^{k+1} \geq \beta_2 \hat{v}_i^k + (1 - \beta_2) (\nabla_i^k)^2$$

by induction we have

$$\hat{v}_i^{k+1} \geq (1 - \beta_2) \sum_{l=0}^k \beta_2^{k-l} (\nabla_i^l)^2. \quad (33)$$

Using (32) and (33), we have

$$\begin{aligned} |h_i^{k+1}|^2 &\leq (1 - \beta_3)^{-1} \left(\sum_{l=0}^k \beta_2^{k-l} (\nabla_i^l)^2 \right) \\ &\leq (1 - \beta_2)^{-1} (1 - \beta_3)^{-1} \hat{v}_i^{k+1}. \end{aligned}$$

From the update (2c), we have

$$\begin{aligned} \|\theta^{k+1} - \theta^k\|^2 &= \alpha_k^2 \sum_{i=1}^p (\epsilon + \hat{v}_i^{k+1})^{-1} |h_i^{k+1}|^2 \\ &\leq \alpha_k^2 p (1 - \beta_2)^{-1} (1 - \beta_3)^{-1} \end{aligned} \quad (34)$$

which completes the proof.

7 Missing Derivations in Section 2.2

The analysis in this part is analogous to that in [12]. We define an auxiliary function as

$$\psi_m(\theta) = \mathcal{L}_m(\theta) - \mathcal{L}_m(\theta^*) - \left\langle \nabla \mathcal{L}_m(\theta^*), \theta - \theta^* \right\rangle$$

where θ^* is a minimizer of \mathcal{L} . Assume that $\nabla \ell(\theta; \xi_m)$ is \bar{L} -Lipschitz continuous for all ξ_m , we have

$$\|\nabla \ell(\theta; \xi_m) - \nabla \ell(\theta^*; \xi_m)\|^2 \leq 2\bar{L} \left(\ell(\theta; \xi_m) - \ell(\theta^*; \xi_m) - \left\langle \nabla \ell(\theta^*; \xi_m), \theta - \theta^* \right\rangle \right).$$

Taking expectation with respect to ξ_m , we can obtain

$$\mathbb{E}_{\xi_m} [\|\nabla \ell(\theta; \xi_m) - \nabla \ell(\theta^*; \xi_m)\|^2] \leq 2\bar{L} \left(\mathcal{L}_m(\theta) - \mathcal{L}_m(\theta^*) - \left\langle \nabla \mathcal{L}_m(\theta^*), \theta - \theta^* \right\rangle \right) = 2\bar{L} \psi_m(\theta).$$

Note that $\nabla \mathcal{L}_m$ is also \bar{L} -Lipschitz continuous and thus

$$\|\nabla \mathcal{L}_m(\theta) - \nabla \mathcal{L}_m(\theta^*)\|^2 \leq 2\bar{L} \left(\mathcal{L}_m(\theta) - \mathcal{L}_m(\theta^*) - \left\langle \nabla \mathcal{L}_m(\theta^*), \theta - \theta^* \right\rangle \right) = 2\bar{L} \psi_m(\theta).$$

7.1 Derivations of (6)

By (21), we can derive that

$$\|\theta_1 + \theta_2\| \leq 2\|\theta_1\|^2 + 2\|\theta_2\|^2$$

which also implies $\|\theta_1\|^2 \geq \frac{1}{2}\|\theta_1 + \theta_2\|^2 - \|\theta_2\|^2$.

As a consequence, we can obtain

$$\begin{aligned} &\mathbb{E} \left[\|\nabla \ell(\theta^k; \xi_m^k) - \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k})\|^2 \right] \\ &\geq \frac{1}{2} \mathbb{E} \left[\left\| (\nabla \ell(\theta^k; \xi_m^k) - \nabla \mathcal{L}_m(\theta^k)) + (\nabla \mathcal{L}_m(\theta^{k-\tau_m^k}) - \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k})) \right\|^2 \right] \\ &\quad - \mathbb{E} \left[\|\nabla \mathcal{L}_m(\theta^k) - \nabla \mathcal{L}_m(\theta^{k-\tau_m^k})\|^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[\|\nabla \ell(\theta^k; \xi_m^k) - \nabla \mathcal{L}_m(\theta^k)\|^2 \right] + \frac{1}{2} \mathbb{E} \left[\|\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \mathcal{L}_m(\theta^{k-\tau_m^k})\|^2 \right] \\ &\quad + \underbrace{\mathbb{E} \left[\left\langle \nabla \ell(\theta^k; \xi_m^k) - \nabla \mathcal{L}_m(\theta^k), \nabla \mathcal{L}_m(\theta^{k-\tau_m^k}) - \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right\rangle \right]}_{I_3} - \mathbb{E} \left[\|\nabla \mathcal{L}_m(\theta^k) - \nabla \mathcal{L}_m(\theta^{k-\tau_m^k})\|^2 \right] \end{aligned}$$

where we used the fact that $I_3 = 0$ to obtain (6), that is

$$I_3 = \mathbb{E} \left[\left\langle \mathbb{E} [\nabla \ell(\theta^k; \xi_m^k) | \Theta^k] - \nabla \mathcal{L}_m(\theta^k), \nabla \mathcal{L}_m(\theta^{k-\tau_m^k}) - \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) \right\rangle \right] = 0.$$

7.2 Derivations of (9)

Recall that

$$\begin{aligned}\tilde{\delta}_m^k - \tilde{\delta}_m^{k-\tau_m^k} &= (\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\tilde{\theta}; \xi_m^k) + \nabla\mathcal{L}_m(\tilde{\theta})) - (\nabla\ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla\ell(\tilde{\theta}; \xi_m^{k-\tau_m^k}) + \nabla\mathcal{L}_m(\tilde{\theta})) \\ &= \underbrace{(\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\tilde{\theta}; \xi_m^k) + \nabla\psi_m(\tilde{\theta}))}_{g_m^k} - \underbrace{(\nabla\ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla\ell(\tilde{\theta}; \xi_m^{k-\tau_m^k}) + \nabla\psi_m(\tilde{\theta}))}_{g_m^{k-\tau_m^k}}.\end{aligned}$$

And by (21), we have $\|\tilde{\delta}_m^k - \tilde{\delta}_m^{k-\tau_m^k}\|^2 \leq 2\|g_m^k\|^2 + 2\|g_m^{k-\tau_m^k}\|^2$. We decompose the first term as

$$\begin{aligned}\mathbb{E}[\|g_m^k\|^2] &\leq 2\mathbb{E}[\|\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k)\|^2] + 2\mathbb{E}[\|\nabla\ell(\tilde{\theta}; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k) - \nabla\psi_m(\tilde{\theta})\|^2] \\ &= 2\mathbb{E}[\mathbb{E}[\|\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k)\|^2 | \Theta^k]] \\ &\quad + 2\mathbb{E}[\|\nabla\ell(\tilde{\theta}; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k) - \mathbb{E}[\nabla\ell(\tilde{\theta}; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k) | \Theta^k]\|^2] \\ &\leq 4\bar{L}\mathbb{E}\psi_m(\theta^k) + 2\mathbb{E}[\|\nabla\ell(\tilde{\theta}; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k)\|^2] \\ &= 4\bar{L}\mathbb{E}\psi_m(\theta^k) + 2\mathbb{E}[\mathbb{E}[\|\nabla\ell(\tilde{\theta}; \xi_m^k) - \nabla\ell(\theta^*; \xi_m^k)\|^2 | \Theta^k]] \\ &\leq 4\bar{L}\mathbb{E}\psi_m(\theta^k) + 4\bar{L}\mathbb{E}\psi_m(\tilde{\theta}).\end{aligned}$$

