

## A. Proof of Lemma 2

We recall the standard Descent Lemma (Nesterov, 2018), i.e.,  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ ,

$$F(\mathbf{y}) \leq F(\mathbf{x}) + \langle \nabla F(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \quad (15)$$

since the global function  $F$  is  $L$ -smooth. Setting  $\mathbf{y} = \bar{\mathbf{x}}_{t+1}$  and  $\mathbf{x} = \bar{\mathbf{x}}_t$  in (15) and using (5), we have:  $\forall t \geq 0$ ,

$$\begin{aligned} F(\bar{\mathbf{x}}_{t+1}) &\leq F(\bar{\mathbf{x}}_t) - \langle \nabla F(\bar{\mathbf{x}}_t), \bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t \rangle + \frac{L}{2} \|\bar{\mathbf{x}}_{t+1} - \bar{\mathbf{x}}_t\|^2 \\ &\leq F(\bar{\mathbf{x}}_t) - \alpha \langle \nabla F(\bar{\mathbf{x}}_t), \bar{\mathbf{v}}_t \rangle + \frac{L\alpha^2}{2} \|\bar{\mathbf{v}}_t\|^2. \end{aligned} \quad (16)$$

Using  $\langle \mathbf{a}, \mathbf{b} \rangle = \frac{1}{2} (\|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2)$ ,  $\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ , in (16) gives: for  $0 < \alpha \leq \frac{1}{2L}$  and  $\forall t \geq 0$ ,

$$\begin{aligned} F(\bar{\mathbf{x}}_{t+1}) &\leq F(\bar{\mathbf{x}}_t) - \frac{\alpha}{2} \|\nabla F(\bar{\mathbf{x}}_t)\|^2 - \left( \frac{\alpha}{2} - \frac{L\alpha^2}{2} \right) \|\bar{\mathbf{v}}_t\|^2 + \frac{\alpha}{2} \|\bar{\mathbf{v}}_t - \nabla F(\bar{\mathbf{x}}_t)\|^2, \\ &\leq F(\bar{\mathbf{x}}_t) - \frac{\alpha}{2} \|\nabla F(\bar{\mathbf{x}}_t)\|^2 - \left( \frac{\alpha}{2} - \frac{L\alpha^2}{2} \right) \|\bar{\mathbf{v}}_t\|^2 + \alpha \|\bar{\mathbf{v}}_t - \bar{\nabla \mathbf{f}}(\mathbf{x}_t)\|^2 + \alpha \|\bar{\nabla \mathbf{f}}(\mathbf{x}_t) - \nabla F(\bar{\mathbf{x}}_t)\|^2, \\ &\stackrel{(i)}{\leq} F(\bar{\mathbf{x}}_t) - \frac{\alpha}{2} \|\nabla F(\bar{\mathbf{x}}_t)\|^2 - \frac{\alpha}{4} \|\bar{\mathbf{v}}_t\|^2 + \alpha \|\bar{\mathbf{v}}_t - \bar{\nabla \mathbf{f}}(\mathbf{x}_t)\|^2 + \frac{\alpha L^2}{n} \|\mathbf{x}_t - \mathbf{Jx}_t\|^2, \end{aligned} \quad (17)$$

where (i) is due to Lemma 1(c) and that  $\frac{L\alpha^2}{2} \leq \frac{\alpha}{4}$  since  $0 < \alpha \leq \frac{1}{2L}$ . Rearranging (17), we have: for  $0 < \alpha \leq \frac{1}{2L}$  and  $\forall t \geq 0$ ,

$$\|\nabla F(\bar{\mathbf{x}}_t)\|^2 \leq \frac{2(F(\bar{\mathbf{x}}_t) - F(\bar{\mathbf{x}}_{t+1}))}{\alpha} - \frac{1}{2} \|\bar{\mathbf{v}}_t\|^2 + 2 \|\bar{\mathbf{v}}_t - \bar{\nabla \mathbf{f}}(\mathbf{x}_t)\|^2 + \frac{2L^2}{n} \|\mathbf{x}_t - \mathbf{Jx}_t\|^2. \quad (18)$$

Taking the telescoping sum of (18) over  $t$  from 0 to  $T$ ,  $\forall T \geq 0$  and using the fact that  $F$  bounded below by  $F^*$  in the resulting inequality finishes the proof.

## B. Proof of Lemma 3

### B.1. Proof of Eq. (7)

We recall that the update of each local stochastic gradient estimator  $\mathbf{v}_t^i$ ,  $\forall t \geq 1$ , in (2) may be written equivalently as follows:

$$\mathbf{v}_t^i = \beta \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) + (1 - \beta) \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \mathbf{v}_{t-1}^i \right),$$

where  $\beta \in (0, 1)$ . We have:  $\forall t \geq 1$  and  $\forall i \in \mathcal{V}$ ,

$$\begin{aligned} \mathbf{v}_t^i - \nabla f_i(\mathbf{x}_t^i) &= \beta \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) + (1 - \beta) \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \mathbf{v}_{t-1}^i \right) - \beta \nabla f_i(\mathbf{x}_t^i) - (1 - \beta) \nabla f_i(\mathbf{x}_t^i) \\ &= \beta \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i) \right) + (1 - \beta) \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_t^i) \right) \\ &= \beta \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i) \right) + (1 - \beta) \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \nabla f_i(\mathbf{x}_{t-1}^i) - \nabla f_i(\mathbf{x}_t^i) \right) \\ &\quad + (1 - \beta) \left( \mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i) \right). \end{aligned} \quad (19)$$

In (19), we observe that  $\forall t \geq 1$  and  $\forall i \in \mathcal{V}$ ,

$$\mathbb{E} \left[ \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i) | \mathcal{F}_t \right] = \mathbf{0}_p, \quad (20)$$

$$\mathbb{E} \left[ \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \nabla f_i(\mathbf{x}_{t-1}^i) - \nabla f_i(\mathbf{x}_t^i) | \mathcal{F}_t \right] = \mathbf{0}_p, \quad (21)$$

by the definition of the filtration  $\mathcal{F}_t$  in (1). Averaging (19) over  $i$  from 1 to  $n$  gives:  $\forall t \geq 0$ ,

$$\begin{aligned} \bar{\mathbf{v}}_t - \bar{\nabla f}(\mathbf{x}_t) &= (1 - \beta) \left( \bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1}) \right) + \underbrace{\beta \cdot \frac{1}{n} \sum_{i=1}^n \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i) \right)}_{=: \mathbf{s}_t} \\ &\quad + (1 - \beta) \cdot \underbrace{\frac{1}{n} \sum_{i=1}^n \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \nabla f_i(\mathbf{x}_{t-1}^i) - \nabla f_i(\mathbf{x}_t^i) \right)}_{=: \mathbf{z}_t}. \end{aligned} \quad (22)$$

