

## A Fokker-Planck Equation and Backward Kolmogorov Equation

In this section, we introduce the Fokker-Planck Equation and the Backward Kolmogorov equation. Fokker-Planck equation addresses the evolution of probability density  $p(\mathbf{x})$  that associates with the SDE. We give the following specific definition.

**Definition A.1** (Fokker–Planck Equation). Let  $p(\mathbf{x}, t)$  be the probability density at time  $t$  of the stochastic differential equation and denote  $p_0(\mathbf{x})$  the initial probability density. Then

$$\partial_t p(\mathbf{x}, t) = \mathcal{L}^* p(\mathbf{x}, t), \quad p(\mathbf{x}, 0) = p_0(\mathbf{x}),$$

where  $\mathcal{L}^*$  is the formal adjoint of  $\mathcal{L}$ .

Fokker-Planck equation gives us a way to find whether there exists a stationary distribution for the SDE. It can be shown [30] that for the stochastic differential equation (1.2), its stationary distribution exists and satisfies

$$\pi(d\mathbf{x}) = \frac{1}{Q} e^{-\beta F_n(\mathbf{x})}, \quad Q = \int e^{-\beta F_n(\mathbf{x})} d\mathbf{x}. \quad (\text{A.1})$$

This is also known as Gibbs measure.

Backward Kolmogorov equation describes the evolution of  $\mathbb{E}[g(\mathbf{X}(t)) | \mathbf{X}(0) = \mathbf{x}]$  with  $g$  being a smooth test function.

**Definition A.2** (Backward Kolmogorov Equation). Let  $\mathbf{X}(t)$  solves the stochastic differential equation (1.2). Let  $u(\mathbf{x}, t) = \mathbb{E}[g(\mathbf{X}(t)) | \mathbf{X}(0) = \mathbf{x}]$ , we have

$$\partial_t u(\mathbf{x}, t) = \mathcal{L} u(\mathbf{x}, t), \quad u(\mathbf{x}, 0) = g(\mathbf{x}).$$

Now consider doing first order Taylor expansion on  $u(\mathbf{x}, t)$ , we have

$$\begin{aligned} u(\mathbf{x}, t) &= u(\mathbf{x}, 0) + \frac{\partial}{\partial t} u(\mathbf{x}, t)|_{t=0} \cdot (t - 0) + O(t^2) \\ &= g(\mathbf{x}) + t \mathcal{L} g(\mathbf{x}) + O(t^2). \end{aligned} \quad (\text{A.2})$$

## B Proof of Corollaries

In this section, we provide the proofs of corollaries for iteration complexity in our main theory section.

*Proof of Corollary 3.4.* To ensure the iterate error converge to  $\epsilon$  precision, we need

$$\Theta e^{-\lambda K \eta} \leq \frac{\epsilon}{2}, \quad \frac{C_\psi \eta}{\beta} \leq \frac{\epsilon}{2}.$$

The second inequality can be easily satisfied with  $\eta = O(\epsilon)$  and the first inequality implies

$$K \geq \frac{1}{\lambda \eta} \log \left( \frac{2\Theta}{\epsilon} \right).$$

Combining with  $\eta = O(\epsilon)$  and  $\Theta = O(d^2/\rho^{d/2})$ , we obtain the iteration complexity

$$K = O\left(\frac{d}{\epsilon \lambda} \cdot \log \left( \frac{1}{\epsilon} \right)\right),$$

which completes the proof. □

*Proof of Corollary 3.7.* To ensure the iterate error of SGLD converging to  $\epsilon$  precision, we require the following inequalities to hold

$$C_1 \sqrt{\beta} \Gamma(M \sqrt{\Gamma} + G) K \eta \left[ \frac{n - B}{B(n - 1)} \right]^{1/4} \leq \frac{\epsilon}{3}, \quad \Theta e^{-\lambda K \eta} \leq \frac{\epsilon}{3}, \quad \frac{C_\psi \eta}{\beta} \leq \frac{\epsilon}{3}.$$

