

Supplemental Materials

In this section, we present the complete proofs of the above lemmas and theorems. In the beginning, we give some useful properties of CoordSGE and RandSGE, respectively.

Lemma 10 ((Liu et al. 2018)). *Suppose that the function $f(x)$ is L -smooth. Let $\hat{\nabla}f(x)$ denote the estimated gradient defined by CoordSGE. Define $f_\mu = \mathbb{E}_{u \sim U[-\mu, \mu]} f(x + ue_j)$, where $U[-\mu, \mu]$ denotes the uniform distribution on the interval $[-\mu, \mu]$. Then for any $x \in \mathbb{R}^d$ we have*

1. f_μ is L -smooth, and $\hat{\nabla}f(x) = \sum_{j=1}^d \frac{\partial f_\mu(x)}{\partial x_j} e_j$.
2. $|f_\mu(x) - f(x)| \leq \frac{L\mu^2}{2}$ and $\left| \frac{\partial f_\mu(x)}{\partial x_j} \right| \leq \frac{L\mu^2}{2}$.
3. $\|\hat{\nabla}f(x) - \nabla f(x)\|^2 \leq \frac{L^2 d^2 \mu^2}{4}$.

Lemma 11. *Assume that the function $f(x)$ is L -smooth. Let $\hat{\nabla}_r f(x)$ denote the estimated gradient defined by RandSGE. Define $f_\mu = \mathbb{E}_{u \sim U_S}[f(x + \mu u)]$, where U is uniform distribution over a d -dimensional unit ball S . Then, we have*

1. For any $x \in \mathbb{R}^d$, $\nabla f_\mu(x) = \mathbb{E}_u[\hat{\nabla}_r f(x, u)]$.
2. $|f_\mu(x) - f(x)| \leq \frac{\mu^2 L}{2}$ and $\|f_\mu(x) - f(x)\| \leq \frac{\mu L d}{2}$ for any $x \in \mathbb{R}^d$.
3. $\mathbb{E}_u \|\hat{\nabla}_r f(x, u) - \hat{\nabla}_r f(y, u)\|^2 \leq 3dL^2 \|x - y\|^2 + \frac{3L^2 d^2 \mu^2}{2}$.

Proof. The proof of items 1 and 2 can be found in (Gao, Jiang, and Zhang 2018). Item 3 is due to Lemma 5 in (Ji et al. 2019). □

Lemma 12. *For a given $x \in \mathbb{R}^d$, let $\bar{x} = \text{Prox}_{\eta h}(x - \eta v)$, then we have for all $w \in \mathbb{R}^d$*

$$\begin{aligned} F(\bar{x}) &\leq F(w) + \langle \nabla f(x) - v, \bar{x} - w \rangle - \frac{1}{\eta} \langle \bar{x} - x, \bar{x} - w \rangle \\ &\quad + \frac{L}{2} \|\bar{x} - x\|^2 + \frac{L}{2} \|w - x\|^2 \end{aligned} \quad (1)$$

Proof. First, we recall the proximal operator

$$\text{Prox}_{\eta h}(x - \eta v) := \arg \min_{y \in \mathbb{R}^d} \left(h(y) + \frac{1}{2\eta} \|y - x\|^2 + \langle v, y \rangle \right) \quad (2)$$

For the nonsmooth function $h(x)$, for all $w \in \mathbb{R}^d$ we have

$$\begin{aligned} h(\bar{x}) &\leq h(w) + \langle p, \bar{x} - w \rangle \\ &= h(w) - \left\langle v + \frac{1}{\eta} (\bar{x} - x), \bar{x} - w \right\rangle \end{aligned} \quad (3)$$

where $p \in \partial h(\bar{x})$ such that $p + \frac{1}{\eta} (\bar{x} - x) + v = 0$ according to the optimality condition of (2), and (3) due to the convexity of h . In addition, since $f(x)$ is L -Lipschitz continuous, we have

$$f(\bar{x}) \leq f(x) + \langle \nabla f(x), \bar{x} - x \rangle + \frac{L}{2} \|\bar{x} - x\|^2 \quad (4)$$

and

$$f(x) \leq f(w) + \langle \nabla f(x), x - w \rangle + \frac{L}{2} \|w - x\|^2 \quad (5)$$

This lemma is obtained by adding (3), (4), (5), and using $F(x) = f(x) + h(x)$. □

Lemma 13. *Let $x_t^s = \text{Prox}_{\eta h}(x_{t-1}^s - \eta \hat{v}_{t-1}^s)$ and \bar{x}_t^s be the proximal projection using full true gradient, i.e., $\bar{x}_t^s = \text{Prox}_{\eta h}(x_{t-1}^s - \eta \nabla f(x_{t-1}^s))$. Then the following inequality holds*

$$\langle \nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \leq \eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\|^2$$

Proof. Based on inequality (3) we obtain

$$h(x_t^s) \leq h(\bar{x}_t^s) - \left\langle \hat{v}_{t-1}^s + \frac{1}{\eta} (x_t^s - x_{t-1}^s), x_t^s - \bar{x}_t^s \right\rangle \quad (6)$$

$$h(\bar{x}_t^s) \leq h(x_t^s) - \left\langle \nabla f(x_{t-1}^s) + \frac{1}{\eta} (\bar{x}_t^s - x_{t-1}^s), \bar{x}_t^s - x_t^s \right\rangle \quad (7)$$

