
Supplementary Material for Paper: Moniqua: Modulo Quantized Communication in Decentralized SGD

A. Overview

This supplementary material contains proof to all the theoretical results. It is organized as follows: In Section B, we analyze how to work with Modulo and quantization, as proofs to Lemma 1 and Lemma 2 in the paper. In Section C, we provably explain why using shared randomness in communication with stochastic rounding can improve performance. In Section D, we illustrate why directly quantizing communication in D-PSGD fails to converge asymptotically, as a proof to Theorem 1. In Section E, we introduce some useful tools of modeling communication as a Markov Chain for the rest of the proof (part of the intuition is illustrated in the paper). We recommend to go through this before getting into Section F to H. Finally we will provide proof to Theorem 2 to 5 from Section F to H.

B. Modulo Operation with Quantization

Proof to Lemma 1.

Proof. Rewrite x and y as

$$\begin{aligned} x &= N_x a + r_x, -\frac{a}{2} \leq r_x < \frac{a}{2} \\ y &= N_y a + r_y, -\frac{a}{2} \leq r_y < \frac{a}{2} \end{aligned}$$

where $N_x, N_y \in \mathbb{Z}$ then,

$$\begin{aligned} \text{LHS} &= (r_x - r_y) \bmod a \\ \text{RHS} &= ((N_x - N_y)a + r_x - r_y) \bmod a = (r_x - r_y) \bmod a = \text{LHS} \end{aligned}$$

Thus we complete the proof. □

Proof to Lemma 2.

Proof. We start from

$$B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) + x = B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) + x - y + y$$

If B_θ is sufficiently large such that $B_\theta \geq 2\theta + 2\delta B_\theta > 2|x - y| + 2\delta B_\theta$, we could put a "mod B_θ " to the first four terms as follows:

$$\begin{aligned} & B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) + x - y + y \\ &= \left(B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) + x - y \right) \bmod B_\theta + y \\ &\stackrel{\text{Lemma 1}}{=} \left[\left(B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) + x \right) \bmod B_\theta - y \bmod B_\theta \right] \bmod B_\theta + y \\ &\stackrel{\text{Lemma 1}}{=} \left\{ \left[B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) \bmod B_\theta - \left(B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) - x \right) \bmod B_\theta \right] \bmod B_\theta - y \bmod B_\theta \right\} \bmod B_\theta + y \end{aligned}$$

Note that the term $\left(B_\theta \left(\frac{x}{B_\theta} \bmod 1\right) - x\right) \bmod B_\theta = 0$, then we can proceed as:

$$\begin{aligned} & \left\{ \left[B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) \bmod B_\theta - \left(B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) - x \right) \bmod B_\theta \right] \bmod B_\theta - y \bmod B_\theta \right\} \bmod B_\theta + y \\ &= \left(B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) \bmod B_\theta - y \bmod B_\theta \right) \bmod B_\theta + y \\ &= \left(B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - y \right) \bmod B_\theta + y \end{aligned}$$

By moving x to the right side we obtain

$$\left| \left(B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - y \right) \bmod B_\theta + y - x \right| = \left| B_\theta \mathcal{Q}_\delta \left(\frac{x}{B_\theta} \bmod 1 \right) - B_\theta \left(\frac{x}{B_\theta} \bmod 1 \right) \right| \leq \delta B_\theta$$

That completes the proof. \square

C. Shared Randomness

In this section, we provide a theoretical explanation why using shared randomness in the stochastic rounding is able to improve the performance. Without the loss of generality, in the following analysis, we let the quantization step associated with stochastic rounding quantizer \mathcal{Q}_δ be $\delta = 1$. For any $z \in \mathbb{R}$ quantized using \mathcal{Q}_δ , let $z_f = z - \lfloor z \rfloor$, the variance of quantization error can be expressed as

$$\mathbb{E} |\mathcal{Q}_\delta(z) - z|^2 = (1 - z_f)(-z_f)^2 + z_f(1 - z_f)^2 = z_f(1 - z_f) \quad (1)$$

Note that in Moniqua, the term associate with quantization error is

$$\mathbb{E} \left\| (\mathbf{q}_{k,j} - \mathbf{x}_{k,j}) - (\mathbf{q}_{k,i} - \mathbf{x}_{k,i}) \right\|^2$$

We now show for $\forall x, y \in \mathbb{R}$

$$\mathbb{E} |(\mathcal{Q}_\delta(x) - x) - (\mathcal{Q}_\delta(y) - y)|^2 = \mathbb{E} |\mathcal{Q}_\delta(y - x) - (y - x)|^2$$

With out the loss of generality, let $x - \lfloor x \rfloor \leq y - \lfloor y \rfloor$. Let $x_f = x - \lfloor x \rfloor$ and $y_f = y - \lfloor y \rfloor$, then

$$\begin{aligned} \lfloor x + u \rfloor &= \lfloor x \rfloor \quad \text{and} \quad \lfloor y + u \rfloor = \lfloor y \rfloor, \text{ with probability } \lceil y \rceil - y \\ \lfloor x + u \rfloor &= \lceil x \rceil \quad \text{and} \quad \lfloor y + u \rfloor = \lceil y \rceil, \text{ with probability } x - \lfloor x \rfloor \\ \lfloor x + u \rfloor &= \lfloor x \rfloor \quad \text{and} \quad \lfloor y + u \rfloor = \lceil y \rceil, \text{ with probability } (\lceil x \rceil - x) - (\lceil y \rceil - y) \end{aligned}$$

Then we have

$$\begin{aligned} & \mathbb{E} |(\mathcal{Q}_\delta(x) - x) - (\mathcal{Q}_\delta(y) - y)|^2 \\ &= \mathbb{E} \left| \left(\delta \left\lfloor \frac{x}{\delta} + u \right\rfloor - x \right) - \left(\delta \left\lfloor \frac{y}{\delta} + u \right\rfloor - y \right) \right|^2 \\ &= (\lceil y \rceil - y)((\lfloor x \rfloor - x) - (\lfloor y \rfloor - y))^2 + (x - \lfloor x \rfloor)((\lceil x \rceil - x) - (\lceil y \rceil - y))^2 \\ & \quad + ((\lceil x \rceil - x) - (\lceil y \rceil - y))((\lfloor x \rfloor - x) - (\lceil y \rceil - y))^2 \\ &= (1 - y_f)(x_f - y_f)^2 + (x_f)(x_f - y_f) + (y_f - x_f)(y_f - x_f - 1)^2 \\ &= (1 - y_f + x_f)(y_f - x_f)^2 + (y_f - x_f)(y_f - x_f - 1)^2 \\ &= (1 - y_f + x_f)(y_f - x_f) \\ &= \mathbb{E} |\mathcal{Q}_\delta(y - x) - (y - x)|^2 \end{aligned}$$

The last equality holds due to equation 1. Next, for $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ let

$$\mathbf{\Delta} = \mathbf{y} - \mathbf{x}$$

$$\mathbf{r} = \mathcal{Q}_\delta(\mathbf{\Delta}) - \mathbf{\Delta}$$

And let r_h denote h -th entry of \mathbf{r} , let $\mathbf{\Delta}_h$ denote h -th entry of $\mathbf{\Delta}$. We obtain

$$\begin{aligned} r_h &= \mathcal{Q}_\delta(\mathbf{\Delta}_h) - \mathbf{\Delta}_h \\ &= \delta \begin{cases} -\frac{\mathbf{\Delta}_h}{\delta} + \lfloor \frac{\mathbf{\Delta}_h}{\delta} \rfloor + 1, & p_t \leq \frac{\mathbf{\Delta}_h}{\delta} - \lfloor \frac{\mathbf{\Delta}_h}{\delta} \rfloor \\ -\frac{\mathbf{\Delta}_h}{\delta} + \lfloor \frac{\mathbf{\Delta}_h}{\delta} \rfloor, & \text{otherwise} \end{cases} \\ &= \delta \begin{cases} -q + 1, & p_t \leq q \\ -q, & \text{otherwise} \end{cases} \end{aligned}$$

where

$$q = \frac{\mathbf{\Delta}_h}{\delta} - \left\lfloor \frac{\mathbf{\Delta}_h}{\delta} \right\rfloor, q \in [0, 1]$$

Based on that, we have

$$\begin{aligned} \mathbb{E}[r_h^2] &\leq \delta^2((-q+1)^2q + (-q)^2(1-q)) \\ &= \delta^2q(1-q) \\ &\leq \delta^2 \min\{q, 1-q\} \end{aligned}$$

Since $\min\{q, 1-q\} \leq \left| \frac{\mathbf{\Delta}_h}{\delta} \right|$, we have

$$\mathbb{E}[r_h^2] \leq \delta^2 \left| \frac{\mathbf{\Delta}_h}{\delta} \right| \leq \delta |\mathbf{\Delta}_h|$$

Summing over the index h yields,

$$\mathbb{E}\|\mathbf{r}\|_2^2 \leq \delta \mathbb{E}\|\mathbf{\Delta}\|_1 \leq \sqrt{d}\delta \mathbb{E}\|\mathbf{\Delta}\|_2$$

Pushing back \mathbf{x} and \mathbf{r} , we have

$$\mathbb{E}\|\mathcal{Q}_\delta(\mathbf{y} - \mathbf{x}) - (\mathbf{y} - \mathbf{x})\|^2 \leq \sqrt{d}\delta \mathbb{E}\|\mathbf{y} - \mathbf{x}\| = \sqrt{d}\delta \mathbb{E}\|\mathbf{x} - \mathbf{y}\|$$

Putting it back we have

$$\mathbb{E}\|(\mathcal{Q}_\delta(\mathbf{x}) - \mathbf{x}) - (\mathcal{Q}_\delta(\mathbf{y}) - \mathbf{y})\|^2 \leq \sqrt{d}\delta \mathbb{E}\|\mathbf{x} - \mathbf{y}\|$$

Now we can see that the error term is bounded by the distance of two quantized tensor, which, in decentralized training, refers to the distance between two models on adjacent workers. In such a way, the error bound can be reduced since the workers are getting close to each other.

D. Why Naive Quantization Fails in D-PSGD (Proof to Theorem 1)

The update rule of naive quantization on D-PSGD is

$$\mathbf{x}_{k+1,i} = \mathbf{x}_{k,i} \mathbf{W}_{ii} + \sum_{j=1, j \neq i}^n \mathcal{Q}_\delta(\mathbf{x}_{k,j}) \mathbf{W}_{ji} - \alpha_k \tilde{\mathbf{g}}_{k,i} = \mathbf{x}_{k,i} + \sum_{j=1, j \neq i}^n (\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,i}) \mathbf{W}_{ji} - \alpha_k \tilde{\mathbf{g}}_{k,i}$$

where α_k is allowed to vary with any policy. Let

$$\begin{aligned} \mathbf{X}_k &= [\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,n}] \in \mathbb{R}^{d \times n} \\ \mathbf{\Omega}_k &= \left[\sum_{j \neq 1} \mathbf{W}_{j1} (\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,1}), \dots, \sum_{j \neq n} \mathbf{W}_{jn} (\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,n}) \right] \in \mathbb{R}^{d \times n} \\ \tilde{\mathbf{G}}_k &= [\tilde{\mathbf{g}}_{k,1}, \dots, \tilde{\mathbf{g}}_{k,n}] \in \mathbb{R}^{d \times n} \end{aligned}$$

by rewriting the update rule, we obtain

$$\mathbf{X}_{k+1} = \mathbf{X}_k + \mathbf{\Omega}_k - \alpha_k \tilde{\mathbf{G}}_k$$

Let $\mathbf{Y}_k = \mathbf{X}_k - \mathbf{x}^* \mathbf{1}^\top$, and considering the fact that $\nabla f(\mathbf{x}) = \mathbf{x} - \delta \mathbf{1}/2 = \mathbf{x} - \mathbf{x}^*$, we can rewrite the update rule as

$$\mathbf{Y}_{k+1} \mathbf{e}_i = \mathbf{Y}_k \mathbf{e}_i + \mathbf{\Omega}_k \mathbf{e}_i - \alpha_k \mathbf{Y}_k \mathbf{e}_i + \alpha_k \left(\tilde{\mathbf{G}}_k - \mathbf{G}_k \right) \mathbf{e}_i$$

where $\left(\tilde{\mathbf{G}}_k - \mathbf{G}_k \right)$ denotes variance in the gradient sampling.

Suppose that by using the update rule of naive quantization, worker i converges to \mathbf{x}^* . Then there must exist a K such that $\forall k \geq K$,

$$\mathbb{E} \|\mathbf{Y}_{k+1} \mathbf{e}_i\|^2 \leq \mathbb{E} \|\mathbf{Y}_k \mathbf{e}_i\|^2 < \frac{\phi^2 \delta^2}{8(1+\phi^2)} \quad (2)$$

Next we show that this assumption lets us derive a contradiction. Firstly, considering the property of linear quantizer,

$$\frac{\delta^2}{4} \leq \mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,i}) - \mathbf{x}^*\|^2 \leq 2\mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,i}) - \mathbf{x}_{k,i}\|^2 + 2\mathbb{E} \|\mathbf{x}_{k,i} - \mathbf{x}^*\|^2$$

As a result

$$\mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,i}) - \mathbf{x}_{k,i}\|^2 \geq \frac{\delta^2}{8} - \frac{\phi^2 \delta^2}{8(1+\phi^2)} = \frac{\delta^2}{8(1+\phi^2)}$$

Since \mathcal{Q}_δ is unbiased, that means $\mathbb{E}[\mathcal{Q}_\delta(\mathbf{x}) - \mathbf{x}] = 0$, then we have

$$\begin{aligned} & \mathbb{E} \|\mathbf{\Omega}_k \mathbf{e}_i\|^2 \\ &= \mathbb{E} \left\| \sum_{j \neq i} \mathbf{W}_{ji} (\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,i}) \right\|^2 \\ &= \sum_{j \in \mathcal{N}_i} \mathbf{W}_{ji}^2 \mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,i}\|^2 + \sum_{m \neq n \neq i} \mathbb{E} \langle (\mathcal{Q}_\delta(\mathbf{x}_{k,m}) - \mathbf{x}_{k,i}) \mathbf{W}_{mi}, (\mathcal{Q}_\delta(\mathbf{x}_{k,n}) - \mathbf{x}_{k,i}) \mathbf{W}_{ni} \rangle \\ &\geq \phi^2 \sum_{j \in \mathcal{N}_i} \mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,i}\|^2 + \sum_{m \neq n \neq i} \mathbb{E} \langle (\mathcal{Q}_\delta(\mathbf{x}_{k,m}) - \mathbf{x}_{k,i}) \mathbf{W}_{mi}, (\mathcal{Q}_\delta(\mathbf{x}_{k,n}) - \mathbf{x}_{k,i}) \mathbf{W}_{ni} \rangle \\ &\stackrel{(*)}{=} \phi^2 \sum_{j \in \mathcal{N}_i} \mathbb{E} \|\mathcal{Q}_\delta(\mathbf{x}_{k,j}) - \mathbf{x}_{k,i}\|^2 \\ &\geq \frac{\phi^2 \delta^2}{8(1+\phi^2)} \end{aligned}$$

where step (*) holds due to unbiased quantizer. Putting it back to the update rule, we obtain

$$\begin{aligned} & \mathbb{E} \|\mathbf{Y}_{k+1} \mathbf{e}_i\|^2 \\ &= \mathbb{E} \left\| \left(\mathbf{Y}_k + \mathbf{\Omega}_k - \alpha_k \mathbf{Y}_k + \alpha_k \left(\tilde{\mathbf{G}}_k - \mathbf{G}_k \right) \right) \mathbf{e}_i \right\|^2 \\ &\stackrel{(*)}{=} \mathbb{E} \|(1 - \alpha_k) \mathbf{Y}_k \mathbf{e}_i\|^2 + \mathbb{E} \|\mathbf{\Omega}_k \mathbf{e}_i\|^2 + \mathbb{E} \left\| \alpha_k \left(\tilde{\mathbf{G}}_k - \mathbf{G}_k \right) \mathbf{e}_i \right\|^2 \\ &\geq \mathbb{E} \|\mathbf{\Omega}_k \mathbf{e}_i\|^2 \\ &\geq \frac{\phi^2 \delta^2}{8(1+\phi^2)} \end{aligned}$$

where cross terms in the (*) step are all 0 due to the unbiased quantizer and unbiased sampling of the gradient. Her we obtain the contradictory that $\frac{\phi^2 \delta^2}{8(1+\phi^2)} \leq \mathbb{E} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 < \frac{\phi^2 \delta^2}{8(1+\phi^2)}$. That being said, for $\forall k, i$

$$\mathbb{E} \|\mathbf{x}_{k,i} - \mathbf{x}^*\|^2 = \mathbb{E} \|\nabla f(\mathbf{x}_{k,i})\|^2 \geq \frac{\phi^2 \delta^2}{8(1+\phi^2)}$$

Thus we complete the proof.

E. A Markov Chain Analysis on the Communication

To better understand how the parallel workers reach consensus over a communication matrix, in this section we use theory from the analysis of Markov Chains to obtain some useful lemmas for proof of Moniqua on D-PSGD and AD-PSGD.

Since the communication matrix \mathbf{W} is doubly stochastic (each row and column sum to 1), it has the same structure as the transition matrix of a Markov Chain with $\frac{1}{n}$ as its the stationary distribution ($\mathbf{W}\frac{\mathbf{1}}{n} = \frac{\mathbf{1}}{n}$). Now let t_{mix} and $d(t)$ denote the mixing time and maximal distance between initial state and stationary distribution as defined in Markov Chain theory.¹

E.1. D-PSGD

In D-PSGD, the communication matrix is fixed during the training. That makes it perfectly aligned with the structure of a Markov Chain. As a result, we obtain the following lemma:

Lemma E.1.

$$\left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 2 \cdot 2^{-\lfloor \frac{t}{t_{\text{mix}}} \rfloor}$$

Proof. For $\forall \mathbf{x} \in \mathbb{R}^d$, let $\mathbf{u} \in \mathbb{R}^d$ be such a vector that every entry of \mathbf{u} is the positive entry of \mathbf{x} and 0 otherwise. Let $\mathbf{v} \in \mathbb{R}^d$ be such a vector that every entry of \mathbf{v} is the absolute value of negative entry of \mathbf{x} and 0 otherwise. The setting above means $\mathbf{x} = \mathbf{u} - \mathbf{v}$. For example,

$$\begin{aligned} \mathbf{x} &= [2, -1]^\top \\ \mathbf{u} &= [2, 0]^\top \\ \mathbf{v} &= [0, 1]^\top \end{aligned}$$

And we have

$$\begin{aligned} & \left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{x} \right\|_1 \\ &= \left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) (\mathbf{u} - \mathbf{v}) \right\|_1 \\ &\leq \left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{u} \right\|_1 + \left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{v} \right\|_1 \\ &= \mathbf{1}^\top \mathbf{u} \left\| \mathbf{W}^t \frac{\mathbf{u}}{\mathbf{1}^\top \mathbf{u}} - \frac{\mathbf{1}}{n} \right\|_1 + \mathbf{1}^\top \mathbf{v} \left\| \mathbf{W}^t \frac{\mathbf{v}}{\mathbf{1}^\top \mathbf{v}} - \frac{\mathbf{1}}{n} \right\|_1 \\ &\leq 2(\mathbf{1}^\top \mathbf{u} + \mathbf{1}^\top \mathbf{v})d(t) \\ &\leq 2d(t) \|\mathbf{x}\|_1 \end{aligned}$$

Considering the definition of L1-norm, we have

$$\left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 = \max_{\|\mathbf{x}\|_1} \frac{\left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{x} \right\|_1}{\|\mathbf{x}\|_1} \leq 2d(t)$$

According to a well-known results on the theory of Markov Chains,² $d(lt_{\text{mix}}) \leq 2^{-l}$ holds for any non-negative integer l , so we have

$$\left\| \mathbf{W}^t \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 2d(t) \leq 2d \left(\frac{t}{t_{\text{mix}}} \cdot t_{\text{mix}} \right) \leq 2d \left(\left\lfloor \frac{t}{t_{\text{mix}}} \right\rfloor t_{\text{mix}} \right) \leq 2 \cdot 2^{-\lfloor \frac{t}{t_{\text{mix}}} \rfloor}$$

That completes the proof. \square

¹Here we are using notation from Chapter 4.5 of *Markov Chains and Mixing Times* (Levin 2009), available at <https://pages.uoregon.edu/dlevin/MARKOV/markovmixing.pdf>

²Again, see *Markov Chains and Mixing Times* for more details.