The third inequality can be easily satisfied with  $\eta = O(\epsilon)$ . For the second inequality, similar as in the proof of Corollary 3.4, we have

$$K\eta \geq \frac{1}{\lambda} \log \left( \frac{3\Theta}{\epsilon} \right).$$

Since  $\epsilon < 1$ , we know that  $\log(1/\epsilon)$  will not go to zero when  $\epsilon$  goes to zero. In fact, if we set  $\eta = O(\epsilon)$  and  $K = O(d/(\lambda\epsilon) \log(1/\epsilon))$ , the first term in (3.2) scales as

$$C_1 \sqrt{\beta} \Gamma (M\sqrt{\Gamma} + G) K \eta \left[ \frac{n-B}{B(n-1)} \right]^{1/4} = O \left( \frac{d^{3/2} K \eta}{B^{1/4}} \right) = O \left( \frac{d^{3/2}}{B^{1/4} \lambda} \log \left( \frac{1}{\epsilon} \right) \right).$$

Therefore, within  $K = O(d/(\epsilon\lambda) \cdot \log(1/\epsilon))$  iterations, the iterate error of SGLD scales as

$$O \left( \frac{d^{3/2}}{B^{1/4} \lambda} \log \left( \frac{1}{\epsilon} \right) + \epsilon \right).$$

□

*Proof of Corollary 3.11.* Similar to previous proofs, in order to achieve an  $\epsilon$ -precision iterate error for SVRG-LD, we require

$$C_1 \Gamma K^{3/4} \eta \left[ \frac{L\beta M^2(n-B)}{B(n-1)} \left( 9\eta(M^2\Gamma + G^2) + \frac{d}{\beta} \right) \right]^{1/4} \leq \frac{\epsilon}{3}, \quad \Theta e^{-\lambda K \eta} \leq \frac{\epsilon}{3}, \quad \frac{C_\psi \eta}{\beta} \leq \frac{\epsilon}{3}.$$

By previous proofs we know that the second and third inequalities imply  $\eta = O(\epsilon)$  and  $K\eta = O(1/\lambda \log(3\Theta/\epsilon))$  respectively. Combining with the first inequality, we have

$$\eta^{1/4} = O \left( \frac{B^{1/4} \epsilon}{(K\eta)^{3/4} d^{5/4} L^{1/4}} \right)$$

Combining with the first inequality, we have

$$\eta = O \left( \min \left\{ \frac{B\epsilon^4}{(K\eta)^3 d^5 L}, \epsilon \right\} \right)$$

Combining the above requirements yields

$$K = O \left( \frac{Ld^5}{B\lambda^4 \epsilon^4} \log^4 \left( \frac{1}{\epsilon} \right) + \frac{1}{\epsilon} \right). \quad (\text{B.1})$$

For gradient complexity, note that for each iteration we need  $B$  stochastic gradient evaluations and we also need in total  $K/L$  full gradient calculations. Therefore, the gradient complexity for SVRG-LD is

$$O(K \cdot B + K/L \cdot n) = \tilde{O} \left( \left( \frac{n}{B} + L \right) \frac{1}{\epsilon^4} + \left( \frac{n}{L} + B \right) \frac{1}{\epsilon} \right) \cdot e^{\tilde{O}(d)}.$$

If we solve for the best  $B$  and  $L$ , we obtain  $B = \sqrt{n}\epsilon^{-3/2}$ ,  $L = \sqrt{n}\epsilon^{3/2}$ . Therefore, we have the optimal gradient complexity for SVRG-LD as

$$\tilde{O} \left( \frac{\sqrt{n}}{\epsilon^{5/2}} \right) \cdot e^{\tilde{O}(d)}.$$

□

## C Proof of Technical Lemmas

In this section, we provide proofs of the technical lemmas used in the proof of our main theory.