By summing (6) and (7), we have

$$\begin{aligned} \frac{1}{\eta} \langle x_t^s - \bar{x}_t^s, x_t^s - \bar{x}_t^s \rangle &\leq \langle \nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\ \frac{1}{\eta} \|x_t^s - \bar{x}_t^s\|^2 &\leq \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\| \|x_t^s - \bar{x}_t^s\| \end{aligned} \quad (8)$$

where (8) holds by Cauchy-Schwarz inequality. Thus, we obtain

$$\|x_t^s - \bar{x}_t^s\| \leq \eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\| \quad (9)$$

Now the proof is complete using Cauchy-Schwarz inequality and (9). \square

Proof of Lemma 1

Proof. We have

$$\begin{aligned} &\mathbb{E} \left[\eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{g}^s) \right\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - \left(\nabla f(x_{t-1}^s) - \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} \hat{\nabla} f_j(\tilde{x}^{s-1}) \right) \right\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1})) + \left(\frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} \hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1}) \right) \right\|^2 \right] \\ &= \eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) + \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\ &\leq 2\eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) + \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left[\|\hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s)\|^2 \right] \end{aligned} \quad (10)$$

$$\begin{aligned} &= 2\eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\hat{\nabla} f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) \right\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] + 2\eta \mathbb{E} \left[\|\hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s)\|^2 \right] \end{aligned} \quad (11)$$

$$\begin{aligned} &= \frac{2\eta}{b^2} \mathbb{E} \left[\sum_{i \in I_b} \left\| ((\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})) - (\hat{\nabla} f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) \right\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] + 2\eta \mathbb{E} \left[\|\hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s)\|^2 \right] \end{aligned} \quad (12)$$

$$\begin{aligned} &\leq \frac{2\eta}{b^2} \mathbb{E} \left[\sum_{i \in I_b} \|\hat{\nabla} f_i(x_{t-1}^s) - \hat{\nabla} f_i(\tilde{x}^{s-1})\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] + 2\eta \mathbb{E} \left[\|\hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s)\|^2 \right] \quad (13) \\ &\leq \frac{3\eta L^2 d^2 \mu^2}{b} + \frac{6\eta L^2}{b} \mathbb{E} \left[\|x_{t-1}^s - \tilde{x}^{s-1}\|^2 \right] + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \end{aligned}$$

$$+ 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2 \quad (14)$$

$$\leq \frac{6\eta L^2}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2 + \frac{3\eta L^2 d^2 \mu^2}{b} \quad (15)$$

$$\leq \frac{6\eta L^2}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} + \eta \frac{L^2 d^2 \mu^2}{2} + \frac{3\eta L^2 d^2 \mu^2}{b} \quad (16)$$

$$\leq \frac{6\eta L^2}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \quad (17)$$

where, recalling that a deterministic gradient estimator is employed and the expectations are taking with respect to I_b and $I_{\mathcal{B}}$. The inequality (10) holds by the Jensen's inequality. (11) and (12) are due to $\mathbb{E}[\|x_1 + x_2 + \dots + x_k\|^2] = \sum_{i=1}^k \mathbb{E}[\|x_i\|^2]$ if x_1, x_2, \dots, x_k are independent and of mean zero. Recall that I_b and $I_{\mathcal{B}}$ are also independent. (13) applies the fact that for any random variable z , $\mathbb{E}[\|z - \mathbb{E}[z]\|^2] \leq \mathbb{E}[\|z\|^2]$. (14) employs following inequality

$$\begin{aligned} \mathbb{E} \left\| \hat{\nabla} f_i(x_t^s) - \hat{\nabla} f_i(\tilde{x}^s) \right\|^2 &= \mathbb{E} \left\| \hat{\nabla} f_i(x_t^s) - \nabla f_i(x_t^s) + \nabla f_i(x_t^s) - \hat{\nabla} f_i(\tilde{x}^s) + \nabla f_i(\tilde{x}^s) - \nabla f_i(\tilde{x}^s) \right\|^2 \\ &\leq 3\mathbb{E} \left\| \hat{\nabla} f_i(x_t^s) - \nabla f_i(x_t^s) \right\|^2 + 3\mathbb{E} \left\| \hat{\nabla} f_i(\tilde{x}^s) - \nabla f_i(\tilde{x}^s) \right\|^2 \\ &\quad + 3\mathbb{E} \left\| \nabla f_i(x_t^s) - \nabla f_i(\tilde{x}^s) \right\|^2 \\ &\leq \frac{3L^2 d^2 \mu^2}{2} + 3L^2 \|x_t^s - \tilde{x}^s\|^2 \end{aligned} \quad (18)$$

where the last inequality used the fact that f_{i,μ_j} is L -smooth. (15) is by Assumption 2 and (16) uses Lemma 10. The proof is now complete. \square