Note that  $\mathbb{E}[\mathbf{s}_t | \mathcal{F}_t] = \mathbb{E}[\mathbf{z}_t | \mathcal{F}_t] = \mathbf{0}_p$  by (20) and (21). In light of (22), we have:  $\forall t \geq 1$ ,

$$\begin{aligned} \mathbb{E} \left[ \|\bar{\mathbf{v}}_t - \bar{\nabla f}(\mathbf{x}_t)\|^2 | \mathcal{F}_t \right] &= (1 - \beta)^2 \|\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1})\|^2 + \mathbb{E} \left[ \|\beta \mathbf{s}_t + (1 - \beta) \mathbf{z}_t\|^2 | \mathcal{F}_t \right] \\ &\quad + 2\mathbb{E} \left[ \left\langle (1 - \beta) (\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1})), \beta \mathbf{s}_t + (1 - \beta) \mathbf{z}_t \right\rangle | \mathcal{F}_t \right] \\ &\stackrel{(i)}{=} (1 - \beta)^2 \|\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1})\|^2 + \mathbb{E} \left[ \|\beta \mathbf{s}_t + (1 - \beta) \mathbf{z}_t\|^2 | \mathcal{F}_t \right] \\ &\leq (1 - \beta)^2 \|\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1})\|^2 + 2\beta^2 \mathbb{E} \left[ \|\mathbf{s}_t\|^2 | \mathcal{F}_t \right] + 2(1 - \beta)^2 \mathbb{E} \left[ \|\mathbf{z}_t\|^2 | \mathcal{F}_t \right], \end{aligned} \quad (23)$$

where (i) is due to

$$\mathbb{E} \left[ \left\langle (1 - \beta) (\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1})), \beta \mathbf{s}_t + (1 - \beta) \mathbf{z}_t \right\rangle | \mathcal{F}_t \right] = 0,$$

since  $\mathbb{E}[\mathbf{s}_t | \mathcal{F}_t] = \mathbb{E}[\mathbf{z}_t | \mathcal{F}_t] = \mathbf{0}_p$  and  $(\bar{\mathbf{v}}_{t-1} - \bar{\nabla f}(\mathbf{x}_{t-1}))$  is  $\mathcal{F}_t$ -measurable. We next bound the second and the third term in (23) respectively. For the second term in (23), we observe that  $\forall t \geq 1$ ,

$$\begin{aligned} \mathbb{E} [\|\mathbf{s}_t\|^2] &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i)\|^2 \right] + \frac{1}{n^2} \sum_{i \neq j} \mathbb{E} \left[ \left\langle \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i), \mathbf{g}_j(\mathbf{x}_t^j, \boldsymbol{\xi}_t^j) - \nabla f_j(\mathbf{x}_t^j) \right\rangle \right] \\ &\stackrel{(i)}{=} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i)\|^2 \right] \\ &\leq \frac{\bar{\nu}^2}{n}. \end{aligned} \quad (24)$$

We note that (i) in (24) uses that whenever  $i \neq j$ ,

$$\begin{aligned} &\mathbb{E} \left[ \left\langle \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i), \mathbf{g}_j(\mathbf{x}_t^j, \boldsymbol{\xi}_t^j) - \nabla f_j(\mathbf{x}_t^j) \right\rangle | \mathcal{F}_t \right] \\ &\stackrel{(ii)}{=} \mathbb{E} \left[ \left\langle \mathbb{E} [\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) | \sigma(\boldsymbol{\xi}_t^j, \mathcal{F}_t)] - \nabla f_i(\mathbf{x}_t^i), \mathbf{g}_j(\mathbf{x}_t^j, \boldsymbol{\xi}_t^j) - \nabla f_j(\mathbf{x}_t^j) \right\rangle | \mathcal{F}_t \right] \\ &\stackrel{(iii)}{=} \mathbb{E} \left[ \left\langle \mathbb{E} [\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) | \mathcal{F}_t] - \nabla f_i(\mathbf{x}_t^i), \mathbf{g}_j(\mathbf{x}_t^j, \boldsymbol{\xi}_t^j) - \nabla f_j(\mathbf{x}_t^j) \right\rangle | \mathcal{F}_t \right] \\ &= 0, \end{aligned} \quad (25)$$

where (ii) is due to the tower property of the conditional expectation and (iii) uses that  $\boldsymbol{\xi}_t^j$  is independent of  $\{\boldsymbol{\xi}_t^i, \mathcal{F}_t\}$  and  $\mathbf{x}_t^i$  is  $\mathcal{F}_t$ -measurable. Towards the third term (23), we define for the ease of exposition,  $\forall t \geq 1$ ,

$$\hat{\nabla}_t^i := \nabla f_i(\mathbf{x}_t^i) - \nabla f_i(\mathbf{x}_{t-1}^i)$$

and recall that  $\mathbb{E} [\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) | \mathcal{F}_t] = \widehat{\nabla}_t^i$ . Observe that  $\forall t \geq 1$ ,

$$\begin{aligned}
\mathbb{E} [\|\mathbf{z}_t\|^2 | \mathcal{F}_t] &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \widehat{\nabla}_t^i) \right\|^2 | \mathcal{F}_t \right] \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \widehat{\nabla}_t^i \right\|^2 | \mathcal{F}_t \right] \\
&\quad + \underbrace{\frac{1}{n^2} \sum_{i \neq j} \mathbb{E} \left[ \langle \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \widehat{\nabla}_t^i, \mathbf{g}_j(\mathbf{x}_t^j, \boldsymbol{\xi}_t^j) - \mathbf{g}_j(\mathbf{x}_{t-1}^j, \boldsymbol{\xi}_t^j) - \widehat{\nabla}_t^j \rangle | \mathcal{F}_t \right]}_{=0} \\
&\stackrel{(i)}{=} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \widehat{\nabla}_t^i \right\|^2 | \mathcal{F}_t \right], \\
&\stackrel{(ii)}{\leq} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) \right\|^2 | \mathcal{F}_t \right], \tag{26}
\end{aligned}$$

where (i) follows from a similar line of arguments as (25) and (ii) uses the conditional variance decomposition, i.e., for any random vector  $\mathbf{a} \in \mathbb{R}^p$  consisted of square-integrable random variables,

$$\mathbb{E} \left[ \left\| \mathbf{a} - \mathbb{E} [\mathbf{a} | \mathcal{F}_t] \right\|^2 | \mathcal{F}_t \right] = \mathbb{E} \left[ \left\| \mathbf{a} \right\|^2 | \mathcal{F}_t \right] - \left\| \mathbb{E} [\mathbf{a} | \mathcal{F}_t] \right\|^2. \tag{27}$$