Additionally, based on standard results in the theory of reversible Markov Chains, we also have³

$$t_{\text{mix}} \leq \log \left(\frac{1}{\frac{1}{4} \cdot \frac{1}{n}} \right) \frac{1}{1 - \rho} \leq \frac{\log(4n)}{1 - \rho}.$$

E.2. AD-PSGD

Note that unlike D-PSGD, here \mathbf{W}_k can be different at each update step and usually each individually have spectral radius $\rho = 1$, so we can't expect to get a bound in terms of a bound on the spectral gap as we did in Theorems 2 and 3. Instead, we require the following condition, which is inspired by the literature on Markov chain Monte Carlo methods: for some constant t_{mix} (here t_{mix} is the same as t_{mix} in the paper) and for any k and any non-negative vector $\boldsymbol{\mu} \in \mathbb{R}^d$ such that $\mathbf{1}^\top \boldsymbol{\mu} = 1$, it must hold that

$$\left\| \left(\prod_{i=1}^{t_{\text{mix}}} \mathbf{W}_{k+i} \right) \boldsymbol{\mu} - \frac{\mathbf{1}}{n} \right\|_1 \leq \frac{1}{2}.$$

We call this constant t_{mix} because it is effectively the *mixing time* of the time-inhomogeneous Markov chain with transition probability matrix \mathbf{W}_k at time k . Note that this condition is more general than those used in previous work on AD-PSGD because it does not require that the \mathbf{W}_k are sampled independently or in an unbiased manner. Based on the above analysis, we can prove the following lemma, which is analogous to the lemma used in the synchronous case.

Lemma E.2. *For any $k \geq 0$ and for any $b \geq a \geq 0$, there exists t_{mix} such that*

$$\left\| \prod_{q=a}^b \mathbf{W}_q \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 2 \cdot 2^{-\lfloor \frac{b-a+1}{t_{\text{mix}}} \rfloor}$$

Proof. Note that for any $\mathbf{x} \in \mathbb{R}^d$, and let \mathbf{u} and \mathbf{v} be two vectors having same definition as in Lemma E.1 with respect to \mathbf{x} , then we have for any k

$$\begin{aligned} & \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{x} \right\|_1 \\ &= \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) (\mathbf{u} - \mathbf{v}) \right\|_1 \\ &\leq \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{u} \right\|_1 + \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{v} \right\|_1 \\ &= \mathbf{1}^\top \mathbf{u} \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \frac{\mathbf{u}}{\mathbf{1}^\top \mathbf{u}} - \frac{\mathbf{1}}{n} \right\|_1 + \mathbf{1}^\top \mathbf{v} \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \frac{\mathbf{v}}{\mathbf{1}^\top \mathbf{v}} - \frac{\mathbf{1}}{n} \right\|_1 \\ &\leq \frac{1}{2} (\mathbf{1}^\top \mathbf{u} + \mathbf{1}^\top \mathbf{v}) \\ &\leq \frac{1}{2} \|\mathbf{x}\|_1 \end{aligned}$$

Considering the definition of the induced ℓ_1 operator norm, we have

$$\left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 = \max_{\mathbf{x}} \frac{\left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{q+k} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \mathbf{x} \right\|_1}{\|\mathbf{x}\|_1} \leq \frac{1}{2}$$

As a result, from the submultiplicativity of the matrix induced norm, we obtain

$$\left\| \prod_{q=a}^b \mathbf{W}_q \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1$$

³Detailed analysis and proofs of this result can be found in chapter 12.2 of *Markov Chains and Mixing Times*.

$$\begin{aligned} &\leq \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{a-1+q} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \cdots \left\| \prod_{q=1}^{t_{\text{mix}}} \mathbf{W}_{\dots+q} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \cdot \left\| \prod_{q=1}^{t_r} \mathbf{W}_{\dots+q} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \\ &\leq 2^{-\lfloor \frac{b-a+1}{t_{\text{mix}}} \rfloor} \left\| \prod_{q=1}^{t_r} \mathbf{W}_{\dots+q} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \end{aligned}$$

where $t_r = (b - a + 1) \bmod t_{\text{mix}}$. Note that

$$\left\| \prod_{q=1}^{t_r} \mathbf{W}_q \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 1 - \frac{1}{n} + (n-1)\frac{1}{n} = 2 - \frac{2}{n} \leq 2$$

Putting it back we obtain

$$\left\| \prod_{q=a}^b \mathbf{W}_{\dots+q} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 2 \cdot 2^{-\lfloor \frac{b-a+1}{t_{\text{mix}}} \rfloor}$$

That completes the proof. \square

Note that in the analysis of Moniqua on AD-PSGD (Section H), we will use this lemma as an assumption.

F. Moniqua on D-PSGD (Proof to Theorem 2 and 3)

F.1. Notations

For convenience, we adopt the following notation

$$\begin{aligned} \mathbf{X}_k &= [\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,n}], & \hat{\mathbf{X}}_k &= [\hat{\mathbf{x}}_{k,1}, \dots, \hat{\mathbf{x}}_{k,n}] \\ \tilde{\mathbf{G}}_k &= [\tilde{\mathbf{g}}_{k,1}, \dots, \tilde{\mathbf{g}}_{k,n}], & \mathbf{G}_k &= [\mathbf{g}_{k,1}, \dots, \mathbf{g}_{k,n}] \\ \bar{\mathbf{X}} &= \mathbf{X} \frac{\mathbf{1}}{n}, \forall \mathbf{X} \in \mathbb{R}^{d \times n}, & \boldsymbol{\Omega}_k &= (\hat{\mathbf{X}}_k - \mathbf{X}_k)(\mathbf{W} - \mathbf{I}) \end{aligned}$$

where $\mathbf{g}_{k,i}$ denotes gradient computed via the whole dataset \mathcal{D}_i and $\mathbf{x}_{k,i}$

From a local view, the update rule on worker i at iteration k can be written as

$$\mathbf{x}_{k+1,i} \leftarrow \mathbf{x}_{k,i} + \sum_{j \in \mathcal{N}_i} (\hat{\mathbf{x}}_{k,j} - \hat{\mathbf{x}}_{k,i}) \mathbf{W}_{ji} - \alpha_k \tilde{\mathbf{g}}_{k,i}$$

which is equivalent to

$$\mathbf{x}_{k+1,i} = \sum_{j=1}^n \mathbf{x}_{k,j} \mathbf{W}_{ji} - \alpha_k \tilde{\mathbf{g}}_{k,i} + \sum_{j=1}^n ((\hat{\mathbf{x}}_{k,j} - \mathbf{x}_{k,j}) - (\hat{\mathbf{x}}_{k,i} - \mathbf{x}_{k,i})) \mathbf{W}_{ji} \quad (3)$$

with a more compact notation, this can be expressed as:

$$\mathbf{X}_{k+1} = \mathbf{X}_k + \hat{\mathbf{X}}_k(\mathbf{W} - \mathbf{I}) - \alpha_k \tilde{\mathbf{G}}_k = \mathbf{X}_k \mathbf{W} - \alpha_k \tilde{\mathbf{G}}_k + (\hat{\mathbf{X}}_k - \mathbf{X}_k)(\mathbf{W} - \mathbf{I}) \quad (4)$$

F.2. Proof to Theorem 2.

Proof. From Lemma F.4 we have

$$\begin{aligned} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq 4(\mathbb{E}f(\mathbf{0}) - \mathbb{E}f^*) + \frac{2\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{8\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{24\varsigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \\ &\quad + \frac{8L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\boldsymbol{\Omega}_k\|_F^2 \end{aligned}$$

Note that

$$\sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2 = \sum_{k=0}^{K-1} \alpha_k \sum_{i=1}^n \mathbb{E} \left\| \sum_{j=1}^n ((\hat{\mathbf{x}}_{k,j} - \mathbf{x}_{k,j}) - (\hat{\mathbf{x}}_{k,i} - \mathbf{x}_{k,i})) \mathbf{W}_{ji} \right\|^2 \stackrel{\text{Lemma F.1, F.3}}{\leq} 4 \sum_{k=0}^{K-1} \alpha_k \delta^2 B_{\theta_k}^2 nd$$

By using Lemma F.3 and by assigning $\delta = \frac{1-\eta\rho}{8C_\alpha^2\eta\log(16n)+2(1-\eta\rho)}$, we obtain

$$\sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2 \leq \frac{G_\infty^2 dn}{C_\alpha^2} \sum_{k=0}^{K-1} \alpha_k^3$$

Pushing it back we obtain

$$\begin{aligned} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq 4(\mathbb{E}f(\mathbf{0}) - \mathbb{E}f^*) + \frac{2\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{8\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{24\zeta^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \\ &\quad + \frac{8G_\infty^2 dL^2}{(1-\rho)^2 C_\alpha^2} \sum_{k=0}^{K-1} \alpha_k^3 \end{aligned}$$

That completes the proof. \square

F.3. Proof to Corollary 1.

Proof. When $\alpha_k = \alpha$, $C_\alpha = \eta = 1$, and we have:

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \leq \frac{4(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2\alpha L}{n} \sigma^2 + \frac{8\alpha^2 L^2 (\sigma^2 + 3\zeta^2)}{(1-\rho)^2} + \frac{8\alpha^2 G_\infty^2 dL^2}{(1-\rho)^2}$$

By setting $\alpha = \frac{1}{\zeta^{\frac{2}{3}} K^{\frac{1}{3}} + \sigma \sqrt{\frac{K}{n}} + 2L}$, we have

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq \frac{8(f(\mathbf{0}) - f^*)L}{K} + \frac{4\sigma(f(\mathbf{0}) - f^* + L/2)}{\sqrt{nK}} + \frac{4\zeta^{\frac{2}{3}}(f(\mathbf{0}) - f^*)}{K^{\frac{2}{3}}} \\ &\quad + \frac{8L^2\sigma^2 n}{(1-\rho)^2(\sigma^2 K + 4nL^2)} + \frac{24L^2\zeta^{\frac{2}{3}}}{(1-\rho)^2 K^{\frac{2}{3}}} + \frac{8G_\infty^2 dnL^2}{(1-\rho)^2(\sigma^2 K + 4nL^2)} \\ &\lesssim \frac{1}{K} + \frac{\sigma}{\sqrt{nK}} + \frac{\zeta^{\frac{2}{3}}}{K^{\frac{2}{3}}} + \frac{\sigma^2 n}{\sigma^2 K + n} + \frac{G_\infty^2 dn}{\sigma^2 K + n} \end{aligned}$$

That completes the proof of Corollary 1. \square

F.4. Lemma for Moniqua on D-PSGD

Lemma F.1. *If $\|\mathbf{x}_{t,i} - \mathbf{x}_{t,j}\|_\infty < \theta_t$, $\forall i, j$ holds at iteration t , then*

$$\left\| \sum_{j=1}^n ((\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}) - (\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i})) \mathbf{W}_{ji} \right\|_\infty \leq \frac{4\delta}{1-2\delta} \theta_t$$

Proof. Let $B_{\theta_t} = \frac{2}{1-2\delta} \theta_t$, based on the algorithm, we obtain

$$\begin{aligned} \hat{\mathbf{x}}_{t,j} &= \left(B_{\theta_t} \mathcal{Q}_\delta \left(\frac{\mathbf{x}_{t,j}}{B_{\theta_t}} \bmod 1 \right) - \mathbf{x}_{t,i} \right) \bmod B_{\theta_t} + \mathbf{x}_{t,i} \\ \hat{\mathbf{x}}_{t,i} &\stackrel{\text{Lemma 2}}{=} B_{\theta_t} \mathcal{Q}_\delta \left(\frac{\mathbf{x}_{t,i}}{B_{\theta_t}} \bmod 1 \right) - B_{\theta_t} \left(\frac{\mathbf{x}_{t,i}}{B_{\theta_t}} \bmod 1 \right) + \mathbf{x}_{t,i} \end{aligned}$$

We start from

$$\begin{aligned} \left\| \sum_{j=1}^n ((\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}) - (\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i})) \mathbf{W}_{ji} \right\|_{\infty} &\leq \sum_{j=1}^n \mathbf{W}_{ji} \|(\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}) - (\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i})\|_{\infty} \\ &\leq \sum_{j=1}^n \mathbf{W}_{ji} \|\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}\|_{\infty} + \sum_{j=1}^n \mathbf{W}_{ji} \|\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i}\|_{\infty} \end{aligned}$$

On the first hand, due to Lemma 2 we obtain

$$\|\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}\|_{\infty} \leq \delta B_{\theta_t}$$

on the other hand,

$$\|\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i}\|_{\infty} = \left\| B_{\theta_t} \mathcal{Q}_{\delta} \left(\frac{\mathbf{x}_{t,i}}{B_{\theta_t}} \bmod 1 \right) - B_{\theta_t} \left(\frac{\mathbf{x}_{t,i}}{B_{\theta_t}} \bmod 1 \right) \right\|_{\infty} \leq \delta B_{\theta_t}$$

Putting it back, we obtain

$$\left\| \sum_{j=1}^n ((\hat{\mathbf{x}}_{t,j} - \mathbf{x}_{t,j}) - (\hat{\mathbf{x}}_{t,i} - \mathbf{x}_{t,i})) \mathbf{W}_{ji} \right\|_{\infty} \leq 2\delta B_{\theta_t} = \frac{4\delta}{1-2\delta} \theta_t$$

which completes the proof. \square

Lemma F.2. For any $\mathbf{X}_t \in \mathbb{R}^{d \times n}$, we have

$$\left\| \sum_{t=0}^{k-1} \mathbf{X}_t \left(\frac{\mathbf{1}\mathbf{1}^{\top}}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F^2 \leq \left(\sum_{t=0}^{k-1} \rho^{k-t-1} \|\mathbf{X}_t\|_F \right)^2$$

Proof.

$$\begin{aligned} \left\| \sum_{t=0}^{k-1} \mathbf{X}_t \left(\frac{\mathbf{1}\mathbf{1}^{\top}}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F^2 &= \left(\left\| \sum_{t=0}^{k-1} \mathbf{X}_t \left(\frac{\mathbf{1}\mathbf{1}^{\top}}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F \right)^2 \\ &\leq \left(\sum_{t=0}^{k-1} \left\| \mathbf{X}_t \left(\frac{\mathbf{1}\mathbf{1}^{\top}}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F \right)^2 \\ &\leq \left(\sum_{t=0}^{k-1} \|\mathbf{X}_t\|_F \left\| \frac{\mathbf{1}\mathbf{1}^{\top}}{n} - \mathbf{W}^{k-t-1} \right\| \right)^2 \\ &\leq \left(\sum_{t=0}^{k-1} \rho^{k-t-1} \|\mathbf{X}_t\|_F \right)^2 \end{aligned}$$

That completes the proof. \square

Lemma F.3. In any iteration $k \geq 0$, and for any two worker i and j , when $\delta = \frac{1-\eta\rho}{8C_{\alpha}^2\eta \log(16n)+2(1-\eta\rho)}$ we have:

$$\|\mathbf{X}_k(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} < \frac{2\alpha_k G_{\infty} C_{\alpha} \eta \log(16n)}{1 - \eta\rho} = \theta_k$$

Proof. We use mathematical induction to prove this:

I. When $k = 0$, $\|\mathbf{X}_0(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} = 0 < \theta_0, \forall i, j$

II. Suppose $\|\mathbf{X}_t(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < \theta_t, \forall t \leq k, \forall i, j$, we obtain

$$\begin{aligned}
 \|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty &= \left\| \sum_{t=0}^k (-\alpha_t \mathbf{G}_t + \Omega_t) \mathbf{W}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\
 &\leq \sum_{t=0}^k \|\alpha_t \mathbf{G}_t\|_{1,\infty} \|\mathbf{W}^{k-t}(\mathbf{e}_i - \mathbf{e}_j)\|_1 + \sum_{t=0}^k \|\Omega_t\|_{1,\infty} \|\mathbf{W}^{k-t}(\mathbf{e}_i - \mathbf{e}_j)\|_1 \\
 &\stackrel{\text{Lemma F.1}}{\leq} \sum_{t=0}^k \alpha_t G_\infty \|\mathbf{W}^{k-t}(\mathbf{e}_i - \mathbf{e}_j)\|_1 + \frac{4\delta}{1-2\delta} \sum_{t=0}^k \theta_t \|\mathbf{W}^{k-t}(\mathbf{e}_i - \mathbf{e}_j)\|_1 \\
 &\leq \alpha_{k+1} G_\infty \sum_{t=0}^k \frac{\alpha_{k-t}}{\alpha_{k+1}} \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 + \frac{4\delta\theta_k}{1-2\delta} \sum_{t=0}^k \frac{\theta_t}{\theta_k} \|\mathbf{W}^{k-t}(\mathbf{e}_i - \mathbf{e}_j)\|_1 \\
 &< \alpha_{k+1} G_\infty C_\alpha \eta \sum_{t=0}^\infty \eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 + \frac{4\delta C_\alpha \theta_k}{1-2\delta} \sum_{t=0}^\infty \eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1
 \end{aligned}$$

For any $t \geq 0$, on one hand

$$\|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 \leq \sqrt{n} \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_2 \leq \sqrt{n} \left\| \mathbf{W}^t \mathbf{e}_i - \frac{\mathbf{1}}{n} \right\| + \sqrt{n} \left\| \mathbf{W}^t \mathbf{e}_j - \frac{\mathbf{1}}{n} \right\| \leq 2\sqrt{n}\rho^t$$

where the last step holds due to the diagonalizability of \mathbf{W} . On the other hand,

$$\|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 \leq \mathbf{1}^\top \mathbf{W}^t \mathbf{e}_i + \mathbf{1}^\top \mathbf{W}^t \mathbf{e}_j = \mathbf{1}^\top \mathbf{e}_i + \mathbf{1}^\top \mathbf{e}_j = 2$$

As a result

$$\eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 \leq \min\{2\sqrt{n}(\eta\rho)^t, 2\}$$

Let $T_0 = \left\lceil \frac{-\log(\sqrt{n})}{\log(\eta\rho)} \right\rceil$, so that $\sqrt{n}(\eta\rho)^{T_0} \leq 1$, then we have

$$\begin{aligned}
 \sum_{t=0}^\infty \eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 &= \sum_{t=0}^{T_0-1} \eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 + \sum_{t=T_0}^\infty \eta^t \|\mathbf{W}^t(\mathbf{e}_i - \mathbf{e}_j)\|_1 \\
 &\leq \sum_{t=0}^{T_0-1} 2 + \sum_{t=0}^\infty 2\sqrt{n}(\eta\rho)^{t+T_0} \\
 &\leq 2 \left\lceil \frac{-\log(\sqrt{n})}{\log(\eta\rho)} \right\rceil + \sum_{t=0}^\infty 2(\sqrt{n}(\eta\rho)^{T_0}) (\eta\rho)^t \\
 &\leq \frac{2\log(\sqrt{n})}{1-\eta\rho} + 2 + \frac{2}{1-\eta\rho} \\
 &\leq \frac{\log(16n)}{1-\eta\rho}
 \end{aligned}$$