### C.1 Proof of Lemma 4.1

Geometric ergodicity of dynamical systems has been studied a lot in the literature [47, 40]. In particular, Roberts and Tweedie [47] proved that even when the diffusion converges exponentially fast to its stationary distribution, the Euler-Maruyama discretization in (2.2) may still lose the convergence properties and examples for Langevin diffusion can be found therein. To further address this problem, [40] built their analysis of ergodicity for SDEs on a *minorization* condition and the existence of a Lyapunov function. In time discretization of dynamics systems, they studied how time-discretization affects the minorization condition and the Lyapunov structure. For the self-containedness of our analysis, we present the minorization condition on a compact set  $\mathcal{C}$  as follows.

**Proposition C.1.** There exist  $t_0 \in \mathbb{R}$  and  $\xi > 0$  such that the Markov process  $\{\mathbf{X}(t)\}$  satisfies

$$\mathbb{P}(\mathbf{X}(t_0) \in A | \mathbf{X}(0) = \mathbf{x}) \geq \xi \nu(A),$$

for any  $A \in \mathcal{B}(\mathbb{R}^d)$ , some fixed compact set  $\mathcal{C} \in \mathcal{B}(\mathbb{R}^d)$ , and  $\mathbf{x} \in \mathcal{C}$ , where  $\mathcal{B}(\mathbb{R}^d)$  denotes the Borel  $\sigma$ -algebra on  $\mathbb{R}^d$  and  $\nu$  is a probability measure with  $\nu(\mathcal{C}^c) = 0$  and  $\nu(\mathcal{C}) = 1$ .

Proposition C.1 does not always hold for a Markov process generated by an arbitrary SDE. However, for Langevin diffusion (1.2) studied in this paper, Mattingly et al. [40] proved that this minorization condition actually holds under the dissipative and smooth assumptions (see Corollary 7.4 in Mattingly et al. [40]). For more explanation on the existence and robustness of the minorization condition under discretization approximations for Langevin diffusion, we refer interested readers to Corollary 7.5 and the proof of Theorem 6.2 in Mattingly et al. [40]. Now we are going to prove Lemma 4.1, which requires the following useful lemmas:

**Lemma C.2.** Let  $V(\mathbf{x}) = C + \|\mathbf{x}\|_2^2$  be a function on  $\mathbb{R}^d$ , where  $C > 0$  is a constant. Denote the expectation with Markov process  $\{\mathbf{X}(t)\}$  starting at  $\mathbf{x}$  by  $\mathbb{E}^\mathbf{x}[\cdot] = \mathbb{E}[\cdot | \mathbf{X}(0) = \mathbf{x}]$ . Under Assumption 3.2, we have

$$\mathbb{E}^\mathbf{x}[V(\mathbf{X}(t))] \leq e^{-2mt} V(\mathbf{x}) + \frac{b + m + d/\beta}{m} (1 - e^{-2mt}),$$

for all  $\mathbf{x} \in \mathbb{R}^d$ .

**Lemma C.3.** (Theorem 7.3 in Mattingly et al. [40]) Under Assumptions 3.1 and 3.2, let  $V(\mathbf{x}) = C_0 + M/2 \|\mathbf{x}\|_2^2$  be an essential quadratic function. The numerical approximation (2.1) (GLD) of Langevin diffusion (1.2) has a unique invariant measure  $\mu$  and for all test function  $g$  such that  $|g| \leq V$ , we have

$$|\mathbb{E}[g(\mathbf{X}_k)] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq C\kappa\rho^{-d/2}(1 + \kappa e^{m\eta}) \exp\left(-\frac{2mk\eta\rho^d}{\log(\kappa)}\right),$$

where  $\rho \in (0, 1)$ ,  $C > 0$  are absolute constants, and  $\kappa = 2M(b + m + d)/m$ .