To proceed from (26), we take its expectation and observe that  $\forall t \geq 1$ ,

$$\begin{aligned}
\mathbb{E} [\|\mathbf{z}_t\|^2] &\leq \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) \right\|^2 \right] \\
&\stackrel{(i)}{\leq} \frac{L^2}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{x}_t^i - \mathbf{x}_{t-1}^i \right\|^2 \right] \\
&= \frac{L^2}{n^2} \mathbb{E} \left[ \left\| \mathbf{x}_t - \mathbf{x}_{t-1} \right\|^2 \right] \\
&= \frac{L^2}{n^2} \mathbb{E} \left[ \left\| \mathbf{x}_t - \mathbf{Jx}_t + \mathbf{Jx}_t - \mathbf{Jx}_{t-1} + \mathbf{Jx}_{t-1} - \mathbf{x}_{t-1} \right\|^2 \right] \\
&\leq \frac{3L^2}{n^2} \mathbb{E} \left[ \left\| \mathbf{x}_t - \mathbf{Jx}_t \right\|^2 + n \left\| \bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-1} \right\|^2 + \left\| \mathbf{x}_{t-1} - \mathbf{Jx}_{t-1} \right\|^2 \right] \\
&\stackrel{(ii)}{=} \frac{3L^2\alpha^2}{n} \mathbb{E} \left[ \left\| \bar{\mathbf{v}}_{t-1} \right\|^2 \right] + \frac{3L^2}{n^2} \left( \mathbb{E} \left[ \left\| \mathbf{x}_t - \mathbf{Jx}_t \right\|^2 + \left\| \mathbf{x}_{t-1} - \mathbf{Jx}_{t-1} \right\|^2 \right] \right), \tag{28}
\end{aligned}$$

where (i) uses the mean-squared smoothness of each  $\mathbf{g}_i$  and (ii) uses the update of  $\bar{\mathbf{x}}_t$  in (5). The proof follows by taking the expectation (23) and then using (24) and (28) in the resulting inequality.

## B.2. Proof of Eq. (8)

We recall from (19) the following relationship:  $\forall t \geq 1$ ,

$$\begin{aligned}
\mathbf{v}_t^i - \nabla f_i(\mathbf{x}_t^i) &= \beta \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i) \right) + (1 - \beta) \left( \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) + \nabla f_i(\mathbf{x}_{t-1}^i) - \nabla f_i(\mathbf{x}_t^i) \right) \\
&\quad + (1 - \beta) \left( \mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i) \right). \tag{29}
\end{aligned}$$

We recall that the conditional expectation of the first and the second term in (29) with respect to  $\mathcal{F}_t$  is zero and that the third term in (29) is  $\mathcal{F}_t$ -measurable. Following a similar procedure in the proof of (23), we have:  $\forall t \geq 1$ ,

$$\begin{aligned} \mathbb{E} [\|\mathbf{v}_t^i - \nabla f_i(\mathbf{x}_t^i)\|^2 | \mathcal{F}_t] &\leq (1-\beta)^2 \|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2 + 2\beta^2 \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i)\|^2 | \mathcal{F}_t] \\ &\quad + 2(1-\beta)^2 \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - (\nabla f_i(\mathbf{x}_t^i) - \nabla f_i(\mathbf{x}_{t-1}^i))\|^2 | \mathcal{F}_t] \\ &\stackrel{(i)}{\leq} (1-\beta)^2 \|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2 + 2\beta^2 \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_t^i)\|^2 | \mathcal{F}_t] \\ &\quad + 2(1-\beta)^2 \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i)\|^2 | \mathcal{F}_t] \end{aligned} \quad (30)$$

where (i) uses the conditional variance decomposition (27). We then take the expectation of (30) with the help of the mean-squared smoothness and the bounded variance of each  $\mathbf{g}_i$  to proceed:  $\forall t \geq 1$ ,

$$\begin{aligned} \mathbb{E} [\|\mathbf{v}_t^i - \nabla f_i(\mathbf{x}_t^i)\|^2] &\leq (1-\beta)^2 \mathbb{E} [\|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2] + 2\beta^2 \nu_i^2 + 2(1-\beta)^2 L^2 \mathbb{E} [\|\mathbf{x}_t^i - \mathbf{x}_{t-1}^i\|^2] \\ &\leq (1-\beta)^2 \mathbb{E} [\|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2] + 2\beta^2 \nu_i^2 \\ &\quad + 6(1-\beta)^2 L^2 \left( \mathbb{E} [\|\mathbf{x}_t^i - \bar{\mathbf{x}}_t\|^2 + \|\bar{\mathbf{x}}_t - \bar{\mathbf{x}}_{t-1}\|^2 + \|\bar{\mathbf{x}}_{t-1} - \mathbf{x}_{t-1}^i\|^2] \right), \\ &= (1-\beta)^2 \mathbb{E} [\|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2] + 2\beta^2 \nu_i^2 + 6(1-\beta)^2 L^2 \alpha^2 \mathbb{E} [\|\bar{\mathbf{x}}_{t-1}\|^2] \\ &\quad + 6(1-\beta)^2 L^2 \mathbb{E} [\|\mathbf{x}_t^i - \bar{\mathbf{x}}_t\|^2 + \|\mathbf{x}_{t-1}^i - \bar{\mathbf{x}}_{t-1}\|^2], \end{aligned} \quad (31)$$

where the last line uses the  $\bar{\mathbf{x}}_t$ -update in (5). Summing up (31) over  $i$  from 1 to  $n$  completes the proof.

