

# Supplemental Material for “Stochastic L-BFGS Revisited: Improved Convergence Rates and Practical Acceleration Strategies”

## S-1 Technical Lemmas

We first provide some technical lemmas that will be used in our subsequent proofs. Lemma S-1 states two classical inequalities of a strongly convex and smooth function, and its proof can be found in Boyd and Vandenberghe (2004, Chapter 9). Lemma S-2 states some results about importance sampling, and can be derived using definition of expectation and Lemma S-1. For a complete proof, see Xiao and Zhang (2014, Section 3). Lemma S-3 states a classical result on the distance between (uniformly) subsampled mean and (global) mean of a finite number of vectors. Its proof can be found in Konečný et al. (2016, Appendix B.A).

**Lemma S-1.** *If a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex and  $L$ -smooth, then for any  $\mathbf{x} \in \mathbb{R}^d$ ,*

$$\frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|^2, \quad (\text{S-1})$$

$$\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}^*\|^2, \quad (\text{S-2})$$

where  $\mathbf{x}^*$  denotes the unique minimizer of  $f$  on  $\mathbb{R}^d$ .

**Lemma S-2.** *Let  $f$  be defined as in (1) and satisfies Assumption 2. Define a distribution  $p$  with support  $[n]$  such that  $p_i = L_i/(n\bar{L})$ , for any  $i \in [n]$ . Then for any  $\mathbf{x} \in \mathbb{R}^d$ ,*

$$\mathbb{E}_{i \sim p} \left[ \frac{1}{np_i} \nabla f_i(\mathbf{x}) \right] = \nabla f(\mathbf{x}), \quad (\text{S-3})$$

$$\mathbb{E}_{i \sim p} \left[ \left\| \frac{1}{np_i} (\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}^*)) \right\|^2 \right] \leq 2\bar{L}(f(\mathbf{x}) - f(\mathbf{x}^*)), \quad (\text{S-4})$$

where  $\mathbf{x}^*$  denotes the unique minimizer of  $f$  on  $\mathbb{R}^d$ .

**Lemma S-3.** *Let  $\{\mathbf{z}_i\}_{i=1}^n \subseteq \mathbb{R}^d$  and define  $\bar{\mathbf{z}} \triangleq 1/n \sum_{i=1}^n \mathbf{z}_i$ . Uniformly sample a random subset  $\mathcal{S}$  of  $[n]$  with size  $b$  without replacement. Then*

$$\mathbb{E}_{\mathcal{S}} \left[ \left\| \frac{1}{n} \sum_{i \in \mathcal{S}} \mathbf{z}_i - \bar{\mathbf{z}} \right\|^2 \right] \leq \frac{n-b}{b(n-1)} \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i\|^2 \right). \quad (\text{S-5})$$

## S-2 Proof of Theorem 1

*Proof.* Fix an outer iteration  $s$  and consider an inner iteration  $t$ . Define  $r \triangleq \lfloor (\sum_{i=0}^{s-1} T_i + t)/L \rfloor$ ,<sup>11</sup> then the iteration in (10) becomes

$$\tilde{\mathbf{x}}_{s,t+1,r} = \tilde{\mathbf{x}}_{s,t,r} - \eta \tilde{\mathbf{v}}_{s,t,r}. \quad (\text{S-6})$$

In addition, from Lemmas 1 and 3, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_{s,t}} [\tilde{\mathbf{v}}_{s,t,r} | \mathcal{F}_{s,t}] &= \mathbf{H}_r^{1/2} \nabla f(\mathbf{x}_{s,t}) = \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}), \\ \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \left\| \tilde{\mathbf{v}}_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) \right\|^2 \middle| \mathcal{F}_{s,t} \right] &\leq \left\| \mathbf{H}^{1/2} \right\|^2 \frac{4\bar{L}}{b} (f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*) + f(\mathbf{x}^s) - f(\mathbf{x}^*)) \\ &\leq \frac{4\Gamma\bar{L}}{b} \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right), \end{aligned} \quad (\text{S-8})$$

where  $r' \triangleq \lfloor sm/L \rfloor$ . Next we express the distance between  $\tilde{\mathbf{x}}_{s,t+1,r}$  and  $\tilde{\mathbf{x}}_r^*$  as

$$\|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 = \|\tilde{\mathbf{x}}_{s,t,r} - \eta \tilde{\mathbf{v}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \quad (\text{S-9})$$

$$= \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 + 2\eta \left( \frac{\eta}{2} \|\tilde{\mathbf{v}}_{s,t,r}\|^2 - \langle \tilde{\mathbf{v}}_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle \right). \quad (\text{S-10})$$

<sup>11</sup>To avoid cluttered notations, we omit showing the dependence of  $r$  on  $s$  and  $t$ .

