

APPENDIX

A PROOF OF CLAIM [1](#)

Proof. Consider n scalar functions that for each $i \in \mathcal{V}$, $f_i(x) = (1/2)(x - a_i)^2$ for some a_i , and complete graph with $\mathbf{W} = (1/n)\mathbf{1}_n\mathbf{1}_n^\top$. Let $a_i = a, \forall i = 1, \dots, n/2$, and $a_i = b, \forall i = n/2 + 1, \dots, n$, and $b - a > 2B + 1$. Let $x_i^0 = a + 0.5, \forall i \in \mathcal{V}$. Then, $\forall i \in \mathcal{V}$, vanilla normalization reduces to

$$\begin{aligned} x_i^1 &= \frac{1}{n} \sum_{r=1}^n \left(x_r^0 - \alpha \text{sign}(x_r^0 - a_r) \right) \\ &= x_r^0 - \frac{\alpha}{n} \sum_{r=1}^n \text{sign}(x_r^0 - a_r) \\ &= x_r^0 - \frac{\alpha}{n} \sum_{r=1}^{n/2} \text{sign}(0.5) - \frac{\alpha}{n} \sum_{r=n/2+1}^n \text{sign}(0.5 - (b - a)) \\ &= x_r^0. \end{aligned}$$

Therefore, $x_r^t = a + 0.5, \forall r \in \mathcal{V}, \forall t \geq 0$. Since the optimal solution to the original problem is $\frac{a+b}{2}$, the optimality gap is $\frac{b-a}{2} - 0.5 \geq B$. \square

Remark 9. Note that the proof above can be further extended to the case where the gradient oracle admits almost surely bounded gradient noise. We can use the noise bound to adapt the choices of a, b, ε such that all signs still get canceled. Similar examples have been used to show divergence results in [Shulgin et al. \(2025\)](#).

B PROOFS OF THEOREMS

Proof structure. The central recursion of our analysis leverages the descent lemma for L smooth functions applied to consecutive network averages of $\{\mathbf{x}_i^{t+1}\}$ and $\{\mathbf{x}_i^t\}$ ([Lemma 3](#)). This recursion involves two coupled error sources: consensus errors between $\{\mathbf{x}_i^t\}$, $\{\mathbf{y}_i^t\}$, and gradient estimation errors. We establish the intricate coupling between these errors through a series of intermediate lemmas. Our proof strategy proceeds as follows: we first derive two key lemmas that bound consensus errors in terms of gradient estimation errors ([Lemma 4](#) and [6](#)), then decompose gradient estimation errors into constituent components, including average gradient estimation errors ([Lemma 7](#)) and stack gradient estimation errors ([Lemma 8](#)), and bound each separately. With all error bounds established, we substitute these into the main recursion and optimize hyperparameter selection to achieve the final order-optimal convergence rates.

B.1 PRELIMINARIES

We define some stacked long vectors,

$$\begin{aligned} F(\mathbf{x}^t) &:= [f_1(\mathbf{x}_1^t), \dots, f_n(\mathbf{x}_n^t)]^\top, \\ \nabla F(\mathbf{x}^t) &:= [\nabla f_1(\mathbf{x}_1^t)^\top, \dots, \nabla f_n(\mathbf{x}_n^t)^\top]^\top, \\ \mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) &:= [\mathbf{g}_1(\mathbf{x}_1^t, \boldsymbol{\xi}_1^t)^\top, \dots, \mathbf{g}_n(\mathbf{x}_n^t, \boldsymbol{\xi}_n^t)^\top]^\top \\ \mathbf{v}^t &:= [(\mathbf{v}_1^t)^\top, \dots, (\mathbf{v}_n^t)^\top]^\top, \\ \mathcal{N}(\mathbf{y}^t) &:= \left[\frac{(\mathbf{y}_1^t)^\top}{\|\mathbf{y}_1^t\|}, \dots, \frac{(\mathbf{y}_n^t)^\top}{\|\mathbf{y}_n^t\|} \right]^\top, \\ \mathbf{x}^t &:= [(\mathbf{x}_1^t)^\top, \dots, (\mathbf{x}_n^t)^\top]^\top. \end{aligned}$$

Then, Algorithm [1](#) can be rewritten in the compact long-vector form:

$$\mathbf{v}^t = \beta \mathbf{v}^{t-1} + (1 - \beta) \mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t); \quad (6)$$

$$\mathbf{y}^t = (\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{y}^{t-1} + \mathbf{v}^t - \mathbf{v}^{t-1}), \quad (7)$$

$$\mathbf{x}^{t+1} = (\mathbf{W} \otimes \mathbf{I}_d)(\mathbf{x}^t - \alpha \mathcal{N}(\mathbf{y}^t)). \quad (8)$$

We define the following averages over network:

$$\bar{\mathbf{v}}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^t, \quad \bar{\mathbf{y}}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^t, \quad \tilde{\mathbf{y}}^t = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{y}_i^t}{\|\mathbf{y}_i^t\|}, \quad \bar{\mathbf{x}}^t = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^t, \quad \bar{\nabla} F(\mathbf{x}^t) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_i^t). \quad (9)$$

From the doubly stochasticity of \mathbf{W} , the global average updates as

$$\bar{\mathbf{x}}^{t+1} = \bar{\mathbf{x}}^t - \frac{\alpha}{n} \sum_{r=1}^n \frac{\mathbf{y}_r^t}{\|\mathbf{y}_r^t\|} = \bar{\mathbf{x}}^t - \alpha \tilde{\mathbf{y}}^t. \quad (10)$$

B.2 INTERMEDIATE LEMMAS

We first present some standard useful relations to be used in our analysis.

Lemma 1. *The following relations hold:*

1. $\bar{\mathbf{y}}^t = \bar{\mathbf{v}}^t$;
2. $\mathbf{W} - \mathbf{1}_n \mathbf{1}_n^\top / n = (\mathbf{W} - \mathbf{1}_n \mathbf{1}_n^\top / n)(\mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n^\top / n) = (\mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n^\top / n)(\mathbf{W} - \mathbf{1}_n \mathbf{1}_n^\top / n)$;
3. $\mathbf{W}^k - \mathbf{1}_n \mathbf{1}_n^\top / n = (\mathbf{W} - \mathbf{1}_n \mathbf{1}_n^\top / n)^k, \forall k \in \mathbb{N}_+$;
4. $(1/\sqrt{n}) \sum_{i=1}^n \|\mathbf{a}_i\| \leq \|\mathbf{a}\| \leq \sum_{i=1}^n \|\mathbf{a}_i\|, \forall \mathbf{a} = [\mathbf{a}_1^\top, \dots, \mathbf{a}_n^\top]^\top \in \mathbb{R}^{nd}$,
5. $\sum_{i=1}^m a_i^p \leq (\sum_{i=1}^m a_i)^p \leq m^{p-1} \sum_{i=1}^m a_i^p, \forall m \in \mathbb{N}_+, \forall a_i \in \mathbb{R}_+$.

We then present a standard decent lemma for L -smooth functions .

Lemma 2 (Decent lemma for L -smooth functions). *Let Assumption [2](#) hold. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, there holds*

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

We next present the main descent lemma on the network average.