## C. Proof of Lemma 5

### C.1. Proof of Lemma 5(a)

Recall the initialization of GT-HSGD that  $\mathbf{v}_{-1} = \mathbf{0}_{np}$ ,  $\mathbf{y}_0 = \mathbf{0}_{np}$ , and  $\mathbf{v}_0^i = \frac{1}{b_0} \sum_{r=1}^{b_0} \mathbf{g}_i(\mathbf{x}_0^i, \boldsymbol{\xi}_{0,r}^i)$ . Using the gradient tracking update (4a) at iteration  $t = 0$ , we have:

$$\begin{aligned} \mathbb{E} [\|\mathbf{y}_1 - \mathbf{J}\mathbf{y}_1\|^2] &= \mathbb{E} [\|\mathbf{W}(\mathbf{y}_0 + \mathbf{v}_0 - \mathbf{v}_{-1}) - \mathbf{J}\mathbf{W}(\mathbf{y}_0 + \mathbf{v}_0 - \mathbf{v}_{-1})\|^2] \\ &\stackrel{(i)}{=} \mathbb{E} [\|(\mathbf{W} - \mathbf{J})\mathbf{v}_0\|^2] \\ &\stackrel{(ii)}{\leq} \lambda^2 \mathbb{E} [\|\mathbf{v}_0 - \nabla \mathbf{f}(\mathbf{x}_0) + \nabla \mathbf{f}(\mathbf{x}_0)\|^2] \\ &= \lambda^2 \sum_{i=1}^n \mathbb{E} [\|\mathbf{v}_0^i - \nabla f_i(\mathbf{x}_0^i)\|^2] + \lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2 \\ &\stackrel{(iii)}{=} \lambda^2 \sum_{i=1}^n \mathbb{E} \left[ \left\| \frac{1}{b_0} \sum_{r=1}^{b_0} (\mathbf{g}_i(\mathbf{x}_0^i, \boldsymbol{\xi}_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)) \right\|^2 \right] + \lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2 \\ &\stackrel{(iv)}{=} \frac{\lambda^2}{b_0^2} \sum_{i=1}^n \sum_{r=1}^{b_0} \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_0^i, \boldsymbol{\xi}_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)\|^2] + \lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2, \end{aligned} \quad (32)$$

where (i) uses  $\mathbf{J}\mathbf{W} = \mathbf{J}$  and the initial condition of  $\mathbf{v}_{-1}$  and  $\mathbf{y}_0$ , (ii) uses  $\|\mathbf{W} - \mathbf{J}\| = \lambda$ , (iii) is due to the initialization of  $\mathbf{v}_0^i$ , and (iv) follows from the fact that  $\{\boldsymbol{\xi}_{0,1}^i, \boldsymbol{\xi}_{0,2}^i, \dots, \boldsymbol{\xi}_{0,b_0}^i\}$ ,  $\forall i \in \mathcal{V}$ , is an independent family of random vectors, by a similar line of arguments in (24) and (25). The proof then follows by using the bounded variance of each  $\mathbf{g}_i$  in (32).

## C.2. Proof of Lemma 5(b)

Following the gradient tracking update (4a), we have:  $\forall t \geq 1$ ,

$$\begin{aligned} \|\mathbf{y}_{t+1} - \mathbf{J}\mathbf{y}_{t+1}\|^2 &= \|\mathbf{W}(\mathbf{y}_t + \mathbf{v}_t - \mathbf{v}_{t-1}) - \mathbf{J}\mathbf{W}(\mathbf{y}_t + \mathbf{v}_t - \mathbf{v}_{t-1})\|^2 \\ &\stackrel{(i)}{=} \|\mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t + (\mathbf{W} - \mathbf{J})(\mathbf{v}_t - \mathbf{v}_{t-1})\|^2 \\ &= \|\mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2 + 2\langle \mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t, (\mathbf{W} - \mathbf{J})(\mathbf{v}_t - \mathbf{v}_{t-1}) \rangle + \|(\mathbf{W} - \mathbf{J})(\mathbf{v}_t - \mathbf{v}_{t-1})\|^2 \\ &\stackrel{(ii)}{\leq} \lambda^2 \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2 + \underbrace{2\langle \mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t, (\mathbf{W} - \mathbf{J})(\mathbf{v}_t - \mathbf{v}_{t-1}) \rangle}_{=: A_t} + \lambda^2 \|\mathbf{v}_t - \mathbf{v}_{t-1}\|^2, \end{aligned} \quad (33)$$

where (i) uses  $\mathbf{J}\mathbf{W} = \mathbf{J}$  and (ii) is due to  $\|\mathbf{W} - \mathbf{J}\| = \lambda$ . In the following, we bound  $A_t$  and the last term in (33) respectively. We recall the update of each local stochastic gradient estimator  $\mathbf{v}_t^i$  in (2):  $\forall t \geq 1$ ,

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

We observe that  $\forall t \geq 1$  and  $\forall i \in \mathcal{V}$ ,

$$\begin{aligned} \mathbf{v}_t^i - \mathbf{v}_{t-1}^i &= \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \beta\mathbf{v}_{t-1}^i - (1 - \beta)\mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) \\ &= \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \beta\mathbf{v}_{t-1}^i + \beta\mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) \\ &= \mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \beta(\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)) + \beta(\mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_{t-1}^i)). \end{aligned} \quad (34)$$

Moreover, we observe from (34) that  $\forall t \geq 1$ ,

$$\mathbb{E}[\mathbf{v}_t - \mathbf{v}_{t-1} | \mathcal{F}_t] = \nabla \mathbf{f}(\mathbf{x}_t) - \nabla \mathbf{f}(\mathbf{x}_{t-1}) - \beta(\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})). \quad (35)$$

Towards  $A_t$ , we have:  $\forall t \geq 1$ ,

$$\begin{aligned} \mathbb{E}[A_t | \mathcal{F}_t] &\stackrel{(i)}{=} 2\left\langle \mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t, (\mathbf{W} - \mathbf{J})\mathbb{E}[\mathbf{v}_t - \mathbf{v}_{t-1} | \mathcal{F}_t] \right\rangle \\ &\stackrel{(ii)}{=} 2\left\langle \mathbf{W}\mathbf{y}_t - \mathbf{J}\mathbf{y}_t, (\mathbf{W} - \mathbf{J})\left(\nabla \mathbf{f}(\mathbf{x}_t) - \nabla \mathbf{f}(\mathbf{x}_{t-1}) - \beta(\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1}))\right) \right\rangle \\ &\stackrel{(iii)}{\leq} 2\lambda \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\| \cdot \lambda \left\| \nabla \mathbf{f}(\mathbf{x}_t) - \nabla \mathbf{f}(\mathbf{x}_{t-1}) - \beta(\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})) \right\| \\ &\stackrel{(iv)}{\leq} \frac{1 - \lambda^2}{2} \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2 + \frac{2\lambda^4}{1 - \lambda^2} \left\| \nabla \mathbf{f}(\mathbf{x}_t) - \nabla \mathbf{f}(\mathbf{x}_{t-1}) - \beta(\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})) \right\|^2, \\ &\stackrel{(v)}{\leq} \frac{1 - \lambda^2}{2} \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2 + \frac{4\lambda^4 L^2}{1 - \lambda^2} \|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2 + \frac{4\lambda^4 \beta^2}{1 - \lambda^2} \|\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})\|^2, \end{aligned} \quad (36)$$