*Proof of Lemma 4.1.* The proof is majorly adapted from that of Theorem 7.3 and Corollary 7.5 in Mattingly et al. [40]. By Assumption 3.1,  $F_n$  is  $M$ -smooth. Thus we have

$$F_n(\mathbf{x}) \leq F_n(\mathbf{y}) + \langle \nabla F_n(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{M}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ . By Lemma D.1 and choosing  $\mathbf{y} = \mathbf{0}$ , this immediately implies that  $F_n(\mathbf{x})$  can always be bounded by a quadratic function  $V(\mathbf{x})$ , i.e.,

$$F_n(\mathbf{x}) \leq \frac{M}{2} V(\mathbf{x}) = \frac{M}{2} (C_0 + \|\mathbf{x}\|_2^2).$$

Therefore  $V(\mathbf{x})$  is an essentially quadratic Lyapunov function such that  $|F_n(\mathbf{x})| \leq MV(\mathbf{x})/2$  for  $\mathbf{x} \in \mathbb{R}^d$ . By Lemma C.2 the Lyapunov function satisfies

$$\mathbb{E}^{\mathbf{x}_0}[V(\mathbf{X}(t))] \leq e^{-2mt} V(\mathbf{x}_0) + \frac{b + m + d/\beta}{m} (1 - e^{-2mt}).$$

According to Corollary 7.5 in Mattingly et al. [40], the Markov chain  $\{\mathbf{X}_k\}_{k=1,2,\dots,K}$  satisfies

$$\mathbb{E}^{\mathbf{x}_0}[MV(\mathbf{X}_1)/2] \leq e^{-2m\eta}[MV(\mathbf{x}_0)/2] + \frac{M(b + m + d/\beta)}{2m}. \quad (\text{C.1})$$

Recall the GLD update formula defined in (2.1)

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k) + \sqrt{2\eta\beta^{-1}} \cdot \epsilon_k.$$

Define  $F'(\mathbf{X}_k) = \beta F_n(\mathbf{X}_k)$  and  $\eta' = \eta/\beta$ , we have

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \eta' \nabla F'(\mathbf{X}_k) + \sqrt{2\eta'} \cdot \epsilon_k. \quad (\text{C.2})$$

This suggests that the result for  $\beta \neq 1$  is equivalent to rescaling  $\eta$  to  $\eta/\beta$  and  $F_n(\cdot)$  to  $\beta F_n(\cdot)$ . Therefore, in the following proof, we will assume that  $\beta = 1$  and then rescale  $\eta$ ,  $F_n(\cdot)$  at last. Similar tricks are used in Raginsky et al. [45], Zhang et al. [53]. Under Assumptions 3.1 and 3.2, it is proved that Euler-Maruyama approximation of Langevin dynamics (1.2) has a unique invariant measure  $\mu$  on  $\mathbb{R}^d$ . Denote  $\mathbf{X}^\mu$  as a random vector which is sampled from measure  $\mu$ . By Lemma C.3, for all test function  $g$  such that  $|g| \leq V$ , it holds that

$$|\mathbb{E}[g(\mathbf{X}_k)] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq C\kappa'\rho^{-d/2}(1 + \kappa'e^{m\eta}) \exp\left(-\frac{2mk\eta\rho^d}{\log(\kappa')}\right),$$

where  $\rho, \delta \in (0, 1), C > 0$  are absolute constants, and  $\kappa' = 2M(b + m + d)/m$ . Take  $F_n$  as the test function and  $\mathbf{X}_0 = \mathbf{0}$ , and by rescaling  $\eta$  and  $F_n(\cdot)$  (dissipative and smoothness parameters), we have

$$|\mathbb{E}[F_n(\mathbf{X}_k)] - \mathbb{E}[F_n(\mathbf{X}^\mu)]| \leq C\kappa\rho^{-d/2}(1 + \kappa e^{m\eta}) \exp\left(-\frac{2mk\eta\rho^d}{\log(\kappa)}\right),$$

where  $\kappa = 2M(b\beta + m\beta + d)/m$ . □

## C.2 Proof of Lemma 4.2

To prove Lemma 4.2, we lay down the following supporting lemma, of which the derivation is inspired and adapted from Chen et al. [12].