As a result, we have

$$\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < \frac{\alpha_{k+1} G_\infty C_\alpha \eta \log(16n)}{1-\eta\rho} + \frac{4\delta C_\alpha}{1-2\delta} \cdot \frac{\log(16n)}{1-\eta\rho} \theta_k$$

with $\delta = \frac{1-\eta\rho}{8C_\alpha^2 \eta \log(16n) + 2(1-\eta\rho)}$,

$$\begin{aligned}
 \|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty &< \frac{\alpha_{k+1} G_\infty C_\alpha \eta \log(16n)}{1-\eta\rho} + \frac{4\delta C_\alpha}{1-2\delta} \cdot \frac{\log(16n)}{1-\eta\rho} \cdot \frac{2\alpha_k G_\infty C_\alpha \eta \log(16n)}{1-\eta\rho} \\
 &\leq \frac{\alpha_{k+1} G_\infty C_\alpha \eta \log(16n)}{1-\eta\rho} + \frac{4\delta C_\alpha}{1-2\delta} \cdot \frac{\log(16n)}{1-\eta\rho} \cdot \frac{2\alpha_{k+1} C_\alpha \eta G_\infty C_\alpha \eta \log(16n)}{1-\eta\rho}
 \end{aligned}$$

$$\leq \frac{2\alpha_{k+1}G_\infty C_\alpha \eta \log(16n)}{1 - \eta\rho} = \theta_{k+1}$$

Combining I and II, we complete the proof. \square

Lemma F.4. *The running average of the gradient norm has the following bound:*

$$\begin{aligned} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq 4(\mathbb{E}f(\mathbf{0}) - \mathbb{E}f^*) + \frac{2\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{8\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{24\varsigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \\ &\quad + \frac{8L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2 \end{aligned}$$

Proof. Let $\mathbf{1}$ denote a n -dimensional vector with all the entries be 1. And we have

$$\bar{\mathbf{X}}_{k+1} = (\mathbf{X}_k \mathbf{W} - \alpha_k \tilde{\mathbf{G}}_k + \Omega_k) \frac{\mathbf{1}}{n} = \bar{\mathbf{X}}_k - \alpha_k \bar{\tilde{\mathbf{G}}}_k + (\hat{\mathbf{X}}_k - \mathbf{X}_k)(\mathbf{W} - \mathbf{I}) \frac{\mathbf{1}}{n} = \bar{\mathbf{X}}_k - \alpha_k \bar{\tilde{\mathbf{G}}}_k$$

And by Taylor Expansion, we have

$$\begin{aligned} \mathbb{E}f(\bar{\mathbf{X}}_{k+1}) &= \mathbb{E}f\left(\frac{(\mathbf{X}_k \mathbf{W} - \alpha_k \tilde{\mathbf{G}}_k + \Omega_k)\mathbf{1}}{n}\right) \\ &= \mathbb{E}f\left(\bar{\mathbf{X}}_k - \alpha_k \bar{\tilde{\mathbf{G}}}_k\right) \\ &\leq \mathbb{E}f(\bar{\mathbf{X}}_k) - \alpha_k \mathbb{E}\langle \nabla f(\bar{\mathbf{X}}_k), \bar{\tilde{\mathbf{G}}}_k \rangle + \frac{\alpha_k^2 L}{2} \mathbb{E} \|\bar{\tilde{\mathbf{G}}}_k\|^2 \end{aligned}$$

And for the last term, we have

$$\begin{aligned} \mathbb{E} \|\bar{\tilde{\mathbf{G}}}_k\|^2 &= \mathbb{E} \left\| \frac{\sum_{i=1}^n \tilde{\mathbf{g}}_{k,i}}{n} \right\|^2 \\ &= \mathbb{E} \left\| \frac{\sum_{i=1}^n \tilde{\mathbf{g}}_{k,i} - \sum_{i=1}^n \mathbf{g}_{k,i}}{n} + \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 \\ &= \mathbb{E} \left\| \frac{\sum_{i=1}^n \tilde{\mathbf{g}}_{k,i} - \sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 + \mathbb{E} \left\| \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 + \mathbb{E} \left\langle \frac{\sum_{i=1}^n \tilde{\mathbf{g}}_{k,i} - \sum_{i=1}^n \mathbf{g}_{k,i}}{n} + \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\rangle \\ &= \mathbb{E} \left\| \frac{\sum_{i=1}^n \tilde{\mathbf{g}}_{k,i} - \sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 + \mathbb{E} \left\| \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 \\ &\stackrel{\text{Assumption 3}}{=} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \|\tilde{\mathbf{g}}_{k,i} - \mathbf{g}_{k,i}\|^2 + \mathbb{E} \left\| \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 \\ &\leq \frac{\sigma^2}{n} + \mathbb{E} \left\| \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 \end{aligned}$$

Putting it back, we obtain

$$\begin{aligned} \mathbb{E}f(\bar{\mathbf{X}}_{k+1}) &\leq \mathbb{E}f(\bar{\mathbf{X}}_k) - \alpha_k \mathbb{E}\langle \nabla f(\bar{\mathbf{X}}_k), \bar{\tilde{\mathbf{G}}}_k \rangle + \frac{\alpha_k^2 L}{2n} \sigma^2 + \frac{\alpha_k^2 L}{2} \mathbb{E} \left\| \frac{\sum_{i=1}^n \mathbf{g}_{k,i}}{n} \right\|^2 \\ &= \mathbb{E}f(\bar{\mathbf{X}}_k) - \frac{\alpha_k - \alpha_k^2 L}{2} \mathbb{E} \|\bar{\tilde{\mathbf{G}}}_k\|^2 - \frac{\alpha_k}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \frac{\alpha_k^2 L}{2n} \sigma^2 + \frac{\alpha_k}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \bar{\tilde{\mathbf{G}}}_k\|^2 \end{aligned}$$

where the last step comes from $2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 = \|\mathbf{a} - \mathbf{b}\|^2$ And

$$\mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \bar{\tilde{\mathbf{G}}}_k\|^2 \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\| \nabla f_i \left(\frac{\sum_{i'=1}^n \mathbf{x}_{k,i'}}{n} \right) - \nabla f_i(\mathbf{x}_{k,i}) \right\|^2$$

$$\begin{aligned}
 &\stackrel{\text{Assumption 1}}{\leq} \frac{L^2}{n} \sum_{i=1}^n \mathbb{E} \left\| \frac{\sum_{i'=1}^n \mathbf{x}_{k,i'}}{n} - \mathbf{x}_{k,i} \right\|^2 \\
 &= \frac{L^2}{n} \sum_{i=1}^n \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2
 \end{aligned}$$

by Lipschitz assumption, we obtain

$$\frac{\alpha_k - \alpha_k^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{\alpha_k}{2} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \mathbb{E} f(\bar{\mathbf{X}}_{k+1}) + \frac{\alpha_k^2 L}{2n} \sigma^2 + \frac{\alpha_k L^2}{2n} \sum_{i=1}^n \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2$$

summing over from $k = 0$ to $K - 1$ on both sides, we have

$$\begin{aligned}
 \sum_{k=0}^{K-1} (\alpha_k - \alpha_k^2 L) \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq 2(\mathbb{E} f(\bar{\mathbf{X}}_0) - \mathbb{E} f(\bar{\mathbf{X}}_K)) + \frac{\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 \\
 &\quad + \frac{L^2}{n} \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2
 \end{aligned}$$

From Lemma F.5, we have

$$\begin{aligned}
 &\sum_{k=0}^{K-1} (\alpha_k - \alpha_k^2 L) \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\
 &\leq 2(\mathbb{E} f(\bar{\mathbf{X}}_0) - \mathbb{E} f(\bar{\mathbf{X}}_K)) + \frac{\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{L^2}{n} \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 \\
 &\leq 2(\mathbb{E} f(\bar{\mathbf{X}}_0) - \mathbb{E} f(\bar{\mathbf{X}}_K)) + \frac{\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{4\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12\zeta^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\
 &\quad + \frac{4L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2
 \end{aligned}$$

Rearrange the terms, we have

$$\begin{aligned}
 \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &\leq 4(\mathbb{E} f(\mathbf{0}) - \mathbb{E} f^*) + \frac{2\sigma^2 L}{n} \sum_{k=0}^{K-1} \alpha_k^2 + \frac{8\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{24\zeta^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \\
 &\quad + \frac{8L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2
 \end{aligned}$$

and that completes the proof \square

Lemma F.5.

$$\begin{aligned}
 \frac{L^2}{n} \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 &\leq \frac{4\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12\zeta^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\
 &\quad + \frac{4L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|\Omega_k\|_F^2
 \end{aligned}$$

Proof.

$$\sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2$$

$$\begin{aligned}
 &= \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \left\| \mathbf{X}_k \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &= \sum_{k=1}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \left\| \left(\mathbf{X}_{k-1} \mathbf{W} - \alpha \tilde{\mathbf{G}}_{k-1} + \boldsymbol{\Omega}_{k-1} \right) \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &\stackrel{\mathbf{x}_{0,i}=0}{=} \sum_{k=1}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \left\| \sum_{t=0}^{k-1} \left(-\alpha_t \tilde{\mathbf{G}}_t + \boldsymbol{\Omega}_t \right) \left(\frac{\mathbf{1}}{n} - \mathbf{W}^{k-t-1} \mathbf{e}_i \right) \right\|^2 \\
 &\leq 2 \sum_{k=1}^{K-1} \alpha_k \sum_{i=1}^n \mathbb{E} \left\| \sum_{t=0}^{k-1} \alpha_t \tilde{\mathbf{G}}_t \left(\frac{\mathbf{1}}{n} - \mathbf{W}^{k-t-1} \mathbf{e}_i \right) \right\|^2 + 2 \sum_{k=1}^{K-1} \alpha_k \sum_{i=1}^n \mathbb{E} \left\| \sum_{t=0}^{k-1} \boldsymbol{\Omega}_t \left(\frac{\mathbf{1}}{n} - \mathbf{W}^{k-t-1} \mathbf{e}_i \right) \right\|^2 \\
 &= 2 \sum_{k=1}^{K-1} \alpha_k \mathbb{E} \left\| \sum_{t=0}^{k-1} \alpha_t \tilde{\mathbf{G}}_t \left(\frac{\mathbf{1}\mathbf{1}^\top}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F^2 + 2 \sum_{k=1}^{K-1} \mathbb{E} \left\| \sum_{t=0}^{k-1} \boldsymbol{\Omega}_t \left(\frac{\mathbf{1}\mathbf{1}^\top}{n} - \mathbf{W}^{k-t-1} \right) \right\|_F^2 \\
 &\stackrel{\text{Lemma F.2}}{\leq} 2 \sum_{k=1}^{K-1} \alpha_k \left(\sum_{t=0}^{k-1} \rho^{k-t-1} \alpha_t \mathbb{E} \left\| \tilde{\mathbf{G}}_t \right\|_F \right)^2 + 2 \sum_{k=1}^{K-1} \alpha_k \left(\sum_{t=0}^{k-1} \rho^{k-t-1} \mathbb{E} \left\| \boldsymbol{\Omega}_t \right\|_F \right)^2 \\
 &\stackrel{\text{Lemma F.7}}{\leq} \frac{2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \tilde{\mathbf{G}}_k \right\|_F^2 + \frac{2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \left\| \boldsymbol{\Omega}_k \right\|_F^2 \\
 &\stackrel{\text{Lemma F.6}}{\leq} \frac{2}{(1-\rho)^2} \left(n\sigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k^3 \mathbb{E} \left\| \bar{\mathbf{X}}_k - \mathbf{x}_{k,i} \right\|^2 + 3n\varsigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3n \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2 \right) \\
 &\quad + \frac{2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \left\| \boldsymbol{\Omega}_k \right\|_F^2
 \end{aligned}$$

Rearrange the terms, we have

$$\begin{aligned}
 \sum_{k=0}^{K-1} \alpha_k \left(1 - \frac{6\alpha_k^2 L^2}{(1-\rho)^2} \right) \sum_{i=1}^n \mathbb{E} \left\| \bar{\mathbf{X}}_k - \mathbf{x}_{k,i} \right\|^2 &\leq \frac{2n\sigma^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{6n\varsigma^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{6n}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2 \\
 &\quad + \frac{2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \left\| \boldsymbol{\Omega}_k \right\|_F^2
 \end{aligned}$$

Let $1 - \frac{6\alpha_k^2 L^2}{(1-\rho)^2} \geq \frac{1}{2}$, we have

$$\begin{aligned}
 \frac{L^2}{n} \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k \mathbb{E} \left\| \bar{\mathbf{X}}_k - \mathbf{x}_{k,i} \right\|^2 &\leq \frac{4\sigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12\varsigma^2 L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 + \frac{12L^2}{(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2 \\
 &\quad + \frac{4L^2}{n(1-\rho)^2} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \left\| \boldsymbol{\Omega}_k \right\|_F^2
 \end{aligned}$$

That completes the proof. \square

Lemma F.6.

$$\sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \tilde{\mathbf{G}}_k \right\|_F^2 \leq n\sigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k^3 \mathbb{E} \left\| \bar{\mathbf{X}}_k - \mathbf{x}_{k,i} \right\|^2 + 3n\varsigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3n \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2$$

Proof. From the property of Frobenius norm, we have

$$\mathbb{E} \left\| \tilde{\mathbf{G}}_k \right\|_F^2 = \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{g}}_{k,i} \right\|^2$$

Since

$$\begin{aligned}\mathbb{E} \|\tilde{\mathbf{g}}_{k,i}\|^2 &= \mathbb{E} \|\tilde{\mathbf{g}}_{k,i} - \mathbf{g}_{k,i}\|^2 + \mathbb{E} \|\mathbf{g}_{k,i}\|^2 \\ &= \sigma^2 + 3\mathbb{E} \|\nabla f_i(\mathbf{x}_{k,i}) - \nabla f_i(\bar{\mathbf{X}}_k)\|^2 + 3\mathbb{E} \|\nabla f_i(\bar{\mathbf{X}}_k) - \nabla f(\bar{\mathbf{X}}_k)\|^2 + 3\mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ &\leq \sigma^2 + 3L^2\mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 + 3\varsigma^2 + 3\mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2\end{aligned}$$

Summing from $k = 0$ to $K - 1$, we obtain

$$\begin{aligned}& \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \|\tilde{\mathbf{G}}_k\|_F^2 \\ &= \sum_{k=0}^{K-1} \alpha_k^3 \sum_{i=1}^n \mathbb{E} \|\tilde{\mathbf{g}}_{k,i}\|^2 \\ &\leq \sum_{k=0}^{K-1} \alpha_k^3 \sum_{i=1}^n \sigma^2 + 3L^2 \sum_{k=0}^{K-1} \alpha_k^3 \sum_{i=1}^n \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 + 3 \sum_{k=0}^{K-1} \alpha_k^3 \sum_{i=1}^n \varsigma^2 + 3 \sum_{k=0}^{K-1} \alpha_k^3 \sum_{i=1}^n \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ &= n\sigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n \alpha_k^3 \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 + 3n\varsigma^2 \sum_{k=0}^{K-1} \alpha_k^3 + 3n \sum_{k=0}^{K-1} \alpha_k^3 \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2\end{aligned}$$

That completes the proof. \square

Lemma F.7. Given $0 \leq \rho < 1$ and T , a positive integer. Also given non-negative sequences $\{a_t\}_{t=1}^\infty$ and $\{b_t\}_{t=1}^\infty$ with $\{a_t\}_{t=1}^\infty$ being non-increasing, the following inequalities holds:

$$\begin{aligned}\sum_{t=1}^k a_t \left(\sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \right) &\leq \frac{T}{1-\rho} \sum_{s=1}^k a_s b_s \\ \sum_{t=1}^k a_t \left(\sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \right)^2 &\leq \frac{T^2}{(1-\rho)^2} \sum_{s=1}^k a_s b_s^2\end{aligned}$$

Proof. Firstly,

$$S_k = \sum_{t=1}^k a_t \left(\sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \right) = \sum_{s=1}^k \sum_{t=s}^k \alpha_t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \leq \sum_{s=1}^k a_s b_s \sum_{t=0}^{T-1} \sum_{m=0}^{\infty} \rho^m \leq \frac{T}{1-\rho} \sum_{s=1}^k a_s b_s$$

further we have

$$\begin{aligned}& \sum_{t=1}^k a_t \left(\sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \right)^2 = \sum_{t=1}^k a_t \sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s \sum_{r=1}^t \rho^{-\lfloor \frac{t-r}{T} \rfloor} b_r = \sum_{t=1}^k a_t \sum_{s=1}^t \sum_{r=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor + \lfloor \frac{t-r}{T} \rfloor} b_s b_r \\ &\leq \sum_{t=1}^k a_t \sum_{s=1}^t \sum_{r=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor + \lfloor \frac{t-r}{T} \rfloor} \frac{b_s^2 + b_r^2}{2} = \sum_{t=1}^k a_t \sum_{s=1}^t \sum_{r=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor + \lfloor \frac{t-r}{T} \rfloor} b_s^2 \\ &\leq \sum_{t=1}^k a_t \sum_{s=1}^t b_s^2 \rho^{-\lfloor \frac{t-s}{T} \rfloor} \sum_{r=1}^t \rho^{-\lfloor \frac{t-r}{T} \rfloor} \leq \sum_{t=1}^k a_t \sum_{s=1}^t b_s^2 \rho^{-\lfloor \frac{t-s}{T} \rfloor} \sum_{r=0}^{T-1} \sum_{m=0}^{\infty} \rho^m \\ &\leq \frac{T}{1-\rho} \sum_{t=1}^k a_t \sum_{s=1}^t \rho^{-\lfloor \frac{t-s}{T} \rfloor} b_s^2 \stackrel{\text{Using } S_k}{\leq} \frac{T^2}{(1-\rho)^2} \sum_{s=1}^k a_s b_s^2\end{aligned}$$

That completes the proof. \square

F.5. Proof to Theorem 3.

Proof. Let $\bar{\rho}$ denote the spectral gap of matrix $\bar{\mathbf{W}}$, it is straightforward to know that $\bar{\rho} = \gamma\rho + (1 - \gamma)$. we first use mathematical induction to prove at iteration $\forall k \leq K$, for any worker i and j , with probability $(1 - \epsilon)^k$

$$\|\mathbf{X}_k(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < \theta = \frac{2\alpha \log(16n)G_\infty}{\gamma(1 - \rho)}$$

where $\gamma = \frac{2}{1 - \rho + \frac{16\delta^2}{(1 - 2\delta)^2} \cdot \frac{32 \log(4n)}{1 - \rho} \log(\frac{1}{\epsilon})}$.