By nonnegativity of ψ_m , we have

$$\begin{aligned}\mathbb{E}[\|g_m^k\|^2] &\leq 4\bar{L} \sum_{m \in \mathcal{M}} \mathbb{E}\psi_m(\theta^k) + 4\bar{L} \sum_{m \in \mathcal{M}} \mathbb{E}\psi_m(\tilde{\theta}) \\ &= 4M\bar{L}(\mathbb{E}\mathcal{L}(\theta^k) - \mathcal{L}(\theta^*)) + 4M\bar{L}(\mathbb{E}\mathcal{L}(\tilde{\theta}) - \mathcal{L}(\theta^*)).\end{aligned}\tag{35}$$

Similarly, we can prove

$$\mathbb{E}[\|g_m^{k-\tau_m^k}\|^2] \leq 4M\bar{L}(\mathbb{E}\mathcal{L}(\theta^{k-\tau_m^k}) - \mathcal{L}(\theta^*)) + 4M\bar{L}(\mathbb{E}\mathcal{L}(\tilde{\theta}) - \mathcal{L}(\theta^*)).\tag{36}$$

Therefore, it follows that

$$\begin{aligned}\mathbb{E}[\|\tilde{\delta}_m^k - \tilde{\delta}_m^{k-\tau_m^k}\|^2] &\leq 8M\bar{L}(\mathbb{E}\mathcal{L}(\theta^k) - \mathcal{L}(\theta^*)) + 8M\bar{L}(\mathbb{E}\mathcal{L}(\theta^{k-\tau_m^k}) - \mathcal{L}(\theta^*)) + 16M\bar{L}(\mathbb{E}\mathcal{L}(\tilde{\theta}) - \mathcal{L}(\theta^*)).\end{aligned}$$

7.3 Derivations of (11)

The LHS of (10) can be written as

$$\begin{aligned}\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^{k-\tau_m^k}; \xi_m^k) &= (\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^{k-\tau_m^k}; \xi_m^k) + \nabla\mathcal{L}_m(\theta^{k-\tau_m^k})) - \nabla\mathcal{L}_m(\theta^{k-\tau_m^k}) \\ &= (\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^{k-\tau_m^k}; \xi_m^k) + \nabla\psi_m(\theta^{k-\tau_m^k})) - \nabla\psi_m(\theta^{k-\tau_m^k}).\end{aligned}$$

Similar to (35), we can obtain

$$\begin{aligned}\mathbb{E}[\|\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^{k-\tau_m^k}; \xi_m^k) + \nabla\psi_m(\theta^{k-\tau_m^k})\|^2] &\leq 4M\bar{L}(\mathbb{E}\mathcal{L}(\theta^k) - \mathcal{L}(\theta^*)) + 4M\bar{L}(\mathbb{E}\mathcal{L}(\theta^{k-\tau_m^k}) - \mathcal{L}(\theta^*)).\end{aligned}$$

Combined with the fact

$$\begin{aligned}\mathbb{E}[\|\nabla\psi_m(\theta^{k-\tau_m^k})\|^2] &= \mathbb{E}[\|\nabla\mathcal{L}_m(\theta^{k-\tau_m^k}) - \nabla\mathcal{L}_m(\theta^*)\|^2] \\ &\leq 2\bar{L}\mathbb{E}\psi_m(\theta^{k-\tau_m^k}) \leq 2M\bar{L}(\mathbb{E}\mathcal{L}(\theta^{k-\tau_m^k}) - \mathcal{L}(\theta^*))\end{aligned}$$

we have

$$\mathbb{E}[\|\nabla\ell(\theta^k; \xi_m^k) - \nabla\ell(\theta^{k-\tau_m^k}; \xi_m^k)\|^2] \leq 8M\bar{L}(\mathbb{E}\mathcal{L}(\theta^k) - \mathcal{L}(\theta^*)) + 12M\bar{L}(\mathbb{E}\mathcal{L}(\theta^{k-\tau_m^k}) - \mathcal{L}(\theta^*)).$$

8 Proof of Lemma 1

Using the smoothness of $\mathcal{L}(\theta)$ in Assumption 1, we have

$$\begin{aligned}\mathcal{L}(\theta^{k+1}) &\leq \mathcal{L}(\theta^k) + \left\langle \nabla \mathcal{L}(\theta^k), \theta^{k+1} - \theta^k \right\rangle + \frac{L}{2} \|\theta^{k+1} - \theta^k\|^2 \\ &= \mathcal{L}(\theta^k) - \alpha_k \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} h^{k+1} \right\rangle + \frac{L}{2} \|\theta^{k+1} - \theta^k\|^2.\end{aligned}\quad (37)$$

We can further decompose the inner product as

$$\begin{aligned}& - \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} h^{k+1} \right\rangle \\ &= - (1 - \beta_1) \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} \nabla^k \right\rangle - \underbrace{\beta_1 \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle}_{I_1^k} \\ & \quad - \underbrace{\left\langle \nabla \mathcal{L}(\theta^k), \left((\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} - (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} \right) h^{k+1} \right\rangle}_{I_2^k}\end{aligned}\quad (38)$$

where we again decompose the first inner product as

$$\begin{aligned}-(1 - \beta_1) \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} \nabla^k \right\rangle &= - (1 - \beta_1) \underbrace{\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \nabla^k \right\rangle}_{I_3^k} \\ & \quad - \underbrace{(1 - \beta_1) \left\langle \nabla \mathcal{L}(\theta^k), \left((\epsilon I + \hat{V}^k)^{-\frac{1}{2}} - (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \right) \nabla^k \right\rangle}_{I_4^k}.\end{aligned}\quad (39)$$

Next, we bound the terms $I_1^k, I_2^k, I_3^k, I_4^k$ separately.

Taking expectation on I_1^k conditioned on Θ^k , we have

$$\begin{aligned}\mathbb{E}[I_1^k \mid \Theta^k] &= -\mathbb{E} \left[\beta_1 \left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle \mid \Theta^k \right] \\ &= -\beta_1 \left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle - \beta_1 \left\langle \nabla \mathcal{L}(\theta^k) - \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle \\ &\stackrel{(a)}{\leq} -\beta_1 \left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle + \alpha_{k-1}^{-1} \beta_1 L \|\theta^k - \theta^{k-1}\|^2 \\ &\stackrel{(b)}{\leq} \beta_1 (I_1^{k-1} + I_2^{k-1} + I_3^{k-1} + I_4^{k-1}) + \alpha_{k-1}^{-1} \beta_1 L \|\theta^k - \theta^{k-1}\|^2\end{aligned}\quad (40)$$

where follows from the L -smoothness of $\mathcal{L}(\theta)$ implied by Assumption 1; and (b) uses again the decomposition (38) and (39).