Defining  $\tilde{\delta}_{s,t,r} \triangleq \tilde{\mathbf{v}}_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r})$ , we now bound

$$\frac{\eta}{2} \|\tilde{\mathbf{v}}_{s,t,r}\|^2 - \langle \tilde{\mathbf{v}}_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle \leq - \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \right) - \langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^* \rangle - \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2. \quad (\text{S-11})$$

This is because

$$\tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) + \frac{\eta}{2} \|\tilde{\mathbf{v}}_{s,t,r}\|^2 - \langle \tilde{\mathbf{v}}_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle + \langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^* \rangle + \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \quad (\text{S-12})$$

$$= \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \langle \tilde{\mathbf{v}}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} + \frac{\eta}{2} \tilde{\mathbf{v}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle + \langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^* \rangle + \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \quad (\text{S-13})$$

$$\stackrel{(a)}{\leq} \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \frac{\Gamma\bar{L}}{2} \eta^2 \|\tilde{\mathbf{v}}_{s,t,r}\|^2 - \langle \tilde{\mathbf{v}}_{s,t,r} - \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^* \rangle + \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \quad (\text{S-14})$$

$$= \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \langle \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}), \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_{s,t,r} \rangle - \frac{\Gamma\bar{L}}{2} \|\eta \tilde{\mathbf{v}}_{s,t,r}\|^2 - \langle \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}), \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle + \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \quad (\text{S-15})$$

$$\stackrel{(b)}{\leq} \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) + \langle \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}), \tilde{\mathbf{x}}_r^* - \tilde{\mathbf{x}}_{s,t,r} \rangle + \frac{\gamma\bar{\mu}}{2} \|\tilde{\mathbf{x}}_r^* - \tilde{\mathbf{x}}_{s,t,r}\|^2 \quad (\text{S-16})$$

$$\stackrel{(c)}{\leq} \tilde{f}_r(\tilde{\mathbf{x}}_r^*), \quad (\text{S-17})$$

where (a) follows from  $\eta \leq 1/(\Gamma\bar{L})$ , (b) follows from the  $(\Gamma\bar{L})$ -smoothness of  $\tilde{f}_r$  and (c) follows from the  $(\gamma\bar{\mu})$ -strong convexity of  $\tilde{f}_r$  (see Lemma 2). Substituting (S-11) into (S-10), we have

$$\|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 \leq (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 - 2\eta \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \right) - 2\eta \langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^* \rangle \quad (\text{S-18})$$

$$= (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 - 2\eta \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \right) + 2\eta^2 \|\tilde{\delta}_{s,t,r}\|^2 - 2\eta \langle \tilde{\delta}_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle. \quad (\text{S-19})$$

Taking expectation with respect to (w.r.t.)  $\mathcal{B}_{s,t}$  and using (S-7) and (S-8), we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 \middle| \mathcal{F}_{s,t} \right] + 2\eta \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_{s,t} \right] \\ & \leq (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 + \frac{8}{b} \Gamma\bar{L}\eta^2 \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right). \end{aligned} \quad (\text{S-20})$$

By relaxing the factor  $1 - \eta\gamma\bar{\mu}$  to 1, we are able to telescope (S-20) over  $t = 0, \dots, m-1$  using the tower property of conditional expectations (Williams, 1991, Chapter 9.7) and obtain

$$\begin{aligned} & \mathbb{E}_{\mathcal{B}_{s,(m-1)}} \left[ \|\tilde{\mathbf{x}}_{s,m,r} - \tilde{\mathbf{x}}_r^*\|^2 \middle| \mathcal{F}_s \right] + 2m\eta \left( 1 - \frac{4}{b} \Gamma\bar{L}\eta \right) \frac{1}{m} \sum_{t=1}^m \mathbb{E}_{\mathcal{B}_{s,(t-1)}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_{s,t-1} \right] \\ & \leq \|\tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^*\|^2 + \frac{8}{b} \Gamma\bar{L}\eta^2 (1+m) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right), \end{aligned} \quad (\text{S-21})$$

where  $\mathcal{B}_{s,(m-1)} \triangleq \{\mathcal{B}_{s,i}\}_{i=0}^{m-1}$ . If we use option II to choose  $\mathbf{x}^{s+1}$  (in line 20), we have  $\tilde{\mathbf{x}}^{s+1,r''} = 1/m \sum_{i=1}^m \tilde{\mathbf{x}}_{s,t,r''}$ , where  $r'' \triangleq \lfloor (s+1)m/L \rfloor$ . Using (7) and Jensen's inequality, we have

$$\frac{1}{m} \sum_{t=1}^m \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \right) = \frac{1}{m} \sum_{t=1}^m \left( \tilde{f}_{r''}(\tilde{\mathbf{x}}_{s,t,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) \right) \geq \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*). \quad (\text{S-22})$$