I. When $k = 0$, $\|\mathbf{X}_0(\mathbf{e}_i - \mathbf{e}_j)\|_\infty = 0 < \theta$

II. Suppose $\|\mathbf{X}_t(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < \theta$ holds for $\forall t \leq k$, then for $k + 1$ we have

$$\begin{aligned} \|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty &= \left\| \left(\mathbf{X}_k \bar{\mathbf{W}} - \alpha \tilde{\mathbf{G}}_k + \gamma \Omega_k \right) (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\ &\stackrel{\mathbf{X}_0=0}{=} \left\| \sum_{t=0}^k \left(-\alpha \tilde{\mathbf{G}}_t + \gamma \Omega_t \right) \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\ &\leq \left\| \sum_{t=0}^k \alpha \tilde{\mathbf{G}}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty + \left\| \sum_{t=0}^k \gamma \Omega_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \end{aligned}$$

We bound these two terms separately. First from Lemma F.3 we know that

$$\sum_{t=0}^{\infty} \left\| \bar{\mathbf{W}}^t (\mathbf{e}_i - \mathbf{e}_j) \right\|_1 < \frac{\log(16n)}{1 - \bar{\rho}} = \frac{\log(16n)}{\gamma(1 - \rho)} \quad (5)$$

then we have for the first term,

$$\begin{aligned} \left\| \sum_{t=0}^k \alpha \tilde{\mathbf{G}}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty &\leq \sum_{t=0}^k \left\| \alpha \tilde{\mathbf{G}}_t \right\|_{1, \infty} \left\| \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_1 \\ &\leq \alpha G_\infty \sum_{t=0}^{\infty} \left\| \bar{\mathbf{W}}^t (\mathbf{e}_i - \mathbf{e}_j) \right\|_1 \\ &< \frac{\alpha \log(16n)G_\infty}{\gamma(1 - \rho)} \end{aligned}$$

Next, we bound the second term. Suppose the infinity norm of the term $\sum_{t=0}^k \gamma \Omega_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j)$ is taken at coordinate h , then we have

$$\begin{aligned} \left\| \sum_{t=0}^k \gamma \Omega_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty &= \gamma \left| \mathbf{e}_h^\top \left(\sum_{t=0}^k \Omega_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right) \right| \\ &= \gamma \left| \sum_{t=0}^k \mathbf{e}_h^\top \left(\Omega_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j) \right) \right| \end{aligned}$$

Let

$$u_t = \sum_{m=0}^t \mathbf{e}_h^\top \left(\Omega_{k-m} \bar{\mathbf{W}}^m (\mathbf{e}_i - \mathbf{e}_j) \right)$$

from the induction hypothesis we know that $\{u_t\}_{t \leq k}$ is a martingale sequence. Note that,

$$\begin{aligned} |u_t - u_{t-1}| &= \left| \mathbf{e}_h^\top \left(\Omega_{k-t} \bar{\mathbf{W}}^t (\mathbf{e}_i - \mathbf{e}_j) \right) \right| \\ &\leq \left\| \Omega_{k-t} \bar{\mathbf{W}}^t (\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \end{aligned}$$

$$\begin{aligned}
 & \text{Equation 5} \\
 & \leq \|\boldsymbol{\Omega}_{k-t}\|_{1,\infty} \min\{2\sqrt{n\bar{\rho}^t}, 2\} \\
 & \leq 2\delta B_\theta \min\{2\sqrt{n\bar{\rho}^t}, 2\}
 \end{aligned}$$

where $B_\theta = \frac{2}{1-2\delta}\theta$, then by using Azuma's inequality we obtain

$$\begin{aligned}
 \mathbb{P}\left[\left|\sum_{t=0}^k \mathbf{e}_h^\top \left(\boldsymbol{\Omega}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j)\right)\right| > a\right] & \leq \exp\left(-\frac{a^2}{8\delta^2 B_\theta^2 \sum_{t=0}^k \min\{2\sqrt{n\bar{\rho}^t}, 2\}^2}\right) \\
 & \leq \exp\left(-\frac{a^2}{32\delta^2 B_\theta^2 \sum_{t=0}^\infty \min\{n\bar{\rho}^{2t}, 1\}}\right)
 \end{aligned}$$

Here we use the induction hypothesis. Similar as before, Let $T_0 = \left\lceil \frac{-\log(n)}{2\log(\bar{\rho})} \right\rceil$, so that $n\bar{\rho}^{2T_0} \leq 1$, then we have

$$\begin{aligned}
 \sum_{t=0}^\infty \min\{n\bar{\rho}^{2t}, 1\} & = \sum_{t=0}^{T_0-1} \min\{n\bar{\rho}^{2t}, 1\} + \sum_{t=T_0}^\infty \min\{n\bar{\rho}^{2t}, 1\} \\
 & < \sum_{t=0}^{T_0-1} 1 + \sum_{t=0}^\infty n\bar{\rho}^{2t+2T_0} \\
 & \leq \left\lceil \frac{-\log(n)}{2\log(\bar{\rho})} \right\rceil + \sum_{t=0}^\infty (n\bar{\rho}^{2T_0}) \bar{\rho}^{2t} \\
 & \leq \frac{\log(n)}{1-\bar{\rho}^2} + 1 + \frac{1}{1-\bar{\rho}^2} \\
 & \leq \frac{\log(4n)}{1-\bar{\rho}^2} \\
 & = \frac{\log(4n)}{\gamma(1-\rho)(2-\gamma(1-\rho))}
 \end{aligned}$$

Putting it back, we obtain

$$\mathbb{P}\left[\left|\sum_{t=0}^k \mathbf{e}_h^\top \left(\boldsymbol{\Omega}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j)\right)\right| > a\right] \leq \exp\left(-\frac{a^2 \gamma(1-\rho)(2-\gamma(1-\rho))}{32\delta^2 B_\theta^2 \log(4n)}\right)$$

In other words, with probability $1 - \epsilon$,

$$\left\|\sum_{t=0}^k \gamma \boldsymbol{\Omega}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j)\right\|_\infty = \gamma \left|\sum_{t=0}^k \mathbf{e}_h^\top \left(\boldsymbol{\Omega}_t \bar{\mathbf{W}}^{k-t} (\mathbf{e}_i - \mathbf{e}_j)\right)\right| \leq \delta B_\theta \sqrt{\frac{32 \log(4n) \gamma}{(1-\rho)(2-\gamma(1-\rho))} \log\left(\frac{1}{\epsilon}\right)}$$

Combine them together, we obtain

$$\begin{aligned}
 \|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty & < \frac{\alpha \log(16n) G_\infty}{\gamma(1-\rho)} + \delta B_\theta \sqrt{\frac{32 \log(4n) \gamma}{(1-\rho)(2-\gamma(1-\rho))} \log\left(\frac{1}{\epsilon}\right)} \\
 & < \frac{\alpha \log(16n) G_\infty}{\gamma(1-\rho)} + \frac{2\delta}{1-2\delta} \theta \sqrt{\frac{32 \log(4n) \gamma}{(1-\rho)(2-\gamma(1-\rho))} \log\left(\frac{1}{\epsilon}\right)}
 \end{aligned}$$

$$\text{Let } \gamma = \frac{2}{1-\rho + \frac{16\delta^2}{(1-2\delta)^2} \cdot \frac{32 \log(4n)}{1-\rho} \log\left(\frac{1}{\epsilon}\right)}$$

$$\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < \frac{\alpha \log(16n) G_\infty}{\gamma(1-\rho)} + \frac{1}{2} \theta \leq \theta$$

Combining I and II, we complete the proof.

We proceed to obtain the convergence rate. From Theorem 2 we have with $\alpha_k = \alpha$

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \leq \frac{4(\mathbb{E}f(\mathbf{0}) - \mathbb{E}f^*)}{\alpha K} + \frac{2\alpha\sigma^2 L}{n} + \frac{8\alpha^2\sigma^2 L^2}{(1-\bar{\rho})^2} + \frac{24\alpha^2\zeta^2 L^2}{(1-\bar{\rho})^2} + \frac{8\alpha L^2}{n(1-\bar{\rho})^2 K} \sum_{k=0}^{K-1} \mathbb{E} \|\gamma \boldsymbol{\Omega}_k\|_F^2$$

Note that with probability $(1-\epsilon)^K$

$$\sum_{k=0}^{K-1} \mathbb{E} \|\gamma \boldsymbol{\Omega}_k\|_F^2 = \gamma^2 \sum_{k=0}^{K-1} \sum_{i=1}^n \mathbb{E} \left\| \sum_{j=1}^n ((\hat{\mathbf{x}}_{k,j} - \mathbf{x}_{k,j}) - (\hat{\mathbf{x}}_{k,i} - \mathbf{x}_{k,i})) \mathbf{W}_{ji} \right\|^2 \stackrel{\text{Lemma F.1}}{\leq} \frac{16\delta^2 \gamma^2}{(1-2\delta)^2} \theta^2 dnK$$

Fit in $\theta = \frac{2\alpha \log(16n)G_\infty}{\gamma(1-\rho)}$, we obtain

$$\sum_{k=0}^{K-1} \mathbb{E} \|\gamma \boldsymbol{\Omega}_k\|_F^2 \leq \frac{64\alpha^2 \delta^2 \log^2(16n)G_\infty^2}{(1-2\delta)^2(1-\rho)^2} dnK$$

Let \mathcal{E} denote the event that the bound θ holds for all $0 \leq t \leq T-1$, then,

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 &= \left[\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \mid \mathcal{E} \right] \mathbb{P}(\mathcal{E}) + \left[\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \mid \neg \mathcal{E} \right] \mathbb{P}(\neg \mathcal{E}) \\ &\leq \frac{4(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2\alpha L}{n} \sigma^2 + \frac{8\alpha^2 L^2 (\sigma^2 + 3\zeta^2)}{(1-\bar{\rho})^2} + \frac{8L^2}{nK(1-\bar{\rho})^2} \sum_{k=1}^{K-1} \mathbb{E} \|\gamma \boldsymbol{\Omega}_k\|_F^2 \\ &\quad + G_\infty^2 d (1 - (1-\epsilon)^K) \\ &\leq \frac{4(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2\alpha L}{n} \sigma^2 + \frac{8\alpha^2 L^2 (\sigma^2 + 3\zeta^2)}{\gamma^2(1-\rho)^2} + \frac{512\alpha^2 \delta^2 L^2 \log^2(16n)G_\infty^2 d}{\gamma^2(1-\rho)^4(1-2\delta)^2} \\ &\quad + G_\infty^2 d (1 - (1-\epsilon)^K) \end{aligned}$$

Assign $\epsilon = \frac{1}{K^2}$ and set $\alpha = \frac{1}{\varsigma^{\frac{2}{3}} K^{\frac{1}{3}} + \sigma \sqrt{\frac{K}{n}} + 2L}$, we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \lesssim \frac{\sigma}{\sqrt{nK}} + \frac{1}{K} + \frac{\varsigma^{\frac{2}{3}} \delta^4 \log^2(n) \log^2(K)}{K^{\frac{2}{3}}(1-2\delta)^4} + \frac{\sigma^2 n \delta^4 \log^2(n) \log^2(K)}{(\sigma^2 K + n)(1-2\delta)^4} + \frac{n \delta^6 \log^4(n) \log^2(K)}{(\sigma^2 K + n)(1-2\delta)^6}$$

That completes the proof \square

G. Moniqua on D^2 (Proof to Theorem 4)

G.1. Setting

We first show the pseudo code in Algorithm 1.

D^2 makes the following assumptions (1-4), and we add the additional assumption (5):

1. **Lipschitzian Gradient:** All the function f_i have L-Lipschitzian gradients.
2. **Communication Matrix:** Communication matrix \mathbf{W} is a symmetric doubly stochastic matrix. Let the eigenvalues of $\mathbf{W} \in \mathbb{R}^{n \times n}$ be $\lambda_1 \geq \dots \geq \lambda_n$. We assume $\lambda_2 < 1, \lambda_n > -\frac{1}{3}$.

3. **Bounded Variance:**

$$\mathbb{E}_{\xi_i \sim \mathcal{D}_i} \left\| \nabla \tilde{f}_i(\mathbf{x}; \xi_i) - \nabla f_i(\mathbf{x}) \right\|^2 \leq \sigma^2, \forall i$$

where $\nabla \tilde{f}_i(\mathbf{x}; \xi_i)$ denotes gradient sample on worker i computed via data sample ξ_i .

4. **Initialization:** All the models are initialized by the same parameters: $\mathbf{x}_{0,i} = \mathbf{x}_0, \forall i$ and with out the loss of generality $\mathbf{x}_0 = 0$.
5. **Gradient magnitude:** The norm of a sampled gradient is bounded by $\|\tilde{\mathbf{g}}_{k,i}\|_\infty \leq G_\infty$ for some constant G_∞ .

Algorithm 1 Moniqua with Variance Reduction on worker i

Require: initial point $\mathbf{x}_{0,i} = \mathbf{x}_0$, step size α , the discrepancy bound B_θ , communication matrix \mathbf{W} , number of iterations K , neighbor list of worker i : \mathcal{N}_i , quantizer Q_δ

- 1: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
- 2: Randomly sample data $\xi_{k,i}$ from local memory
- 3: Compute a local stochastic gradient based on $\xi_{k,i}$ and current weight $\mathbf{x}_{k,i}$: $\tilde{\mathbf{g}}_{k,i}$
- 4: **if** $k = 0$ **then**
- 5: Update local weight: $\mathbf{x}_{k+\frac{1}{2},i} \leftarrow \mathbf{x}_{k,i} - \alpha\tilde{\mathbf{g}}_{k,i}$
- 6: **else**
- 7: Update local weight: $\mathbf{x}_{k+\frac{1}{2},i} \leftarrow 2\mathbf{x}_{k,i} - \mathbf{x}_{k-1,i} - \alpha\tilde{\mathbf{g}}_{k,i} + \alpha\tilde{\mathbf{g}}_{k-1,i}$
- 8: **end if**
- 9: Send modulo-ed model to neighbors: $\mathbf{q}_{k+\frac{1}{2},i} \leftarrow Q_\delta \left(\frac{\mathbf{x}_{k+\frac{1}{2},i}}{B_\theta} \bmod 1 \right)$
- 10: Compute local biased term $\hat{\mathbf{x}}_{k+\frac{1}{2},i}$ as:

$$\hat{\mathbf{x}}_{k+\frac{1}{2},i} = \mathbf{q}_{k+\frac{1}{2},i} B_\theta - \mathbf{x}_{k+\frac{1}{2},i} \bmod B_\theta + \mathbf{x}_{k+\frac{1}{2},i}$$

- 11: Recover model received from worker j as:

$$\hat{\mathbf{x}}_{k+\frac{1}{2},j} = (\mathbf{q}_{k+\frac{1}{2},j} B_\theta - \mathbf{x}_{k+\frac{1}{2},j}) \bmod B_\theta + \mathbf{x}_{k+\frac{1}{2},j}$$

- 12: Average with neighboring workers: $\mathbf{x}_{k+1,i} \leftarrow \mathbf{x}_{k+\frac{1}{2},i} + \sum_{j \in \mathcal{N}_i} (\hat{\mathbf{x}}_{k+\frac{1}{2},j} - \hat{\mathbf{x}}_{k+\frac{1}{2},i}) \mathbf{W}_{ji}$
 - 13: **end for**
 - 14: **return** $\bar{\mathbf{X}}_K = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{K,i}$
-

G.2. Proof to Theorem 4

Proof. From a local view, define $\mathbf{x}_{-1} = \tilde{\mathbf{g}}_{-1} = 0$, the update rule of Moniqua on D^2 on worker i in iteration k can be written as

$$\begin{aligned} \mathbf{x}_{k+\frac{1}{2},i} &= 2\mathbf{x}_{k,i} - \mathbf{x}_{k-1,i} - \alpha\tilde{\mathbf{g}}_{k,i} + \alpha\tilde{\mathbf{g}}_{k-1,i} \\ \mathbf{x}_{k+1,i} &= \sum_{j=1}^n \mathbf{x}_{k+\frac{1}{2},j} \mathbf{W}_{ji} + \sum_{j=1}^n \left((\hat{\mathbf{x}}_{k+\frac{1}{2},j} - \mathbf{x}_{k+\frac{1}{2},j}) - (\hat{\mathbf{x}}_{k+\frac{1}{2},i} - \mathbf{x}_{k+\frac{1}{2},i}) \right) \mathbf{W}_{ji} \end{aligned}$$

For a more compact expression,

$$\begin{aligned} \mathbf{X}_{k+\frac{1}{2}} &= 2\mathbf{X}_k - \mathbf{X}_{k-1} - \alpha\tilde{\mathbf{G}}_k + \alpha\tilde{\mathbf{G}}_{k-1} \\ \mathbf{X}_{k+1} &= \mathbf{X}_{k+\frac{1}{2}} \mathbf{W} + (\hat{\mathbf{X}}_{k+\frac{1}{2}} - \mathbf{X}_{k+\frac{1}{2}})(\mathbf{W} - \mathbf{I}) \end{aligned}$$

Define

$$\mathbf{\Omega}_k = (\hat{\mathbf{X}}_{k+\frac{1}{2}} - \mathbf{X}_{k+\frac{1}{2}})(\mathbf{W} - \mathbf{I})$$

Since \mathbf{W} is symmetric, it can be diagonalized as $\mathbf{W} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^\top$, where the i -th column of \mathbf{P} and $\mathbf{\Lambda}$ are \mathbf{W} 's i -th eigenvector and eigenvalue, respectively. And we obtain

$$\mathbf{X}_{k+1} = 2\mathbf{X}_k \mathbf{P} \mathbf{\Lambda} \mathbf{P}^\top - \mathbf{X}_{k-1} \mathbf{P} \mathbf{\Lambda} \mathbf{P}^\top - \alpha\tilde{\mathbf{G}}_k \mathbf{P} \mathbf{\Lambda} \mathbf{P}^\top + \alpha\tilde{\mathbf{G}}_{k-1} \mathbf{P} \mathbf{\Lambda} \mathbf{P}^\top + \mathbf{\Omega}_k$$

and

$$\mathbf{X}_{k+1} \mathbf{P} = 2\mathbf{X}_k \mathbf{P} \mathbf{\Lambda} - \mathbf{X}_{k-1} \mathbf{P} \mathbf{\Lambda} - \alpha\tilde{\mathbf{G}}_k \mathbf{P} \mathbf{\Lambda} + \alpha\tilde{\mathbf{G}}_{k-1} \mathbf{P} \mathbf{\Lambda} + \mathbf{\Omega}_k \mathbf{P}$$

Denote $\mathbf{Y}_k = \mathbf{X}_k \mathbf{P}$, $\mathbf{H}(\mathbf{X}_k; \xi_k) = \tilde{\mathbf{G}}_k \mathbf{P}$, and denote $\mathbf{y}_{k,i}$, $\mathbf{h}_{k,i}$ and $\mathbf{r}_{k,i}$ as the i -th column of \mathbf{Y}_k , \mathbf{H}_k and $\mathbf{\Omega}_k \mathbf{P}$, respectively. Then we have

$$\mathbf{y}_{k+1,i} = \lambda_i (2\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i} - \alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}) + \mathbf{r}_{k,i}$$

From Lemma G.5 (Constants C_1, C_2, C_3 and C_4 are defined in the Lemma G.1. Constants D_1 and D_2 are defined in Lemma G.5) we get

$$\begin{aligned} & \left(1 - \frac{3C_1\alpha^2L^2}{C_4}\right) \mathbb{E} \|\nabla f(\mathbf{0})\| + \left(1 - \alpha L - 3\frac{C_2}{C_4}\alpha^4L^4\right) \frac{1}{K} \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ & \leq \frac{2(f(0) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1\alpha^2L^2(\sigma^2 + \varsigma_0^2)}{C_4K} + 6\frac{C_2}{C_4}\alpha^2\sigma^2L^2 + 3\frac{C_2}{nC_4}\alpha^4\sigma^2L^4 + \frac{C_3L^2}{C_4} \left(\frac{6D_1n + 8}{6D_2n + 1}\right)^2 \alpha^2 G_\infty^2 d \end{aligned}$$

Let $\alpha = \frac{1}{\sigma\sqrt{K/n+2L}}$, we have

$$\begin{aligned} & \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ & \leq \frac{2(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1\alpha^2L^2(\sigma^2 + \varsigma_0^2)}{C_4K} + 6\frac{C_2}{C_4}\alpha^2\sigma^2L^2 + 3\frac{C_2}{nC_4}\alpha^4\sigma^2L^4 + \left(\frac{6D_1n + 8}{6D_2n + 1}\right)^2 \frac{C_3L^2}{C_4} G_\infty^2 d \alpha^2 \\ & \leq \frac{4(f(\mathbf{0}) - f^*)L}{K} + \frac{2\sigma(f(\mathbf{0}) - f^* + L/2)}{\sqrt{nK}} + \frac{3C_1L^2(\sigma^2 + \varsigma_0^2)n}{C_4(\sigma^2K^2 + 4nL^2K)} + \frac{6C_2L^2\sigma^2n}{C_4(\sigma^2K + 4nL^2)} \\ & \quad + \frac{3C_2n\sigma^2L^2}{C_4(\sigma^4K^2 + 16n^2L^4)} + \left(\frac{6D_1n + 8}{6D_2n + 1}\right)^2 \frac{C_3G_\infty^2dL^2n}{C_4(\sigma^2K + 4nL^2)} \\ & \lesssim \frac{1}{K} + \frac{\sigma}{\sqrt{nK}} + \frac{(\sigma^2 + \varsigma_0^2)n}{\sigma^2K^2 + nK} + \frac{\sigma^2n}{\sigma^2K + n} + \frac{\sigma^2n}{\sigma^4K^2 + n^2} + \frac{G_\infty^2dn}{\sigma^2K + n} \\ & \lesssim \frac{1}{K} + \frac{\sigma}{\sqrt{nK}} + \frac{\sigma^2n}{\sigma^2K + n} + \frac{G_\infty^2dn}{\sigma^2K + n} \end{aligned}$$