where (i) is due to the  $\mathcal{F}_t$ -measurability of  $\mathbf{y}_t$ , (ii) uses (35), (iii) is due to the Cauchy-Schwarz inequality and  $\|\mathbf{W} - \mathbf{J}\| = \lambda$ , (iv) uses the elementary inequality that  $2ab \leq \eta a^2 + b^2/\eta$ , with  $\eta = \frac{1-\lambda^2}{2\lambda^2}$  for any  $a, b \in \mathbb{R}$ , and (v) holds since each  $f_i$  is  $L$ -smooth. Next, towards the last term in (33), we take the expectation of (34) to obtain:  $\forall t \geq 1$  and  $\forall i \in \mathcal{V}$ ,

$$\begin{aligned} \mathbb{E}\left[\|\mathbf{v}_t^i - \mathbf{v}_{t-1}^i\|^2\right] &\leq 3\mathbb{E}\left[\|\mathbf{g}_i(\mathbf{x}_t^i, \boldsymbol{\xi}_t^i) - \mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i)\|^2\right] + 3\beta^2\mathbb{E}\left[\|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2\right] \\ &\quad + 3\beta^2\mathbb{E}\left[\|\mathbf{g}_i(\mathbf{x}_{t-1}^i, \boldsymbol{\xi}_t^i) - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2\right] \\ &\leq 3L^2\mathbb{E}\left[\|\mathbf{x}_t^i - \mathbf{x}_{t-1}^i\|^2\right] + 3\beta^2\mathbb{E}\left[\|\mathbf{v}_{t-1}^i - \nabla f_i(\mathbf{x}_{t-1}^i)\|^2\right] + 3\beta^2\nu_i^2, \end{aligned} \quad (37)$$

where (37) is due to the mean-squared smoothness and the bounded variance of each  $\mathbf{g}_i$ . Summing up (37) over  $i$  from 1 to  $n$  gives an upper bound on the last term in (33):  $\forall t \geq 1$ ,

$$\lambda^2\mathbb{E}\left[\|\mathbf{v}_t - \mathbf{v}_{t-1}\|^2\right] \leq 3\lambda^2 L^2\mathbb{E}\left[\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2\right] + 3\lambda^2\beta^2\mathbb{E}\left[\|\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})\|^2\right] + 3\lambda^2 n\beta^2\nu^2. \quad (38)$$

We now use (36) and (38) in (33) to obtain:  $\forall t \geq 1$ ,

$$\begin{aligned}\mathbb{E} [\|\mathbf{y}_{t+1} - \mathbf{J}\mathbf{y}_{t+1}\|^2] &\leq \frac{1+\lambda^2}{2} \mathbb{E} [\|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2] + \frac{7\lambda^2 L^2}{1-\lambda^2} \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2] \\ &\quad + \frac{7\lambda^2 \beta^2}{1-\lambda^2} \mathbb{E} [\|\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})\|^2] + 3\lambda^2 n \beta^2 \bar{\nu}^2.\end{aligned}\quad (39)$$

Towards the second term in (39), we use (10) to obtain:  $\forall t \geq 1$ ,

$$\begin{aligned}\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2 &= \|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t + \mathbf{J}\mathbf{x}_t - \mathbf{J}\mathbf{x}_{t-1} + \mathbf{J}\mathbf{x}_{t-1} - \mathbf{x}_{t-1}\|^2 \\ &\stackrel{(i)}{\leq} 3\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2 + 3n\alpha^2 \|\bar{\mathbf{v}}_{t-1}\|^2 + 3\|\mathbf{x}_{t-1} - \mathbf{J}\mathbf{x}_{t-1}\|^2 \\ &\leq 6\lambda^2 \alpha^2 \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2 + 3n\alpha^2 \|\bar{\mathbf{v}}_{t-1}\|^2 + 9\|\mathbf{x}_{t-1} - \mathbf{J}\mathbf{x}_{t-1}\|^2,\end{aligned}\quad (40)$$

where (i) uses the  $\bar{\mathbf{x}}_t$ -update in (5). Finally, we use (40) in (39) to obtain:  $\forall t \geq 1$ ,

$$\begin{aligned}\mathbb{E} [\|\mathbf{y}_{t+1} - \mathbf{J}\mathbf{y}_{t+1}\|^2] &\leq \left( \frac{1+\lambda^2}{2} + \frac{42\lambda^4 L^2 \alpha^2}{1-\lambda^2} \right) \mathbb{E} [\|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2] + \frac{21\lambda^2 n L^2 \alpha^2}{1-\lambda^2} \mathbb{E} [\|\bar{\mathbf{v}}_{t-1}\|^2] \\ &\quad + \frac{63\lambda^2 L^2}{1-\lambda^2} \mathbb{E} [\|\mathbf{x}_{t-1} - \mathbf{J}\mathbf{x}_{t-1}\|^2] + \frac{7\lambda^2 \beta^2}{1-\lambda^2} \mathbb{E} [\|\mathbf{v}_{t-1} - \nabla \mathbf{f}(\mathbf{x}_{t-1})\|^2] + 3\lambda^2 n \beta^2 \bar{\nu}^2.\end{aligned}$$

The proof is completed by the fact that  $\frac{1+\lambda^2}{2} + \frac{42\lambda^4 L^2 \alpha^2}{1-\lambda^2} \leq \frac{3+\lambda^2}{4}$  if  $0 < \alpha \leq \frac{1-\lambda^2}{2\sqrt{42\lambda^2 L}}$ .

## D. Proof of Lemma 6

### D.1. Proof of Eq. (11)

We recursively apply the inequality on  $V_t$  from  $t$  to 0 to obtain:  $\forall t \geq 1$ ,

$$\begin{aligned}V_t &\leq qV_{t-1} + qR_{t-1} + Q_t + C \\ &\leq q^2 V_{t-2} + (q^2 R_{t-2} + qR_{t-1}) + (qQ_{t-1} + Q_t) + (qC + C) \\ &\quad \dots \\ &\leq q^t V_0 + \sum_{i=0}^{t-1} q^{t-i} R_i + \sum_{i=1}^t q^{t-i} Q_i + C \sum_{i=0}^{t-1} q^i.\end{aligned}\quad (41)$$

Summing up (41) over  $t$  from 1 to  $T$  gives:  $\forall T \geq 1$ ,

$$\begin{aligned}\sum_{t=0}^T V_t &\leq V_0 \sum_{t=0}^T q^t + \sum_{t=1}^T \sum_{i=0}^{t-1} q^{t-i} R_i + \sum_{t=1}^T \sum_{i=1}^t q^{t-i} Q_i + C \sum_{t=1}^T \sum_{i=0}^{t-1} q^i \\ &\leq V_0 \sum_{t=0}^{\infty} q^t + \sum_{t=0}^{T-1} \left( \sum_{i=0}^{\infty} q^i \right) R_t + \sum_{t=1}^T \left( \sum_{i=0}^{\infty} q^i \right) Q_t + C \sum_{t=1}^T \sum_{i=0}^{\infty} q^i,\end{aligned}$$

and the proof follows by  $\sum_{i=0}^{\infty} q^i = (1-q)^{-1}$ .