Therefore, we have

$$\frac{1}{m} \sum_{t=1}^m \mathbb{E}_{\mathcal{B}_{s,(t-1)}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_{s,t-1} \right] \geq \mathbb{E}_{\mathcal{B}_{s,(m-1)}} \left[ \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) \middle| \mathcal{F}_s \right]. \quad (\text{S-23})$$

Alternatively, if we use option I to determine  $\mathbf{x}^{s+1}$  (in line 19), we still have (S-23) (with inequality replaced by equality). If we further use (S-1) in Lemma S-1 to bound  $\|\tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^*\|^2$  in (S-21), we have

$$\begin{aligned} & 2m\eta \left(1 - \frac{4}{b}\Gamma\bar{L}\eta\right) \mathbb{E}_{\mathcal{B}_{s,(m-1)}} \left[ \left| \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) \right| \mathcal{F}_s \right] \\ & \leq \left( \frac{8}{b}\Gamma\bar{L}\eta^2(1+m) + \frac{2}{\gamma\bar{\mu}} \right) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right). \end{aligned} \quad (\text{S-24})$$

Using (7) again and rearranging, we have

$$\mathbb{E} \left[ f(\mathbf{x}^{s+1}) - f(\mathbf{x}^*) \mid \mathcal{F}_s \right] \leq \rho(f(\mathbf{x}^s) - f(\mathbf{x}^*)). \quad (\text{S-25})$$

Taking expectation on both sides and we complete the proof.  $\square$

### S-3 Proof of Lemma 3

First from (S-4) we immediately have  $\mathbb{E}_{\mathcal{B}_{s,t}} [\mathbf{v}_{s,t} \mid \mathcal{F}_{s,t}] = \nabla f(\mathbf{x}_{s,t})$ . Then

$$\mathbb{E}_{\mathcal{B}_{s,t}} \left[ \|\mathbf{v}_{s,t} - \nabla f(\mathbf{x}_{s,t})\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-26})$$

$$= \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \left\| \frac{1}{b} \sum_{j=1}^b (1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^s)) + \nabla f(\mathbf{x}^s) - \nabla f(\mathbf{x}_{s,t})) \right\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-27})$$

$$\stackrel{(a)}{=} \frac{1}{b^2} \sum_{j=1}^b \mathbb{E}_{i_j} \left[ \left\| 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^s)) + \nabla f(\mathbf{x}^s) - \nabla f(\mathbf{x}_{s,t}) \right\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-28})$$

$$\leq \frac{1}{b^2} \sum_{j=1}^b \mathbb{E}_{i_j} \left[ \left\| 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^s)) \right\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-29})$$

$$= \frac{1}{b^2} \sum_{j=1}^b \mathbb{E}_{i_j} \left[ \left\| 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^*)) + 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}^*) - \nabla f_{i_j}(\mathbf{x}^s)) \right\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-30})$$

$$\leq \frac{2}{b^2} \sum_{j=1}^b \mathbb{E}_{i_j} \left[ \left\| 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^*)) \right\|^2 \mid \mathcal{F}_{s,t} \right] + \mathbb{E}_{i_j} \left[ \left\| 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}^s) - \nabla f_{i_j}(\mathbf{x}^*)) \right\|^2 \mid \mathcal{F}_{s,t} \right] \quad (\text{S-31})$$

$$\leq \frac{2}{b^2} \sum_{j=1}^b 2\bar{L}(f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*)) + 2\bar{L}(f(\mathbf{x}^s) - f(\mathbf{x}^*)) \quad (\text{S-32})$$

$$= \frac{4\bar{L}}{b} (f(\mathbf{x}_{s,t}) - f(\mathbf{x}^*) + f(\mathbf{x}^s) - f(\mathbf{x}^*)), \quad (\text{S-33})$$

where (a) follows from the independence of  $i_j$  and  $i_{j'}$  when  $j \neq j'$  and

$$\mathbb{E}_{i_j} \left[ 1/(np_{i_j}) (\nabla f_{i_j}(\mathbf{x}_{s,t}) - \nabla f_{i_j}(\mathbf{x}^s)) + \nabla f(\mathbf{x}^s) - \nabla f(\mathbf{x}_{s,t}) \mid \mathcal{F}_{s,t} \right] = 0, \quad (\text{S-34})$$

(b) follows from (S-3) and  $\mathbb{E} \left[ \|\mathbf{a} - \mathbb{E}[\mathbf{a}]\|^2 \mid \mathcal{G} \right] \leq \mathbb{E} \left[ \|\mathbf{a}\|^2 \mid \mathcal{G} \right]$  almost surely, for any  $\sigma$ -algebra  $\mathcal{G}$ , (c) follows from  $\|\mathbf{a} + \mathbf{b}\|^2 \leq 2\|\mathbf{a}\|^2 + 2\|\mathbf{b}\|^2$ , and (d) follows from (S-4).