That completes the proof. □

G.3. Lemma for D^2

Lemma G.1. Define

$$\begin{aligned} D_1 &= \max \left\{ |v_n| + \frac{2|\lambda_n|}{1 - |v_n|}, \sqrt{\frac{\lambda_2}{1 - \lambda_2}} + \frac{2\lambda_2}{1 - \lambda_2} \right\} \\ D_2 &= \max \left\{ \frac{2}{1 - |v_n|}, \frac{2}{\sqrt{1 - \lambda_2}} \right\} \\ v_n &= \lambda_n - \sqrt{\lambda_n^2 - \lambda_n} \end{aligned}$$

Let $\delta = \frac{1}{12nD_2+2}$, and we have for $\forall i, j$

$$\left\| \mathbf{x}_{k+\frac{1}{2}}(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty < \theta = (6D_1n + 8)\alpha G_\infty$$

Proof. We use mathematical induction to prove this:

I. When $k = 0$,

$$\left\| \mathbf{X}_{0+\frac{1}{2}}(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty = \left\| -\alpha \tilde{\mathbf{G}}_0(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \leq \alpha \left\| \tilde{\mathbf{G}}_0 \right\|_{1,\infty} \|\mathbf{e}_i - \mathbf{e}_j\|_1 < 2\alpha G_\infty \leq (6D_1n + 8)\alpha G_\infty$$

II. Suppose for $k \geq 0, \forall t \leq k$, we have $\left\| \mathbf{X}_{t+\frac{1}{2}}(\mathbf{e}_i - \mathbf{e}_j) \right\| < (6D_1n + 8)\alpha G_\infty$, then for $\forall i, j$

$$\begin{aligned} & \left\| \mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\ & \leq \left\| \mathbf{X}_{k+1} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|_\infty + \left\| \mathbf{X}_{k+1} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_j \right) \right\|_\infty \end{aligned}$$

$$\begin{aligned}
 &= \left\| \left\| \mathbf{X}_{k+1} \mathbf{P} \mathbf{P}^\top \mathbf{e}_i - \mathbf{X}_{k+1} \mathbf{P} \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \mathbf{P}^\top \mathbf{e}_i \right\|_\infty + \left\| \left\| \mathbf{X}_{k+1} \mathbf{P} \mathbf{P}^\top \mathbf{e}_j - \mathbf{X}_{k+1} \mathbf{P} \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \mathbf{P}^\top \mathbf{e}_j \right\|_\infty \right\|_\infty \\
 &\leq \left\| \left\| \mathbf{X}_{k+1} \mathbf{P} \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right\|_{1,\infty} \left\| \mathbf{P}^\top \mathbf{e}_i \right\|_1 + \left\| \left\| \mathbf{X}_{k+1} \mathbf{P} \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right\|_{1,\infty} \left\| \mathbf{P}^\top \mathbf{e}_j \right\|_1 \right\|_\infty \\
 &\leq 2\sqrt{n} \left\| \left\| \mathbf{X}_{k+1} \mathbf{P} \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right\|_{1,\infty} \right\|_\infty
 \end{aligned}$$

From the update rule, we have

$$\mathbf{y}_{k+1,i} = \lambda_i(2\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i} - \alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}) + \mathbf{r}_{k,i} = \lambda_i(2\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}) + \lambda_i\boldsymbol{\beta}_{k,i} + \mathbf{r}_{k,i}$$

where $\boldsymbol{\beta}_{k,i} = -\alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}$, for all \mathbf{y}_i with $-\frac{1}{3} < \lambda_i < 0$, from Lemma G.3 we have

$$\mathbf{y}_{k+1,i} = \mathbf{y}_{1,i} \left(\frac{u_i^{k+1} - v_i^{k+1}}{u_i - v_i} \right) + \sum_{s=1}^k (\lambda_i \boldsymbol{\beta}_{s,i} + \mathbf{r}_{s,i}) \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i}$$

where $u_i = \lambda_i + \sqrt{\lambda_i^2 - \lambda_i}$ and $v_i = \lambda_i - \sqrt{\lambda_i^2 - \lambda_i}$, we obtain

$$\|\mathbf{y}_{k+1,i}\|_\infty \leq \|\mathbf{y}_{1,i}\|_\infty \left| \frac{u_i^{k+1} - v_i^{k+1}}{u_i - v_i} \right| + |\lambda_i| \sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\|_\infty \left| \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i} \right| + \sum_{s=1}^k \|\mathbf{r}_{s,i}\|_\infty \left| \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i} \right|$$

Since

$$\left| \frac{u_i^{n+1} - v_i^{n+1}}{u_i - v_i} \right| \leq |v_i|^n \left| \frac{u_i \left(\frac{u_i}{v_i} \right)^n - v_i}{u_i - v_i} \right| \leq |v_i|^n$$

We obtain

$$\|\mathbf{y}_{k+1,i}\|_\infty \leq \|\mathbf{y}_{1,i}\|_\infty |v_i|^k + |\lambda_i| \sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\|_\infty |v_i|^{k-s} + \sum_{s=1}^k \|\mathbf{r}_{s,i}\|_\infty |v_i|^{k-s}$$

For $\boldsymbol{\beta}_{s,i}$, we have

$$\begin{aligned}
 \|\boldsymbol{\beta}_{s,i}\|_\infty &= \|-\alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}\|_\infty \leq 2\alpha(\|\mathbf{h}_{k,i}\|_\infty + \|\mathbf{h}_{k-1,i}\|_\infty) \\
 &\leq 2\alpha(\|\mathbf{G}_k\|_{1,\infty} \|\mathbf{P}\mathbf{e}_i\|_1 + \|\mathbf{G}_{k-1}\|_{1,\infty} \|\mathbf{P}\mathbf{e}_i\|_1) \\
 &\leq 2\alpha\sqrt{n}G_\infty
 \end{aligned}$$

For $\mathbf{r}_{s,i}$, we have

$$\|\mathbf{r}_{k,i}\|_\infty = \|\boldsymbol{\Omega}_k \mathbf{P}\mathbf{e}_i\|_\infty \leq \|\boldsymbol{\Omega}_k\|_{1,\infty} \|\mathbf{P}\mathbf{e}_i\|_1 \leq 2\sqrt{n}\delta B_\theta$$

when $\lambda_i < 0$, we have

$$\|\mathbf{y}_{k+1,i}\|_\infty \leq \|\mathbf{y}_{1,i}\|_\infty |v_i|^k + |\lambda_i| \sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\|_\infty |v_i|^{k-s} + \sum_{s=1}^k \|\mathbf{r}_{s,i}\|_\infty |v_i|^{k-s}$$

$$\begin{aligned}
 &\leq \|\mathbf{y}_{1,i}\|_\infty |v_n|^k + |\lambda_n| \sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\|_\infty |v_n|^{k-s} + \sum_{s=1}^k \|\mathbf{r}_{s,i}\|_\infty |v_n|^{k-s} \\
 &\leq \alpha\sqrt{n}G_\infty |v_n|^k + 2\alpha\sqrt{n}G_\infty |\lambda_n| \sum_{s=1}^\infty |v_n|^{k-s} + 2\sqrt{n}\delta B_\theta \sum_{s=1}^\infty |v_n|^{k-s} \\
 &\leq \alpha\sqrt{n}G_\infty |v_n| + \frac{2\alpha\sqrt{n}G_\infty |\lambda_n|}{1 - |v_n|} + \frac{2\sqrt{n}\delta B_\theta}{1 - |v_n|}
 \end{aligned}$$

where $v_n = \lambda_n - \sqrt{\lambda_n^2 - \lambda_n}$.

On the other hand, when $0 \leq \lambda_i < 1$, from Lemma G.3 we have

$$\mathbf{y}_{k+1,i} \sin \phi_i = \mathbf{y}_{1,i} \lambda_i^{\frac{k}{2}} \sin[(t+1)\phi_i] + \lambda_i \sum_{s=1}^k \boldsymbol{\beta}_{s,i} \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i] + \sum_{s=1}^k \mathbf{r}_{s,i} \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i]$$

By taking norm, we get

$$\begin{aligned}
 \|\mathbf{y}_{k+1,i}\|_\infty |\sin \phi_i| &= \|\mathbf{y}_{1,i}\|_\infty \lambda_i^{\frac{k}{2}} |\sin[(t+1)\phi_i]| + \lambda_i \sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\|_\infty |\lambda_i^{\frac{k-s}{2}}| |\sin[(k+1-s)\phi_i]| \\
 &\quad + \sum_{s=1}^k \|\mathbf{r}_{s,i}\|_\infty |\lambda_i^{\frac{k-s}{2}}| |\sin[(k+1-s)\phi_i]| \\
 &< \|\mathbf{y}_{1,i}\|_\infty \lambda_2^{\frac{k}{2}} + 2\alpha\sqrt{n}G_\infty \lambda_2 \sum_{s=1}^\infty \lambda_2^{\frac{s}{2}} + 2\sqrt{n}\delta B_\theta \sum_{s=1}^\infty \lambda_2^{\frac{s}{2}} \\
 &\leq \alpha\sqrt{n}G_\infty \sqrt{\lambda_2} + \frac{2\alpha\sqrt{n}G_\infty \lambda_2 + 2\sqrt{n}\delta B_\theta}{\sqrt{1 - \lambda_2}}
 \end{aligned}$$

Since $|\sin \phi_i| \geq \sqrt{1 - \lambda_2}$, putting it back, we get

$$\|\mathbf{y}_{k+1,i}\|_\infty < \alpha\sqrt{n}G_\infty \sqrt{\frac{\lambda_2}{1 - \lambda_2}} + \frac{2\alpha\sqrt{n}G_\infty \lambda_2 + 2\sqrt{n}\delta B_\theta}{1 - \lambda_2}$$

So there exists D_1, D_2

$$\begin{aligned}
 D_1 &= \max \left\{ |v_n| + \frac{2|\lambda_n|}{1 - |v_n|}, \sqrt{\frac{\lambda_2}{1 - \lambda_2}} + \frac{2\lambda_2}{1 - \lambda_2} \right\} \\
 D_2 &= \max \left\{ \frac{2}{1 - |v_n|}, \frac{2}{\sqrt{1 - \lambda_2}} \right\}
 \end{aligned}$$

such that

$$\|\mathbf{y}_{k+1,i}\|_\infty < D_1 \alpha\sqrt{n}G_\infty + D_2 \sqrt{n}\delta B_\theta$$

Putting it back we have $\forall i, j$

$$\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty < D_1 \alpha n G_\infty + D_2 n \delta B_\theta$$

As a result

$$\begin{aligned}
 &\left\| \mathbf{X}_{k+1+\frac{1}{2}}(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\
 &= \left\| (2\mathbf{X}_{k+1} - \mathbf{X}_k - \alpha\tilde{\mathbf{G}}_{k+1} + \alpha\tilde{\mathbf{G}}_k)(\mathbf{e}_i - \mathbf{e}_j) \right\|_\infty \\
 &\leq 2\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_\infty + \|\mathbf{X}_k(\mathbf{e}_i - \mathbf{e}_j)\|_\infty + \alpha \left\| \tilde{\mathbf{G}}_{k+1} \right\|_{1,\infty} \|\mathbf{e}_i - \mathbf{e}_j\|_1 + \alpha \left\| \tilde{\mathbf{G}}_k \right\|_{1,\infty} \|\mathbf{e}_i - \mathbf{e}_j\|_1 \\
 &< 3(D_1 \alpha n G_\infty + D_2 n \delta B_\theta) + 4\alpha G_\infty
 \end{aligned}$$

$$\leq (6D_1n + 8)\alpha G_\infty$$

The last step is because $\delta = \frac{1}{12nD_2+2}$

Combining I and II we complete the proof. \square

Lemma G.2. *By defining*

$$\begin{aligned} C_1 &= \max \left\{ \frac{3}{1 - |v_n|^2}, \frac{3}{(1 - \lambda_2)^2} \right\} \\ C_2 &= \max \left\{ \frac{3\lambda_n^2}{(1 - |v_n|)^2}, \frac{3\lambda_2^2}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \\ C_3 &= \max \left\{ \frac{3}{(1 - |v_n|)^2}, \frac{3}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \end{aligned}$$

we have

$$\begin{aligned} & (1 - 12C_2\alpha^2L^2) \sum_{i=1}^n \sum_{k=1}^K \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 \\ & \leq 3C_1\alpha^2n\sigma^2 + 3C_1\alpha^2n\zeta_0^2 + 3C_1\alpha^2n\mathbb{E} \|\nabla f(\mathbf{0})\| + 6C_2\alpha^2n\sigma^2K + 3C_2\alpha^4\sigma^2L^2K \\ & \quad + 3C_2\alpha^4nL^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2 \end{aligned}$$

Proof.

$$\begin{aligned} \sum_{i=1}^n \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 &= \sum_{i=1}^n \left\| \mathbf{X}_k \left(\mathbf{e}_i - \frac{\mathbf{1}}{n} \right) \right\|^2 \\ &= \left\| \mathbf{X}_k \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_F^2 \\ &= \left\| \mathbf{X}_k \mathbf{P} \mathbf{P}^\top - \mathbf{X}_k \mathbf{v}_1 \mathbf{v}_1^\top \right\|_F^2 \\ &\stackrel{\text{Lemma G.4}}{=} \left\| \mathbf{X}_k \mathbf{P} \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right\|_F^2 \\ &= \sum_{i=2}^n \|\mathbf{y}_{k,i}\|^2 \end{aligned}$$

From the update rule, we obtain,

$$\mathbf{y}_{k+1,i} = \lambda_i(2\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i} - \alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}) + \mathbf{r}_{k,i} = \lambda_i(2\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}) + \lambda_i\boldsymbol{\beta}_{k,i} + \mathbf{r}_{k,i}$$

where $\boldsymbol{\beta}_{k,i} = -\alpha\mathbf{h}_{k,i} + \alpha\mathbf{h}_{k-1,i}$, for all \mathbf{y}_i with $-\frac{1}{3} < \lambda_i < 0$, from Lemma G.3 we have

$$\mathbf{y}_{k+1,i} = \mathbf{y}_{1,i} \left(\frac{u_i^{k+1} - v_i^{k+1}}{u_i - v_i} \right) + \sum_{s=1}^k (\lambda_i \boldsymbol{\beta}_{s,i} + \mathbf{r}_{k,i}) \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i}$$

where $u_i = \lambda_i + \sqrt{\lambda_i^2 - \lambda_i}$ and $v_i = \lambda_i - \sqrt{\lambda_i^2 - \lambda_i}$, we obtain

$$\|\mathbf{y}_{k+1,i}\|^2 \leq 3 \|\mathbf{y}_{1,i}\|^2 \left(\frac{u_i^{k+1} - v_i^{k+1}}{u_i - v_i} \right)^2 + 3\lambda_i^2 \left(\sum_{s=1}^k \|\boldsymbol{\beta}_{s,i}\| \left| \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i} \right| \right)^2$$

$$+ 3 \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| \left| \frac{u_i^{k-s+1} - v_i^{k-s+1}}{u_i - v_i} \right| \right)^2$$

Since

$$\left| \frac{u_i^{n+1} - v_i^{n+1}}{u_i - v_i} \right| \leq |v_i|^n \left| \frac{u_i \left(\frac{u_i}{v_i} \right)^n - v_i}{u_i - v_i} \right| \leq |v_i|^n$$

We obtain

$$\|\mathbf{y}_{k+1,i}\|^2 \leq 3 \|\mathbf{y}_{1,i}\|^2 |v_i|^{2t} + 3\lambda_i^2 \left(\sum_{s=1}^k \|\beta_{s,i}\| |v_i|^{k-s} \right)^2 + 3 \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| |v_i|^{k-s} \right)^2$$

Summing over from $k = 0$ to $t = K - 1$, we obtain

$$\begin{aligned} \sum_{k=0}^{K-1} \|\mathbf{y}_{k+1,i}\|^2 &= \sum_{k=1}^K \|\mathbf{y}_{k,i}\|^2 \\ &\leq 3 \|\mathbf{y}_{1,i}\|^2 \sum_{k=0}^{K-1} |v_i|^{2k} + 3\lambda_i^2 \sum_{k=1}^{K-1} \left(\sum_{s=1}^k \|\beta_{s,i}\| |v_i|^{k-s} \right)^2 + 3 \sum_{k=1}^{K-1} \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| |v_i|^{k-s} \right)^2 \\ &\leq \frac{3 \|\mathbf{y}_{1,i}\|^2}{1 - |v_i|^2} + \frac{3\lambda_i^2}{(1 - |v_i|)^2} \sum_{k=1}^{K-1} \|\beta_{k,i}\|^2 + \frac{3}{(1 - |v_i|)^2} \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2 \\ &\leq \frac{3 \|\mathbf{y}_{1,i}\|^2}{1 - |v_n|^2} + \frac{3\lambda_n^2}{(1 - |v_n|)^2} \sum_{k=1}^{K-1} \|\beta_{k,i}\|^2 + \frac{3}{(1 - |v_n|)^2} \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2 \end{aligned}$$

where $v_n = \lambda_n - \sqrt{\lambda_n^2 - \lambda_n}$.