Lemma 3 (Decent lemma for network average). *Let Assumption [2](#) hold. Let $\boldsymbol{\epsilon}^t = \bar{\mathbf{y}}^t - \nabla f(\bar{\mathbf{x}}^t)$. We have*

$$\sum_{t=0}^{T-1} \alpha \|\nabla f(\bar{\mathbf{x}}^t)\| \leq f(\bar{\mathbf{x}}^0) - f_* + \sum_{t=0}^{T-1} 2\alpha \|\boldsymbol{\epsilon}^t\| + \sum_{t=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^n \|\bar{\mathbf{y}}^t - \mathbf{y}_i^t\| + \sum_{t=0}^{T-1} \frac{L}{2} \alpha^2.$$

Proof. Since $\|\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t\| = \alpha \|\tilde{\mathbf{y}}^t\| = \alpha$, applying Lemma [2](#) on $\bar{\mathbf{x}}^{t+1}, \bar{\mathbf{x}}^t$ gives that

$$\begin{aligned} f(\bar{\mathbf{x}}^{t+1}) &\leq f(\bar{\mathbf{x}}^t) + \nabla f(\bar{\mathbf{x}}^t)^\top (\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t) + \frac{L}{2} \|\bar{\mathbf{x}}^{t+1} - \bar{\mathbf{x}}^t\|^2 \\ &\stackrel{(i)}{\leq} f(\bar{\mathbf{x}}^t) - \alpha (\bar{\mathbf{y}}^t - \boldsymbol{\epsilon}^t)^\top \tilde{\mathbf{y}}^t + \frac{L}{2} \alpha^2 \\ &\stackrel{(ii)}{\leq} f(\bar{\mathbf{x}}^t) - \alpha (\bar{\mathbf{y}}^t)^\top \tilde{\mathbf{y}}^t + \alpha \|\boldsymbol{\epsilon}^t\| + \frac{L}{2} \alpha^2, \end{aligned} \quad (11)$$

where we used the definitions (9)(10) in (i), and used Cauchy-Schwartz inequality followed by $\|\tilde{\mathbf{y}}^t\| \leq 1$ in (ii). Next,

$$\begin{aligned}
-(\bar{\mathbf{y}}^t)^\top \tilde{\mathbf{y}}^t &= -(\mathbf{y}^t)^\top \left[\frac{\tilde{\mathbf{y}}^t}{\|\tilde{\mathbf{y}}^t\|} + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^t \left(\frac{1}{\|\mathbf{y}_i^t\|} - \frac{1}{\|\tilde{\mathbf{y}}^t\|} \right) \right] \\
&\leq -\|\tilde{\mathbf{y}}^t\| + \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^t \left(\frac{\|\tilde{\mathbf{y}}^t\|}{\|\mathbf{y}_i^t\|} - 1 \right) \\
&\stackrel{(i)}{\leq} -\|\nabla f(\bar{\mathbf{x}}^t)\| + \|\epsilon^t\| + \frac{1}{n} \sum_{i=1}^n \|\|\tilde{\mathbf{y}}^t\| - \|\mathbf{y}_i^t\|\| \\
&\stackrel{(ii)}{\leq} -\|\nabla f(\bar{\mathbf{x}}^t)\| + \|\epsilon^t\| + \frac{1}{n} \sum_{i=1}^n \|\tilde{\mathbf{y}}^t - \mathbf{y}_i^t\|, \tag{12}
\end{aligned}$$

where we used $\|\tilde{\mathbf{y}}^t\| = \|\nabla f(\bar{\mathbf{x}}^t) + \epsilon^t\| \geq \|\nabla f(\bar{\mathbf{x}}^t)\| - \|\epsilon^t\|$, and Cauchy-Schwartz inequality in (i), and $\|\|\mathbf{a}\| - \|\mathbf{b}\|\| \leq \|\mathbf{a} - \mathbf{b}\|$ for any $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ in (ii). Plugging in (12) into (11), and summing over $t = 0, \dots, T-1$, we have

$$f(\bar{\mathbf{x}}^T) \leq f(\bar{\mathbf{x}}^0) - \sum_{t=0}^{T-1} \alpha \|\nabla f(\bar{\mathbf{x}}^t)\| + \sum_{t=0}^{T-1} 2\alpha \|\epsilon^t\| + \sum_{t=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^n \|\tilde{\mathbf{y}}^t - \mathbf{y}_i^t\| + \sum_{t=0}^{T-1} \frac{L}{2} \alpha^2.$$

Using $f(\bar{\mathbf{x}}^T) \geq f_*$ and rearranging terms above give the desired result. \square

With Lemma 3, it remains to bound the gradient estimation error $\|\epsilon^t\|$ and the consensus error $\mathbf{y}_i^t - \bar{\mathbf{y}}^t$. Let us decompose the gradient estimation error as follows:

$$\epsilon^t = \bar{\mathbf{y}}^t - \nabla f(\bar{\mathbf{x}}^t) = \bar{\mathbf{v}}^t - \nabla f(\bar{\mathbf{x}}^t) = \underbrace{\bar{\mathbf{v}}^t - \bar{\nabla} F(\mathbf{x}^t)}_{:= \epsilon_1^t \in \mathbb{R}^d} + \underbrace{\bar{\nabla} F(\mathbf{x}^t) - \nabla f(\bar{\mathbf{x}}^t)}_{:= \epsilon_2^t \in \mathbb{R}^d}. \tag{13}$$

It is clear that ϵ_1^t is the gradient estimation error, and ϵ_2^t , exploiting the smoothness property in 2, can be bounded by the consensus error $\mathbf{x}_i^t - \bar{\mathbf{x}}^t$. Since the consensus error is also used in bounding ϵ_1^t , we need to first bound the consensus errors $\mathbf{x}_i^t - \bar{\mathbf{x}}^t$ and $\mathbf{y}_i^t - \bar{\mathbf{y}}^t$.

Lemma 4 (Consensus errors of $\{\mathbf{x}_i^t\}$). *We have for all $t = 0, \dots, T$,*

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i^t - \bar{\mathbf{x}}^t\| \leq \frac{\alpha\lambda}{1-\lambda}. \tag{14}$$

Proof. Using the relation 4 in Lemma 1 we have

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i^t - \bar{\mathbf{x}}^t\| \leq \frac{1}{\sqrt{n}} \|\mathbf{x}^t - \mathbf{1}_n \otimes \bar{\mathbf{x}}^t\|. \tag{15}$$

From the compact form update in (8), we have

$$\mathbf{x}^t = (\mathbf{W} \otimes I_d) \mathbf{x}^0 - \alpha \sum_{k=0}^{t-1} (\mathbf{W} \otimes I_d)^{t-k} \mathcal{N}(\mathbf{y}^k).$$