### D.2. Proof of Eq. (12)

We recursively apply the inequality on  $V_t$  from  $t+1$  to 1 to obtain:  $\forall t \geq 1$ ,

$$\begin{aligned}V_{t+1} &\leq qV_t + R_{t-1} + C \\ &\leq q^2 V_{t-1} + (qR_{t-2} + R_{t-1}) + (qC + C) \\ &\quad \dots \\ &\leq q^t V_1 + \sum_{i=0}^{t-1} q^{t-1-i} R_i + C \sum_{i=0}^{t-1} q^i.\end{aligned}\quad (42)$$

We sum up (42) over  $t$  from 1 to  $T - 1$  to obtain:  $\forall T \geq 2$ ,

$$\begin{aligned} \sum_{t=0}^{T-1} V_{t+1} &\leq V_1 \sum_{t=0}^{T-1} q^t + \sum_{t=1}^{T-1} \sum_{i=0}^{t-1} q^{t-1-i} R_i + C \sum_{t=1}^{T-1} \sum_{i=0}^{t-1} q^i \\ &\leq V_1 \sum_{t=0}^{\infty} q^t + \sum_{t=0}^{T-2} \left( \sum_{i=0}^{\infty} q^i \right) R_t + C \sum_{t=1}^{T-1} \sum_{i=0}^{\infty} q^i, \end{aligned}$$

and the proof follows by  $\sum_{i=0}^{\infty} q^i = (1 - q)^{-1}$ .

## E. Proof of Lemma 7

### E.1. Proof of Eq. (13)

We first observe that  $\frac{1}{1-(1-\beta)^2} \leq \frac{1}{\beta}$  for  $\beta \in (0, 1)$ . Applying (11) to (7) gives:  $\forall T \geq 1$ ,

$$\begin{aligned} &\sum_{t=0}^T \mathbb{E} \left[ \|\bar{\mathbf{v}}_t - \nabla \bar{\mathbf{f}}(\mathbf{x}_t)\|^2 \right] \\ &\leq \frac{\mathbb{E} [\|\bar{\mathbf{v}}_0 - \nabla \bar{\mathbf{f}}(\mathbf{x}_0)\|^2]}{\beta} + \frac{6L^2\alpha^2}{n\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{6L^2}{n^2\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{J}\mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{2\beta\bar{\nu}^2 T}{n} \\ &\leq \frac{\mathbb{E} [\|\bar{\mathbf{v}}_0 - \nabla \bar{\mathbf{f}}(\mathbf{x}_0)\|^2]}{\beta} + \frac{6L^2\alpha^2}{n\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{12L^2}{n^2\beta} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{2\beta\bar{\nu}^2 T}{n}. \end{aligned} \quad (43)$$

Towards the first term in (43), we observe that

$$\begin{aligned} \mathbb{E} [\|\bar{\mathbf{v}}_0 - \nabla \bar{\mathbf{f}}(\mathbf{x}_0)\|^2] &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \frac{1}{b_0} \sum_{r=1}^{b_0} (\mathbf{g}_i(\mathbf{x}_0^i, \xi_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)) \right\|^2 \right] \\ &\stackrel{(i)}{=} \frac{1}{n^2 b_0^2} \sum_{i=1}^n \sum_{r=1}^{b_0} \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_0^i, \xi_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)\|^2] \leq \frac{\bar{\nu}^2}{nb_0}, \end{aligned} \quad (44)$$

where (i) follows from a similar line of arguments in (25). Then (13) follows from using (44) in (43).

### E.2. Proof of Eq. (14)

We apply (11) to (8) to obtain:  $\forall T \geq 1$ ,

$$\begin{aligned} &\sum_{t=0}^T \mathbb{E} [\|\mathbf{v}_t - \nabla \mathbf{f}(\mathbf{x}_t)\|^2] \\ &\leq \frac{\mathbb{E} [\|\mathbf{v}_0 - \nabla \mathbf{f}(\mathbf{x}_0)\|^2]}{\beta} + \frac{6nL^2\alpha^2}{\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{6L^2}{\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{J}\mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + 2n\beta\bar{\nu}^2 T \\ &\leq \frac{\mathbb{E} [\|\mathbf{v}_0 - \nabla \mathbf{f}(\mathbf{x}_0)\|^2]}{\beta} + \frac{6nL^2\alpha^2}{\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{12L^2}{\beta} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + 2n\beta\bar{\nu}^2 T. \end{aligned} \quad (45)$$

In (45), we observe that

$$\begin{aligned} \mathbb{E} [\|\mathbf{v}_0 - \nabla \mathbf{f}(\mathbf{x}_0)\|^2] &= \sum_{i=1}^n \mathbb{E} \left[ \left\| \frac{1}{b_0} \sum_{r=1}^{b_0} (\mathbf{g}_i(\mathbf{x}_0^i, \xi_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)) \right\|^2 \right] \\ &\stackrel{(i)}{=} \frac{1}{b_0^2} \sum_{i=1}^n \sum_{r=1}^{b_0} \mathbb{E} [\|\mathbf{g}_i(\mathbf{x}_0^i, \xi_{0,r}^i) - \nabla f_i(\mathbf{x}_0^i)\|^2] \leq \frac{n\bar{\nu}^2}{b_0}, \end{aligned} \quad (46)$$

where (i) follows from a similar line of arguments in (25). Then (14) follows from using (46) in (45).