On the other hand, when $0 \leq \lambda_i < 1$, from Lemma G.3 we have

$$\mathbf{y}_{k+1,i} \sin \phi_i = \mathbf{y}_{1,i} \lambda_i^{\frac{k}{2}} \sin[(t+1)\phi_i] + \lambda_i \sum_{s=1}^k \beta_{s,i} \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i] + \sum_{s=1}^k \mathbf{r}_{s,i} \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i]$$

And we have

$$\begin{aligned} \|\mathbf{y}_{k+1,i}\|^2 \sin^2 \phi_i &\leq 3 \|\mathbf{y}_{1,i}\|^2 \lambda_i^k \sin^2[(t+1)\phi_i] + 3\lambda_i^2 \left(\sum_{s=1}^k \|\beta_{s,i}\| \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i] \right)^2 \\ &\quad + 3 \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| \lambda_i^{\frac{k-s}{2}} \sin[(k+1-s)\phi_i] \right)^2 \\ &\leq 3 \|\mathbf{y}_{1,i}\|^2 \lambda_i^k + 3\lambda_i^2 \left(\sum_{s=1}^k \|\beta_{s,i}\| \lambda_i^{\frac{k-s}{2}} \right)^2 + 3 \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| \lambda_i^{\frac{k-s}{2}} \right)^2 \end{aligned}$$

Summing from $k = 0$ to $K - 1$, we have

$$\begin{aligned} \sum_{k=0}^{K-1} \|\mathbf{y}_{k+1,i}\|^2 \sin^2 \phi_i &= \sum_{k=1}^K \|\mathbf{y}_{k,i}\|^2 \sin^2 \phi_i \\ &\leq 3 \|\mathbf{y}_{1,i}\|^2 \sum_{k=0}^{K-1} \lambda_i^k + 3\lambda_i^2 \sum_{k=1}^{K-1} \left(\sum_{s=1}^k \|\beta_{s,i}\| \lambda_i^{\frac{k-s}{2}} \right)^2 + 3 \sum_{k=1}^{K-1} \left(\sum_{s=1}^k \|\mathbf{r}_{s,i}\| \lambda_i^{\frac{k-s}{2}} \right)^2 \\ &\leq \frac{3 \|\mathbf{y}_{1,i}\|^2}{1 - \lambda_i} + \frac{3\lambda_i^2}{(1 - \sqrt{\lambda_i})^2} \sum_{k=1}^{K-1} \|\beta_{k,i}\|^2 + \frac{3}{(1 - \sqrt{\lambda_i})^2} \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2 \end{aligned}$$

Since $\sin^2 \phi_i = 1 - \lambda_i$, we have

$$\begin{aligned} \sum_{k=1}^K \|\mathbf{y}_{k,i}\|^2 &\leq \frac{3 \|\mathbf{y}_{1,i}\|^2}{(1 - \lambda_i)^2} + \frac{3\lambda_i^2}{(1 - \sqrt{\lambda_i})^2(1 - \lambda_i)} \sum_{k=1}^{K-1} \|\boldsymbol{\beta}_{k,i}\|^2 + \frac{3}{(1 - \sqrt{\lambda_i})^2(1 - \lambda_i)} \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2 \\ &\leq \frac{3 \|\mathbf{y}_{1,i}\|^2}{(1 - \lambda_2)^2} + \frac{3\lambda_2^2}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \sum_{k=1}^{K-1} \|\boldsymbol{\beta}_{k,i}\|^2 + \frac{3}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2 \end{aligned}$$

So there exists C_1, C_2, C_3

$$\begin{aligned} C_1 &= \max \left\{ \frac{3}{1 - |v_n|^2}, \frac{3}{(1 - \lambda_2)^2} \right\} \\ C_2 &= \max \left\{ \frac{3\lambda_n^2}{(1 - |v_n|)^2}, \frac{3\lambda_2^2}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \\ C_3 &= \max \left\{ \frac{3}{(1 - |v_n|)^2}, \frac{3}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \end{aligned}$$

$$\sum_{k=1}^K \|\mathbf{y}_{k,i}\|^2 \leq C_1 \|\mathbf{y}_{1,i}\|^2 + C_2 \sum_{k=1}^{K-1} \|\boldsymbol{\beta}_{k,i}\|^2 + C_3 \sum_{k=1}^{K-1} \|\mathbf{r}_{k,i}\|^2$$

By taking expectation we have

$$\sum_{k=1}^K \mathbb{E} \|\mathbf{y}_{k,i}\|^2 \leq C_1 \mathbb{E} \|\mathbf{y}_{1,i}\|^2 + C_2 \sum_{k=1}^{K-1} \mathbb{E} \|\boldsymbol{\beta}_{k,i}\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\mathbf{r}_{k,i}\|^2$$

We next analyze $\beta_{k,i}$:

$$\begin{aligned} &\sum_{i=2}^n \mathbb{E} \|\boldsymbol{\beta}_{k,i}\|^2 \\ &= \alpha^2 \sum_{i=2}^n \mathbb{E} \|\mathbf{h}_{k,i} - \mathbf{h}_{k-1,i}\|^2 \\ &= \alpha^2 \sum_{i=2}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_k \mathbf{P} \mathbf{e}_i - \tilde{\mathbf{G}}_{k-1} \mathbf{P} \mathbf{e}_i \right\|^2 \\ &\leq \alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_k \mathbf{P} \mathbf{e}_i - \tilde{\mathbf{G}}_{k-1} \mathbf{P} \mathbf{e}_i \right\|^2 \\ &\leq \alpha^2 \mathbb{E} \left\| \tilde{\mathbf{G}}_k \mathbf{P} - \tilde{\mathbf{G}}_{k-1} \mathbf{P} \right\|_F^2 \\ &\stackrel{\text{Lemma G.4}}{\leq} \alpha^2 \mathbb{E} \left\| \tilde{\mathbf{G}}_k - \tilde{\mathbf{G}}_{k-1} \right\|_F^2 \\ &= \alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_k \mathbf{e}_i - \tilde{\mathbf{G}}_{k-1} \mathbf{e}_i \right\|^2 \\ &\leq 3\alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_k \mathbf{e}_i - \mathbf{G}_k \mathbf{e}_i \right\|^2 + 3\alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_{k-1} \mathbf{e}_i - \mathbf{G}_{k-1} \mathbf{e}_i \right\|^2 \\ &\quad + 3\alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \mathbf{G}_k \mathbf{e}_i - \mathbf{G}_{k-1} \mathbf{e}_i \right\|^2 \\ &\leq 6\alpha^2 n \sigma^2 + 3\alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \mathbf{G}_k \mathbf{e}_i - \mathbf{G}_{k-1} \mathbf{e}_i \right\|^2 \end{aligned}$$

$$\begin{aligned}
 &\leq 6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \sum_{i=1}^n \mathbb{E} \|\mathbf{x}_{k,i} - \mathbf{x}_{k-1,i}\|^2 \\
 &\leq 6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \sum_{i=1}^n \mathbb{E} \left\| \mathbf{Y}_k \mathbf{P}^\top \mathbf{e}_i - \mathbf{Y}_{k-1} \mathbf{P}^\top \mathbf{e}_i \right\|^2 \\
 &\leq 6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \mathbb{E} \left\| \mathbf{Y}_k \mathbf{P}^\top - \mathbf{Y}_{k-1} \mathbf{P}^\top \right\|_F^2 \\
 &\stackrel{\text{Lemma G.4}}{\leq} 6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \mathbb{E} \|\mathbf{Y}_k - \mathbf{Y}_{k-1}\|_F^2 \\
 &\leq 6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \sum_{i=1}^n \mathbb{E} \|\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}\|^2
 \end{aligned}$$

Putting it back, we have

$$\begin{aligned}
 &\sum_{i=2}^n \sum_{k=1}^K \mathbb{E} \|\mathbf{y}_{k,i}\|^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + C_2 \sum_{i=2}^n \sum_{k=1}^{K-1} \mathbb{E} \|\beta_{k,i}\|^2 + C_3 \sum_{k=1}^{K-1} \sum_{i=2}^n \mathbb{E} \|\mathbf{r}_{k,i}\|^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + C_2 \sum_{k=1}^{K-1} \left(6\alpha^2 n\sigma^2 + 3\alpha^2 L^2 \sum_{i=1}^n \mathbb{E} \|\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}\|^2 \right) + C_3 \sum_{k=1}^{K-1} \sum_{i=2}^n \mathbb{E} \|\mathbf{r}_{k,i}\|^2 \\
 &\stackrel{\text{Lemma G.4}}{\leq} C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + 6C_2 \alpha^2 n\sigma^2 K + 3C_2 \alpha^2 L^2 \sum_{k=1}^{K-1} \sum_{i=1}^n \mathbb{E} \|\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2
 \end{aligned}$$

Since

$$\begin{aligned}
 &\mathbb{E} \|\mathbf{y}_{k,1} - \mathbf{y}_{k-1,1}\|^2 = \mathbb{E} \|\mathbf{X}_k \mathbf{P} \mathbf{e}_1 - \mathbf{X}_{k-1} \mathbf{P} \mathbf{e}_1\|^2 = \mathbb{E} \|\mathbf{X}_k \mathbf{v}_1 - \mathbf{X}_{k-1} \mathbf{v}_1\|^2 \\
 &= \mathbb{E} \left\| \mathbf{X}_k \frac{1}{\sqrt{n}} \mathbf{1} - \mathbf{X}_{k-1} \frac{1}{\sqrt{n}} \mathbf{1} \right\|^2 = n \mathbb{E} \|\bar{\mathbf{X}}_k - \bar{\mathbf{X}}_{k-1}\|^2 = n\alpha^2 \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \\
 &\leq n\alpha^2 \mathbb{E} \|\bar{\mathbf{G}}_k - \bar{\mathbf{G}}_k\|^2 + n\alpha^2 \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \leq n\alpha^2 \frac{\sigma^2}{n} + n\alpha^2 \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \\
 &= \alpha^2 \sigma^2 + n\alpha^2 \mathbb{E} \|\bar{\mathbf{G}}_k\|^2
 \end{aligned}$$

Putting it back, and we obtain

$$\begin{aligned}
 &\sum_{i=2}^n \sum_{k=1}^K \mathbb{E} \|\mathbf{y}_{k,i}\|^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + 6C_2 \alpha^2 n\sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K + 3C_2 \alpha^4 nL^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \\
 &\quad + 3C_2 \alpha^2 L^2 \sum_{k=1}^{K-1} \sum_{i=2}^n \mathbb{E} \|\mathbf{y}_{k,i} - \mathbf{y}_{k-1,i}\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + 6C_2 \alpha^2 n\sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K + 3C_2 \alpha^4 nL^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \\
 &\quad + 6C_2 \alpha^2 L^2 \sum_{k=1}^{K-1} \sum_{i=2}^n \mathbb{E} \left(\|\mathbf{y}_{k,i}\|^2 + \|\mathbf{y}_{k-1,i}\|^2 \right) + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2
 \end{aligned}$$

$$\begin{aligned}
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + 6C_2 \alpha^2 n \sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K + 3C_2 \alpha^4 n L^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 \\
 &\quad + 12C_2 \alpha^2 L^2 \sum_{k=1}^{K-1} \sum_{i=2}^n \mathbb{E} \|\mathbf{y}_{k,i}\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\boldsymbol{\Omega}_k\|_F^2
 \end{aligned}$$

Rearrange the terms, we get

$$\begin{aligned}
 &(1 - 12C_2 \alpha^2 L^2) \sum_{i=2}^n \sum_{k=1}^K \mathbb{E} \|\mathbf{y}_{k,i}\|^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{Y}_1\|_F^2 + 6C_2 \alpha^2 n \sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K + 3C_2 \alpha^4 n L^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\boldsymbol{\Omega}_k\|_F^2 \\
 &\leq C_1 \mathbb{E} \|\mathbf{X}_1\|_F^2 + 6C_2 \alpha^2 n \sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K + 3C_2 \alpha^4 n L^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\boldsymbol{\Omega}_k\|_F^2
 \end{aligned}$$

Considering

$$\begin{aligned}
 \mathbb{E} \|\mathbf{X}_1\|_F^2 &= \alpha^2 \mathbb{E} \|\tilde{\mathbf{G}}_0\|_F^2 \\
 &= \alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_{0,i} - \mathbf{G}_{0,i} + \mathbf{G}_{0,i} - \nabla f(\mathbf{0}) + \nabla f(\mathbf{0}) \right\|^2 \\
 &\leq 3\alpha^2 \sum_{i=1}^n \mathbb{E} \left\| \tilde{\mathbf{G}}_{0,i} - \mathbf{G}_{0,i} \right\|^2 + 3\alpha^2 \sum_{i=1}^n \mathbb{E} \|\mathbf{G}_{0,i} - \nabla f(\mathbf{0})\|^2 + 3\alpha^2 \sum_{i=1}^n \mathbb{E} \|\nabla f(\mathbf{0})\|^2 \\
 &\leq 3\alpha^2 n \sigma^2 + 3\alpha^2 n \varsigma_0^2 + 3\alpha^2 n \mathbb{E} \|\nabla f(\mathbf{0})\|
 \end{aligned}$$

We finally get

$$\begin{aligned}
 &(1 - 12C_2 \alpha^2 L^2) \sum_{i=2}^n \sum_{k=1}^K \mathbb{E} \|\mathbf{y}_{k,i}\|^2 \\
 &= (1 - 12C_2 \alpha^2 L^2) \sum_{i=1}^n \sum_{k=1}^K \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2 \\
 &\leq 3C_1 \alpha^2 n \sigma^2 + 3C_1 \alpha^2 n \varsigma_0^2 + 3C_1 \alpha^2 n \mathbb{E} \|\nabla f(\mathbf{0})\| + 6C_2 \alpha^2 n \sigma^2 K + 3C_2 \alpha^4 \sigma^2 L^2 K \\
 &\quad + 3C_2 \alpha^4 n L^2 \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + C_3 \sum_{k=1}^{K-1} \mathbb{E} \|\boldsymbol{\Omega}_k\|_F^2
 \end{aligned}$$

That completes the proof. \square

Lemma G.3. Given $\rho \in (-\frac{1}{3}, 0) \cup (0, 1)$, for any two sequence $\{a_t\}_{t=1}^\infty$, $\{b_t\}_{t=1}^\infty$ and $\{c_t\}_{t=1}^\infty$ that satisfying

$$\begin{aligned}
 a_0 &= b_0 = 0, \\
 a_{t+1} &= \rho(2a_t - a_{t-1}) + b_t - b_{t-1} + c_t, \forall t \geq 1
 \end{aligned}$$

we have

$$a_{t+1} = a_1 \left(\frac{u^{t+1} - v^{t+1}}{u - v} \right) + \sum_{s=1}^t (b_s - b_{s-1} + c_s) \left(\frac{u^{t-s+1} - v^{t-s+1}}{u - v} \right), \forall t \geq 0$$

where

$$u = \rho + \sqrt{\rho^2 - \rho}, v = \rho - \sqrt{\rho^2 - \rho}$$

Moreover, if $0 < \rho < 1$, we have

$$a_{t+1} = a_1 \rho^{\frac{t}{2}} \frac{\sin[(t+1)\phi]}{\sin \phi} + \sum_{s=1}^t (b_s - b_{s-1} + c_s) \rho^{\frac{t-s}{2}} \frac{\sin[(t-s+1)\phi]}{\sin \phi}$$

where

$$\phi = \arccos(\sqrt{\rho})$$

Proof. when $t \geq 1$, we have

$$a_{t+1} = 2\rho a_t - \rho a_{t-1} + b_t - b_{t-1} + c_t$$

since,

$$u = \rho + \sqrt{\rho^2 - \rho}, v = \rho - \sqrt{\rho^2 - \rho}$$

we obtain

$$a_{t+1} - ua_t = (a_t - ua_{t-1})v + b_t - b_{t-1} + c_t$$

Recursively we have

$$\begin{aligned} a_{t+1} - ua_t &= (a_t - ua_{t-1})v + b_t - b_{t-1} + c_t \\ &= (a_{t-1} - ua_{t-2})v^2 + (b_{t-1} - b_{t-2} + c_{t-1})v + b_t - b_{t-1} + c_t \\ &= (a_1 - ua_0)v^t + \sum_{s=1}^t (b_s - b_{s-1} + c_s)v^{t-s} \\ &= a_1 v^t + \sum_{s=1}^t (b_s - b_{s-1} + c_s)v^{t-s} \end{aligned}$$

Dividing both sides by u^{t+1} , we have

$$\begin{aligned} \frac{a_{t+1}}{u^{t+1}} &= \frac{a_t}{u^t} + u^{-(t+1)} \left(a_1 v^t + \sum_{s=1}^t (b_s - b_{s-1} + c_s)v^{t-s} \right) \\ &= \frac{a_{t-1}}{u^{t-1}} + u^{-t} \left(a_1 v^{t-1} + \sum_{s=1}^{t-1} (b_s - b_{s-1} + c_s)v^{t-1-s} \right) \\ &\quad + u^{-(t+1)} \left(a_1 v^t + \sum_{s=1}^t (b_s - b_{s-1} + c_s)v^{t-s} \right) \\ &= \frac{a_1}{u} + \sum_{k=1}^t u^{-k-1} \left(a_1 v^k + \sum_{s=1}^k (b_s - b_{s-1} + c_s)v^{k-s} \right) \end{aligned}$$

Multiplying both sides by u^{t+1}

$$\begin{aligned} a_{t+1} &= a_1 u^t + \sum_{k=1}^t u^{t-k} \left(a_1 v^k + \sum_{s=1}^k (b_s - b_{s-1} + c_s)v^{k-s} \right) \\ &= a_1 u^t \left(1 + \sum_{k=1}^t \left(\frac{v}{u} \right)^k \right) + u^t \sum_{k=1}^t \sum_{s=1}^k (b_s - b_{s-1} + c_s)v^{-s} \left(\frac{v}{u} \right)^k \\ &= a_1 u^t \sum_{k=0}^t \left(\frac{v}{u} \right)^k + u^t \sum_{s=1}^t \sum_{k=s}^t (b_s - b_{s-1} + c_s)v^{-s} \left(\frac{v}{u} \right)^k \\ &= a_1 u^t \left(\frac{1 - \left(\frac{v}{u} \right)^{t+1}}{1 - \frac{v}{u}} \right) + u^t \sum_{s=1}^t (b_s - b_{s-1} + c_s)v^{-s} \left(\frac{v}{u} \right)^s \frac{1 - \left(\frac{v}{u} \right)^{t-s+1}}{1 - \frac{v}{u}} \\ &= a_1 \left(\frac{u^{t+1} - v^{t+1}}{u - v} \right) + \sum_{s=1}^t (b_s - b_{s-1} + c_s) \frac{u^{t-s+1} - v^{t-s+1}}{u - v} \end{aligned}$$

Note that when $0 < \rho < 1$, both u and v are complex numbers, we have

$$u = \sqrt{\rho}e^{i\phi}, v = \sqrt{\rho}e^{-i\phi}$$

where $\phi = \arccos \sqrt{\rho}$. And under this context, we have

$$a_{t+1} = a_1 \rho^{\frac{t}{2}} \frac{\sin[(t+1)\phi]}{\sin \phi} + \sum_{s=1}^t (b_s - b_{s-1} + c_s) \rho^{\frac{t-s}{2}} \frac{\sin[(t-s+1)\phi]}{\sin \phi}$$

That completes the proof. □

Lemma G.4. For any matrix $\mathbf{X} \in \mathbb{R}^{N \times n}$, we have

$$\begin{aligned} \sum_{i=2}^n \|\mathbf{X}\mathbf{v}_i\|^2 &\leq \sum_{i=1}^n \|\mathbf{X}\mathbf{v}_i\|^2 = \|\mathbf{X}\|_F^2 \\ \sum_{i=1}^n \|\mathbf{X}\mathbf{P}^\top \mathbf{e}_i\|^2 &= \|\mathbf{X}\mathbf{P}^\top\|_F^2 = \|\mathbf{X}\|_F^2 \end{aligned}$$

Proof.

$$\sum_{i=2}^n \|\mathbf{X}\mathbf{v}_i\|^2 \leq \sum_{i=1}^n \|\mathbf{X}\mathbf{v}_i\|^2 = \|\mathbf{X}_t \mathbf{P}\|_F^2 = \text{Tr}(\mathbf{X}_t \mathbf{P} \mathbf{P}^\top \mathbf{X}_t^\top) = \text{Tr}(\mathbf{X}_t \mathbf{X}_t^\top) = \|\mathbf{X}_t\|_F^2$$

And similarly,

$$\sum_{i=1}^n \|\mathbf{X}\mathbf{P}^\top \mathbf{e}_i\|^2 = \|\mathbf{X}\mathbf{P}^\top\|_F^2 = \text{Tr}(\mathbf{X}_t \mathbf{P}^\top \mathbf{P} \mathbf{X}_t^\top) = \text{Tr}(\mathbf{X}_t \mathbf{X}_t^\top) = \|\mathbf{X}_t\|_F^2$$

That completes the proof. □

Lemma G.5. If we run Algorithm 1 for K iterations the following inequality holds:

$$\begin{aligned} &\left(1 - \frac{3C_1\alpha^2 L^2}{C_4}\right) \mathbb{E} \|\nabla f(\mathbf{0})\| + \left(1 - \alpha L - 3\frac{C_2}{C_4}\alpha^4 L^4\right) \frac{1}{K} \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ &\leq \frac{2(f(0) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1\alpha^2 L^2(\sigma^2 + \varsigma_0^2)}{C_4 K} + 6\frac{C_2}{C_4}\alpha^2 \sigma^2 L^2 + 3\frac{C_2}{nC_4}\alpha^4 \sigma^2 L^4 \\ &+ \frac{C_3 L^2}{C_4} \left(\frac{6D_1 n + 8}{6D_2 n + 1}\right)^2 \alpha^2 G_\infty^2 d \end{aligned}$$

where

$$\begin{aligned} C_1 &= \max \left\{ \frac{3}{1 - |v_n|^2}, \frac{3}{(1 - \lambda_2)^2} \right\} \\ C_2 &= \max \left\{ \frac{3\lambda_n^2}{(1 - |v_n|)^2}, \frac{3\lambda_2^2}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \\ C_3 &= \max \left\{ \frac{3}{(1 - |v_n|)^2}, \frac{3}{(1 - \sqrt{\lambda_2})^2(1 - \lambda_2)} \right\} \\ C_4 &= 1 - 12C_2\alpha^2 L^2 \end{aligned}$$