It follows that

$$\begin{aligned}
& \|\mathbf{x}^t - \mathbf{1}_n \otimes \bar{\mathbf{x}}^t\| \\
&= \|(\mathbf{I}_{nd} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d)\mathbf{x}^t\| \\
&\stackrel{(i)}{=} \alpha \left\| \sum_{k=0}^{t-1} (\mathbf{I}_{nd} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) (\mathbf{W} \otimes \mathbf{I}_d)^{t-k} \mathcal{N}(\mathbf{y}^k) \right\| \\
&\leq \alpha \sum_{k=0}^{t-1} \|\mathbf{W}^{t-k} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top\|_2 \|\mathcal{N}(\mathbf{y}^k)\| \\
&\stackrel{(ii)}{\leq} \alpha \sum_{k=0}^{t-1} \|\mathbf{W} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top\|_2^{t-k} \|\mathcal{N}(\mathbf{y}^k)\| \\
&\leq \alpha\sqrt{n} \sum_{k=0}^{t-1} \lambda^{t-k} \tag{16}
\end{aligned}$$

$$\stackrel{(iii)}{\leq} \frac{\alpha\sqrt{n}\lambda}{1-\lambda}. \tag{17}$$

where we used the double stochasticity of \mathbf{W} and $\mathbf{x}_i^0 = \bar{\mathbf{x}}^0, \forall i \in [n]$ in (i), the relation 3 in Lemma 1 in (ii), and Assumption 4 in (iii). Substituting (17) into (15) gives the desired bound in (14). \square

Before proceeding to bound consensus errors for $\{\mathbf{y}_i^t\}$, we present the following bound on vector-valued martingale difference sequence from Liu & Zhou (2025).

Lemma 5. *Given a sequence of random vectors $\mathbf{d}_t \in \mathbb{R}^d, \forall t$ such that $\mathbb{E}[\mathbf{d}_t | \mathcal{F}_{t-1}] = \mathbf{0}$ where $\mathcal{F}_t = \sigma(\mathbf{d}_1, \dots, \mathbf{d}_t)$ is the natural filtration, then for any $p \in [1, 2]$, there is*

$$\mathbb{E} \left[\left\| \sum_{t=1}^T \mathbf{d}_t \right\|^p \right] \leq 2\sqrt{2} \mathbb{E} \left[\left(\sum_{t=1}^T \|\mathbf{d}_t\|^p \right)^{\frac{1}{p}} \right], \forall T \geq 0.$$

Lemma 6 (Consensus errors for $\{\mathbf{y}_i^t\}$). *We have for all $t = 0, \dots, T$,*

$$\begin{aligned}
& \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^n \|\mathbf{y}_i^t - \bar{\mathbf{y}}^t\| \right] \\
&\leq 2\sqrt{2}n^{\frac{1}{2}} \left(\frac{1}{\beta} - 1 \right) \left(\sum_{k=0}^t \lambda^{(t-k+1)p} \right)^{\frac{1}{p}} \sigma + \frac{1}{\sqrt{n}} \left(\frac{1}{\beta} - 1 \right) \sum_{k=0}^t \lambda^{t-k+1} \mathbb{E} [\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|].
\end{aligned}$$

Proof. Similar to (15), we have

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_i^t - \bar{\mathbf{y}}^t\| \leq \frac{1}{\sqrt{n}} \|\mathbf{y}^t - \mathbf{1}_n \otimes \bar{\mathbf{y}}^t\|. \tag{18}$$

Following from (7),

$$\begin{aligned}
& \mathbf{y}^t - \mathbf{1}_n \otimes \bar{\mathbf{y}}^t \tag{19} \\
&\stackrel{(7)}{=} (\mathbf{I}_{nd} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) (\mathbf{W} \otimes \mathbf{I}_d) (\mathbf{y}^{t-1} + \mathbf{v}^t - \mathbf{v}^{t-1}) \\
&= (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) \mathbf{y}^{t-1} + (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) (\mathbf{v}^t - \mathbf{v}^{t-1}) \\
&\stackrel{(i)}{=} (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) (\mathbf{I}_{nd} - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) \mathbf{y}^{t-1} + (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d) (\mathbf{v}^t - \mathbf{v}^{t-1}) \\
&\stackrel{(ii)}{=} \sum_{k=0}^t (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \otimes \mathbf{I}_d)^{t-k+1} (\mathbf{v}^t - \mathbf{v}^{t-1}), \tag{20}
\end{aligned}$$

where we used relation 3 in Lemma 1 in (i) and used $\mathbf{y}_i^0 = \mathbf{0}_d, \forall i \in [n]$ in (ii). From the update in (6), we have

$$\mathbf{v}^t - \mathbf{v}^{t-1} = (\beta - 1)\mathbf{v}^{t-1} + (1 - \beta)\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) = (1 - \beta)(\mathbf{v}^t - \mathbf{v}^{t-1}) + (1 - \beta)(\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \mathbf{v}^t).$$

Then, there holds,

$$\mathbf{v}^t - \mathbf{v}^{t-1} = \left(\frac{1}{\beta} - 1\right)(\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \mathbf{v}^t) = \left(\frac{1}{\beta} - 1\right)(\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \nabla F(\mathbf{x}^t) + \nabla F(\mathbf{x}^t) - \mathbf{v}^t). \quad (21)$$

Putting the relation above into (20) and applying (20) recursively, from $\mathbf{y}_i^0 = \bar{\mathbf{y}}^0$, we have

$$\begin{aligned} \|\mathbf{y}^t - \mathbf{1}_n \otimes \bar{\mathbf{y}}^t\| &\leq \left(\frac{1}{\beta} - 1\right) \left\| \sum_{k=0}^t (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d)^{t-k+1} (\mathbf{g}(\mathbf{x}^k, \boldsymbol{\xi}^k) - \nabla F(\mathbf{x}^k)) \right\| \\ &\quad + \left(\frac{1}{\beta} - 1\right) \left\| \sum_{k=0}^t (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d)^{t-k+1} (\nabla F(\mathbf{x}^k) - \mathbf{v}^k) \right\| \end{aligned} \quad (22)$$

We note that the first half of the right hand side above can be addressed by Lemma 5:

$$\begin{aligned} &\mathbb{E} \left[\left\| \sum_{k=0}^t (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d)^{t-k+1} (\mathbf{g}(\mathbf{x}^k, \boldsymbol{\xi}^k) - \nabla F(\mathbf{x}^k)) \right\|^p \right] \\ &\leq 2\sqrt{2} \mathbb{E} \left[\left(\sum_{k=0}^t \lambda^{(t-k+1)p} \|\mathbf{g}(\mathbf{x}^k, \boldsymbol{\xi}^k) - \nabla F(\mathbf{x}^k)\|^p \right)^{\frac{1}{p}} \right]. \end{aligned} \quad (23)$$