### S-4 Proof of Lemma 4

For any  $\mathbf{x} \in \mathbb{R}^d$  and  $\mathcal{T} \subseteq [n]$  with cardinality  $b_{\mathbf{H}}$ ,

$$\bar{\mu}_{b_{\mathbf{H}}} \mathbf{I} \preceq \nabla^2 f_{\mathcal{T}}(\mathbf{x}) \preceq \bar{L}_{b_{\mathbf{H}}} \mathbf{I}, \quad (\text{S-35})$$

Define

$$\mathbf{V}_k \triangleq \mathbf{I} - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \quad \text{and} \quad \mathbf{Q}_k \triangleq \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}. \quad (\text{S-36})$$

Fix any  $r \in \mathbb{N}$ . Since

$$\mathbf{y}_k = \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_k, \quad (\text{S-37})$$

we have

$$\mathbf{V}_k = \mathbf{I} - \frac{\nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_k} \quad (\text{S-38})$$

$$= \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r)^{1/2} \left( \mathbf{I} - \frac{\tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^T}{\tilde{\mathbf{s}}_k^T \tilde{\mathbf{s}}_k} \right) \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r)^{-1/2}, \quad (\text{S-39})$$

where  $\tilde{\mathbf{s}}_k \triangleq \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r)^{1/2} \mathbf{s}_k$ . Hence

$$\|\mathbf{V}_k\| \leq \left\| \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r)^{1/2} \right\| \left\| \mathbf{I} - \frac{\tilde{\mathbf{s}}_k \tilde{\mathbf{s}}_k^T}{\|\tilde{\mathbf{s}}_k\|^2} \right\| \left\| \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r)^{-1/2} \right\| \quad (\text{S-40})$$

$$\leq \bar{L}_{b\text{H}}^{1/2} \bar{\mu}_{b\text{H}}^{-1/2} = \kappa_{b\text{H}}^{1/2}. \quad (\text{S-41})$$

Similarly,

$$\|\mathbf{Q}_k\| = \frac{\|\mathbf{s}_k\|^2}{\mathbf{s}_k^T \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \mathbf{s}_k} \leq \frac{1}{\bar{\mu}_{b\text{H}}}. \quad (\text{S-42})$$

By (5), we have

$$\left\| \mathbf{H}_r^{(k)} \right\| \leq \|\mathbf{V}_k\|^2 \left\| \mathbf{H}_r^{(k-1)} \right\| + \|\mathbf{Q}_k\| \leq \kappa_{b\text{H}} \left\| \mathbf{H}_r^{(k-1)} \right\| + \frac{1}{\bar{\mu}_{b\text{H}}}. \quad (\text{S-43})$$

Now, apply (S-43) repeatedly over  $k = r - M' + 1, \dots, r$  and we have

$$\|\mathbf{H}_r\| = \left\| \mathbf{H}_r^{(r)} \right\| \leq \kappa_{b\text{H}}^{M'} \left\| \mathbf{H}_r^{(r-M')} \right\| + \frac{1}{\bar{\mu}_{b\text{H}}} \sum_{i=0}^{M'-1} \kappa_{b\text{H}}^i \quad (\text{S-44})$$

$$\stackrel{\text{(a)}}{\leq} \frac{1}{\bar{\mu}_{b\text{H}}} \left( \kappa_{b\text{H}}^{M'} + \frac{\kappa_{b\text{H}}^{M'} - 1}{\kappa_{b\text{H}} - 1} \right) \quad (\text{S-45})$$

$$\leq \frac{1}{\bar{\mu}_{b\text{H}}} \kappa_{b\text{H}}^{M'} \left( 1 + \frac{1}{\kappa_{b\text{H}} - 1} \right) \quad (\text{S-46})$$

$$\stackrel{\text{(b)}}{\leq} \frac{\kappa_{b\text{H}}^{M+1}}{\bar{\mu}_{b\text{H}} (\kappa_{b\text{H}} - 1)}, \quad (\text{S-47})$$

where (a) follows from

$$\mathbf{H}_r^{(r-M')} = \frac{\mathbf{s}_r^T \mathbf{y}_r}{\mathbf{y}_r^T \mathbf{y}_r} \mathbf{I} \quad (\text{S-48})$$

and

$$\frac{\mathbf{s}_r^T \mathbf{y}_r}{\mathbf{y}_r^T \mathbf{y}_r} = \frac{\tilde{\mathbf{s}}_r^T \tilde{\mathbf{s}}_r}{\tilde{\mathbf{s}}_r^T \nabla^2 f_{\mathcal{T}_r}(\bar{\mathbf{x}}_r) \tilde{\mathbf{s}}_r} \leq \frac{1}{\bar{\mu}_{b\text{H}}}, \quad (\text{S-49})$$