Taking expectation on I_2^k over all the randomness, we have

$$\begin{aligned}\mathbb{E}[I_2^k] &= \mathbb{E} \left[- \left\langle \nabla \mathcal{L}(\theta^k), \left((\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} - (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} \right) h^{k+1} \right\rangle \right] \\ &= \mathbb{E} \left[\sum_{i=1}^p \nabla_i \mathcal{L}(\theta^k) h_i^{k+1} \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\ &\stackrel{(d)}{\leq} \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\| \|h^{k+1}\| \sum_{i=1}^p \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\ &\stackrel{(e)}{\leq} \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right]\end{aligned}\quad (41)$$

where (d) follows from the Cauchy-Schwarz inequality and (e) is due to Assumption 2.

Regarding I_3^k , we will bound separately in Lemma 2.

Taking expectation on I_4^k over all the randomness, we have

$$\begin{aligned}
 \mathbb{E}[I_4^k] &= \mathbb{E} \left[- (1 - \beta_1) \left\langle \nabla \mathcal{L}(\theta^k), \left((\epsilon I + \hat{V}^k)^{-\frac{1}{2}} - (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \right) \nabla^k \right\rangle \right] \\
 &= - (1 - \beta_1) \mathbb{E} \left[\sum_{i=1}^p \nabla_i \mathcal{L}(\theta^k) \nabla_i^k \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} \right) \right] \\
 &\leq (1 - \beta_1) \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\| \|\nabla^k\| \sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} \right) \right] \\
 &\leq (1 - \beta_1) \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} \right) \right]. \tag{42}
 \end{aligned}$$

Taking expectation on (37) over all the randomness, and plugging (40), (41), and (42), we have

$$\begin{aligned}
 \mathbb{E}[\mathcal{L}(\theta^{k+1})] - \mathbb{E}[\mathcal{L}(\theta^k)] &\leq -\alpha_k \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} h^{k+1} \right\rangle \right] + \frac{L}{2} \mathbb{E} \left[\|\theta^{k+1} - \theta^k\|^2 \right] \\
 &= \alpha_k \mathbb{E} [I_1^k + I_2^k + I_3^k + I_4^k] + \frac{L}{2} \mathbb{E} \left[\|\theta^{k+1} - \theta^k\|^2 \right] \\
 &\leq -\alpha_k (1 - \beta_1) \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \nabla^k \right\rangle \right] \\
 &\quad - \alpha_k \beta_1 \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle \right] \\
 &\quad + \alpha_k \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\
 &\quad + \alpha_k (1 - \beta_1) \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} \right) \right] \\
 &\quad + \left(\frac{L}{2} + \alpha_k \alpha_{k-1}^{-1} \beta_1 L \right) \mathbb{E} \left[\|\theta^{k+1} - \theta^k\|^2 \right]. \tag{43}
 \end{aligned}$$

Since $(\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} \leq (\epsilon + \hat{v}_i^{k-1})^{-\frac{1}{2}}$, we have

$$\begin{aligned}
 &\sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) + (1 - \beta_1) \sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^k)^{-\frac{1}{2}} \right) \right] \\
 &\leq (2 - \beta_1) \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right]. \tag{44}
 \end{aligned}$$

Plugging (44) into (43) leads to the statement of Lemma 1.

9 Proof of Lemma 2

We first analyze the inner produce under CADA2 and then CADA1.

First recall that $\bar{\nabla}^k = \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \ell(\theta^k; \xi_m^k)$. Using the law of total probability implies that

$$\begin{aligned}
 \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \bar{\nabla}^k \right\rangle \right] &= \mathbb{E} \left[\mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \bar{\nabla}^k \right\rangle \mid \Theta^k \right] \right] \\
 &= \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \mathbb{E} [\bar{\nabla}^k \mid \Theta^k] \right\rangle \right] \\
 &= \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right]. \tag{45}
 \end{aligned}$$

Taking expectation on $\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \nabla^k \rangle$ over all randomness, we have

$$\begin{aligned}
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \nabla^k \right\rangle \right] \\
 &= - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \bar{\nabla}^k \right\rangle \right] \\
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \frac{1}{M} \sum_{m \in \mathcal{M}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right] \\
 & \stackrel{(a)}{=} - \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] \\
 & - \frac{1}{M} \sum_{m \in \mathcal{M}} \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right]
 \end{aligned} \tag{46}$$

where (a) uses (45).

Decomposing the inner product, for the CADA2 rule (10), we have

$$\begin{aligned}
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right] \\
 &= - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) \right) \right\rangle \right] \\
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right] \\
 & \stackrel{(b)}{\leq} \frac{L\epsilon^{-\frac{1}{2}}}{12\alpha_k} \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] + 6DL\alpha_k \epsilon^{-\frac{1}{2}} \sigma_m^2 \\
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right]
 \end{aligned} \tag{47}$$

where (b) follows from Lemma 4.

Using the Young's inequality, we can bound the last inner product in (47) as

$$\begin{aligned}
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right] \\
 & \leq \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{1}{2} \mathbb{E} \left[\left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\|^2 \right] \\
 & \stackrel{(g)}{\leq} \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{1}{2} \mathbb{E} \left[\left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left\| \nabla \ell(\theta^{k-\tau_m^k}; \xi_m^k) - \nabla \ell(\theta^k; \xi_m^k) \right\|^2 \right] \right] \\
 & \stackrel{(h)}{\leq} \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{c}{2d_{\max}} \mathbb{E} \left[\left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left\| \sum_{d=1}^{d_{\max}} \|\theta^{k+1-d} - \theta^{k-d}\|^2 \right\| \right] \right] \\
 & \stackrel{(i)}{\leq} \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{c\epsilon^{-\frac{1}{2}}}{2d_{\max}} \sum_{d=1}^D \mathbb{E} \left[\|\theta^{k+1-d} - \theta^{k-d}\|^2 \right]
 \end{aligned} \tag{48}$$

where (g) follows from the Cauchy-Schwarz inequality, and (h) uses the adaptive communication condition (10) in CADA2, and (i) follows since \hat{V}^{k-D} is entry-wise nonnegative and $\|\theta^{k+1-d} - \theta^{k-d}\|^2$ is nonnegative.

Similarly for CADA1's condition (7), we have

$$\begin{aligned}
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\theta^{k-\tau_m^k}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\theta^k; \xi_m^k) \right) \right\rangle \right] \\
 &= - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\nabla \ell(\tilde{\theta}; \xi_m^{k-\tau_m^k}) - \nabla \ell(\tilde{\theta}; \xi_m^k) \right) \right\rangle \right] \\
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\tilde{\delta}_m^{k-\tau_m^k} - \tilde{\delta}_m^k \right) \right\rangle \right] \\
 & \stackrel{(j)}{\leq} \frac{L\epsilon^{-\frac{1}{2}}}{12\alpha_k} \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] + 6DL\alpha_k \epsilon^{-\frac{1}{2}} \sigma_m^2 - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\tilde{\delta}_m^{k-\tau_m^k} - \tilde{\delta}_m^k \right) \right\rangle \right]
 \end{aligned} \tag{49}$$

where (j) follows from Lemma 4 since $\tilde{\theta}$ is a snapshot among $\{\theta^k, \dots, \theta^{k-D}\}$.