We observe that

$$\begin{aligned} &2\sqrt{2} \mathbb{E} \left[\left(\sum_{k=0}^t \lambda^{(t-k+1)p} \|\mathbf{g}(\mathbf{x}^k, \boldsymbol{\xi}^k) - \nabla F(\mathbf{x}^k)\|^p \right)^{\frac{1}{p}} \mid \mathcal{F}_{t-1} \right] \\ &\leq 2\sqrt{2} \mathbb{E} \left[\left(\sum_{k=0}^t \lambda^{(t-k+1)p} \left(\sum_{i=1}^n \|\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k)\|^p \right)^{\frac{1}{p}} \mid \mathcal{F}_{t-1} \right) \right] \\ &\stackrel{(i)}{\leq} 2\sqrt{2} \mathbb{E} \left[\left(\sum_{k=0}^t \sum_{i=1}^n \lambda^{(t-k+1)p} n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k)\|^p \right)^{\frac{1}{p}} \mid \mathcal{F}_{t-1} \right] \\ &\stackrel{(ii)}{\leq} 2\sqrt{2} \left(\mathbb{E} \left[\sum_{i=1}^n \lambda^p n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p \mid \mathcal{F}_{t-1} \right] \right. \\ &\quad \left. + \sum_{k=0}^{t-1} \sum_{i=1}^n \lambda^{(t-k+1)p} n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k)\|^p \right)^{\frac{1}{p}} \\ &\stackrel{(iii)}{\leq} 2\sqrt{2} \left(\lambda^p n^p \sigma^p + \sum_{k=0}^{t-1} \sum_{i=1}^n \lambda^{(t-k+1)p} n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k)\|^p \right)^{\frac{1}{p}}, \end{aligned} \quad (24)$$

where we used relation 5 from Lemma 1 in (i), Jensen's inequality in (ii), and Assumption 3 in (iii). From (23), taking expectations on both sides of (24), and applying the above arguments recursively from \mathcal{F}_{t-2} to \mathcal{F}_0 , we have

$$\mathbb{E} \left[\left\| \sum_{k=0}^t (\mathbf{W} \otimes \mathbf{I}_d - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \otimes \mathbf{I}_d)^{t-k+1} (\mathbf{g}(\mathbf{x}^k, \boldsymbol{\xi}^k) - \nabla F(\mathbf{x}^k)) \right\|^p \right] \leq 2\sqrt{2} \left(\sum_{k=0}^t \lambda^{(t-k+1)p} \right)^{\frac{1}{p}} n \sigma.$$

Therefore, using the relation above, and (18), (22), we reach the desired relation. \square

We then bound average gradient estimation errors $\boldsymbol{\epsilon}_1^t = \bar{\mathbf{v}}^t - \bar{\nabla} F(\mathbf{x}^t)$.

Lemma 7 (Average gradient estimation errors). *For all $t = 0, \dots, T$, we have*

$$\mathbb{E} [\|\bar{\mathbf{v}}^t - \bar{\nabla} F(\mathbf{x}^t)\|] \leq \beta^{t+1} \|\nabla f(\bar{\mathbf{x}}^0)\| + \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}} \left(\sum_{k=0}^t \beta^{(t-k)p} (1 - \beta)^p \right)^{\frac{1}{p}} \sigma + \sum_{k=0}^t \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha \right) L.$$

Proof. Following from the step 4 in Algorithm [1](#) $\forall i \in [n]$,

$$\mathbf{v}_i^t - \nabla f_i(\mathbf{x}_i^t) = \beta(\mathbf{v}_i^{t-1} - \nabla f_i(\mathbf{x}_i^{t-1})) + (1 - \beta)(\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)) + \beta(\nabla f_i(\mathbf{x}_i^{t-1}) - \nabla f_i(\mathbf{x}_i^t)). \quad (25)$$

Averaging the above relation over $i = 1, \dots, n$ leads to that:

$$\begin{aligned} \boldsymbol{\epsilon}_1^t &= \bar{\mathbf{v}}^t - \bar{\nabla} F(\mathbf{x}^t) \\ &= \beta(\bar{\mathbf{v}}^{t-1} - \bar{\nabla} F(\mathbf{x}^{t-1})) + (1 - \beta) \cdot \underbrace{\frac{1}{n} \sum_{i=1}^n (\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t))}_{:= \mathbf{s}^t \in \mathbb{R}^d} + \beta \cdot \underbrace{\frac{1}{n} \sum_{i=1}^n (\nabla f_i(\mathbf{x}_i^{t-1}) - \nabla f_i(\mathbf{x}_i^t))}_{:= \mathbf{z}^t \in \mathbb{R}^d} \\ &= \beta^{t+1} \boldsymbol{\epsilon}_1^{-1} + \sum_{k=0}^t \beta^{t-k} (1 - \beta) \mathbf{s}^k + \sum_{k=0}^t \beta^{t-k+1} \mathbf{z}^k. \end{aligned}$$

Taking Euclidean norms on both sides gives that

$$\|\boldsymbol{\epsilon}_1^t\| \leq \beta^{t+1} \|\boldsymbol{\epsilon}_1^{-1}\| + \left\| \sum_{k=0}^t \beta^{t-k} (1 - \beta) \mathbf{s}^k \right\| + \left\| \sum_{k=0}^t \beta^{t-k+1} \mathbf{z}^k \right\|. \quad (26)$$

We now bound the terms on the right hand side of [\(26\)](#) one by one. First,

$$\|\boldsymbol{\epsilon}_1^{-1}\| = \|\bar{\mathbf{v}}^{-1} - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\bar{\mathbf{x}}^0)\| = \|\nabla f(\bar{\mathbf{x}}^0)\|. \quad (27)$$

Second, notice that $\{\beta^{t-k}(1 - \beta)(\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\}$ is a martingale difference sequence that falls into the pursuit of Lemma [5](#) and thus we obtain

$$\begin{aligned} &\mathbb{E} \left[\left\| \sum_{k=0}^t \beta^{t-k} (1 - \beta) \mathbf{s}^k \right\| \right] \\ &= \frac{1}{n} \mathbb{E} \left[\left\| \sum_{k=0}^t \sum_{i=1}^n \beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k)) \right\| \right] \\ &\leq \frac{2\sqrt{2}}{n} \mathbb{E} \left[\left(\sum_{k=0}^t \sum_{i=1}^n \|\beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\|^p \right)^{\frac{1}{p}} \right]. \end{aligned} \quad (28)$$

Note that

$$\begin{aligned} &\frac{2\sqrt{2}}{n} \mathbb{E} \left[\left(\sum_{k=0}^t \sum_{i=1}^n \|\beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\|^p \right)^{\frac{1}{p}} \mid \mathcal{F}_{t-1} \right] \\ &\stackrel{(i)}{\leq} \frac{2\sqrt{2}}{n} \left(\mathbb{E} \left[\sum_{k=0}^t \sum_{i=1}^n \|\beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\|^p \mid \mathcal{F}_{t-1} \right] \right)^{\frac{1}{p}} \\ &\leq \frac{2\sqrt{2}}{n} \left(\mathbb{E} \left[\sum_{i=1}^n (1 - \beta)^p \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p \mid \mathcal{F}_{t-1} \right] \right. \\ &\quad \left. + \sum_{k=0}^{t-1} \sum_{i=1}^n \|\beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\|^p \right)^{\frac{1}{p}} \\ &\stackrel{(ii)}{\leq} \frac{2\sqrt{2}}{n} \left(n(1 - \beta)^p \sigma^p + \sum_{k=0}^{t-1} \sum_{i=1}^n \|\beta^{t-k} (1 - \beta) (\mathbf{g}_i(\mathbf{x}_i^k, \boldsymbol{\xi}_i^k) - \nabla f_i(\mathbf{x}_i^k))\|^p \right)^{\frac{1}{p}}, \end{aligned} \quad (29)$$

where we used Jensen's inequality in (i) and Assumption [3](#) in (ii). From [\(28\)](#), taking expectations on [\(29\)](#), and recursively applying the preceding arguments from \mathcal{F}_{t-2} to \mathcal{F}_0 , we have

$$\mathbb{E} \left[\left\| \sum_{k=0}^t \beta^{t-k} (1 - \beta) \mathbf{s}^k \right\| \right] \leq \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}} \left(\sum_{k=0}^t \beta^{(t-k)p} (1 - \beta)^p \right)^{\frac{1}{p}} \sigma. \quad (30)$$