Proof of Lemma 2

Proof. We have

$$\begin{aligned} &\mathbb{E} \left[\eta \left\| \nabla f(x_{t-1}^s) - \tilde{v}_{t-1} \right\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{g}^s) \right\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - \left(\nabla f(x_{t-1}^s) - \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} \hat{\nabla} f_j(\tilde{x}^{s-1}) \right) \right\|^2 \right] \\ &= \mathbb{E} \left[\eta \left\| \frac{1}{b} \sum_{i \in I_b} (\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1})) + \left(\frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} \hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1}) \right) \right\|^2 \right] \\ &= \eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\nabla f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) + \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\ &\leq 2\eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\hat{\nabla} f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) + \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2 \quad (19) \end{aligned}$$

$$\begin{aligned} &= 2\eta \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in I_b} ((\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\hat{\nabla} f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1}))) \right\|^2 \right] \\ &\quad + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2 \quad (20) \end{aligned}$$

$$\begin{aligned}
&= \frac{2\eta}{b^2} \mathbb{E} \left[\sum_{i \in I_b} \left\| \left((\hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1})) - (\hat{\nabla} f(x_{t-1}^s) - \hat{\nabla} f(\tilde{x}^{s-1})) \right) \right\|^2 \right] \\
&\quad + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2
\end{aligned} \tag{21}$$

$$\begin{aligned}
&\leq \frac{2\eta}{b^2} \mathbb{E} \left[\sum_{i \in I_b} \left\| \hat{\nabla}_r f_i(x_{t-1}^s) - \hat{\nabla}_r f_i(\tilde{x}^{s-1}) \right\|^2 \right] + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\
&\quad + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2
\end{aligned} \tag{22}$$

$$\begin{aligned}
&\leq \frac{3\eta L^2 d^2 \mu^2}{b} + \frac{6\eta L^2 d}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2\eta \mathbb{E} \left[\left\| \frac{1}{\mathcal{B}} \sum_{j \in I_{\mathcal{B}}} (\hat{\nabla} f_j(\tilde{x}^{s-1}) - \hat{\nabla} f(\tilde{x}^{s-1})) \right\|^2 \right] \\
&\quad + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2
\end{aligned} \tag{23}$$

$$\begin{aligned}
&\leq \frac{6\eta L^2 d}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} + 2\eta \mathbb{E} \left\| \hat{\nabla} f(x_{t-1}^s) - \nabla f(x_{t-1}^s) \right\|^2 \\
&\quad + \frac{3\eta L^2 d^2 \mu^2}{b}
\end{aligned} \tag{24}$$

$$\begin{aligned}
&\leq \frac{6\eta L^2 d}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} \\
&\quad + \eta \frac{L^2 d^2 \mu^2}{2} + \frac{3\eta L^2 d^2 \mu^2}{b}
\end{aligned} \tag{25}$$

$$\leq \frac{6\eta L^2 d}{b} \mathbb{E} \left[\left\| x_{t-1}^s - \tilde{x}^{s-1} \right\|^2 \right] + 2 \frac{I(\mathcal{B} < n) \eta \sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \tag{26}$$

where, the expectations are taking with respect to I_b and $I_{\mathcal{B}}$ and random directions $\{u_i\}$ in (2). The inequality (19) holds by the Jensen's inequality. (20) and (21) are based on $\mathbb{E}[\|x_1 + x_2 + \dots + x_k\|^2] = \sum_{i=1}^k \mathbb{E}[\|x_i\|^2]$ if x_1, x_2, \dots, x_k are independent and of mean zero. Note that I_b and $I_{\mathcal{B}}$ are also independent. (22) uses the fact that $\mathbb{E}[\|x - \mathbb{E}[x]\|^2] \leq \mathbb{E}[\|x\|^2]$, for any random variable x . (23) holds due to Lemma 11. (24) is by Assumption 2 and (25) is by Lemma 10. (26) uses $b \geq 1$. The proof is now complete. \square

In the sequel we frequently use the following inequality

$$\|x - z\|^2 \leq (1 + \frac{1}{\alpha}) \|x - y\|^2 + (1 + \alpha) \|y - z\|^2, \forall \alpha > 0 \tag{27}$$

Proof of Theorem 3

Proof. Now, we apply Lemma 12 to prove Theorem 3. Let $x_t^s = \text{Prox}_{\eta h}(x_{t-1}^s - \eta \hat{\nabla}_t^s)$ and $\bar{x}_t^s = \text{Prox}_{\eta h}(x_{t-1}^s - \eta \nabla f(x_{t-1}^s))$. By letting $x^+ = x_t^s$, $x = x_{t-1}^s$, $v = \hat{\nabla}_t^s$ and $z = \bar{x}_t^s$ in (1), we have

$$\begin{aligned}
F(x_t^s) &\leq F(\bar{x}_t^s) + \langle \nabla f(x_{t-1}^s) - \hat{\nabla}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle - \frac{1}{\eta} \langle x_t^s - x_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
&\quad + \frac{L}{2} \|x_t^s - x_{t-1}^s\|^2 + \frac{L}{2} \|\bar{x}_t^s - x_{t-1}^s\|^2.
\end{aligned} \tag{28}$$