And the last product in (49) is bounded by

$$\begin{aligned}
 & - \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}} \left(\tilde{\delta}_m^{k-\tau_m^k} - \tilde{\delta}_m^k \right) \right\rangle \right] \\
 & \leq \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{c}{2} \mathbb{E} \left[\left\| (\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2d_{\max}}} \left\| \sum_{d=1}^{d_{\max}} \|\theta^{k+1-d} - \theta^{k-d}\|^2 \right\| \right\|^2 \right] \\
 & \stackrel{(i)}{\leq} \frac{1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + \frac{c\epsilon^{-\frac{1}{2}}}{2d_{\max}} \sum_{d=1}^D \mathbb{E} \left[\|\theta^{k+1-d} - \theta^{k-d}\|^2 \right]. \tag{50}
 \end{aligned}$$

Combining (46)-(50) leads to the desired statement for CADA1 and CADA2.

10 Proof of Lemma 3

For notational brevity, we re-write the Lyapunov function (14) as

$$\begin{aligned}
 \mathcal{V}^k & := \mathcal{L}(\theta^k) - \mathcal{L}(\theta^*) - c_k \left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle \\
 & \quad + b_k \sum_{d=0}^D \sum_{i=1}^p (\epsilon + \hat{v}_i^{k-d})^{-\frac{1}{2}} + \sum_{d=1}^D \rho_d \|\theta^{k+1-d} - \theta^{k-d}\|^2 \tag{51}
 \end{aligned}$$

where $\{c_k\}$ are some positive constants.

Therefore, taking expectation on the difference of \mathcal{V}^k and \mathcal{V}^{k+1} in (51), we have (with $\rho_{D+1} = 0$)

$$\begin{aligned}
 \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] & = \mathbb{E}[\mathcal{L}(\theta^{k+1})] - \mathbb{E}[\mathcal{L}(\theta^k)] - c_{k+1} \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^k), (\epsilon I + \hat{V}^{k+1})^{-\frac{1}{2}} h^{k+1} \right\rangle \right] \\
 & \quad + c_k \mathbb{E} \left[\left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle \right] \\
 & \quad + b_{k+1} \sum_{d=0}^D \sum_{i=1}^p (\epsilon + \hat{v}_i^{k+1-d})^{-\frac{1}{2}} - b_k \sum_{d=0}^D \sum_{i=1}^p (\epsilon + \hat{v}_i^{k-d})^{-\frac{1}{2}} \\
 & \quad + \rho_1 \mathbb{E} [\|\theta^{k+1} - \theta^k\|^2] + \sum_{d=1}^D (\rho_{d+1} - \rho_d) \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \\
 & \stackrel{(a)}{\leq} (\alpha_k + c_{k+1}) \mathbb{E} [I_1^k + I_2^k + I_3^k + I_4^k] - c_k \mathbb{E} [I_1^{k-1} + I_2^{k-1} + I_3^{k-1} + I_4^{k-1}] \\
 & \quad + b_{k+1} \sum_{i=1}^p \mathbb{E} [(\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}}] - b_k \sum_{i=1}^p \mathbb{E} [(\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}}] \\
 & \quad + \sum_{d=1}^D (b_{k+1} - b_k) \sum_{i=1}^p \mathbb{E} [(\epsilon + \hat{v}_i^{k+1-d})^{-\frac{1}{2}}] + \left(\frac{L}{2} + \rho_1 \right) \mathbb{E} [\|\theta^{k+1} - \theta^k\|^2] \\
 & \quad + \sum_{d=1}^D (\rho_{d+1} - \rho_d) \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \tag{52}
 \end{aligned}$$

where (a) uses the smoothness in Assumption 1 and the definition of $I_1^k, I_2^k, I_3^k, I_4^k$ in (38) and (39).

Note that we can bound $(\alpha_k + c_{k+1}) \mathbb{E} [I_1^k + I_2^k + I_3^k + I_4^k]$ the same as (38) in the proof of Lemma 1. In addition, Lemma 2 implies that

$$\begin{aligned}
 \mathbb{E}[I_3^k] & \leq -\frac{1-\beta_1}{2} \mathbb{E} \left[\left\| \nabla \mathcal{L}(\theta^k) \right\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] \\
 & \quad + (1-\beta_1) \epsilon^{-\frac{1}{2}} \left(\frac{L}{12\alpha_k} + \frac{c}{2d_{\max}} \right) \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] + (1-\beta_1) \frac{6DL\alpha_k \epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2. \tag{53}
 \end{aligned}$$

Therefore, plugging Lemma 1 with α_k replaced by $\alpha_k + c_{k+1}$ into (52), together with (53), leads to

$$\begin{aligned}
 \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] &\leq -(\alpha_k + c_{k+1}) \left(\frac{1 - \beta_1}{2} \right) \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] \\
 &\quad + (\alpha_k + c_{k+1})(1 - \beta_1) \epsilon^{-\frac{1}{2}} \left(\frac{L}{12\alpha_k} + \frac{c}{2d_{\max}} \right) \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \\
 &\quad + (\alpha_k + c_{k+1})(1 - \beta_1) \frac{6DL\alpha_k \epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 \\
 &\quad + ((\alpha_k + c_{k+1})\beta_1 - c_k) \mathbb{E} [I_1^{k-1} + I_2^{k-1} + I_3^{k-1} + I_4^{k-1}] \\
 &\quad + (\alpha_k + c_{k+1})(2 - \beta_1) \sigma^2 \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\
 &\quad + b_{k+1} \sum_{i=1}^p \mathbb{E} \left[(\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right] - b_k \sum_{i=1}^p \mathbb{E} \left[(\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} \right] \\
 &\quad + \sum_{d=1}^D (b_{k+1} - b_k) \sum_{i=1}^p \mathbb{E} \left[(\epsilon + \hat{v}_i^{k+1-d})^{-\frac{1}{2}} \right] \\
 &\quad + \sum_{d=1}^D (\rho_{d+1} - \rho_d) \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \\
 &\quad + \left(\frac{L}{2} + \rho_1 + (\alpha_k + c_{k+1})\alpha_{k-1}^{-1}\beta_1 L \right) \mathbb{E} [\|\theta^{k+1} - \theta^k\|^2]. \tag{54}
 \end{aligned}$$