Third,

$$\begin{aligned}
& \left\| \sum_{k=0}^t \beta^{t-k+1} \mathbf{z}^k \right\| \\
& \leq \sum_{k=0}^t \beta^{t-k+1} \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_i(\mathbf{x}_i^{k-1}) - \nabla f_i(\mathbf{x}_i^k)) \right\| \\
& \leq \sum_{k=0}^t \beta^{t-k+1} \left(\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}_i^{k-1}) - \nabla f_i(\bar{\mathbf{x}}^{k-1})\| + \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\bar{\mathbf{x}}^{k-1}) - \nabla f_i(\bar{\mathbf{x}}^k)\| \right. \\
& \quad \left. + \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\bar{\mathbf{x}}^k) - \nabla f_i(\mathbf{x}_i^k)\| \right) \\
& \stackrel{(i)}{\leq} \sum_{k=0}^t \beta^{t-k+1} \left(L \cdot \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i^{k-1} - \bar{\mathbf{x}}^{k-1}\| + L \cdot \frac{1}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}^{k-1} - \bar{\mathbf{x}}^k\| + L \cdot \frac{1}{n} \sum_{i=1}^n \|\bar{\mathbf{x}}^k - \mathbf{x}_i^k\| \right) \\
& \stackrel{(ii)}{\leq} \sum_{k=0}^t \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha \right) L.
\end{aligned} \tag{31}$$

where in (i) we used Assumption 2 and in (ii) we used (14) in Lemma 4. Putting relations (27) (30) (31) together leads to the final bound for this lemma. \square

We next bound the stacked gradient estimation errors.

Lemma 8 (Stacked gradient estimation errors). *For all $t = 0, \dots, T$, we have*

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{v}^t - \nabla F(\mathbf{x}^t)\|] \\
& \leq \beta^{t+1} \|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\| + 2\sqrt{2} \left(\sum_{k=0}^t \beta^{(t-k)p} (1-\beta)^p \right)^{\frac{1}{p}} n\sigma + n \sum_{k=0}^t \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha \right) L.
\end{aligned}$$

Proof. Define $\tilde{\boldsymbol{\epsilon}}_1^t := \mathbf{v}^t - \nabla F(\mathbf{x}^t) \in \mathbb{R}^{nd}$. Similar to (25), we have

$$\begin{aligned}
\mathbf{v}^t - \nabla F(\mathbf{x}^t) &= \beta(\mathbf{v}^{t-1} - \nabla F(\mathbf{x}^{t-1})) + (1-\beta) \underbrace{(\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \nabla F(\mathbf{x}^t))}_{:= \tilde{\mathbf{s}}^t \in \mathbb{R}^{nd}} + \beta \underbrace{(\nabla F(\mathbf{x}^{t-1}) - \nabla F(\mathbf{x}^t))}_{:= \tilde{\mathbf{z}}^t \in \mathbb{R}^{nd}} \\
&= \beta^{t+1} \tilde{\boldsymbol{\epsilon}}_1^{-1} + \sum_{k=0}^t \beta^{t-k} (1-\beta) \tilde{\mathbf{s}}^k + \sum_{k=0}^t \beta^{t-k+1} \tilde{\mathbf{z}}^k.
\end{aligned}$$

Taking Euclidean norms on both sides gives that

$$\|\tilde{\boldsymbol{\epsilon}}_1^t\| \leq \beta^{t+1} \|\tilde{\boldsymbol{\epsilon}}_1^{-1}\| + \left\| \sum_{k=0}^t \beta^{t-k} (1-\beta) \tilde{\mathbf{s}}^k \right\| + \left\| \sum_{k=0}^t \beta^{t-k+1} \tilde{\mathbf{z}}^k \right\|.$$

Similar to the analysis in Lemma 7, we bound the right hand side above term by term. First,

$$\|\tilde{\boldsymbol{\epsilon}}_1^{-1}\| = \|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|.$$

Second, notice also that $\{\beta^{t-k} (1-\beta) \tilde{\mathbf{s}}^k\}$ is a martingale difference sequence and can be dealt with using Lemma 5. We have

$$\begin{aligned}
& \mathbb{E} \left[\left\| \sum_{k=0}^t \beta^{t-k} (1-\beta) \tilde{\mathbf{s}}^k \right\| \right] \\
&= \mathbb{E} \left[\left\| \sum_{k=0}^t \beta^{t-k} (1-\beta) (\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^{t,b}) - \nabla F(\mathbf{x}^t)) \right\| \right] \\
&\leq 2\sqrt{2} \mathbb{E} \left[\left(\sum_{k=0}^t \beta^{(t-k)p} (1-\beta)^p \|\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \nabla F(\mathbf{x}^t)\|^p \right)^{\frac{1}{p}} \right].
\end{aligned} \tag{32}$$

In addition,

$$\begin{aligned}
& 2\sqrt{2}\mathbb{E}\left[\left(\sum_{k=0}^t \beta^{(t-k)p}(1-\beta)^p \|\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \nabla F(\mathbf{x}^t)\|^p\right)^{\frac{1}{p}} \mid \mathcal{F}_{t-1}\right] \\
& \stackrel{(i)}{\leq} 2\sqrt{2}\left(\mathbb{E}\left[\sum_{k=0}^t \beta^{(t-k)p}(1-\beta)^p \|\mathbf{g}(\mathbf{x}^t, \boldsymbol{\xi}^t) - \nabla F(\mathbf{x}^t)\|^p \mid \mathcal{F}_{t-1}\right]\right)^{\frac{1}{p}} \\
& \stackrel{(ii)}{\leq} 2\sqrt{2}\left(\mathbb{E}\left[\sum_{k=0}^t \beta^{(t-k)p}(1-\beta)^p \left(\sum_{i=1}^n \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|\right)^p \mid \mathcal{F}_{t-1}\right]\right)^{\frac{1}{p}} \\
& \stackrel{(iii)}{\leq} 2\sqrt{2}\left(\mathbb{E}\left[\sum_{k=0}^t \sum_{i=1}^n \beta^{(t-k)p}(1-\beta)^p n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p \mid \mathcal{F}_{t-1}\right]\right)^{\frac{1}{p}} \tag{33} \\
& = 2\sqrt{2}\left(\mathbb{E}\left[\sum_{i=1}^n (1-\beta)^p n^{p-1} \|\nabla \mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p \mid \mathcal{F}_{t-1}\right]\right. \\
& \quad \left. + \sum_{k=0}^{t-1} \sum_{i=1}^n \beta^{(t-k)p}(1-\beta)^p n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p\right)^{\frac{1}{p}} \\
& \leq 2\sqrt{2}\left((1-\beta)^p n^p \sigma^p + \sum_{k=0}^{t-1} \sum_{i=1}^n \beta^{(t-k)p}(1-\beta)^p n^{p-1} \|\mathbf{g}_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \nabla f_i(\mathbf{x}_i^t)\|^p\right)^{\frac{1}{p}},
\end{aligned}$$