and (b) follows from  $\kappa_{b\text{H}} \geq 1$  and  $M' \leq M$ . To show  $\gamma = 1/(M+1)\bar{L}_{b\text{H}}$ , it suffices to show

$$\|\mathbf{B}_r\| \leq (M+1)\bar{L}_{b\text{H}}, \quad (\text{S-50})$$

where  $\mathbf{B}_r \triangleq \mathbf{H}_r^{-1}$ . To begin with, we rewrite (5) as Nocedal and Wright (2006, Chapter 6)

$$\mathbf{B}_r^{(k)} = \mathbf{B}_r^{(k-1)} - \frac{\mathbf{B}_r^{(k-1)} \mathbf{s}_k \mathbf{s}_k^T \mathbf{B}_r^{(k-1)}}{\mathbf{s}_k^T \mathbf{B}_r^{(k-1)} \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k}, \quad (\text{S-51})$$

where  $\mathbf{B}_r^{(k)} \triangleq \left(\mathbf{H}_r^{(k)}\right)^{-1}$ , for each  $k \in [r - M' + 1 : r]$ . Since the update rule (5) preserves positive definiteness of  $\mathbf{H}_r^{(k)}$  (Nocedal and Wright, 2006), we have that  $\mathbf{B}_r^{(k)} \succ 0$  for all  $k \in [r - M' + 1 : r]$ . Define  $\widehat{\mathbf{s}}_k \triangleq \left(\mathbf{B}_r^{(k)}\right)^{1/2} \mathbf{s}_k$ , then

$$\left\| \mathbf{B}_r^{(k-1)} - \frac{\mathbf{B}_r^{(k-1)} \mathbf{s}_k \mathbf{s}_k^T \mathbf{B}_r^{(k-1)}}{\mathbf{s}_k^T \mathbf{B}_r^{(k-1)} \mathbf{s}_k} \right\| \leq \left\| \left(\mathbf{B}_r^{(k-1)}\right)^{1/2} \right\|^2 \left\| \mathbf{I} - \frac{\widehat{\mathbf{s}}_k \widehat{\mathbf{s}}_k^T}{\|\widehat{\mathbf{s}}_k\|^2} \right\| = \left\| \mathbf{B}_r^{(k-1)} \right\|. \quad (\text{S-52})$$

By (S-51), we have that

$$\left\| \mathbf{B}_r^{(k)} \right\| \leq \left\| \mathbf{B}_r^{(k-1)} \right\| + \left\| \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right\| \leq \left\| \mathbf{B}_r^{(k-1)} \right\| + \bar{L}_{b_H}, \quad (\text{S-53})$$

since

$$\left\| \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right\| = \frac{\|\mathbf{y}_k\|^2}{\mathbf{s}_k^T \mathbf{y}_k} = \frac{\widetilde{\mathbf{s}}_k^T \nabla^2 f_{T_r}(\bar{\mathbf{x}}_r) \widetilde{\mathbf{s}}_k}{\widetilde{\mathbf{s}}_k^T \widetilde{\mathbf{s}}_k} \leq \bar{L}_{b_H}. \quad (\text{S-54})$$

Then we have

$$\|\mathbf{B}_r\| = \left\| \mathbf{B}_r^{(r)} \right\| \leq \left\| \mathbf{B}_r^{(r-M')} \right\| + M' \bar{L}_{b_H} \stackrel{(a)}{\leq} (M+1) \bar{L}_{b_H}, \quad (\text{S-55})$$

where (a) follows from (S-48) and (S-54).

## S-5 Proof of Proposition 2

Our proof leverages a refined telescoping of (S-20). For each  $t \in [m]$ , we multiply both sides of (S-20) by  $(1 - \eta\gamma\bar{\mu})^{m-t}$  and obtain

$$\begin{aligned} & (1 - \eta\gamma\bar{\mu})^{m-t} \left( \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \middle| \mathcal{F}_{s,t} \right] + 2\eta \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_{s,t} \right] \right) \\ & \leq (1 - \eta\gamma\bar{\mu})^{m-t} \left( (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t-1,r} - \tilde{\mathbf{x}}_r^*\|^2 + \frac{8}{b} \Gamma \bar{L} \eta^2 \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t-1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right) \right). \end{aligned} \quad (\text{S-56})$$

Telescope (S-56) over  $t = 1, \dots, m$  and we have

$$\begin{aligned} & 2c\eta \left( 1 - \frac{4}{b} \frac{\Gamma \bar{L} \eta}{1 - \eta\gamma\bar{\mu}} \right) \frac{1}{c} \sum_{t=1}^m (1 - \eta\gamma\bar{\mu})^{m-t} \mathbb{E}_{\mathcal{B}_{s,(t)}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_s \right] \\ & \leq (1 - \eta\gamma\bar{\mu})^m \left\| \tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^* \right\|^2 + \frac{8}{b} \Gamma \bar{L} \eta^2 \left( (1 - \eta\gamma\bar{\mu})^m + c \right) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right). \end{aligned} \quad (\text{S-57})$$