Besides, by letting $x^+ = \bar{x}_t^s$, $x = x_{t-1}^s$, $v = \nabla f(x_{t-1}^s)$ and $z = x = x_{t-1}^s$ in (1), we have

$$\begin{aligned}
F(\bar{x}_t^s) &\leq F(x_{t-1}^s) - \frac{1}{\eta} \langle \bar{x}_t^s - x_{t-1}^s, \bar{x}_t^s - x_{t-1}^s \rangle + \frac{L}{2} \|\bar{x}_t^s - x_{t-1}^s\|^2 \\
&= F(x_{t-1}^s) - \left(\frac{1}{\eta} - \frac{L}{2} \right) \|\bar{x}_t^s - x_{t-1}^s\|^2.
\end{aligned} \tag{29}$$

Combining (28) and (29) we have

$$\begin{aligned}
F(x_t^s) &\leq F(x_{t-1}^s) + \frac{L}{2} \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \langle \nabla f(x_{t-1}^s) - \hat{\nabla}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
&\quad - \frac{1}{\eta} \langle x_t^s - x_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
&= F(x_{t-1}^s) + \frac{L}{2} \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \langle \nabla f(x_{t-1}^s) - \hat{\nabla}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle
\end{aligned}$$

$$\begin{aligned}
& -\frac{1}{2\eta} \left(\|x_t^s - x_{t-1}^s\|^2 + \|x_t^s - \bar{x}_t^s\|^2 - \|\bar{x}_t^s - x_{t-1}^s\|^2 \right) \\
& = F(x_{t-1}^s) - \left(\frac{1}{2\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{2\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \langle \nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
& \quad - \frac{1}{2\eta} \|x_t^s - \bar{x}_t^s\|^2 \\
& \leq F(x_{t-1}^s) - \left(\frac{1}{2\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{2\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \langle \nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
& \quad - \frac{1}{8\eta} \|x_t^s - x_{t-1}^s\|^2 + \frac{1}{6\eta} \|\bar{x}_t^s - x_{t-1}^s\|^2 \tag{30}
\end{aligned}$$

$$\begin{aligned}
& = F(x_{t-1}^s) - \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{3\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \langle \nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s, x_t^s - \bar{x}_t^s \rangle \\
& \leq F(x_{t-1}^s) - \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{3\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\|^2 \tag{31}
\end{aligned}$$

where the second inequality uses (27) with $\alpha = 3$ and the last inequality holds due to the Lemma 13.

Note that $x_t^s = \text{Prox}_{\eta h}(x_{t-1}^s - \eta \hat{v}_{t-1}^s)$ is the iterated form in our algorithm. By taking the expectation with respect to all random variables in (31) we obtain

$$\mathbb{E}[F(x_t^s)] \leq \mathbb{E} \left[F(x_{t-1}^s) - \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{3\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 + \eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\|^2 \right] \tag{32}$$

In (32), we further bound $\eta \|\nabla f(x_{t-1}^s) - \hat{v}_{t-1}^s\|^2$ using Lemma 1 to obtain

$$\begin{aligned}
& \mathbb{E}[F(x_t^s)] \\
& \leq \mathbb{E} \left[F(x_{t-1}^s) - \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{1}{3\eta} - L \right) \|\bar{x}_t^s - x_{t-1}^s\|^2 \right] \\
& \quad + \frac{6\eta L^2}{b} \mathbb{E} \|x_{t-1}^s - \bar{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \\
& = \mathbb{E} \left[F(x_{t-1}^s) - \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - x_{t-1}^s\|^2 - \left(\frac{\eta}{3} - L\eta^2 \right) \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad + \frac{6\eta L^2}{b} \mathbb{E} \|x_{t-1}^s - \bar{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \tag{33}
\end{aligned}$$

$$\begin{aligned}
& \leq \mathbb{E} \left[F(x_{t-1}^s) - \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \bar{x}^{s-1}\|^2 - \left(\frac{\eta}{3} - L\eta^2 \right) \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad + \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \bar{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \tag{34}
\end{aligned}$$

where recalling $\bar{x}_t^s := \text{Prox}_{\eta h}(x_{t-1}^s - \eta \nabla f(x_{t-1}^s))$, (33) is based on the definition of gradient mapping $g_\eta(x_{t-1}^s)$. (34) uses (27) by choosing $\alpha = 2t - 1$.

Taking a telescopic sum for $t = 1, 2, \dots, m$ in epoch s from (34) and recalling that $x_m^s = \bar{x}^s$ and $x_0^s = \bar{x}^{s-1}$, we obtain

$$\begin{aligned}
& \mathbb{E}[F(\bar{x}^s)] \\
& \leq \mathbb{E} \left[F(\bar{x}^{s-1}) - \sum_{t=1}^m \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \bar{x}^{s-1}\|^2 - \left(\frac{\eta}{3} - L\eta^2 \right) \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad + \sum_{t=1}^m \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \bar{x}^{s-1}\|^2 \\
& \quad + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2} \\
& \leq \mathbb{E} \left[F(\bar{x}^{s-1}) - \sum_{t=1}^{m-1} \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \bar{x}^{s-1}\|^2 - \left(\frac{\eta}{3} - L\eta^2 \right) \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad + \sum_{t=2}^m \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \bar{x}^{s-1}\|^2
\end{aligned}$$