where in (i) we used Jensen's inequality, and in (ii), (iii) we used relations 4, and 5 in Lemma [1](#) respectively. Based on [\(32\)](#), taking expectations on both sides of [\(33\)](#), and applying the above arguments from \mathcal{F}_{t-2} to \mathcal{F}_0 , we obtain

$$\mathbb{E}\left[\left\|\sum_{k=0}^t \beta^{t-k}(1-\beta)\tilde{\mathbf{s}}^k\right\|\right] \leq 2\sqrt{2}\left(\sum_{k=0}^t \beta^{(t-k)p}(1-\beta)^p\right)^{\frac{1}{p}} n\sigma.$$

Third,

$$\begin{aligned}
& \left\|\sum_{k=0}^t \beta^{t-k+1}\tilde{\mathbf{z}}^k\right\| \\
& \leq \sum_{k=0}^t \beta^{t-k+1} \|\nabla F(\mathbf{x}^{k-1}) - \nabla F(\mathbf{x}^k)\| \\
& \leq \sum_{k=0}^t \sum_{i=1}^n \beta^{t-k+1} \|\nabla f_i(\mathbf{x}_i^{k-1}) - \nabla f_i(\mathbf{x}_i^k)\| \\
& \leq \sum_{k=0}^t \sum_{i=1}^n \beta^{t-k+1} \left(\|\nabla f_i(\mathbf{x}_i^{k-1}) - \nabla f_i(\bar{\mathbf{x}}^{k-1})\| + \|\nabla f_i(\bar{\mathbf{x}}^{k-1}) - \nabla f_i(\bar{\mathbf{x}}^k)\| + \|\nabla f_i(\mathbf{x}_i^k) - \nabla f_i(\bar{\mathbf{x}}^k)\|\right) \\
& \stackrel{(i)}{\leq} n \sum_{k=0}^t \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha\right) L.
\end{aligned}$$

where (i) follows from similar arguments in [\(30\)](#). \square

Now we are ready to prove our main theorems.

Proof of Theorem [7](#) We observe that

$$\begin{aligned}
\frac{1}{n} \sum_{t=0}^{T-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f(\mathbf{x}_i^t)\|] & \leq \frac{1}{n} \sum_{t=0}^{T-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f(\mathbf{x}_i^t) - \nabla f(\bar{\mathbf{x}}^t)\| + \|\nabla f(\bar{\mathbf{x}}^t)\|] \\
& \stackrel{(14)}{\leq} T \cdot \frac{\alpha\lambda L}{1-\lambda} + \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\bar{\mathbf{x}}^t)\|].
\end{aligned} \tag{34}$$

From Lemmas 3, (13), and Lemma 6,

$$\begin{aligned}
& \sum_{t=0}^{T-1} \alpha \|\nabla f(\bar{\mathbf{x}}^t)\| \\
& \leq f(\bar{\mathbf{x}}^0) - f_* + \sum_{t=0}^{T-1} 2\alpha (\|\epsilon_1^t\| + \|\bar{\nabla} F(\mathbf{x}^t) - \nabla f(\bar{\mathbf{x}}^t)\|) + \sum_{t=0}^{T-1} \frac{\alpha}{n} \sum_{i=1}^n \|\bar{\mathbf{y}}^t - \mathbf{y}_i^t\| + \sum_{t=0}^{T-1} \frac{L}{2} \alpha^2 \\
& \leq f(\bar{\mathbf{x}}^0) - f_* \\
& \quad + \sum_{t=0}^{T-1} 2\alpha \left[\beta^{t+1} \|\nabla f(\bar{\mathbf{x}}^0)\| + \frac{2\sqrt{2}}{n^{1-\frac{1}{p}}} \left(\sum_{k=0}^t \beta^{(t-k)p} (1-\beta)^p \right)^{\frac{1}{p}} \sigma + \sum_{k=0}^t \beta^{t-k+1} \left(\frac{2\alpha\lambda}{1-\lambda} + \alpha \right) L + \frac{\alpha\lambda L}{1-\lambda} \right] \\
& \quad + \sum_{t=0}^{T-1} \alpha \left[2\sqrt{2}n^{\frac{1}{2}} \left(\frac{1}{\beta} - 1 \right) \left(\sum_{k=0}^t \lambda^{(t-k+1)p} \right)^{\frac{1}{p}} \sigma + \frac{1}{\sqrt{n}} \left(\frac{1}{\beta} - 1 \right) \sum_{k=0}^t \lambda^{t-k+1} \mathbb{E}[\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|] \right] \\
& \quad + \frac{1}{2} \alpha^2 LT \\
& \stackrel{(i)}{\leq} f(\bar{\mathbf{x}}^0) - f_* + 2\|\nabla f(\bar{\mathbf{x}}^0)\| \cdot \frac{\alpha}{1-\beta} + \frac{4\sqrt{2}\sigma}{n^{1-\frac{1}{p}}} \cdot \alpha(1-\beta)^{1-\frac{1}{p}} T + \frac{4L}{1-\lambda} \cdot \frac{\alpha^2 T}{1-\beta} + \frac{2L}{1-\lambda} \cdot \alpha^2 T \\
& \quad + \frac{2\sqrt{2}\sigma n^{\frac{1}{2}}}{(1-\lambda)^{\frac{1}{p}}} \cdot \left(\frac{1}{\beta} - 1 \right) \alpha T + \frac{1}{2} L \cdot \alpha^2 T \\
& \quad + \frac{1}{\sqrt{n}} \left(\frac{1}{\beta} - 1 \right) \alpha \sum_{t=0}^{T-1} \sum_{k=0}^t \lambda^{t-k+1} \left(\beta^{t+1} \|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\| + 2\sqrt{2}n\sigma(1-\beta)^{1-\frac{1}{p}} + \frac{2nL}{1-\lambda} \cdot \frac{\alpha\beta}{1-\beta} \right)
\end{aligned}$$

where in (i) we used $\beta \leq 1, \lambda < 1$ and Lemma 8. Denote $f(\bar{\mathbf{x}}^0) - f_* = \Delta_0$. Dividing αT from the above relation on both sides, and putting it into (34), then rearranging terms leads to