Proof. Since

$$\begin{aligned} \bar{\mathbf{X}}_{k+1} &= (2\mathbf{X}_k - \mathbf{X}_{k-1} - \alpha\tilde{\mathbf{G}}_k + \alpha\tilde{\mathbf{G}}_{k-1})\mathbf{W}\frac{1}{n} + (\hat{\mathbf{X}}_{k+\frac{1}{2}} - \mathbf{X}_{k+\frac{1}{2}})(\mathbf{W} - \mathbf{I})\frac{1}{n} \\ &= 2\bar{\mathbf{X}}_k - \bar{\mathbf{X}}_{k-1} - \alpha\bar{\mathbf{G}}_k + \alpha\bar{\mathbf{G}}_{k-1} \end{aligned}$$

and we have

$$\begin{aligned}\bar{\mathbf{X}}_{k+1} - \bar{\mathbf{X}}_k &= \bar{\mathbf{X}}_k - \bar{\mathbf{X}}_{k-1} - \alpha \bar{\mathbf{G}}_k + \alpha \bar{\mathbf{G}}_{k-1} \\ &= \bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_0 - \alpha \sum_{t=1}^k (\bar{\mathbf{G}}_t - \bar{\mathbf{G}}_{t-1}) \\ &= -\alpha \bar{\mathbf{G}}_k\end{aligned}$$

Note that the update of the averaged model is exactly the same as D-PSGD, thus we can reuse the result from D-PSGD for D^2 as follows:

$$\frac{1-\alpha L}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \leq \frac{2(f(0) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{L^2}{nK} \sum_{k=0}^{K-1} \sum_{i=1}^n \mathbb{E} \|\bar{\mathbf{X}}_k - \mathbf{x}_{k,i}\|^2$$

From Lemma G.2 we obtain

$$\begin{aligned}& \frac{1-\alpha L}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ & \leq \frac{2(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1 \alpha^2 L^2 (\sigma^2 + \varsigma_0^2 + \mathbb{E} \|\nabla f(\mathbf{0})\|)}{C_4 K} + 6 \frac{C_2}{C_4} \alpha^2 \sigma^2 L^2 + 3 \frac{C_2}{nC_4} \alpha^4 \sigma^2 L^4 \\ & \quad + 3 \frac{C_2}{C_4} \alpha^4 L^4 \frac{1}{K} \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{C_3 L^2}{C_4 n K} \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2\end{aligned}$$

Rearrange the terms, we get

$$\begin{aligned}& \left(1 - \frac{3C_1 \alpha^2 L^2}{C_4}\right) \mathbb{E} \|\nabla f(\mathbf{0})\| + \left(1 - \alpha L - 3 \frac{C_2}{C_4} \alpha^4 L^4\right) \frac{1}{K} \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ & \leq \frac{2(f(0) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1 \alpha^2 L^2 (\sigma^2 + \varsigma_0^2)}{C_4 K} + 6 \frac{C_2}{C_4} \alpha^2 \sigma^2 L^2 + 3 \frac{C_2}{nC_4} \alpha^4 \sigma^2 L^4 + \frac{C_3 L^2}{C_4 n K} \sum_{k=1}^{K-1} \mathbb{E} \|\Omega_k\|_F^2\end{aligned}$$

Similar to the case in D-PSGD, we have

$$\begin{aligned}\sum_{k=0}^{K-1} \mathbb{E} \|\Omega_k\|_F^2 &= \sum_{k=0}^{K-1} \sum_{i=1}^n \mathbb{E} \left\| \sum_{j=1}^n \left((\hat{\mathbf{x}}_{k+\frac{1}{2},j} - \mathbf{x}_{k+\frac{1}{2},j}) - (\hat{\mathbf{x}}_{k+\frac{1}{2},i} - \mathbf{x}_{k+\frac{1}{2},i}) \right) \mathbf{W}_{ji} \right\|^2 \\ &\stackrel{\text{Lemma F.1}}{\leq} 4 \sum_{k=0}^{K-1} \sum_{i=1}^n \delta^2 B_\theta^2 d \leq \left(\frac{6D_1 n + 8}{6D_2 n + 1} \right)^2 \alpha^2 G_\infty^2 d n K\end{aligned}$$

Putting it back, we obtain

$$\begin{aligned}& \left(1 - \frac{3C_1 \alpha^2 L^2}{C_4}\right) \mathbb{E} \|\nabla f(\mathbf{0})\| + \left(1 - \alpha L - 3 \frac{C_2}{C_4} \alpha^4 L^4\right) \frac{1}{K} \sum_{k=1}^{K-1} \mathbb{E} \|\bar{\mathbf{G}}_k\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 \\ & \leq \frac{2(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{\alpha L}{n} \sigma^2 + \frac{3C_1 \alpha^2 L^2 (\sigma^2 + \varsigma_0^2)}{C_4 K} + 6 \frac{C_2}{C_4} \alpha^2 \sigma^2 L^2 + 3 \frac{C_2}{nC_4} \alpha^4 \sigma^2 L^4 + \frac{C_3 L^2}{C_4} \left(\frac{6D_1 n + 8}{6D_2 n + 1} \right)^2 \alpha^2 G_\infty^2 d\end{aligned}$$

That completes the proof. \square

H. Moniqua on AD-PSGD (Proof to Theorem 5)

H.1. Definition and Notation

In the original analysis of AD-PSGD, to better capture the nature of workers computing at different speed, the objective function is expressed as

$$f(\mathbf{x}) = \sum_{i=1}^n p_i f_i(\mathbf{x})$$

Algorithm 2 Moniqua with Asynchronous Communication

Require: initial point $\mathbf{x}_{0,i} = \mathbf{x}_0$, step size α , the discrepancy bound B_θ , number of iterations K , quantization function \mathcal{Q}_δ , initial random seed

- 1: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
- 2: worker i_k is updating the gradient while during this iteration the global communication behaviour is written in the form of \mathbf{W}_k .
- 3: Compute a local stochastic gradient with model delayed by τ_k : $\tilde{\mathbf{g}}_{k-\tau_k, i_k}$
- 4: Send modulo-ed model to one randomly selected neighbor j_k : $\mathbf{q}_{k, i_k} \leftarrow \mathcal{Q}_\delta \left(\frac{\mathbf{x}_{k, i_k}}{B_\theta} \bmod 1 \right)$
- 5: Compute local biased term $\hat{\mathbf{x}}_{k, i_k}$ as:

$$\hat{\mathbf{x}}_{k, i_k} = \mathbf{q}_{k, i_k} B_\theta - \mathbf{x}_{k, i_k} \bmod B_\theta + \mathbf{x}_{k, i_k}$$

- 6: Randomly select one neighbor j_k and recover its model as:

$$\hat{\mathbf{x}}_{k, j_k} = (\mathbf{q}_{k, j_k} B_\theta - \mathbf{x}_{k, i_k}) \bmod B_\theta + \mathbf{x}_{k, i_k}$$

- 7: Average with neighboring workers: $\mathbf{x}_{k, i_k} \leftarrow \mathbf{x}_{k, i_k} + \sum_{j \in \mathcal{N}_i} (\hat{\mathbf{x}}_{k, j_k} - \hat{\mathbf{x}}_{k, i_k}) \mathbf{W}^{ji}$
 - 8: Update the local weight with local gradient: $\mathbf{x}_{k+1, i_k} \leftarrow \mathbf{x}_{k, i_k} - \alpha \tilde{\mathbf{g}}_{k-\tau_k, i_k}$
 - 9: **end for**
 - 10: **return** $\bar{\mathbf{X}}_K = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{K, i}$
-

where p_i is a parameter denoting the speed of i -th worker gradient updates. In the rest of the proof, we denote $p = \max_i \{p_i\}$

For simplicity, we also define the following terms

$$\begin{aligned} \nabla F(\mathbf{X}_k) &= n [p_1 \mathbf{g}_{k,1}, \dots, p_n \mathbf{g}_{k,n}] \in \mathbb{R}^{d \times n} \\ \nabla \tilde{F}(\mathbf{X}_k) &= n [p_1 \tilde{\mathbf{g}}_{k,1}, \dots, p_n \tilde{\mathbf{g}}_{k,n}] \in \mathbb{R}^{d \times n} \\ \tilde{\mathbf{G}}_k &= [\dots, \tilde{\mathbf{g}}_{k, i_k}, \dots] \\ \mathbf{G}_k &= [\dots, \mathbf{g}_{k, i_k}, \dots] \\ \Lambda_a^b &= \frac{\mathbf{1}\mathbf{1}^\top}{n} - \prod_{q=a}^b \mathbf{W}_q \end{aligned}$$

H.2. Setting

The pseudo code can be found in Algorithm 2. We makes the following assumptions:

1. **Lipschitzian Gradient:** All the function f_i have L-Lipschitzian gradients.
2. **Communication Matrix**⁴: The communication matrix \mathbf{W}_k is doubly stochastic for any $k \geq 0$ and for any $b \geq a \geq 0$, there exists t_{mix} such that

$$\left\| \prod_{q=a}^b \mathbf{W}_q \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right) \right\|_1 \leq 2 \cdot 2^{-\lfloor \frac{b-a+1}{t_{\text{mix}}} \rfloor}$$

3. **Bounded Variance:**

$$\begin{aligned} \mathbb{E}_{\xi_i \sim \mathcal{D}_i} \left\| \nabla \tilde{f}_i(\mathbf{x}; \xi_i) - \nabla f_i(\mathbf{x}) \right\|^2 &\leq \sigma^2, \forall i \\ \mathbb{E}_{i \sim \{1, \dots, n\}} \left\| \nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x}) \right\|^2 &\leq \zeta^2, \forall i \end{aligned}$$

where $\nabla \tilde{f}_i(\mathbf{x}; \xi_i)$ denotes gradient sample on worker i computed via data sample ξ_i .

⁴Please refer to Section E for more details

4. **Bounded Staleness:** There exists T such that $\tau_k \leq T, \forall k$

5. **Gradient magnitude:** The norm of a sampled gradient is bounded by $\|\tilde{\mathbf{g}}_{k,i}\|_\infty \leq G_\infty$ for some constant G_∞ .

H.3. Proof to Theorem 5.

Proof. We start from

$$\begin{aligned}
 & \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \left(1 - \frac{2\alpha L}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\
 \stackrel{\text{Lemma H.4}}{\leq} & \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{(\sigma^2 + 6\zeta^2)\alpha L}{n} + \left(2L^2 + \frac{12\alpha L^3}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i\right) \right\|^2 \\
 & + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k})\mathbf{1}}{n} \right\|^2 \\
 \stackrel{\text{Lemma H.5}}{\leq} & \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{(\sigma^2 + 6\zeta^2)\alpha L}{n} + \frac{2\alpha^2 T^2 (\sigma^2 + 6\zeta^2) L^2}{n^2} + \frac{4\alpha^2 T^2 L^2}{n^2 K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \\
 & + \left(2L^2 + \frac{12\alpha L^3}{n} + \frac{24L^4 \alpha^2 T^2}{n^2}\right) \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i\right) \right\|^2 \\
 \stackrel{\text{Lemma H.3}}{\leq} & \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{(\sigma^2 + 6\zeta^2)\alpha L}{n} + \frac{2\alpha^2 T^2 (\sigma^2 + 6\zeta^2) L^2}{n^2} + \frac{4\alpha^2 T^2 L^2}{n^2 K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \\
 & + \frac{128\alpha^2 t_{\text{mix}}^2 L^2}{A_1} \left((\sigma^2 + 6\zeta^2)p + \frac{2p}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + G_\infty^2 d \right)
 \end{aligned}$$

where $A_1 = 1 - 192p\alpha^2 t_{\text{mix}}^2 L^2$ as defined in Lemma H.3.

Rearrange the terms, we get

$$\begin{aligned}
 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 & \leq \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{(\sigma^2 + 6\zeta^2)\alpha L}{n} + \frac{2\alpha^2 T^2 (\sigma^2 + 6\zeta^2) L^2}{n^2} \\
 & + \frac{128p\alpha^2 t_{\text{mix}}^2 L^2}{A_1} (\sigma^2 + 6\zeta^2) + \frac{128\alpha^2 t_{\text{mix}}^2 L^2}{A_1} G_\infty^2 d
 \end{aligned}$$

By setting $\alpha = \frac{n}{2L + \sqrt{K(\sigma^2 + 6\zeta^2)}}$

$$\begin{aligned}
 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 & \lesssim \frac{1}{K} + \frac{\sqrt{\sigma^2 + 6\zeta^2}}{\sqrt{K}} + \frac{pt_{\text{mix}}^2 (\sigma^2 + 6\zeta^2) n^2}{(\sigma^2 + 6\zeta^2) K + 4L^2} + \frac{n^2 t_{\text{mix}}^2 G_\infty^2 d}{(\sigma^2 + 6\zeta^2) K + 4L^2} \\
 & \lesssim \frac{1}{K} + \frac{\sqrt{\sigma^2 + 6\zeta^2}}{\sqrt{K}} + \frac{(\sigma^2 + 6\zeta^2) t_{\text{mix}}^2 n^2}{(\sigma^2 + 6\zeta^2) K + 1} + \frac{n^2 t_{\text{mix}}^2 G_\infty^2 d}{(\sigma^2 + 6\zeta^2) K + 1}
 \end{aligned}$$

□

H.4. Lemma for Moniqua on AD-PSGD

Lemma H.1.

$$\mathbb{E} \left\| \tilde{\mathbf{G}}_{k-\tau_k} \frac{\mathbf{1}}{n} \right\|^2 \leq \frac{\sigma^2}{n^2} + \frac{1}{n^2} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2, \forall k \geq 0.$$

Proof.

$$\begin{aligned}
 \mathbb{E} \left\| \tilde{\mathbf{G}}_{k-\tau_k} \frac{\mathbf{1}}{n} \right\|^2 &\leq \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\tilde{\mathbf{g}}_{k-\tau_k, i}}{n} \right\|^2 \\
 &= \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\tilde{\mathbf{g}}_{k-\tau_k, i} - \mathbf{g}_{k-\tau_k, i}}{n} \right\|^2 + \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\mathbf{g}_{k-\tau_k, i}}{n} \right\|^2 \\
 &\leq \frac{\sigma^2}{n^2} + \frac{1}{n^2} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} \right\|^2
 \end{aligned}$$

□

Lemma H.2.

$$\sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} \right\|^2 \leq 12L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 6\varsigma^2 + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2, \forall k \geq 0.$$

Proof.

$$\begin{aligned}
 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} \right\|^2 &= \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} - \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} + \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \\
 &\leq 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} - \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \\
 &= 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} - \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2
 \end{aligned}$$

And

$$\begin{aligned}
 &\sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} - \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \\
 &\leq 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i} - \nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) - \sum_{j=1}^n p_j \nabla f_j(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\quad + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} - \sum_{j=1}^n p_j \nabla f_j(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\leq 3L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{x}_{k-\tau_k, i} - \bar{\mathbf{X}}_{k-\tau_k} \right\|^2 + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) - \sum_{j=1}^n p_j \nabla f_j(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\quad + 3\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} - \sum_{j=1}^n p_j \nabla f_j(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\leq 3L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 3 \sum_{i=1}^n p_i \mathbb{E} \left\| \nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) - \nabla f(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\quad + 3 \sum_{j=1}^n p_j \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, j} - \nabla f_j(\bar{\mathbf{X}}_{k-\tau_k}) \right\|^2 \\
 &\leq 6L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 3\varsigma^2
 \end{aligned}$$

That completes the proof. \square

Lemma H.3. Let $A_1 = 1 - 192p\alpha^2 t_{\text{mix}}^2 L^2$,

$$\sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \leq \frac{32\alpha^2 t_{\text{mix}}^2}{A_1} \left((\sigma^2 + 6\varsigma^2) pK + 2p \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + G_\infty^2 dK \right)$$

Proof.

$$\begin{aligned} & \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_k \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\ &= \sum_{i=1}^n p_i \mathbb{E} \left\| \left(\mathbf{X}_{k-1} \mathbf{W}_{k-1} - \alpha \tilde{\mathbf{G}}_{k-1-\tau_{k-1}} + \boldsymbol{\Omega}_{k-1} \right) \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\ &\stackrel{X_0 = \mathbf{0}}{=} \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{t=0}^{k-1} \left(-\alpha \tilde{\mathbf{G}}_{t-\tau_t} + \boldsymbol{\Omega}_t \right) \boldsymbol{\Lambda}_{t+1}^{k-1} \mathbf{e}_i \right\|^2 \\ &\leq 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{t=0}^{k-1} \alpha \tilde{\mathbf{G}}_{t-\tau_t} \boldsymbol{\Lambda}_{t+1}^{k-1} \mathbf{e}_i \right\|^2 + 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{t=0}^{k-1} \boldsymbol{\Omega}_t \boldsymbol{\Lambda}_{t+1}^{k-1} \mathbf{e}_i \right\|^2 \end{aligned}$$

Now for the first term, we have

$$\begin{aligned} 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{t=0}^{k-1} \alpha \tilde{\mathbf{G}}_{t-\tau_t} \boldsymbol{\Lambda}_{t+1}^{k-1} \mathbf{e}_i \right\|^2 &\leq 2p\alpha^2 \mathbb{E} \left\| \sum_{t=0}^{k-1} \tilde{\mathbf{G}}_{t-\tau_t} \boldsymbol{\Lambda}_{t+1}^{k-1} \right\|_F^2 \\ &\leq 2p\alpha^2 \mathbb{E} \left(\sum_{t=0}^{k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F \left\| \boldsymbol{\Lambda}_{t+1}^{k-1} \right\| \right)^2 \\ &\leq 2p\alpha^2 \mathbb{E} \left(\sum_{t=0}^{k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F \left\| \boldsymbol{\Lambda}_{t+1}^{k-1} \right\|_1 \right)^2 \\ &\leq 8p\alpha^2 \mathbb{E} \left(\sum_{t=0}^{k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 \end{aligned}$$

Now we replace k with $k - \tau_k$, that is

$$\sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \leq 8p\alpha^2 \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 + 2 \sum_{i=1}^n p_i \mathbb{E} \left\| \sum_{t=0}^{k-\tau_k-1} \boldsymbol{\Omega}_t \boldsymbol{\Lambda}_{t+1}^{k-\tau_k-1} \mathbf{e}_i \right\|^2$$

Summing from $k = 0$ to $K - 1$ on both sides, we obtain

$$\begin{aligned} & \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\ &\leq 8p\alpha^2 \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 \\ &\quad + 2 \sum_{i=1}^n p_i \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{t=0}^{k-\tau_k-1} \boldsymbol{\Omega}_t \boldsymbol{\Lambda}_{t+1}^{k-\tau_k-1} \mathbf{e}_i \right\|^2 \end{aligned}$$

$$\begin{aligned}
 &\leq 8p\alpha^2 \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 \\
 &\quad + 2 \sum_{i=1}^n p_i \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \|\Omega_t\|_{1,2} \left\| \Lambda_{t+1}^{k-\tau_k-1} \right\|_1 \|e_i\|_1 \right)^2 \\
 &\leq 8p\alpha^2 \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 \\
 &\quad + 8 \sum_{i=1}^n p_i \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \|\Omega_t\|_{1,2} 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 \\
 &\stackrel{\text{Lemma H.6}}{\leq} 8p\alpha^2 \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 + 32t_{\text{mix}}^2 \sum_{i=1}^n p_i \sum_{k=0}^{K-1} \mathbb{E} \|\Omega_k\|_{1,2}^2 \\
 &\leq 8p\alpha^2 \sum_{k=0}^{K-1} \mathbb{E} \left(\sum_{t=0}^{k-\tau_k-1} \left\| \tilde{\mathbf{G}}_{t-\tau_t} \right\|_F 2^{-\lfloor \frac{k-\tau_k-t-1}{t_{\text{mix}}} \rfloor} \right)^2 + 128\delta^2 B_\theta^2 dt_{\text{mix}}^2 K \\
 &\stackrel{\text{Lemma H.6}}{\leq} 32p\alpha^2 t_{\text{mix}}^2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \tilde{\mathbf{G}}_{k-\tau_k} \right\|_F^2 + 128\delta^2 B_\theta^2 dt_{\text{mix}}^2 K
 \end{aligned}$$