Select $\alpha_k \leq \alpha_{k-1}$ and $c_k := \sum_{j=k}^{\infty} \alpha_j \beta_1^{j-k+1} \leq (1 - \beta_1)^{-1} \alpha_k$ so that $(\alpha_k + c_{k+1})\beta_1 = c_k$ and

$$\begin{aligned}
 (\alpha_k + c_{k+1})(1 - \beta_1) &\leq (\alpha_k + (1 - \beta_1)^{-1} \alpha_{k+1})(1 - \beta_1) \\
 &\leq \alpha_k (1 + (1 - \beta_1)^{-1})(1 - \beta_1) = \alpha_k (2 - \beta_1).
 \end{aligned}$$

In addition, select b_k to ensure that $b_{k+1} \leq b_k$. Then it follows from (54) that

$$\begin{aligned}
 \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] &\leq -\frac{\alpha_k(1 - \beta_1)}{2} \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|_{(\epsilon I + \hat{V}^{k-D})^{-\frac{1}{2}}}^2 \right] + (2 - \beta_1) \alpha_k^2 \frac{6DL\epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 \\
 &\quad + (2 - \beta_1) \alpha_k \epsilon^{-\frac{1}{2}} \left(\frac{L}{12\alpha_k} + \frac{c}{2d_{\max}} \right) \sum_{d=1}^D \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \\
 &\quad + \left(\frac{(2 - \beta_1)^2}{(1 - \beta_1)} \alpha_k \sigma^2 - b_k \right) \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\
 &\quad + \left(\frac{L}{2} + \rho_1 + (1 - \beta_1)^{-1} L \right) \mathbb{E} [\|\theta^{k+1} - \theta^k\|^2] \\
 &\quad + \sum_{d=1}^D (\rho_{d+1} - \rho_d) \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2] \tag{55}
 \end{aligned}$$

where we have also used the fact that $-(\alpha_k + c_{k+1}) \left(\frac{1 - \beta_1}{2} \right) \leq -\frac{\alpha_k(1 - \beta_1)}{2}$ since $c_{k+1} \geq 0$.

If we choose $\alpha_k \leq \frac{1}{L}$ for $k = 1, 2, \dots, K$, then it follows from (55) that

$$\begin{aligned}
 & \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] \\
 & \leq -\frac{\alpha_k(1-\beta_1)}{2} \left(\epsilon + \frac{\sigma^2}{1-\beta_2} \right)^{-\frac{1}{2}} \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|^2 \right] + (2-\beta_1) \frac{6\alpha_k^2 DL \epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 \\
 & \quad + \underbrace{\left(\frac{(2-\beta_1)^2}{(1-\beta_1)} \alpha_k \sigma^2 - b_k \right)}_{A^k} \mathbb{E} \left[\sum_{i=1}^p \left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}} \right) \right] \\
 & \quad + \left(\frac{L}{2} + \rho_1 + (1-\beta_1)^{-1} L \right) \mathbb{E} [\|\theta^{k+1} - \theta^k\|^2] \\
 & \quad + \sum_{d=1}^D \underbrace{\left((2-\beta_1) \epsilon^{-\frac{1}{2}} \left(\frac{L}{12} + \frac{c\alpha_k}{2d_{\max}} \right) + \rho_{d+1} - \rho_d \right)}_{B_d^k} \mathbb{E} [\|\theta^{k+1-d} - \theta^{k-d}\|^2]. \tag{56}
 \end{aligned}$$

To ensure $A^k \leq 0$ and $B_d^k \leq 0$, it is sufficient to choose $\{b_k\}$ and $\{\rho_d\}$ satisfying (with $\rho_{D+1} = 0$)

$$\frac{(2-\beta_1)^2}{(1-\beta_1)} \alpha_k \sigma^2 - b_k \leq 0, \quad k = 1, \dots, K \tag{57}$$

$$(2-\beta_1) \epsilon^{-\frac{1}{2}} \left(\frac{L}{12} + \frac{c\alpha_k}{2d_{\max}} \right) + \rho_{d+1} - \rho_d \leq 0, \quad d = 1, \dots, D. \tag{58}$$

Solve this system of linear equations and get

$$b_k = \frac{(2-\beta_1)^2}{(1-\beta_1)L} \sigma^2, \quad k = 1, \dots, K \tag{59}$$

$$\rho_d = (2-\beta_1) \epsilon^{-\frac{1}{2}} \left(\frac{L}{12} + \frac{c}{2Ld_{\max}} \right) (D-d+1), \quad d = 1, \dots, D \tag{60}$$

plugging which into (56) leads to the conclusion of Lemma 3.

11 Proof of Theorem 1

From the definition of \mathcal{V}^k , we have for any k , that

$$\begin{aligned}
 \mathbb{E}[\mathcal{V}^k] & \geq \mathcal{L}(\theta^k) - \mathcal{L}(\theta^*) - c_k \left\langle \nabla \mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\rangle + \sum_{d=1}^D \rho_d \|\theta^{k+1-d} - \theta^{k-d}\|^2 \\
 & \geq -|c_k| \|\nabla \mathcal{L}(\theta^{k-1})\| \left\| (\epsilon I + \hat{V}^k)^{-\frac{1}{2}} h^k \right\| \\
 & \geq -(1-\beta_1)^{-1} \alpha_k \sigma^2 \epsilon^{-\frac{1}{2}} \tag{61}
 \end{aligned}$$

where we use Assumption 2 and Lemma 5.

By taking summation on (56) over $k = 0, \dots, K-1$, it follows from that

$$\begin{aligned}
 & \frac{\alpha(1-\beta_1)}{2} \left(\epsilon + \frac{\sigma^2}{1-\beta_2} \right)^{-\frac{1}{2}} \frac{1}{K} \sum_{k=1}^K \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|^2 \right] \\
 & \leq \frac{\mathbb{E}[\mathcal{V}^1] - \mathbb{E}[\mathcal{V}^{K+1}]}{K} + (2-\beta_1) \frac{6\alpha^2 DL \epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 + \frac{(2-\beta_1)^2}{(1-\beta_1)} \sigma^2 p D \epsilon^{-\frac{1}{2}} \frac{\alpha}{K} \\
 & \quad + \left(\frac{L}{2} + \rho_1 + (1-\beta_1)^{-1} L \right) \frac{1}{K} \sum_{k=1}^K \mathbb{E} \left[\|\theta^{k+1} - \theta^k\|^2 \right] \\
 & \stackrel{(a)}{\leq} \frac{\mathbb{E}[\mathcal{V}^1]}{K} + (2-\beta_1) \frac{6\alpha^2 DL \epsilon^{-\frac{1}{2}}}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 + (1-\beta_1)^{-1} \sigma^2 \epsilon^{-\frac{1}{2}} \frac{\alpha}{K} + \frac{(2-\beta_1)^2}{(1-\beta_1)} \sigma^2 p D \epsilon^{-\frac{1}{2}} \frac{\alpha}{K} \\
 & \quad + \left(\frac{L}{2} + \rho_1 + (1-\beta_1)^{-1} L \right) p (1-\beta_2)^{-1} (1-\beta_3)^{-1} \alpha^2
 \end{aligned} \tag{62}$$

where (a) follows from (61) and Lemma 6.