Now we consider using option IV to choose  $\mathbf{x}^{s+1}$ . Since  $0 < \beta \leq 1 - \eta\gamma\bar{\mu}$ , using (7) and Jensen's inequality in a similar manner as in (S-22) and taking expectations, we have

$$\frac{1}{c} \sum_{t=1}^m (1 - \eta\gamma\bar{\mu})^{m-t} \mathbb{E}_{\mathcal{B}_{s,(t)}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_s \right] \geq \mathbb{E}_{\mathcal{B}_{s,(m)}} \left[ \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) \middle| \mathcal{F}_s \right]. \quad (\text{S-58})$$

Alternatively, if  $\mathbf{x}^{s+1}$  is determined using option III, the definition of distribution  $q$  still yields (S-58). Finally, using (S-1) in Lemma S-1 to bound  $\|\tilde{\mathbf{x}}^{s,r'} - \tilde{\mathbf{x}}_{r'}^*\|^2$  in (S-57), we have

$$\begin{aligned} & 2c'\eta \left( 1 - \frac{4}{b} \frac{\Gamma \bar{L} \eta}{1 - \eta\gamma\bar{\mu}} \right) \mathbb{E}_{\mathcal{B}_{s,(m)}} \left[ \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}_{r''}(\tilde{\mathbf{x}}_{r''}^*) \middle| \mathcal{F}_s \right] \\ & \leq \left( \frac{8}{b} \Gamma \bar{L} \eta^2 (1 + c') + \frac{2}{\gamma\bar{\mu}} \right) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right). \end{aligned} \quad (\text{S-59})$$

Taking expectation on both sides, rearranging and using (7), we reach (31).

## S-6 Proof of Proposition 3

The subsampled gradient strategy essentially introduces errors in  $\{\mathbf{g}_s\}_{s \geq 0}$ . Therefore, we explicit model this error by  $\mathbf{e}_s \triangleq \tilde{\mathbf{g}}_s - \mathbf{g}_s$ . We first bound the second moment of  $\mathbf{e}_s$ . Since  $\mathbf{g}_s = \nabla f(\mathbf{x}^s)$ , by Lemma S-3 and  $\tilde{b}_s \geq nS^2\alpha_s/(S^2\alpha_s + (n-1)\xi^2\bar{\rho}^{2s})$ , we have for any  $s \in (S]$ ,

$$\mathbb{E}_{\tilde{\mathcal{B}}_s} [\|\mathbf{e}_s\|^2] \leq \frac{n - \tilde{b}_s}{\tilde{b}_s(n-1)} \alpha_s \leq \frac{\xi^2}{S^2} \bar{\rho}^{2s}. \quad (\text{S-60})$$

As a result,

$$\mathbb{E}_{\tilde{\mathcal{B}}_s} [\|\mathbf{e}_s\|] \leq \frac{\xi}{S} \bar{\rho}^s. \quad (\text{S-61})$$

The introduction of random sets  $\{\tilde{\mathcal{B}}_s\}_{s \geq 0}$  requires us to redefine the filtration  $\{\mathcal{F}_{s,t}\}_{s \geq 0, t \in (m-1]}$  as

$$\mathcal{F}_{s,t} \triangleq \sigma \left( \{\tau_j\}_{j=0}^{s-1} \cup \{\tilde{\mathcal{B}}_j\}_{j \in (s)} \cup \{\mathcal{B}_{i,j}\}_{i \in (s-1), j \in (m-1)} \cup \{\mathcal{B}_{s,j}\}_{j=0}^{t-1} \cup \{\mathcal{T}_j\}_{j=0}^{\lfloor (sm+t)/L \rfloor} \right). \quad (\text{S-62})$$

As usual, we define  $\mathcal{F}_s \triangleq \mathcal{F}_{s,0}$ . In the sequel, we follow the naming convention in coordinate transformation framework as in Section 4.1. In particular, we define  $\tilde{\mathbf{e}}_{s,r'} \triangleq \mathbf{H}_{r'}^{1/2} \mathbf{e}_s$ . Now, define  $\tilde{\mathbf{v}}'_{s,t,r} \triangleq \tilde{\mathbf{v}}_{s,t,r} + \tilde{\mathbf{e}}_{s,r'}$ , then from (S-7) and (S-8), we have

$$\mathbb{E}_{\mathcal{B}_{s,t}} [\tilde{\mathbf{v}}'_{s,t,r} | \mathcal{F}_{s,t}] = \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) + \tilde{\mathbf{e}}_{s,r'}, \quad (\text{S-63})$$

$$\mathbb{E}_{\mathcal{B}_{s,t}} \left[ \left\| \tilde{\mathbf{v}}'_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) \right\|^2 \middle| \mathcal{F}_{s,t} \right] \leq \frac{4\Gamma\bar{L}}{b} \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right) + \|\tilde{\mathbf{e}}_{s,r'}\|^2, \quad (\text{S-64})$$