**Lemma C.4.** Under Assumptions 3.1 and 3.2, the Markov chain  $\{\mathbf{X}_k\}_{k=1}^K$  generated by Algorithm 1 satisfies

$$\left| \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[F_n(\mathbf{X}_k) | \mathbf{X}_0 = \mathbf{x}] - \bar{F} \right| \leq C_\psi \left( \frac{\beta}{\eta K} + \frac{\eta}{\beta} \right),$$

where  $\bar{F} = \int F_n(\mathbf{x})\pi(d\mathbf{x})$  with  $\pi$  being the Gibbs measure for the Langevin diffusion (1.2).

*Proof of Lemma 4.2.* By definition we have

$$|\mathbb{E}[F_n(\mathbf{X}^\mu)] - \mathbb{E}[F_n(\mathbf{X}^\pi)]| = \left| \int F_n(\mathbf{x})\mu(d\mathbf{x}) - \int F_n(\mathbf{x})\pi(d\mathbf{x}) \right|. \quad (\text{C.3})$$

For simplicity, we denote the average  $\int F_n(\mathbf{x})\pi(d\mathbf{x})$  as  $\bar{F}_n$ . Since  $\mu$  is the ergodic limit of the Markov chain generated by the GLD process, for a given test function  $F_n$ , we have

$$\int F_n(\mathbf{x})\mu(d\mathbf{x}) = \int \mathbb{E}[F_n(\mathbf{X}_k) | \mathbf{X}_0 = \mathbf{x}] \cdot \mu(d\mathbf{x}).$$

Since  $\mu$  and  $\pi$  are two invariant measures, we consider the case where  $K \rightarrow \infty$ . Take average over  $K$  steps  $\{\mathbf{X}_k\}_{k=0}^{K-1}$  we have

$$\int F_n(\mathbf{x})\mu(d\mathbf{x}) = \lim_{K \rightarrow \infty} \int \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[F_n(\mathbf{X}_k) | \mathbf{X}_0 = \mathbf{x}] \cdot \mu(d\mathbf{x}). \quad (\text{C.4})$$

Submitting (C.4) back into (C.3) yields

$$\begin{aligned} |\mathbb{E}[F_n(\mathbf{X}^\mu)] - \mathbb{E}[F_n(\mathbf{X}^\pi)]| &= \lim_{K \rightarrow \infty} \left| \int \left[ \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[F_n(\mathbf{X}_k) | \mathbf{X}_0 = \mathbf{x}] - \bar{F} \right] \cdot \mu(d\mathbf{x}) \right| \\ &\leq \lim_{K \rightarrow \infty} \int \left| \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[F_n(\mathbf{X}_k) | \mathbf{X}_0 = \mathbf{x}] - \bar{F} \right| \cdot \mu(d\mathbf{x}). \end{aligned} \quad (\text{C.5})$$

Apply Lemma C.4 with  $g$  chosen as  $F_n$  we further bound (C.5) by

$$\begin{aligned} |\mathbb{E}[F_n(\mathbf{X}^\mu)] - \mathbb{E}[F_n(\mathbf{X}^\pi)]| &\leq C_\psi \cdot \lim_{K \rightarrow \infty} \int \left( \frac{\beta}{\eta K} + \frac{\eta}{\beta} \right) \cdot \mu(d\mathbf{x}) \\ &= C_\psi \cdot \lim_{K \rightarrow \infty} \left( \frac{\beta}{\eta K} + \frac{\eta}{\beta} \right) \\ &= \frac{C_\psi \eta}{\beta}. \end{aligned}$$

□

### C.3 Proof of Lemma 4.4

Lemma 4.4 gives the upper bound of function value gap between the GLD iterates and the SGLD iterates. To bound the difference between  $F_n(\mathbf{X}_K)$  and  $F_n(\mathbf{Y}_K)$ , we need the following lemmas.