$$\begin{aligned}
& \frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^n \mathbb{E}[\|\nabla f(\mathbf{x}_i^t)\|] \\
& \leq \frac{\Delta_0}{\alpha T} + \frac{2\|\nabla f(\bar{\mathbf{x}}^0)\|}{(1-\beta)T} + 4\sqrt{2}\sigma \cdot \frac{(1-\beta)^{1-\frac{1}{p}}}{n^{1-\frac{1}{p}}} + \frac{4L}{1-\lambda} \cdot \frac{\alpha}{1-\beta} + \left(\frac{3L}{1-\lambda} + \frac{L}{2} \right) \alpha \\
& \quad + \frac{2\sqrt{2}\sigma}{(1-\lambda)^{\frac{1}{p}}} \cdot n^{\frac{1}{2}} \left(\frac{1}{\beta} - 1 \right) + \frac{\|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|}{1-\lambda} \cdot \frac{\frac{1}{\beta} - 1}{n^{\frac{1}{2}}} + \frac{2\sqrt{2}\sigma}{1-\lambda} \cdot n^{\frac{1}{2}} \left(\frac{1}{\beta} - 1 \right) (1-\beta)^{1-\frac{1}{p}} \\
& \quad + \frac{2L}{(1-\lambda)^2} \cdot n^{\frac{1}{2}} \alpha \\
& \stackrel{(i)}{\leq} \frac{\Delta_0}{\alpha T} + \frac{2\|\nabla f(\bar{\mathbf{x}}^0)\|}{(1-\beta)T} + 4\sqrt{2}\sigma \cdot \frac{(1-\beta)^{1-\frac{1}{p}}}{n^{1-\frac{1}{p}}} + \frac{4L}{1-\lambda} \cdot \frac{\alpha}{1-\beta} + \frac{3.5L}{1-\lambda} \alpha \\
& \quad + \frac{20\sqrt{2}\sigma}{(1-\lambda)^{\frac{1}{p}}} \cdot n^{\frac{1}{2}} (1-\beta) + \frac{10\|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|}{1-\lambda} \cdot \frac{1-\beta}{n^{\frac{1}{2}}} + \frac{20\sqrt{2}\sigma}{1-\lambda} \cdot n^{\frac{1}{2}} (1-\beta)^{2-\frac{1}{p}} + \frac{2L}{(1-\lambda)^2} \cdot n^{\frac{1}{2}} \alpha \\
& \stackrel{(ii)}{\leq} O\left(\frac{\Delta_0}{T} + \frac{2\|\nabla f(\bar{\mathbf{x}}^0)\|}{(1-\beta)T} + 4\sqrt{2}\sigma \cdot \frac{(1-\beta)^{1-\frac{1}{p}}}{n^{1-\frac{1}{p}}} + \sqrt{\frac{L\Delta_0}{(1-\lambda)(1-\beta)T}} + \sqrt{\frac{3.5L\Delta_0}{(1-\lambda)T}} \right. \\
& \quad \left. + \frac{20\sqrt{2}\sigma}{(1-\lambda)^{\frac{1}{p}}} \cdot n^{\frac{1}{2}} (1-\beta) + \frac{10\|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|}{1-\lambda} \cdot \frac{1-\beta}{n^{\frac{1}{2}}} + \frac{\sigma}{1-\lambda} \cdot n^{\frac{1}{2}} (1-\beta)^{2-\frac{1}{p}} + \sqrt{\frac{n^{\frac{1}{2}}L\Delta_0}{(1-\lambda)^2 T}} \right) \\
& \stackrel{(iii)}{\leq} O\left(\frac{\Delta_0}{T} + \frac{\|\nabla f(\bar{\mathbf{x}}^0)\|}{T^{\frac{2p-2}{3p-2}}} + \frac{\sigma}{n^{1-\frac{1}{p}} T^{\frac{p-1}{3p-2}}} + \sqrt{\frac{L\Delta_0}{(1-\lambda)T^{\frac{2p-2}{3p-2}}}} + \sqrt{\frac{3.5L\Delta_0}{(1-\lambda)T}} \right. \\
& \quad \left. + \frac{\sigma n^{\frac{1}{2}}}{(1-\lambda)^{\frac{1}{p}} T^{\frac{p}{3p-2}}} + \frac{\|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|}{(1-\lambda)n^{\frac{1}{2}} T^{\frac{p}{3p-2}}} + \frac{\sigma}{1-\lambda} \frac{n^{\frac{1}{2}}}{T^{\frac{2p-1}{3p-2}}} + \sqrt{\frac{n^{\frac{1}{2}}L\Delta_0}{(1-\lambda)^2 T}} \right)
\end{aligned} \tag{35}$$

where in (i) we take $\beta \geq 1/10$, in (ii) we used

$$\alpha = \min \left(1, \sqrt{\frac{\Delta_0(1-\beta)(1-\lambda)}{4LT}}, \sqrt{\frac{\Delta_0(1-\lambda)}{3.5LT}}, \sqrt{\frac{(1-\lambda)^2\Delta_0}{2n^{\frac{1}{2}}LT}} \right), \quad (36)$$

and in (iii) we used $1 - \beta = \frac{1}{T^{\frac{p}{3p-2}}}$. \square

Proof of Theorem 2 Note that (35)(ii) still holds under the same choice of α in (36) and $\beta \geq 1/10$. Continuing with $1 - \beta = 1/\sqrt{T}$, we have

$$\begin{aligned} & \frac{1}{nT} \sum_{t=0}^{T-1} \sum_{i=1}^n \mathbb{E} [\|\nabla f(\mathbf{x}_i^t)\|] \\ & \leq O \left(\frac{\Delta_0}{T} + \frac{\|\nabla f(\bar{\mathbf{x}}^0)\|}{\sqrt{T}} + \frac{\sigma}{n^{1-\frac{1}{p}}} \cdot \frac{1}{T^{\frac{p-1}{2p}}} + \frac{1}{T^{\frac{1}{4}}} \sqrt{\frac{L\Delta_0}{1-\lambda}} + \sqrt{\frac{3.5L\Delta_0}{(1-\lambda)T}} + \frac{\sigma n^{\frac{1}{2}}}{(1-\lambda)^{\frac{1}{p}}} \frac{1}{\sqrt{T}} + \right. \\ & \quad \left. \frac{\|\nabla F(\mathbf{1}_n \otimes \bar{\mathbf{x}}^0)\|}{(1-\lambda)n^{\frac{1}{2}}} \cdot \frac{1}{\sqrt{T}} + \frac{\sigma n^{\frac{1}{2}}}{1-\lambda} \cdot \frac{1}{T^{\frac{2p-1}{2p}}} + \sqrt{\frac{n^{\frac{1}{2}}L\Delta_0}{(1-\lambda)^2T}} \right). \end{aligned}$$

Rearranging above terms leads to the desired upper bound. \square

C ADDITIONAL EXPERIMENT DETAILS

C.1 BASELINE DESCRIPTIONS

Please see Table 2 for detailed descriptions of baselines.

C.2 ADDITIONAL DETAILS FOR SYNTHETIC EXPERIMENTS

Loss function. Let $(\mathbf{X}_{i,k}, \mathbf{y}_{i,k})$ denote the k -th sample of sub-dataset $(\mathbf{X}_i, \mathbf{y}_i)$ on node i . The loss function of the considered nonconvex linear regression model on this sample is $\ell(\mathbf{y}_{i,k} - \mathbf{X}_{i,k}\mathbf{w}_i^t)$, where the

$$\ell(r) = \begin{cases} \frac{c^2}{6} \left(1 - \left[1 - \left(\frac{r}{c} \right)^2 \right]^3 \right) & \text{if } |r| \leq c, \\ \frac{c^2}{6} & \text{otherwise} \end{cases},$$

and we use the suggested value $c = 4.6851$ in the robust statistics literature.