Specifically, if we choose a constant stepsize $\alpha := \frac{\eta}{\sqrt{K}}$, where $\eta > 0$ is a constant, and define

$$\tilde{C}_1 := (2-\beta_1) 6DL \epsilon^{-\frac{1}{2}} \tag{63}$$

and

$$\tilde{C}_2 := (1-\beta_1)^{-1} \epsilon^{-\frac{1}{2}} + \frac{(2-\beta_1)^2}{(1-\beta_1)} D \epsilon^{-\frac{1}{2}} \tag{64}$$

and

$$\tilde{C}_3 := \left(\frac{L}{2} + \rho_1 + (1-\beta_1)^{-1} L \right) (1-\beta_2)^{-1} (1-\beta_3)^{-1} \tag{65}$$

and

$$\tilde{C}_4 := \frac{1}{2} (1-\beta_1) \left(\epsilon + \frac{\sigma^2}{1-\beta_2} \right)^{-\frac{1}{2}} \tag{66}$$

we can obtain from (62) that

$$\begin{aligned}
 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla \mathcal{L}(\theta^k)\|^2 \right] & \leq \frac{\frac{\mathcal{L}(\theta^0) - \mathcal{L}(\theta^*)}{K} + \frac{\tilde{C}_1}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 \alpha^2 + \tilde{C}_2 p \sigma^2 \frac{\alpha}{K} + \tilde{C}_3 p \alpha^2}{\alpha \tilde{C}_4} \\
 & \leq \frac{\mathcal{L}(\theta^0) - \mathcal{L}(\theta^*)}{K \alpha \tilde{C}_4} + \frac{\tilde{C}_1 \alpha}{\tilde{C}_4 M} \sum_{m \in \mathcal{M}} \sigma_m^2 + \tilde{C}_2 p \frac{\sigma^2}{K \tilde{C}_4} + \frac{\tilde{C}_3 p \alpha}{\tilde{C}_4} \\
 & = \frac{(\mathcal{L}(\theta^0) - \mathcal{L}(\theta^*)) C_4}{\sqrt{K} \eta} + \frac{C_1 \eta}{\sqrt{K} M} \sum_{m \in \mathcal{M}} \sigma_m^2 + \frac{C_2 p \sigma^2}{K} + \frac{C_3 p \eta}{\sqrt{K}}
 \end{aligned}$$

where we define $C_1 := \tilde{C}_1 / \tilde{C}_4$, $C_2 := \tilde{C}_2 / \tilde{C}_4$, $C_3 := \tilde{C}_3 / \tilde{C}_4$, and $C_4 := 1 / \tilde{C}_4$.

12 Proof of Theorem 2

By the PL-condition of $\mathcal{L}(\theta)$, we have

$$\begin{aligned}
 & -\frac{\alpha_k(1-\beta_1)}{2}\left(\epsilon + \frac{\sigma^2}{1-\beta_2}\right)^{-\frac{1}{2}}\mathbb{E}\left[\|\nabla\mathcal{L}(\theta^k)\|^2\right] \\
 & \leq -\alpha_k\mu(1-\beta_1)\left(\epsilon + \frac{\sigma^2}{1-\beta_2}\right)^{-\frac{1}{2}}\mathbb{E}\left[\mathcal{L}(\theta^k) - \mathcal{L}(\theta^*)\right] \\
 & \stackrel{(a)}{\leq} -2\alpha_k\mu\tilde{C}_4\left(\mathbb{E}[\mathcal{V}^k] + c_k\left\langle\nabla\mathcal{L}(\theta^{k-1}), (\epsilon I + \hat{V}^k)^{-\frac{1}{2}}h^k\right\rangle - b_k\sum_{d=0}^D\sum_{i=1}^p(\epsilon + \hat{v}_i^{k-d})^{-\frac{1}{2}} - \sum_{d=1}^D\rho_d\|\theta^{k+1-d} - \theta^{k-d}\|^2\right) \\
 & \stackrel{(b)}{\leq} -2\alpha_k\mu\tilde{C}_4\mathbb{E}[\mathcal{V}^k] + 2\alpha_k^2\mu\tilde{C}_4(1-\beta_1)^{-1}\sigma^2\epsilon^{-\frac{1}{2}} + 2\alpha_k\mu\tilde{C}_4b_k\sum_{d=0}^D\sum_{i=1}^p\mathbb{E}\left[(\epsilon + \hat{v}_i^{k-d})^{-\frac{1}{2}}\right] \\
 & \quad + 2\alpha_k\mu\tilde{C}_4\sum_{d=1}^D\rho_d\mathbb{E}[\|\theta^{k+1-d} - \theta^{k-d}\|^2]
 \end{aligned} \tag{67}$$

where (a) uses the definition of \tilde{C}_4 in (66), and (b) uses Assumption 2 and Lemma 5.

Plugging (67) into (55), we have

$$\begin{aligned}
 & \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] \\
 & \leq -2\alpha_k\mu\tilde{C}_4\mathbb{E}[\mathcal{V}^k] + (2-\beta_1)\frac{6\alpha_k^2DL\epsilon^{-\frac{1}{2}}}{M}\sum_{m\in\mathcal{M}}\sigma_m^2 \\
 & \quad + \frac{(2-\beta_1)^2}{(1-\beta_1)}\alpha_k\sigma^2\mathbb{E}\left[\sum_{i=1}^p\left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}}\right)\right] \\
 & \quad + b_{k+1}\sum_{i=1}^p\mathbb{E}\left[(\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}}\right] - (b_k - 2\alpha_k\mu\tilde{C}_4b_k)\sum_{i=1}^p\mathbb{E}\left[(\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}}\right] \\
 & \quad + \sum_{d=1}^D(b_{k+1} - b_k + 2\alpha_k\mu\tilde{C}_4b_k)\sum_{i=1}^p\mathbb{E}\left[(\epsilon + \hat{v}_i^{k+1-d})^{-\frac{1}{2}}\right] \\
 & \quad + \left(\frac{L}{2} + \rho_1 + (1-\beta_1)^{-1}L\right)p(1-\beta_2)^{-1}(1-\beta_3)^{-1}\alpha_k^2 + 2\alpha_k^2\mu\tilde{C}_4(1-\beta_1)^{-1}\sigma^2\epsilon^{-\frac{1}{2}} \\
 & \quad + \sum_{d=1}^D\left((2-\beta_1)\epsilon^{-\frac{1}{2}}\left(\frac{L}{12} + \frac{c\alpha_k}{2d_{\max}}\right) + \rho_{d+1} - \rho_d + 2\alpha_k\mu\tilde{C}_4\rho_d\right)\mathbb{E}[\|\theta^{k+1-d} - \theta^{k-d}\|^2].
 \end{aligned} \tag{68}$$