Following a similar argument as in the proof of Theorem 1, we have

$$\|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 \leq (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 - 2\eta \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \right) + 2\eta^2 \left\| \tilde{\delta}'_{s,t,r} \right\|^2 - 2\eta \left\langle \tilde{\delta}'_{s,t,r}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \right\rangle, \quad (\text{S-65})$$

where  $\tilde{\delta}'_{s,t,r} \triangleq \tilde{\mathbf{v}}'_{s,t,r} - \nabla \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r})$ . Taking expectation w.r.t.  $\mathcal{B}_{s,t}$  and using (S-63) and (S-64), we have

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \|\tilde{\mathbf{x}}_{s,t+1,r} - \tilde{\mathbf{x}}_r^*\|^2 \middle| \mathcal{F}_{s,t} \right] + 2\eta \mathbb{E}_{\mathcal{B}_{s,t}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t+1,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_{s,t} \right] &\leq (1 - \eta\gamma\bar{\mu}) \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\|^2 \\ &+ \frac{8}{b} \Gamma\bar{L}\eta^2 \left( \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) + \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right) + 2\eta^2 \|\tilde{\mathbf{e}}_{s,r'}\|^2 - 2\eta \langle \tilde{\mathbf{e}}_{s,r'}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle. \end{aligned} \quad (\text{S-66})$$

Using Cauchy-Schwartz inequality, we have

$$-2\eta \langle \tilde{\mathbf{e}}_{s,r'}, \tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^* \rangle \leq 2\eta \|\tilde{\mathbf{e}}_{s,r'}\| \|\tilde{\mathbf{x}}_{s,t,r} - \tilde{\mathbf{x}}_r^*\| \leq 2\eta \|\tilde{\mathbf{e}}_{s,r'}\| \left\| \mathbf{H}_r^{-1/2} \right\| \|\mathbf{x}_{s,t} - \mathbf{x}^*\| \leq 2\eta B\gamma^{-1/2} \|\tilde{\mathbf{e}}_{s,r'}\|. \quad (\text{S-67})$$

Substituting (S-67) into (S-66), and using the telescoping techniques in Section S-5, we have

$$\begin{aligned} 2c\eta \left( 1 - \frac{4}{b} \frac{\Gamma\bar{L}\eta}{1 - \eta\gamma\bar{\mu}} \right) \frac{1}{c} \sum_{t=1}^m (1 - \eta\gamma\bar{\mu})^{m-t} \mathbb{E}_{\mathcal{B}_{s,(t)}} \left[ \tilde{f}_r(\tilde{\mathbf{x}}_{s,t,r}) - \tilde{f}_r(\tilde{\mathbf{x}}_r^*) \middle| \mathcal{F}_s \right] &\leq (1 - \eta\gamma\bar{\mu})^m \|\tilde{\mathbf{x}}^{s,r'} - \mathbf{x}_r^*\|^2 \\ &+ \frac{8}{b} \Gamma\bar{L}\eta^2 ((1 - \eta\gamma\bar{\mu})^m + c) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right) + 2\eta ((1 - \eta\gamma\bar{\mu})^m + c) \left( B\gamma^{-1/2} \|\tilde{\mathbf{e}}_{s,r'}\| + \eta \|\tilde{\mathbf{e}}_{s,r'}\|^2 \right). \end{aligned} \quad (\text{S-68})$$

With either option III or IV and (S-1), we have

$$\begin{aligned} 2c'\eta \left( 1 - \frac{4}{b} \frac{\Gamma\bar{L}\eta}{1 - \eta\gamma\bar{\mu}} \right) \mathbb{E}_{\mathcal{B}_{s,(m)}} \left[ \tilde{f}_{r''}(\tilde{\mathbf{x}}^{s+1,r''}) - \tilde{f}(\tilde{\mathbf{x}}_{r''}^*) \middle| \mathcal{F}_s \right] \\ \leq \left( \frac{8}{b} \Gamma\bar{L}\eta^2 (1 + c') + \frac{2}{\gamma\bar{\mu}} \right) \left( \tilde{f}_{r'}(\tilde{\mathbf{x}}^{s,r'}) - \tilde{f}_{r'}(\tilde{\mathbf{x}}_{r'}^*) \right) + 2\eta (1 + c') \left( B\gamma^{-1/2} \|\tilde{\mathbf{e}}_{s,r'}\| + \eta \|\tilde{\mathbf{e}}_{s,r'}\|^2 \right). \end{aligned} \quad (\text{S-69})$$