$$\begin{aligned}
& + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2} \\
& = \mathbb{E} \left[F(\tilde{x}^{s-1}) - \left(\frac{\eta}{3} - L\eta^2 \right) \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad - \sum_{t=1}^{m-1} \left(\left(\frac{1}{2t} - \frac{1}{2t+1} \right) \left(\frac{5}{8\eta} - \frac{L}{2} \right) - \frac{6\eta L^2}{b} \right) \mathbb{E} \|x_t^s - \tilde{x}^{s-1}\|^2 \\
& \quad + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2} \\
& \leq \mathbb{E} \left[F(\tilde{x}^{s-1}) - \left(\frac{\eta}{3} - L\eta^2 \right) \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] \\
& \quad - \sum_{t=1}^{m-1} \left(\frac{1}{6t^2} \left(\frac{5}{8\eta} - \frac{L}{2} \right) - \frac{6\eta L^2}{b} \right) \mathbb{E} \|x_t^s - \tilde{x}^{s-1}\|^2 \\
& \quad + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2} \\
& \leq \mathbb{E} \left[F(\tilde{x}^{s-1}) - \frac{\eta}{6} \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2}
\end{aligned} \tag{35}$$

$$\begin{aligned}
& \leq \mathbb{E} \left[F(\tilde{x}^{s-1}) - \frac{\eta}{6} \sum_{t=1}^m \|g_\eta(x_{t-1}^s)\|^2 \right] + \sum_{t=1}^m \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \sum_{t=1}^m \eta \frac{7L^2 d^2 \mu^2}{2}
\end{aligned} \tag{36}$$

where (35) holds since norm is always non-negative and $x_0^s = \tilde{x}^{s-1}$. In (36) we have used the fact that $(\frac{1}{6t^2}(\frac{5}{8\eta} - \frac{L}{2}) - \frac{6\eta L^2}{b}) \geq 0$ for all $1 \leq t \leq m$ and $\frac{\eta}{6} \leq \frac{\eta}{3} - L\eta^2$ since $\eta = \min\{\frac{1}{8L}, \frac{\sqrt{b}}{12mL}\}$. Telescoping the sum for $s = 1, 2, \dots, S$ in (36), we obtain

$$\begin{aligned}
0 & \leq \mathbb{E}[F(\tilde{x}^S) - F(x^*)] \\
& \leq \mathbb{E} \left[F(\tilde{x}^0) - F(x^*) - \sum_{s=1}^S \sum_{t=1}^m \frac{\eta}{6} \|g_\eta(x_{t-1}^s)\|^2 + \sum_{s=1}^S \sum_{t=1}^m \left(\frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \right) \right]
\end{aligned}$$

Thus, we have

$$\mathbb{E}[\|g_\eta(\hat{x})\|^2] \leq \frac{6(F(x_0) - F(x^*))}{\eta Sm} + \frac{I(\mathcal{B} < n)12\sigma^2}{\mathcal{B}} + 21L^2 d^2 \mu^2 \tag{37}$$

where (37) holds since we choose \hat{x} uniformly randomly from $\{x_{t-1}^s\}_{t \in [m], s \in [S]}$. \square

Proof of Corollary 4

Proof. Using Theorem 3, we have $\eta = \min\{\frac{1}{8L}, \frac{\sqrt{b}}{12mL}\}$

$$\begin{aligned}
\mathbb{E}[\|g_\eta(\hat{x})\|^2] & \leq \frac{6(F(x_0) - F(x^*))}{\eta Sm} \\
& \quad + \frac{I(\mathcal{B} < n)12\sigma^2}{\mathcal{B}} + 21L^2 d^2 \mu^2 = 3\epsilon
\end{aligned} \tag{38}$$

Now we derive the total number of iterations $T = Sm = \frac{6(F(x_0) - F(x^*))}{\epsilon\eta}$. Since $\mu \leq \frac{\sqrt{\epsilon}}{5dL}$, and for $\mathcal{B} = n$, the second term in the bound (38) is 0, the proof is completed as the number of SZO call equals to $Sn + Smb = 6(F(x_0) - F(x^*))(\frac{n}{\epsilon\eta m} + \frac{b}{\epsilon\eta})$. If $\mathcal{B} < n$ the number of SZO calls equal to $d(S\mathcal{B} + Smb) = 6d(F(x_0) - F(x^*))(\frac{\mathcal{B}}{\epsilon\eta m} + \frac{b}{\epsilon\eta})$ by noting that $\frac{I(\mathcal{B} < n)12\sigma^2}{\mathcal{B}} \leq \epsilon$ due to $\mathcal{B} \geq 12\sigma^2/\epsilon$. The second part of corollary is obtained by setting $m = \sqrt{b}$ in the first part. \square