## F. Proof of Lemma 8

We recall that  $\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\| = 0$ , since it is assumed without generality that  $\mathbf{x}_0^i = \mathbf{x}_0^j$  for any  $i, j \in \mathcal{V}$ . Applying (11) to (9) yields:  $\forall T \geq 1$ ,

$$\sum_{t=0}^T \|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2 \leq \frac{4\lambda^2\alpha^2}{(1-\lambda^2)^2} \sum_{t=1}^T \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2. \quad (47)$$

To further bound  $\sum_{t=1}^T \|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2$ , we apply (12) in Lemma 5(b) to obtain: if  $0 < \alpha \leq \frac{1-\lambda^2}{2\sqrt{42}\lambda^2 L}$ , then  $\forall T \geq 2$ ,

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E} [\|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2] \\ & \leq \frac{4\mathbb{E} [\|\mathbf{y}_1 - \mathbf{J}\mathbf{y}_1\|^2]}{1-\lambda^2} + \frac{84\lambda^2 n L^2 \alpha^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{252\lambda^2 L^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ & \quad + \frac{28\lambda^2 \beta^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\mathbf{v}_t - \nabla \mathbf{f}(\mathbf{x}_t)\|^2] + \frac{12\lambda^2 n \beta^2 \bar{\nu}^2 T}{1-\lambda^2} \\ & \leq \frac{84\lambda^2 n L^2 \alpha^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{252\lambda^2 L^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ & \quad + \frac{28\lambda^2 \beta^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\mathbf{v}_t - \nabla \mathbf{f}(\mathbf{x}_t)\|^2] + \frac{12\lambda^2 n \beta^2 \bar{\nu}^2 T}{1-\lambda^2} + \frac{4\lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2}{1-\lambda^2} + \frac{4\lambda^2 n \bar{\nu}^2}{(1-\lambda^2)b_0}, \end{aligned} \quad (48)$$

where the last inequality is due to Lemma 5(a). To proceed, we use (14), an upper bound on  $\sum_t \mathbb{E} [\|\mathbf{v}_t - \nabla \mathbf{f}(\mathbf{x}_t)\|^2]$ , in (48) to obtain: if  $0 < \alpha \leq \frac{1-\lambda^2}{2\sqrt{42}\lambda^2 L}$  and  $\beta \in (0, 1)$ , then  $\forall T \geq 2$ ,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [\|\mathbf{y}_t - \mathbf{J}\mathbf{y}_t\|^2] & \leq \frac{252\lambda^2 n L^2 \alpha^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{588\lambda^2 L^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-1} \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ & \quad + \frac{28\lambda^2 n \beta \bar{\nu}^2}{(1-\lambda^2)^2 b_0} + \frac{56\lambda^2 n \beta^3 \bar{\nu}^2 T}{(1-\lambda^2)^2} + \frac{12\lambda^2 n \beta^2 \bar{\nu}^2 T}{1-\lambda^2} + \frac{4\lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2}{1-\lambda^2} + \frac{4\lambda^2 n \bar{\nu}^2}{(1-\lambda^2)b_0} \\ & = \frac{252\lambda^2 n L^2 \alpha^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{588\lambda^2 L^2}{(1-\lambda^2)^2} \sum_{t=0}^{T-1} \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ & \quad + \left( \frac{7\beta}{1-\lambda^2} + 1 \right) \frac{4\lambda^2 n \bar{\nu}^2}{(1-\lambda^2)b_0} + \left( \frac{14\beta}{1-\lambda^2} + 3 \right) \frac{4\lambda^2 n \beta^2 \bar{\nu}^2 T}{1-\lambda^2} + \frac{4\lambda^2 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2}{1-\lambda^2}. \end{aligned} \quad (49)$$

Finally, we use (49) in (47) to obtain:  $\forall T \geq 2$ ,

$$\begin{aligned} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] & \leq \frac{1008\lambda^4 n L^2 \alpha^4}{(1-\lambda^2)^4} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{2352\lambda^4 L^2 \alpha^2}{(1-\lambda^2)^4} \sum_{t=0}^{T-1} \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ & \quad + \left( \frac{7\beta}{1-\lambda^2} + 1 \right) \frac{16\lambda^4 n \bar{\nu}^2 \alpha^2}{(1-\lambda^2)^3 b_0} + \left( \frac{14\beta}{1-\lambda^2} + 3 \right) \frac{16\lambda^4 n \beta^2 \bar{\nu}^2 \alpha^2 T}{(1-\lambda^2)^3} + \frac{16\lambda^4 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2 \alpha^2}{(1-\lambda^2)^3}, \end{aligned}$$

which may be written equivalently as

$$\begin{aligned} \left( 1 - \frac{2352\lambda^4 L^2 \alpha^2}{(1-\lambda^2)^4} \right) \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] & \leq \frac{1008\lambda^4 n L^2 \alpha^4}{(1-\lambda^2)^4} \sum_{t=0}^{T-2} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \left( \frac{7\beta}{1-\lambda^2} + 1 \right) \frac{16\lambda^4 n \bar{\nu}^2 \alpha^2}{(1-\lambda^2)^3 b_0} \\ & \quad + \left( \frac{14\beta}{1-\lambda^2} + 3 \right) \frac{16\lambda^4 n \beta^2 \bar{\nu}^2 \alpha^2 T}{(1-\lambda^2)^3} + \frac{16\lambda^4 \|\nabla \mathbf{f}(\mathbf{x}_0)\|^2 \alpha^2}{(1-\lambda^2)^3}. \end{aligned} \quad (50)$$

We observe in (50) that  $\frac{2352\lambda^4 L^2 \alpha^2}{(1-\lambda^2)^4} \leq \frac{1}{2}$  if  $0 < \alpha \leq \frac{(1-\lambda^2)^2}{70\lambda^2 L}$ , and the proof follows.

## G. Proof of Theorem 1

For the ease of presentation, we denote  $\Delta_0 := F(\bar{\mathbf{x}}_0) - F^*$  in the following. We apply (13) to Lemma 2 to obtain: if  $0 < \alpha \leq \frac{1}{2L}$ , then  $\forall T \geq 1$ ,

$$\begin{aligned} \sum_{t=0}^T \mathbb{E} [\|\nabla F(\bar{\mathbf{x}}_t)\|^2] &\leq \frac{2\Delta_0}{\alpha} - \frac{1}{2} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{2L^2}{n} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ &\quad + \frac{2\bar{\nu}^2}{\beta b_0 n} + \frac{12L^2\alpha^2}{n\beta} \sum_{t=0}^{T-1} \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{24L^2}{n^2\beta} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{4\beta\bar{\nu}^2 T}{n} \\ &\leq \frac{2\Delta_0}{\alpha} - \frac{1}{4} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{2L^2}{n} \left(1 + \frac{12}{n\beta}\right) \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] \\ &\quad + \frac{2\bar{\nu}^2}{\beta b_0 n} + \frac{4\beta\bar{\nu}^2 T}{n} - \left(\frac{1}{4} - \frac{12L^2\alpha^2}{n\beta}\right) \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2]. \end{aligned} \quad (51)$$

Therefore, if  $0 < \alpha < \frac{1}{4\sqrt{3}L}$  and  $\frac{48L^2\alpha^2}{n} \leq \beta < 1$ , i.e.,  $\frac{1}{4} - \frac{12L^2\alpha^2}{n\beta} \geq 0$ , we may drop the last term in (51) to obtain:  $\forall T \geq 1$ ,

$$\sum_{t=0}^T \mathbb{E} [\|\nabla F(\bar{\mathbf{x}}_t)\|^2] \leq \frac{2\Delta_0}{\alpha} - \frac{1}{4} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{2L^2}{n} \left(1 + \frac{12}{n\beta}\right) \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{2\bar{\nu}^2}{\beta b_0 n} + \frac{4\beta\bar{\nu}^2 T}{n}. \quad (52)$$