Note that for the first term, we have

$$\begin{aligned}
 &\sum_{k=0}^{K-1} \mathbb{E} \left\| \tilde{\mathbf{G}}_{k-\tau_k} \right\|_F^2 \\
 &= \sum_{k=0}^{K-1} \mathbb{E} \left\| \tilde{\mathbf{g}}_{k-\tau_k, i_k} \right\|^2 \\
 &= \sum_{k=0}^{K-1} \mathbb{E} \left\| \tilde{\mathbf{g}}_{k-\tau_k, i_k} - \mathbf{g}_{k-\tau_k, i_k} \right\|^2 + \sum_{k=0}^{K-1} \mathbb{E} \left\| \mathbf{g}_{k-\tau_k, i_k} \right\|^2 \\
 &\leq \sigma^2 K + \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{t-\tau_t, i} \right\|^2 \\
 &\leq (\sigma^2 + 6\varsigma^2) K + 12L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2
 \end{aligned}$$

Putting these two terms back, we obtain

$$\begin{aligned}
 &\sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &\leq 32p\alpha^2 t_{\text{mix}}^2 \left((\sigma^2 + 6\varsigma^2) K + 12L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \right) \\
 &\quad + 128\delta^2 B_\theta^2 dt_{\text{mix}}^2 K
 \end{aligned}$$

Rearrange the terms, we obtain

$$\begin{aligned}
 &(1 - 192p\alpha^2 t_{\text{mix}}^2 L^2) \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &\leq 32p\alpha^2 t_{\text{mix}}^2 \left((\sigma^2 + 6\varsigma^2) K + 2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \right) + 128\delta^2 B_\theta^2 t_{\text{mix}}^2 K
 \end{aligned}$$

$$\stackrel{\text{Lemma H.7}}{\leq} 32\alpha^2 t_{\text{mix}}^2 \left((\sigma^2 + 6\varsigma^2)pK + 2p \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + G_\infty^2 dK \right)$$

Let $A_1 = 1 - 192p\alpha^2 t_{\text{mix}}^2 L^2$, we obtain

$$\sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \leq \frac{32\alpha^2 t_{\text{mix}}^2}{A_1} \left((\sigma^2 + 6\varsigma^2)pK + 2p \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 + G_\infty^2 dK \right)$$

□

Lemma H.4.

$$\begin{aligned} & \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \left(1 - \frac{2\alpha L}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ & \leq \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k})\mathbf{1}}{n} \right\|^2 \\ & \quad + \left(2L^2 + \frac{12\alpha L^3}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + \frac{(\sigma^2 + 6\varsigma^2)\alpha L}{n} \end{aligned}$$

Proof. We start from $f(\bar{\mathbf{X}}_{k+1})$ Since

$$\bar{\mathbf{X}}_{k+1} = \mathbf{X}_k \mathbf{W}_k \frac{\mathbf{1}}{n} + (\hat{\mathbf{X}}_k - \mathbf{X}_k)(\mathbf{W}_k - \mathbf{I}) \frac{\mathbf{1}}{n} - \alpha \bar{\mathbf{G}}_{k-\tau_k} = \bar{\mathbf{X}}_k - \alpha \bar{\mathbf{G}}_{k-\tau_k}$$

Then from Taylor Expansion, we have

$$\begin{aligned} & \mathbb{E} f(\bar{\mathbf{X}}_{k+1}) \\ & = \mathbb{E} f\left(\bar{\mathbf{X}}_k - \alpha \bar{\mathbf{G}}_{k-\tau_k}\right) \\ & \leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{G}}_{k-\tau_k} \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-\tau_k}\|^2 \\ & = \mathbb{E} f(\bar{\mathbf{X}}_k) - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{G}}_{k-\tau_k} \rangle - \alpha \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \bar{\mathbf{G}}_{k-\tau_k} - \bar{\mathbf{G}}_{k-\tau_k} \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \|\bar{\mathbf{G}}_{k-\tau_k}\|^2 \\ & = \mathbb{E} f(\bar{\mathbf{X}}_k) - \frac{\alpha}{n} \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \rangle + \frac{\alpha^2 L}{2} \mathbb{E} \left\| \frac{\tilde{\mathbf{g}}_{k-\tau_k, i_k}}{n} \right\|^2 \\ & \leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \frac{\alpha}{n} \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \rangle \\ & \quad + \frac{\alpha^2 L}{2} \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\tilde{\mathbf{g}}_{k-\tau_k, i_k} - \mathbf{g}_{k-\tau_k, i_k}}{n} \right\|^2 + \frac{\alpha^2 L}{2} \sum_{i=1}^n p_i \mathbb{E} \left\| \frac{\mathbf{g}_{k-\tau_k, i}}{n} \right\|^2 \\ & \leq \mathbb{E} f(\bar{\mathbf{X}}_k) - \frac{\alpha}{n} \mathbb{E} \langle \nabla f(\bar{\mathbf{X}}_k), \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \rangle + \frac{\alpha^2 L \sigma^2}{2n^2} + \frac{\alpha^2 L}{2n^2} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2 \\ & = \mathbb{E} f(\bar{\mathbf{X}}_k) + \frac{\alpha}{2n} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 - \frac{\alpha}{2n} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 - \frac{\alpha}{2n} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ & \quad + \frac{\alpha^2 L \sigma^2}{2n^2} + \frac{\alpha^2 L}{2n^2} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2 \end{aligned}$$

Rearrange these terms, we can get

$$\frac{\alpha}{2n} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \frac{\alpha}{2n} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2$$

$$\begin{aligned} &\leq \mathbb{E}f(\bar{\mathbf{X}}_k) - \mathbb{E}f(\bar{\mathbf{X}}_{k+1}) + \frac{\alpha}{2n} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \right\|^2 \\ &\quad + \frac{\alpha^2 L \sigma^2}{2n^2} + \frac{\alpha^2 L}{2n^2} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2 \end{aligned}$$

Summing over $k = 0$ to $K - 1$ on both sides, we can get

$$\begin{aligned} &\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ &\leq \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 + \frac{\alpha L \sigma^2}{n} + \frac{\alpha L}{nK} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2 \end{aligned}$$

For $\sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2$, we have

$$\begin{aligned} &\sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ &\leq 2 \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla f(\bar{\mathbf{X}}_{k-\tau_k})\|^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_{k-\tau_k}) - \nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ &= 2 \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla f(\bar{\mathbf{X}}_{k-\tau_k})\|^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i (\nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) - \mathbf{g}_{k-\tau_k, i}) \right\|^2 \\ &\leq 2 \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k) - \nabla f(\bar{\mathbf{X}}_{k-\tau_k})\|^2 + 2 \sum_{k=0}^{K-1} \mathbb{E} \sum_{i=1}^n p_i \|\nabla f_i(\bar{\mathbf{X}}_{k-\tau_k}) - \mathbf{g}_{k-\tau_k, i}\|^2 \\ &\leq 2L^2 \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \mathbf{1}}{n} \right\|^2 + 2L^2 \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - e_i \right) \right\|^2 \end{aligned}$$

Putting it back, we have

$$\begin{aligned} &\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla f(\bar{\mathbf{X}}_k)\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\nabla \bar{F}(\mathbf{X}_{k-\tau_k})\|^2 \\ &\leq \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \mathbf{1}}{n} \right\|^2 \\ &\quad + \frac{2L^2}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - e_i \right) \right\|^2 + \frac{\alpha L \sigma^2}{n} + \frac{\alpha L}{nK} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \|\mathbf{g}_{k-\tau_k, i}\|^2 \\ &\stackrel{\text{Lemma H.2}}{\leq} \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \mathbf{1}}{n} \right\|^2 \\ &\quad + \frac{2L^2}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - e_i \right) \right\|^2 + \frac{\alpha L \sigma^2}{n} \\ &\quad + \frac{\alpha L}{nK} \sum_{k=0}^{K-1} \left(12L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - e_i \right) \right\|^2 + 6\varsigma^2 + 2 \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 \right) \\ &= \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \mathbf{1}}{n} \right\|^2 \end{aligned}$$

$$\begin{aligned}
 & + \left(2L^2 + \frac{12\alpha L^3}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 & + \frac{(\sigma^2 + 6\varsigma^2)\alpha L}{n} + \frac{2\alpha L}{nK} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2
 \end{aligned}$$

Note that

$$\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2 = \mathbb{E} \left\| \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \right\|^2$$

Moving it to the left side, we finally get

$$\begin{aligned}
 & \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla f(\bar{\mathbf{X}}_k) \right\|^2 + \left(1 - \frac{2\alpha L}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \nabla \bar{F}(\mathbf{X}_{k-\tau_k}) \right\|^2 \\
 & \leq \frac{2n(f(\mathbf{0}) - f^*)}{\alpha K} + \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \frac{(\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \mathbf{1}}{n} \right\|^2 \\
 & + \left(2L^2 + \frac{12\alpha L^3}{n}\right) \frac{1}{K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + \frac{(\sigma^2 + 6\varsigma^2)\alpha L}{n}
 \end{aligned}$$

That completes the proof. \square

Lemma H.5. For all $k \geq 0$, we have

$$\begin{aligned}
 & \frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| (\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \frac{\mathbf{1}}{n} \right\|^2 \\
 & \leq \frac{2\alpha^2 T^2 (\sigma^2 + 6\varsigma^2) L^2}{n^2} + \frac{24L^4 \alpha^2 T^2}{n^2 K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 & + \frac{4\alpha^2 T^2 L^2}{n^2 K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k, i} \right\|^2
 \end{aligned}$$

Proof. From Lemma H.4, we know the fact

$$\bar{\mathbf{X}}_{k+1} = \mathbf{X}_k \mathbf{W}_k \frac{\mathbf{1}}{n} + (\hat{\mathbf{X}}_k - \mathbf{X}_k) (\mathbf{W}_k - \mathbf{I}) \frac{\mathbf{1}}{n} - \alpha \bar{\mathbf{G}}_{k-\tau_k} = \bar{\mathbf{X}}_k - \alpha \bar{\mathbf{G}}_{k-\tau_k}$$

As a result

$$\begin{aligned}
 & \sum_{k=0}^{K-1} \mathbb{E} \left\| (\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \frac{\mathbf{1}}{n} \right\|^2 \\
 & = \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{t=1}^{\tau_k} \alpha \tilde{\mathbf{G}}_{k-t} \frac{\mathbf{1}}{n} \right\|^2 \\
 & \leq \alpha^2 \sum_{k=0}^{K-1} \tau_k \sum_{t=1}^{\tau_k} \mathbb{E} \left\| \tilde{\mathbf{G}}_{k-t} \frac{\mathbf{1}}{n} \right\|^2 \\
 & \leq \alpha^2 \sum_{k=0}^{K-1} \tau_k \sum_{t=1}^{\tau_k} \left(\frac{\sigma^2}{n^2} + \frac{1}{n^2} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-t, i} \right\|^2 \right) \\
 & \leq \frac{\alpha^2 T^2 \sigma^2 K}{n^2} + \frac{\alpha^2 T}{n^2} \sum_{k=0}^{K-1} \sum_{t=1}^{\tau_k} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{g}_{k-t, i} \right\|^2
 \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{\alpha^2 T^2 \sigma^2 K}{n^2} + \frac{\alpha^2 T}{n^2} \sum_{k=0}^{K-1} \sum_{t=1}^{\tau_k} \left(12L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-t} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 6\zeta^2 + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-t,i} \right\|^2 \right) \\
 &\leq \frac{\alpha^2 T^2 \sigma^2 K}{n^2} + \frac{\alpha^2 T^2}{n^2} \sum_{k=0}^{K-1} \left(12L^2 \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 + 6\zeta^2 + 2\mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k,i} \right\|^2 \right) \\
 &= \frac{\alpha^2 T^2 (\sigma^2 + 6\zeta^2) K}{n^2} + \frac{12L^2 \alpha^2 T^2}{n^2} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &\quad + \frac{2\alpha^2 T^2}{n^2} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k,i} \right\|^2
 \end{aligned}$$

And we get

$$\begin{aligned}
 &\frac{2L^2}{K} \sum_{k=0}^{K-1} \mathbb{E} \left\| (\mathbf{X}_k - \mathbf{X}_{k-\tau_k}) \frac{\mathbf{1}}{n} \right\|^2 \\
 &\leq \frac{2\alpha^2 T^2 (\sigma^2 + 6\zeta^2) L^2}{n^2} + \frac{24L^4 \alpha^2 T^2}{n^2 K} \sum_{k=0}^{K-1} \sum_{i=1}^n p_i \mathbb{E} \left\| \mathbf{X}_{k-\tau_k} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|^2 \\
 &\quad + \frac{4\alpha^2 T^2 L^2}{n^2 K} \sum_{k=0}^{K-1} \mathbb{E} \left\| \sum_{i=1}^n p_i \mathbf{g}_{k-\tau_k,i} \right\|^2
 \end{aligned}$$

That completes the proof. \square

Lemma H.6. Given non-negative sequences $\{a_t\}_{t=1}^\infty$, $\{b_t\}_{t=1}^\infty$ and $\{\tau_t\}_{t=1}^\infty$ and a positive number T that satisfying

$$a_t = \sum_{s=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} b_s$$

with $0 \leq \rho < 1$, we have

$$\begin{aligned}
 S_k &= \sum_{t=1}^k a_t \leq \frac{(2-\rho)T}{1-\rho} \sum_{s=1}^k b_s \\
 D_k &= \sum_{t=1}^k a_t^2 \leq \frac{(2-\rho)T^2}{(1-\rho)^2} \sum_{s=1}^k b_s^2
 \end{aligned}$$

Proof.

$$\begin{aligned}
 S_k &= \sum_{t=1}^k a_t = \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} b_s \leq \sum_{t=1}^k \sum_{s=1}^t \rho^{\max(\lfloor \frac{t-\tau_t-s}{T} \rfloor, 0)} b_s = \sum_{s=1}^k \sum_{t=s}^k \rho^{\max(\lfloor \frac{t-\tau_t-s}{T} \rfloor, 0)} b_s \\
 &= \sum_{s=1}^k \sum_{t=0}^{k-\tau_k-s} \rho^{\lfloor \frac{t}{T} \rfloor} b_s + \sum_{s=1}^k \sum_{t=1}^{\tau_k} \rho^0 b_s \leq \sum_{s=1}^k \left(\sum_{t=0}^{T-1} \sum_{m=0}^\infty \rho^m \right) b_s + \tau_k \sum_{s=1}^k b_s \leq \left(T + \frac{T}{1-\rho} \right) \sum_{s=1}^k b_s \\
 D_k &= \sum_{t=1}^k a_t^2 = \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} b_s \sum_{r=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-r}{T} \rfloor} b_r = \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \sum_{r=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor + \lfloor \frac{t-\tau_t-r}{T} \rfloor} b_s b_r \\
 &\leq \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \sum_{r=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor + \lfloor \frac{t-\tau_t-r}{T} \rfloor} \frac{b_s^2 + b_r^2}{2} = \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \sum_{r=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor + \lfloor \frac{t-\tau_t-r}{T} \rfloor} b_s^2
 \end{aligned}$$

$$\begin{aligned}
 &\leq \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} b_s^2 \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} \sum_{r=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-r}{T} \rfloor} \leq \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} b_s^2 \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} \sum_{r=0}^{T-1} \sum_{m=0}^{\infty} \rho^m \\
 cs6 &\leq \frac{T}{1-\rho} \sum_{t=1}^k \sum_{s=1}^{t-\tau_t} \rho^{\lfloor \frac{t-\tau_t-s}{T} \rfloor} b_s^2 \stackrel{\text{Using } S_k}{\leq} \frac{(2-\rho)T^2}{(1-\rho)^2} \sum_{s=1}^k b_s^2
 \end{aligned}$$

□

Lemma H.7. for $\forall i, j$ and $\forall k \geq 0$, we have

$$\|\mathbf{X}_k(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} < \theta = 16t_{\text{mix}}\alpha G_{\infty}$$

Proof. We use mathematical induction to prove this.

I. First, for $k = 0$, we have

$$\|\mathbf{X}_k(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} = 0 < \theta = 16t_{\text{mix}}\alpha G_{\infty}$$

II. Suppose for $k \geq 0$, we have $\|\mathbf{X}_t(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} < \theta, \forall t \leq k$, then we have

$$\begin{aligned}
 &\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_{\infty} \\
 &\leq \left\| \mathbf{X}_{k+1} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|_{\infty} + \left\| \mathbf{X}_{k+1} \left(\frac{\mathbf{1}}{n} - \mathbf{e}_j \right) \right\|_{\infty} \\
 &\leq \left\| \mathbf{X}_{k+1} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right) \right\|_{1,\infty} \|\mathbf{e}_i\|_1 + \left\| \mathbf{X}_{k+1} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right) \right\|_{1,\infty} \|\mathbf{e}_j\|_1 \\
 &= 2 \left\| \mathbf{X}_{k+1} \left(\mathbf{I} - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right) \right\|_{1,\infty} \\
 &\leq 2 \left\| \left(\mathbf{X}_k \mathbf{W}_k - \alpha \tilde{\mathbf{G}}_{k-\tau_k} + \mathbf{\Omega}_k \right) \left(\frac{\mathbf{1}}{n} - \mathbf{e}_i \right) \right\|_{1,\infty} \\
 &= 2 \left\| \sum_{t=0}^k \left(-\alpha \tilde{\mathbf{G}}_{t-\tau_t} + \mathbf{\Omega}_t \right) \left(\prod_{q=t+1}^k \mathbf{W}_q - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right) \right\|_{1,\infty} \\
 &\leq 2 \sum_{t=0}^k \left\| \left(-\alpha \tilde{\mathbf{G}}_{t-\tau_t} + \mathbf{\Omega}_t \right) \left(\prod_{q=t+1}^k \mathbf{W}_q - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right) \right\|_{1,\infty} \\
 &\leq 2 \sum_{t=0}^k \left\| -\alpha \tilde{\mathbf{G}}_{t-\tau_t} + \mathbf{\Omega}_t \right\|_{1,\infty} \left\| \prod_{q=t+1}^k \mathbf{W}_q - \frac{\mathbf{1}\mathbf{1}^{\top}}{n} \right\|_1 \\
 &\leq 4(\alpha G_{\infty} + 2\delta B_{\theta}) \sum_{t=0}^k 2^{-\lfloor (k-t)/t_{\text{mix}} \rfloor} \\
 &< 4(\alpha G_{\infty} + 2\delta B_{\theta}) \sum_{t=0}^{t_{\text{mix}}-1} \sum_{r=0}^{\infty} 2^{-r} \\
 &\leq 8(\alpha G_{\infty} + 2\delta B_{\theta}) t_{\text{mix}}
 \end{aligned}$$

Put in $\delta = \frac{1}{64t_{\text{mix}}+2}$, we obtain

$$\|\mathbf{X}_{k+1}(\mathbf{e}_i - \mathbf{e}_j)\|_2 < 8(\alpha G_{\infty} + 2\delta B_{\theta}) t_{\text{mix}} = 8t_{\text{mix}}\alpha G_{\infty} + 8t_{\text{mix}}\alpha G_{\infty} = 16t_{\text{mix}}\alpha G_{\infty}$$

Combining I and II and we complete the proof. □