If we choose b_k to ensure that $b_{k+1} \leq (1 - 2\alpha_k\mu\tilde{C}_4)b_k$, then we can obtain from (68) that

$$\begin{aligned}
 & \mathbb{E}[\mathcal{V}^{k+1}] - \mathbb{E}[\mathcal{V}^k] \\
 & \leq -2\alpha_k\mu\tilde{C}_4\mathbb{E}[\mathcal{V}^k] + \frac{\tilde{C}_1}{M}\sum_{m\in\mathcal{M}}\sigma_m^2\alpha_k^2 + \tilde{C}_3p\alpha_k^2 + 2\mu\tilde{C}_4(1-\beta_1)^{-1}\sigma^2\epsilon^{-\frac{1}{2}}\alpha_k^2 \\
 & \quad + \left(\frac{(2-\beta_1)^2}{(1-\beta_1)}\alpha_k\sigma^2 - (1 - 2\alpha_k\mu\tilde{C}_4)b_k\right)\mathbb{E}\left[\sum_{i=1}^p\left((\epsilon + \hat{v}_i^{k-D})^{-\frac{1}{2}} - (\epsilon + \hat{v}_i^{k+1})^{-\frac{1}{2}}\right)\right] \\
 & \quad + \sum_{d=1}^D\left((2-\beta_1)\epsilon^{-\frac{1}{2}}\left(\frac{L}{12} + \frac{c\alpha_k}{2d_{\max}}\right) + \rho_{d+1} - \rho_d + 2\alpha_k\mu\tilde{C}_4\rho_d\right)\mathbb{E}[\|\theta^{k+1-d} - \theta^{k-d}\|^2].
 \end{aligned} \tag{69}$$

If $\alpha_k \leq \frac{1}{L}$, we choose parameters $\{b_k, \rho_d\}$ to guarantee that

$$\frac{(2-\beta_1)^2}{(1-\beta_1)L}\sigma^2 - \left(1 - \frac{2\mu\tilde{C}_4}{L}\right)b_k \leq 0, \quad \forall k \tag{70}$$

$$(2-\beta_1)\left(\frac{L}{12} + \frac{c}{2Ld_{\max}}\right)\epsilon^{-\frac{1}{2}} + \rho_{d+1} - \left(1 - \frac{2\mu\tilde{C}_4}{L}\right)\rho_d \leq 0, \quad d = 1, \dots, D \tag{71}$$

and choose $\beta_1, \beta_2, \epsilon$ to ensure that $1 - \frac{2\mu\tilde{C}_4}{L} \geq 0$.

Then we have

$$\begin{aligned} \mathbb{E}[\mathcal{V}^{k+1}] &\leq \left(1 - 2\alpha_k\mu\tilde{C}_4\right) \mathbb{E}[\mathcal{V}^k] + \underbrace{\left(\frac{\tilde{C}_1}{M} \sum_{m \in \mathcal{M}} \sigma_m^2 + \tilde{C}_3 p + 2\mu\tilde{C}_4(1 - \beta_1)^{-1} \sigma^2 \epsilon^{-\frac{1}{2}}\right)}_{\tilde{C}_5} \alpha_k^2 \\ &\leq \prod_{j=0}^k (1 - 2\alpha_j\mu\tilde{C}_4) \mathbb{E}[\mathcal{V}^0] + \sum_{j=0}^k \alpha_j^2 \prod_{i=j+1}^k (1 - 2\alpha_i\mu\tilde{C}_4) \tilde{C}_5. \end{aligned} \quad (72)$$

If we choose $\alpha_k = \frac{1}{\mu(k+K_0)\tilde{C}_4} \leq \frac{1}{L}$, where K_0 is a sufficiently large constant to ensure that α_k satisfies the aforementioned conditions, then we have

$$\begin{aligned} \mathbb{E}[\mathcal{V}^K] &\leq \mathbb{E}[\mathcal{V}^0] \prod_{k=0}^{K-1} (1 - 2\alpha_k\mu\tilde{C}_4) + \tilde{C}_5 \sum_{k=0}^{K-1} \alpha_k^2 \prod_{j=k+1}^{K-1} (1 - 2\alpha_j\mu\tilde{C}_4) \\ &\leq \mathbb{E}[\mathcal{V}^0] \prod_{k=0}^{K-1} \frac{k+K_0-2}{k+K_0} + \frac{\tilde{C}_5}{\mu^2\tilde{C}_4^2} \sum_{k=0}^{K-1} \frac{1}{(k+K_0)^2} \prod_{j=k+1}^{K-1} \frac{j+K_0-2}{j+K_0} \\ &\leq \frac{(K_0-2)(K_0-1)}{(K+K_0-2)(K+K_0-1)} \mathbb{E}[\mathcal{V}^0] + \frac{\tilde{C}_5}{\mu^2\tilde{C}_4^2} \sum_{k=0}^{K-1} \frac{(k+K_0-1)}{(k+K_0)(K+K_0-2)(K+K_0-2)} \\ &\leq \frac{(K_0-1)^2}{(K+K_0-1)^2} \mathbb{E}[\mathcal{V}^0] + \frac{\tilde{C}_5 K}{\mu^2\tilde{C}_4^2 (K+K_0-1)^2} \\ &= \frac{(K_0-1)^2}{(K+K_0-1)^2} (\mathcal{L}(\theta^0) - \mathcal{L}(\theta^*)) + \frac{\tilde{C}_5 K}{\mu^2\tilde{C}_4^2 (K+K_0-2)^2} \end{aligned}$$

from which the proof is complete.

13 Additional Numerical Results

13.1 Logistic regression.

Data pre-processing. For *ijcnn1* and *covtype* datasets, they are imported from the popular library LIBSVM¹ without further preprocessing. For *MNIST*, we normalize the data and subtract the mean. We uniformly partition *ijcnn1* dataset with 91,701 samples and *MNIST* dataset with 60,000 samples into $M = 10$ workers. To simulate the heterogeneous setting, we partition *covtype* dataset with 581,012 samples randomly into $M = 20$ workers with different number of samples per worker.

For *covtype*, we fix the batch ratio to be 0.001 uniformly across all workers; and for *ijcnn1* and *MNIST*, we fix the batch ratio to be 0.01 uniformly across all workers.

Choice of hyperparameters. For the logistic regression task, the hyperparameters in each algorithm are chosen by hand to roughly optimize the training loss performance of each algorithm. We list the values of parameters used in each test in Tables 1-2.

Algorithm	stepsize α	momentum weight β	averaging interval H/D
FedAdam	$\alpha_l = 100 \alpha_s = 0.02$	0.9	$H = 10$
Local momentum	0.1	0.9	$H = 10$
ADAM	0.005	$\beta_1 = 0.9 \beta_2 = 0.999$	/
CADA1&2	0.005	$\beta_1 = 0.9 \beta_2 = 0.999$	$D = 100, d_{\max} = 10$
Stochastic LAG	0.1	/	$d_{\max} = 10$

Table 1: Choice of parameters in *covtype*.

¹<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

Algorithm	stepsize α	momentum weight β	averaging interval H/D
FedAdam	$\alpha_l = 100 \alpha_s = 0.03$	0.9	$H = 10$
Local momentum	0.1	0.9	$H = 20$
ADAM	0.01	$\beta_1 = 0.9 \beta_2 = 0.999$	/
CADA	0.01	$\beta_1 = 0.9 \beta_2 = 0.999$	$D = 100, d_{\max} = 10$
Stochastic LAG	0.1	/	$d_{\max} = 10$

Table 2: Choice of parameters in *ijcnn1*.