Moreover, we observe:  $\forall T \geq 1$ ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E} [\|\nabla F(\mathbf{x}_t^i)\|^2] &\leq \frac{2}{n} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E} [\|\nabla F(\mathbf{x}_t^i) - \nabla F(\bar{\mathbf{x}}_t)\|^2 + \|\nabla F(\bar{\mathbf{x}}_t)\|^2] \\ &= \frac{2L^2}{n} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + 2 \sum_{t=0}^T \mathbb{E} [\|\nabla F(\bar{\mathbf{x}}_t)\|^2], \end{aligned} \quad (53)$$

where the last line uses the  $L$ -smoothness of  $F$ . Using (52) in (53) yields: if  $0 < \alpha < \frac{1}{4\sqrt{3}L}$  and  $48L^2\alpha^2/n \leq \beta < 1$ , then

$$\frac{1}{n} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E} [\|\nabla F(\mathbf{x}_t^i)\|^2] \leq \frac{4\Delta_0}{\alpha} - \frac{1}{2} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{6L^2}{n} \left(1 + \frac{8}{n\beta}\right) \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{4\bar{\nu}^2}{\beta b_0 n} + \frac{8\beta\bar{\nu}^2 T}{n}. \quad (54)$$

According to (54), if  $0 < \alpha < \frac{1}{4\sqrt{3}L}$  and  $\beta = 48L^2\alpha^2/n$ , we have:  $\forall T \geq 1$ ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E} [\|\nabla F(\mathbf{x}_t^i)\|^2] &\leq \frac{4\Delta_0}{\alpha} - \frac{1}{2} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{6L^2}{n} \left(1 + \frac{1}{6L^2\alpha^2}\right) \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{4\bar{\nu}^2}{\beta b_0 n} + \frac{8\beta\bar{\nu}^2 T}{n} \\ &\leq \underbrace{\frac{4\Delta_0}{\alpha} - \frac{1}{2} \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2]}_{=: \Phi_T} + \frac{2}{n\alpha^2} \sum_{t=0}^T \mathbb{E} [\|\mathbf{x}_t - \mathbf{J}\mathbf{x}_t\|^2] + \frac{4\bar{\nu}^2}{\beta b_0 n} + \frac{8\beta\bar{\nu}^2 T}{n}, \end{aligned} \quad (55)$$

where the last line is due to  $6L^2\alpha^2 < 1/8$ . To simplify  $\Phi_T$ , we use Lemma 8 to obtain: if  $0 < \alpha \leq \frac{(1-\lambda^2)^2}{70\lambda^2 L}$  then  $\forall T \geq 2$ ,

$$\begin{aligned} \Phi_T &\leq -\frac{1}{2} \left(1 - \frac{8064\lambda^4 L^2 \alpha^2}{(1-\lambda^2)^4}\right) \sum_{t=0}^T \mathbb{E} [\|\bar{\mathbf{v}}_t\|^2] + \frac{64\lambda^4}{(1-\lambda^2)^3} \frac{\|\nabla f(\mathbf{x}_0)\|^2}{n} \\ &\quad + \left(\frac{7\beta}{1-\lambda^2} + 1\right) \frac{64\lambda^4 \bar{\nu}^2}{(1-\lambda^2)^3 b_0} + \left(\frac{14\beta}{1-\lambda^2} + 3\right) \frac{64\lambda^4 \beta^2 \bar{\nu}^2 T}{(1-\lambda^2)^3}. \end{aligned} \quad (56)$$

In (56), we observe that if  $0 < \alpha \leq \frac{(1-\lambda^2)^2}{90\lambda^2L}$ , then  $1 - \frac{8064L^4\alpha^2}{(1-\lambda^2)^4} \geq 0$  and thus the first term in (56) may be dropped; moreover, if  $0 < \alpha \leq \frac{\sqrt{n(1-\lambda^2)}}{26\lambda L}$ , then  $\beta = \frac{48L^2\alpha^2}{n} \leq \frac{1-\lambda^2}{14\lambda^2}$ . Hence, if  $0 < \alpha \leq \min \left\{ \frac{(1-\lambda^2)^2}{90\lambda^2}, \frac{\sqrt{n(1-\lambda^2)}}{26\lambda} \right\} \frac{1}{L}$ , then (56) reduces to:  $\forall T \geq 2$ ,

$$\Phi_T \leq \frac{64\lambda^4}{(1-\lambda^2)^3} \frac{\|\nabla f(\mathbf{x}_0)\|^2}{n} + \frac{96\lambda^2\bar{\nu}^2}{(1-\lambda^2)^3 b_0} + \frac{256\lambda^2\beta^2\bar{\nu}^2 T}{(1-\lambda^2)^3}. \quad (57)$$

Finally, we use (57) in (55) to obtain: if  $0 < \alpha < \min \left\{ \frac{1}{4\sqrt{3}}, \frac{(1-\lambda^2)^2}{90\lambda^2}, \frac{\sqrt{n(1-\lambda^2)}}{26\lambda} \right\} \frac{1}{L}$ , we have:  $\forall T \geq 2$ ,

$$\begin{aligned} \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E} \left[ \|\nabla F(\mathbf{x}_t^i)\|^2 \right] &\leq \frac{4\Delta_0}{\alpha T} + \frac{4\bar{\nu}^2}{\beta b_0 n T} + \frac{8\beta\bar{\nu}^2}{n} \\ &\quad + \frac{64\lambda^4}{(1-\lambda^2)^3 T} \frac{\|\nabla f(\mathbf{x}_0)\|^2}{n} + \frac{96\lambda^2\bar{\nu}^2}{(1-\lambda^2)^3 b_0 T} + \frac{256\lambda^2\beta^2\bar{\nu}^2}{(1-\lambda^2)^3}. \end{aligned} \quad (58)$$

The proof follows by (58) and that  $\mathbb{E}[\|\nabla F(\tilde{\mathbf{x}}_T)\|^2] = \frac{1}{n(T+1)} \sum_{i=1}^n \sum_{t=0}^T \mathbb{E}[\|\nabla F(\mathbf{x}_t^i)\|^2]$  since  $\tilde{\mathbf{x}}_T$  is chosen uniformly at random from  $\{\mathbf{x}_t^i : \forall i \in \mathcal{V}, 0 \leq t \leq T\}$ .