Hyperparameter tuning. Please see Table 3 for hyperparameter searching ranges for this experiment.

Hardware. We ran this experiment on Mac OS X 15.3, CPU M4 10 Cores, RAM 16GB.

C.3 ADDITIONAL DETAILS FOR DECENTRALIZED TRAINING OF TRANSFORMERS

Transformer architecture. We consider the following decoder-only Transformer model (GPT): vocabulary size is 10208, context length is 64, embedding size is 128, number of attention heads is 4, number of attention layers is 2, the linear projection dimension within attention block is 512, and LayerNorm is applied after the 2nd attention block. The total number of parameters of this model is 3018240.

Hyperparameter tuning. See Table 4 for our grid search range for algorithm hyperparameters.

Hardware. We simulate the distributed training on one NVIDIA H100 GPU, using PyTorch 3.2 with CUDA 12. The total hyperparameter search and training procedure took around 100 GPU hours.

Table 2: Summary of Baseline Methods

Method	Parallel update on node i	Hyper-parameters
DSGD	$\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}(\mathbf{x}_r^t - \alpha g_r(\mathbf{x}_r^t, \boldsymbol{\xi}_r^t))$	α : constant stepsize
DSGD-GClip	$\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}\mathbf{x}_r^t - \alpha \text{clip}(g_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t), \tau)$	α, τ : stepsize α , and ℓ_2 clipping levels τ
DSGD-CClip	$\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}\mathbf{x}_r^t - \alpha \text{clip}(g_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t), \tau)$	α, τ : stepsize α , and component-wise clipping levels τ
DSGD-Clip	$\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}\mathbf{x}_r^t - \alpha_t \text{clip}(g_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t), \tau_t)$	α, τ : stepsize $\alpha_t = \alpha/(t+1)$, and ℓ_2 clipping levels $\tau_t = \tau(t+1)^{2/5}$
GT-DSGD	$\mathbf{y}_i^{t+1} = \sum_{r=1}^n w_{ir}(\mathbf{y}_r^t + g_r(\mathbf{x}_r^t, \boldsymbol{\xi}_r^t) - g_r(\mathbf{x}_r^{t-1}, \boldsymbol{\xi}_r^{t-1}))$ $\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}(\mathbf{x}_i^t - \alpha \mathbf{y}_i^{t+1})$	α : constant stepsize
GT-Adam	$\mathbf{m}_i^{t+1} = \beta_1 \mathbf{m}_i^t + (1 - \beta_1) \mathbf{s}_i^t$ $\mathbf{v}_i^{t+1} = \min(\beta_2 \mathbf{v}_i^t + (1 - \beta_2) \mathbf{s}_i^t \odot \mathbf{s}_i^t, G)$ $\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}\mathbf{x}_r^t - \alpha \frac{\mathbf{m}_i^{t+1}}{\sqrt{\mathbf{v}_i^{t+1} + \epsilon}}$ $\mathbf{g}_i^{t+1} = \nabla f_i(\mathbf{x}_i^{t+1})$ $\mathbf{s}_i^{t+1} = \sum_{r=1}^n w_{ir}\mathbf{s}_r^t + \mathbf{g}_i^{t+1} - \mathbf{g}_i^t$	α, G : constant stepsize α , and upper bound G , stabilization factor ϵ
QG-DSGDm	$\mathbf{m}_i^{t+1} = \beta \hat{\mathbf{m}}_i^t + g_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t)$ $\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}(\mathbf{x}_i^t - \eta \mathbf{m}_i^{t+1})$ $\mathbf{d}_i^t = (\mathbf{x}_i^{t+1} - \mathbf{x}_i^t)/\eta$ $\hat{\mathbf{m}}_i^{t+1} = \mu \hat{\mathbf{m}}_i^t + (1 - \mu) \mathbf{d}_i^t$	η, β, μ : constant stepsize η , momentum parameters β, μ
SClip-EF-Network	$\mathbf{m}_i^{t+1} = \beta_t \mathbf{m}_i^t + (1 - \beta_t) \Psi_t(g_i(\mathbf{x}_i^t, \boldsymbol{\xi}_i^t) - \mathbf{m}_i^t)$ $\mathbf{x}_i^{t+1} = \sum_{r=1}^n w_{ir}(\mathbf{x}_r^t - \alpha_t \mathbf{m}_r^{t+1})$	$c_\varphi, \tau, \alpha, \beta$: Component-wise smooth clipping operator: $\Psi_t(y) = \frac{c_\varphi}{\sqrt{t+1}} \frac{y}{\sqrt{y^2 + \tau(t+1)^{3/5}}}$, stepsize $\alpha_t = \alpha/(t+1)^{1/5}$, momentum stepsize $\beta_t = \beta/\sqrt{t+1}$.

Table 3: Hyperparameter grid search in synthetic experiments

Method	Hyperparameter search set
DSGD	$\alpha \in \{10^{-5}, 5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}$
DSGD-Clip	$\alpha \in \{10^{-5}, 5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, \tau \in \{10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10, 50, 10^2\}$
GT-DSGD	$\alpha \in \{10^{-5}, 5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}$
GT-NSGDm	$\alpha \in \{10^{-5}, 5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, \beta \in \{0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.99\}$
SCLIP-EF-Network	$\alpha \in \{10^{-3}, 10^{-2}, 0.1, 1, 10, 30\}, \beta \in \{10^{-2}, 0.1, 0.5, 0.8, 0.99\}, c_\varphi \in \{1, 5, 10, 20, 30, 50\}, \tau \in \{0.1, 1, 10, 50, 100\}$

Table 4: Hyperparameter grid search in decentralized training of Transformers

Method	Hyperparameter search set
DSGD	$\alpha \in \{10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1\}$
DSGD-GClip	$\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}, \tau \in \{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$
DSGD-CCLIP	$\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}, \tau \in \{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$
DSGD-Clip	$\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}, \tau \in \{10^{-3}, 10^{-2}, 10^{-1}, 1, 10, 10^2\}$
GT-DSGD	$\alpha \in \{10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1\}$
GT-Adam	$\alpha \in \{5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, G \in \{10^{-3}, 10^{-2}, 10^{-1}, 1, 10\}, \epsilon = 10^{-8}$
QG-DSGDm	$\eta \in \{5 * 10^{-5}, 10^{-4}, 5 * 10^{-4}, 10^{-3}, 5 * 10^{-3}, 10^{-2}, 5 * 10^{-2}, 10^{-1}, 0.5, 1, 5, 10\}, \beta = \mu \in \{0.01, 0.2, 0.4, 0.6, 0.8, 0.99\}$
GT-NSGDm	$\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10\}, \beta \in \{0.01, 0.2, 0.4, 0.6, 0.8, 0.99\}$
SCLIP-EF-Network	$\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2\}, \beta \in \{0.01, 0.4, 0.8, 0.99\}, c_\varphi \in \{0.1, 1, 10, 10^2\}, \tau \in \{0.01, 0.1, 1, 10\}$