Taking expectation on both sides and using (7), we have

$$\mathbb{E} [f(\mathbf{x}^{s+1}) - f(\mathbf{x}^*)] \leq \bar{\rho} \mathbb{E} [f(\mathbf{x}^s) - f(\mathbf{x}^*)] + \left(1 + \frac{1}{c'}\right) \frac{b}{b - 4\Gamma\bar{L}\eta/(1 - \eta\gamma\bar{\mu})} \left( B\gamma^{-1/2} \mathbb{E} [\|\tilde{\mathbf{e}}_{s,r'}\|] + \eta \mathbb{E} [\|\tilde{\mathbf{e}}_{s,r'}\|^2] \right). \quad (\text{S-70})$$

From (S-60) and (S-61), we have

$$\mathbb{E} [\|\tilde{\mathbf{e}}_{s,r'}\|] \leq \frac{\Gamma^{1/2}\xi}{S} \bar{\rho}^s, \quad (\text{S-71})$$

$$\mathbb{E} [\|\tilde{\mathbf{e}}_{s,r'}\|^2] \leq \frac{\Gamma\xi^2}{S^2} \bar{\rho}^{2s}. \quad (\text{S-72})$$

Substituting (S-71) and (S-72) into (S-70), we have

$$\begin{aligned} \mathbb{E} [f(\mathbf{x}^{s+1}) - f(\mathbf{x}^*)] &\leq \bar{\rho} \mathbb{E} [f(\mathbf{x}^s) - f(\mathbf{x}^*)] + \left(1 + \frac{1}{c'}\right) \frac{b}{b - 4\Gamma\bar{L}\eta/(1 - \eta\gamma\bar{\mu})} \left( \frac{\kappa_{\text{H}}^{1/2} B \xi}{S} \bar{\rho}^s + \frac{\eta \Gamma \xi^2}{S^2} \bar{\rho}^{2s} \right) \\ &\leq \bar{\rho} \mathbb{E} [f(\mathbf{x}^s) - f(\mathbf{x}^*)] + \left(1 + \frac{1}{c'}\right) \frac{b}{b - 4\Gamma\bar{L}\eta/(1 - \eta\gamma\bar{\mu})} \left( \kappa_{\text{H}}^{1/2} B + \eta \Gamma \xi \right) \frac{\xi}{S} \bar{\rho}^s \end{aligned} \quad (\text{S-73})$$

Applying (S-73) recursively and we complete the proof.

## S-7 Future Work and Open Problems

For future work, we propose to pursue from the following two directions:

- Based on our novel coordinate transformation framework, *proximal* and momentum-based stochastic L-BFGS algorithms can be developed and analyzed. The proximal variant can greatly enable our algorithm to handle composite (convex) objective function, thereby greatly extending the applications of our method. The accelerated variant can potentially improve the linear convergence rate of our algorithm, and thus reduce the total computational complexity.
- The linear convergence of the strategy in Section 6.3 can be analyzed theoretically. In Figure 4, we observed the linear convergence empirically. This indicates that linear convergence of Algorithm 1 may still hold under such a strategy.

Besides future work, there also exists an open problem. Although we have improved the linear convergence rate and computational complexity of Algorithm 1 as compared to those in Moritz et al. (2016) and Gower et al. (2016), it seems our improved complexity in (29) is still inferior to the computational complexity of SVRG. In SVRG, the complexity is  $O((n + \kappa)d \log(1/\epsilon))$  (Xiao and Zhang, 2014), so our complexity (29) has an additional multiplicative factor  $\kappa_{\text{H}}$ . This contradicts the experimental results in Moritz et al. (2016), Gower et al. (2016) and Section 7.2, where stochastic L-BFGS-type algorithms have been repeatedly shown to outperform their first-order counterparts. A careful analysis will reveal that the additional  $\kappa_{\text{H}}$  arises from the uniform spectral bound of the metric matrices  $\{\mathbf{H}_r\}_{r \geq 0}$ . This uniform bound is effectively a worst-case bound, and does not reflect the *local curvature information* contained in recent iterates at any time  $(s, t)$ . Since the judicious use of curvature information serves as a very important reason for the fast convergence of stochastic quasi-Newton algorithms, such information should also be reflected in theoretical analysis as well (possibly in an adaptive spectral bound for  $\{\mathbf{H}_r\}_{r \geq 0}$ ). We also believe an effective adaptive bound can potentially improve the complexity result (29). *In short, how to obtain a (computational) complexity bound of stochastic L-BFGS algorithm that is better (or at least as good as) that of SVRG will be an interesting problem to consider.*