Corollary 14. We set the batch size $\mathcal{B} = \min\{12\sigma^2/\epsilon, n\}$ and the smoothing parameter $\mu \leq \frac{\sqrt{\epsilon}}{5dL}$. Suppose \hat{x} returned by Algorithm 1 is an ϵ -accurate solution for problem (1). Recalling that CoordSGE and RandSGE require $O(d)$ and $O(1)$ function queries respectively, the number of SZO calls is at most

$$\begin{aligned}
(dS\mathcal{B} + Smb) & = 6(F(x_0) - F(x^*)) \left(\frac{\mathcal{B}d}{\epsilon\eta m} + \frac{b}{\epsilon\eta} \right) \\
& = O\left(\frac{\mathcal{B}d}{\epsilon\eta m} + \frac{b}{\epsilon\eta} \right)
\end{aligned} \tag{39}$$

and the number of PO calls is equal to $T = Sm = \frac{6(F(x_0) - F(x^*))}{\epsilon\eta} = O\left(\frac{1}{\epsilon\eta}\right)$. In particular, by setting $m = \sqrt{b}$ and $\eta = \frac{1}{12L\sqrt{d}}$, the number of ZO calls is at most

$$\begin{aligned} & 72L(F(x_0) - F(x^*)) \left(\frac{\mathcal{B}d\sqrt{d}}{\epsilon\sqrt{b}} + \frac{b\sqrt{d}}{\epsilon} \right) \\ & = O\left(s_n \frac{d\sqrt{d}}{\epsilon\sqrt{b}} + \frac{b\sqrt{d}}{\epsilon}\right) \end{aligned} \quad (40)$$

where $s_n = \min\{n, \frac{1}{\epsilon}\}$ and the number of PO queries is equal to $T = Sm = S\sqrt{b} = \frac{72L\sqrt{d}(F(x_0) - F(x^*))}{\epsilon} = O\left(\frac{\sqrt{d}}{\epsilon}\right)$.

Proof of Theorem 7

Proof. We start by recalling inequality (33) from the proof of Theorem 3, i.e.,

$$\begin{aligned} & \mathbb{E}[F(x_t^s)] \\ & \leq \mathbb{E} \left[F(x_{t-1}^s) - \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \left(\frac{\eta}{3} - L\eta^2 \right) \|g_\eta(x_{t-1}^s)\|^2 \right] \\ & \quad + \left(\frac{2\eta L^2 d}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{L^2 d^2 \mu^2}{2} \\ & \leq \mathbb{E} \left[F(x_{t-1}^s) - \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \frac{\eta}{6} \|g_\eta(x_{t-1}^s)\|^2 \right] \\ & \quad + \left(\frac{6\eta L^2 d}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \end{aligned} \quad (41)$$

where in (41) inequality we applied $\eta L \leq \frac{1}{6}$. Moreover, substituting PL inequality, i.e.,

$$\|g_\eta(x)\|^2 \geq 2\lambda(F(x) - F^*) \quad (42)$$

into (41), we obtain

$$\begin{aligned} & \mathbb{E}[F(x_t^s)] \\ & \leq \mathbb{E} \left[F(x_{t-1}^s) - \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \lambda \frac{\eta}{3} (F(x_{t-1}^s) - F^*) \right] \\ & \quad + \left(\frac{6\eta L^2 d}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \end{aligned} \quad (43)$$

Thus, we have

$$\begin{aligned} & \mathbb{E}[F(x_t^s)] \\ & \leq \mathbb{E} \left[\left(1 - \lambda \frac{\eta}{3} \right) (F(x_{t-1}^s) - F^*) - \frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 \right] \\ & \quad + \left(\frac{6\eta L^2 d}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \mathbb{E} \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 + \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \eta \frac{7L^2 d^2 \mu^2}{2} \end{aligned} \quad (44)$$

Let $\beta := 1 - \lambda \frac{\eta}{3}$ and $\Psi_t^s := \frac{\mathbb{E}[F(x_t^s) - F^*]}{\beta^t}$. Combining these definitions with (44), we have

$$\begin{aligned} & \Psi_t^s \\ & \leq \Psi_{t-1}^s - \frac{1}{\beta^t} \mathbb{E} \left[\frac{1}{2t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \left(\frac{6\eta L^2 d}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 \right] \\ & \quad + \frac{1}{\beta^t} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1}{\beta^t} \eta \frac{7L^2 d^2 \mu^2}{2} \end{aligned} \quad (45)$$

Similar to the proof of Theorem 3, summing (45) for $t = 1, 2, \dots, m$ in epoch s and remembering that $x_m^s = \tilde{x}^s$ and $x_0^s = \tilde{x}^{s-1}$, we have