13.2 Training neural networks.

For training neural networks, we use the cross-entropy loss but with different neural network models.

Neural network models. For *MNIST* dataset, we use a convolutional neural network with two convolution-ELUmaxpooling layers (ELU is a smoothed ReLU) followed by two fully-connected layers. The first convolution layer is $5 \times 5 \times 20$ with padding, and the second layer is $5 \times 5 \times 50$ with padding. The output of second layer is followed by two fully connected layers with one being 800×500 and the other being 500×10 . The output goes through a softmax function. For *CIFAR10* dataset, we use the popular neural network architecture *ResNet20*² which has 20 and roughly 0.27 million parameters. We do not use a pre-trained model.

Data pre-processing. We uniformly partition *MNIST* and *CIFAR10* datasets into $M = 10$ workers. For *MNIST*, we use the raw data without preprocessing. The minibatch size per worker is 12. For *CIFAR10*, in addition to normalizing the data and subtracting the mean, we randomly flip and crop part of the original image every time it is used for training. This is a standard technique of data augmentation to avoid over-fitting. The minibatch size for *CIFAR10* is 50 per worker.

Choice of hyperparameters. For *MNIST* dataset which is relatively easy, the hyperparameters in each algorithm are chosen by hand to optimize the performance of each algorithm. We list the values of parameters used in each test in Table 3.

Algorithm	stepsize α	momentum weight β	averaging interval H/D
FedAdam	$\alpha_l = 0.1 \alpha_s = 0.001$	0.9	$H = 8$
Local momentum	0.001	0.9	$H = 8$
ADAM	0.0005	$\beta_1 = 0.9 \beta_2 = 0.999$	/
CADA1&2	0.0005	$\beta_1 = 0.9 \beta_2 = 0.999$	$D = 50, d_{\max} = 10$
Stochastic LAG	0.1	/	$d_{\max} = 10$

Table 3: Choice of parameters in multi-class *MNIST*.

For *CIFAR10* dataset, we search the best values of hyperparameters from the following search grid on a per-algorithm basis to optimize the testing accuracy versus the number of communication rounds. The chosen values of parameter are listed in Table 4.

FedAdam: $\alpha_s \in \{0.1, 0.01, 0.001\}$; $\alpha_l \in \{1, 0.5, 0.1\}$; $H \in \{1, 4, 6, 8, 16\}$.

Local momentum: $\alpha \in \{0.1, 0.01, 0.001\}$; $H \in \{1, 4, 6, 8, 16\}$.

CADA1: $\alpha \in \{0.1, 0.01, 0.001\}$; $c \in \{0.05, 0.1, 0.3, 0.6, 0.9, 1.2, 1.5, 1.8\}$.

CADA2: $\alpha \in \{0.1, 0.01, 0.001\}$; $c \in \{0.05, 0.1, 0.3, 0.6, 0.9, 1.2, 1.5, 1.8\}$.

LAG: $\alpha \in \{0.1, 0.01, 0.001\}$; $c \in \{0.05, 0.1, 0.3, 0.6, 0.9, 1.2, 1.5, 1.8\}$.

²https://github.com/akamaster/pytorch_resnet_cifar10

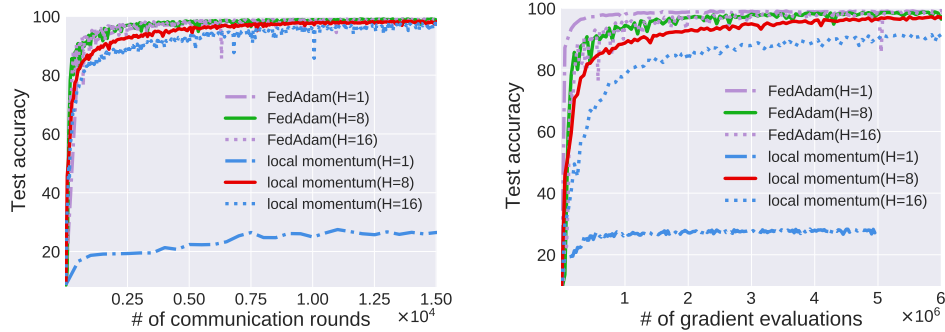


Figure 6: Performance of FedAdam and local momentum on *MNIST* under different averaging interval H .

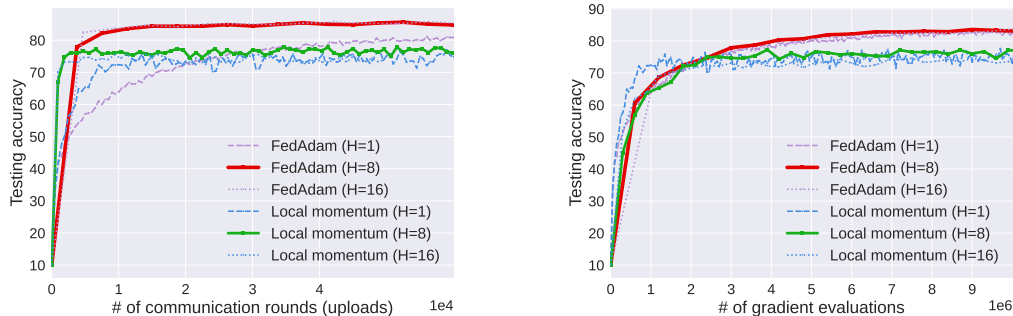


Figure 7: Performance of FedAdam and local momentum on *CIFAR10* under different averaging interval H .

Algorithm	stepsize α	momentum weight β	averaging interval H/D
FedAdam	$\alpha_l = 0.1 \ \alpha_s = 0.1$	0.9	$H = 8$
Local momentum	0.1	0.9	$H = 8$
CADA1	0.1	$\beta_1 = 0.9 \ \beta_2 = 0.99$	$D = 50, \ d_{\max} = 2$
CADA2	0.1	$\beta_1 = 0.9 \ \beta_2 = 0.99$	$D = 50, \ d_{\max} = 2$
Stochastic LAG	0.1	/	$d_{\max} = 2$

Table 4: Choice of parameters in *CIFAR10*.

Additional results. In addition to the results presented in the main paper, we report a new set of simulations on the performance of local update based algorithms under different averaging interval H . Since algorithms under $H = 4, 6$ do not perform as good as $H = 8$, we only plot $H = 1, 8, 16$ in Figures 6 and 7 to ease the comparison. Figure 6 compares the performance of FedAdam and local momentum on *MNIST* dataset under different averaging interval H . Figure 7 compares the performance of FedAdam and local momentum on *CIFAR10* dataset under different H .

Figure 7 compares the performance of FedAdam and local momentum on *CIFAR10* dataset under different averaging interval H . FedAdam and local momentum under a larger averaging interval H have faster convergence speed at the initial stage, but they reach slightly lower testing accuracy. This reduced test accuracy is common among local SGD-type methods, which has also been studied theoretically; see e.g., [14].