$$\mathbb{E}[F(\tilde{x}^s) - F^*]$$

$$\begin{aligned}
&\leq \beta^m \mathbb{E} \left[(F(\tilde{x}^{s-1}) - F^*) \right] + \beta^m \sum_{t=1}^m \frac{1}{\beta^t} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \beta^m \sum_{t=1}^m \frac{1}{\beta^t} \eta \frac{7L^2 d^2 \mu^2}{2} \\
&\quad - \beta^m \mathbb{E} \left[\sum_{t=1}^m \frac{1}{2t\beta^t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \sum_{t=1}^m \frac{1}{\beta^t} \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 \right] \\
&\leq \beta^m \mathbb{E} \left[(F(\tilde{x}^{s-1}) - F^*) \right] + \frac{1-\beta^m}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1-\beta^m}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&\quad - \beta^m \mathbb{E} \left[\sum_{t=1}^m \frac{1}{2t\beta^t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 - \sum_{t=1}^m \frac{1}{\beta^t} \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 \right] \\
&\leq \beta^m \mathbb{E} \left[(F(\tilde{x}^{s-1}) - F^*) \right] + \frac{1-\beta^m}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1-\beta^m}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&\quad - \beta^m \mathbb{E} \left[\sum_{t=1}^{m-1} \frac{1}{2t\beta^t} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 \right] \\
&\quad + \beta^m \mathbb{E} \left[\sum_{t=2}^m \frac{1}{\beta^t} \left(\frac{6\eta L^2}{b} + \frac{1}{2t-1} \left(\frac{5}{8\eta} - \frac{L}{2} \right) \right) \|x_{t-1}^s - \tilde{x}^{s-1}\|^2 \right] \tag{46} \\
&\leq \beta^m \mathbb{E} \left[(F(\tilde{x}^{s-1}) - F^*) \right] + \frac{1-\beta^m}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1-\beta^m}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&\quad - \beta^m \mathbb{E} \left[\sum_{t=1}^{m-1} \frac{1}{\beta^{t+1}} \left(\left(\frac{\beta}{2t} - \frac{1}{2t+1} \right) \left(\frac{5}{8\eta} - \frac{L}{2} \right) - \frac{6\eta L^2}{b} \right) \|x_t^s - \tilde{x}^{s-1}\|^2 \right] \\
&\leq \beta^m \mathbb{E} \left[(F(\tilde{x}^{s-1}) - F^*) \right] + \frac{1-\beta^m}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1-\beta^m}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \tag{47}
\end{aligned}$$

where (46) since $\|\cdot\|^2$ always is non-negative and $x_0^s = \tilde{x}^{s-1}$. (47) holds since it is sufficient to show $\left(\frac{\beta}{2t} - \frac{1}{2t+1} \right) \left(\frac{5}{8\eta} - \frac{L}{2} \right) - \frac{6\eta L^2}{b} \geq 0$, for all $t = 1, 2, \dots, m$. It is easy to see that this inequality is valid due to $\eta \leq \min\{\frac{1}{8L}, \frac{\sqrt{\gamma b}}{12mL}\}$, where $\gamma = 1 - \frac{2\lambda\eta}{3}m - \frac{\lambda\eta}{3} > 0$. Similarly, let $\tilde{\beta} = \beta^m$ and $\tilde{\Psi}^s = \frac{\mathbb{E}[F(\tilde{x}^s) - F^*]}{\tilde{\beta}^s}$. Substituting these definitions into (47), we have

$$\tilde{\Psi}^s \leq \tilde{\Psi}^{s-1} + \frac{1}{\tilde{\beta}^s} \frac{1-\tilde{\beta}}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1}{\tilde{\beta}^s} \frac{1-\tilde{\beta}}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \tag{48}$$

Taking a telescopic sum from (48) for all epochs $1 \leq s \leq S$, we obtain

$$\begin{aligned}
\mathbb{E}[F(\tilde{x}^S) - F^*] &\leq \tilde{\beta}^S \mathbb{E}[F(\tilde{x}^0) - F^*] + \tilde{\beta}^S \sum_{s=1}^S \frac{1}{\tilde{\beta}^s} \frac{1-\tilde{\beta}}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \tilde{\beta}^S \sum_{s=1}^S \frac{1}{\tilde{\beta}^s} \frac{1-\tilde{\beta}}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&= \beta^{Sm} \mathbb{E}[F(\tilde{x}^0) - F^*] + \frac{1-\tilde{\beta}^S}{1-\tilde{\beta}} \frac{1-\tilde{\beta}}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1-\tilde{\beta}^S}{1-\tilde{\beta}} \frac{1-\tilde{\beta}}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&\leq \beta^{Sm} \mathbb{E}[F(\tilde{x}^0) - F^*] + \frac{1}{1-\beta} \frac{2I(\mathcal{B} < n)\eta\sigma^2}{\mathcal{B}} + \frac{1}{1-\beta} \frac{7\eta L^2 d^2 \mu^2}{2} \\
&= \left(1 - \frac{\lambda\eta}{3} \right)^{Sm} \mathbb{E}[F(\tilde{x}^0) - F^*] + \frac{6I(\mathcal{B} < n)\sigma^2}{\lambda\mathcal{B}} + \frac{21L^2 d^2 \mu^2}{2\lambda} \tag{49}
\end{aligned}$$

where in (49) we recall that $\beta = 1 - \frac{\lambda\eta}{3}$. \square

Proof of Corollary 8

Proof. From Theorem 7, we have

$$\begin{aligned}
\mathbb{E}[F(\tilde{x}^S) - F^*] &\leq \left(1 - \frac{\lambda\eta}{3} \right)^{Sm} \mathbb{E}[F(\tilde{x}^0) - F^*] \\
&\quad + \frac{6I(\mathcal{B} < n)\sigma^2}{\lambda\mathcal{B}} + \frac{21L^2 d^2 \mu^2}{2\lambda} = 3\epsilon
\end{aligned}$$

which gives the total number of iterations $T = Sm = O(\frac{1}{\lambda\eta} \log \frac{1}{\epsilon})$ and is equal to the number of PO calls. Since $\mu \leq \frac{\sqrt{\lambda\epsilon}}{4Ld}$, we have $\frac{21L^2 d^2 \mu^2}{2\lambda} \leq \epsilon$. The number of SZO calls is equal to $d(S\mathcal{B} + Sm\mathcal{B}) = O(\frac{\mathcal{B}d}{\lambda\eta m} \log \frac{1}{\epsilon} + \frac{bd}{\lambda\eta} \log \frac{1}{\epsilon})$. Note that if $\mathcal{B} < n$ then $\frac{6I(\mathcal{B} < n)\sigma^2}{\lambda\mathcal{B}} \leq \epsilon$ since

$\mathcal{B} \geq 6\sigma^2/\lambda\epsilon$. With $m = \sqrt{b}$ and $\eta = \frac{\sqrt{y}}{12L}$, the number of PO queries to $T = Sm = O(\frac{1}{\lambda\sqrt{y}} \log \frac{1}{\epsilon})$ and the number of SZO calls is equal to $d(S\mathcal{B} + Smb) = O(\frac{\mathcal{B}d}{\lambda\sqrt{y}m} \log \frac{1}{\epsilon} + \frac{bd}{\lambda\sqrt{y}} \log \frac{1}{\epsilon})$. \square

Corollary 15

Corollary 15. Suppose the final iteration point \tilde{x}^S in Algorithm 1 satisfies $\mathbb{E}[F(\tilde{x}^S) - F^*] \leq \epsilon$ under PL condition. Under Assumptions 1 and 2, we set batch size $\mathcal{B} = \min\{\frac{6\sigma^2}{\lambda\epsilon}, n\}$ and the smoothing parameter $\mu \leq \frac{\sqrt{\lambda\epsilon}}{4Ld}$. The number of SZO calls is bounded by

$$(S\mathcal{B}d + Smb) = O(\frac{s_n d}{\lambda\eta m} \log \frac{1}{\epsilon} + \frac{b}{\lambda\eta} \log \frac{1}{\epsilon})$$

where $s_n = \min\{n, \frac{1}{\lambda\epsilon}\}$. The number of PO calls is equal to the total number of iterations T which is given by

$$T = Sm = O(\frac{1}{\lambda\eta} \log \frac{1}{\epsilon})$$

In particular, given the setting $m = \sqrt{b}$ and $\eta = \frac{\sqrt{y}}{12L\sqrt{d}}$, the number of SZO calls simplifies to $(S\mathcal{B}d + Smb) = O(\frac{\mathcal{B}d\sqrt{d}}{\lambda\sqrt{y}m} \log \frac{1}{\epsilon} + \frac{b\sqrt{d}}{\lambda\sqrt{y}} \log \frac{1}{\epsilon})$.

Proof. From Theorem 7, we have

$$\begin{aligned} \mathbb{E}[F(\tilde{x}^S) - F^*] &\leq \left(1 - \frac{\lambda\eta}{3}\right)^{Sm} \mathbb{E}[F(\tilde{x}^0) - F^*] \\ &\quad + \frac{6I(\mathcal{B} < n)\sigma^2}{\lambda\mathcal{B}} + \frac{21L^2d^2\mu^2}{2\lambda} = 3\epsilon \end{aligned} \quad (50)$$

which gives the total number of iterations $T = Sm = O(\frac{1}{\lambda\eta} \log \frac{1}{\epsilon})$ and equals to the number of PO calls. Since $\mu \leq \frac{\sqrt{\lambda\epsilon}}{4Ld}$, we have $\frac{21L^2d^2\mu^2}{2\lambda} \leq \epsilon$. The number of SZO calls equals to $(S\mathcal{B}d + Smb) = O(\frac{\mathcal{B}d}{\lambda\eta m} \log \frac{1}{\epsilon} + \frac{b}{\lambda\eta} \log \frac{1}{\epsilon})$. Note that if $\mathcal{B} < n$ then $\frac{6I(\mathcal{B} < n)\sigma^2}{\lambda\mathcal{B}} \leq \epsilon$ since $\mathcal{B} \geq 6\sigma^2/\lambda\epsilon$. With $m = \sqrt{b}$ and $\eta = \frac{\sqrt{y}}{12L\sqrt{d}}$, the number of PO calls equals to $T = Sm = O(\frac{\sqrt{d}}{\lambda\sqrt{y}} \log \frac{1}{\epsilon})$ and the number of SZO calls equals to $(S\mathcal{B}d + Smb) = O(\frac{\mathcal{B}d\sqrt{d}}{\lambda\sqrt{y}m} \log \frac{1}{\epsilon} + \frac{b\sqrt{d}}{\lambda\sqrt{y}} \log \frac{1}{\epsilon})$. \square

References

- Gao, X.; Jiang, B.; and Zhang, S. 2018. On the information-adaptive variants of the admm: an iteration complexity perspective. *Journal of Scientific Computing* 76(1):327–363.
- Ji, K.; Wang, Z.; Zhou, Y.; and Liang, Y. 2019. Improved zeroth-order variance reduced algorithms and analysis for nonconvex optimization. In *International Conference on Machine Learning*, 3100–3109.
- Liu, S.; Kailkhura, B.; Chen, P.-Y.; Ting, P.; Chang, S.; and Amini, L. 2018. Zeroth-order stochastic variance reduction for nonconvex optimization. In *Advances in Neural Information Processing Systems*, 3727–3737.