**Lemma C.5.** Under Assumptions 3.1 and 3.2, for any  $\mathbf{x} \in \mathbb{R}^d$ , it holds that

$$\mathbb{E} \left\| \nabla F_n(\mathbf{x}) - \frac{1}{B} \sum_{i \in I_k} \nabla f_i(\mathbf{x}) \right\|_2^2 \leq \frac{4(n-B)(M\|\mathbf{x}\|_2 + G)^2}{B(n-1)},$$

where  $B = |I_k|$  is the mini-batch size and  $G = \max_{i=1, \dots, n} \{\|\nabla f_i(\mathbf{x}^*)\|_2\} + bM/m$ .

The following lemma describes the  $L^2$  bound for discrete processes  $\mathbf{X}_k$  (GLD),  $\mathbf{Y}_k$  (SGLD) and  $\mathbf{Z}_k$  (SVRG-LD). Note that for SGLD, similar result is also presented in Raginsky et al. [45].

**Lemma C.6.** Under Assumptions 3.1 and 3.2, for sufficiently small step size  $\eta$ , suppose the initial points of Algorithms 1, 2 and 3 are chosen at  $\mathbf{0}$ , then the  $L^2$  bound of the GLD process (2.1), SGLD process (2.2) and SVRG-LD process (2.3) can be uniformly bounded by

$$\max\{\mathbb{E}[\|\mathbf{X}_k\|_2^2], \mathbb{E}[\|\mathbf{Y}_k\|_2^2], \mathbb{E}[\|\mathbf{Z}_k\|_2^2]\} \leq \Gamma \quad \text{where} \quad \Gamma := 2 \left( 1 + \frac{1}{m} \right) \left( b + 2G^2 + \frac{d}{\beta} \right),$$

for any  $k = 0, 1, \dots, K$ , where  $G = \max_{i=1, \dots, n} \{\|\nabla f_i(\mathbf{x}^*)\|_2\} + bM/m$ .

The following lemma gives out the upper bound for the exponential  $L^2$  bound of  $\mathbf{X}_k$ .

**Lemma C.7.** Under Assumptions 3.1 and 3.2, for sufficiently small step size  $\eta < 1$  and the inverse temperature satisfying  $\beta \geq \max\{2/(m - M^2\eta), 4\eta\}$ , it holds that

$$\log \mathbb{E}[\exp(\|\mathbf{X}_k\|_2^2)] \leq \|\mathbf{X}_0\|_2^2 + \frac{2\beta(b + G^2) + 2d}{\beta - 4\eta} k\eta.$$

**Lemma C.8.** [44, 45] For any two probability density functions  $\mu, \nu$  with bounded second moments, let  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $C^1$  function such that

$$\|\nabla g(\mathbf{x})\|_2 \leq C_1 \|\mathbf{x}\|_2 + C_2, \forall \mathbf{x} \in \mathbb{R}^d$$

for some constants  $C_1, C_2 \geq 0$ . Then

$$\left| \int_{\mathbb{R}^d} g(\mathbf{x}) d\mu - \int_{\mathbb{R}^d} g(\mathbf{x}) d\nu \right| \leq (C_1\sigma + C_2) \mathcal{W}_2(\mu, \nu),$$

where  $\mathcal{W}_2$  is the 2-Wasserstein distance and  $\sigma^2 = \max \left\{ \int_{\mathbb{R}^d} \|\mathbf{x}\|_2^2 \mu(d\mathbf{x}), \int_{\mathbb{R}^d} \|\mathbf{x}\|_2^2 \nu(d\mathbf{x}) \right\}$ .

**Lemma C.9.** (Corollary 2.3 in Bolley and Villani [5]) Let  $\nu$  be a probability measure on  $\mathbb{R}^d$ . Assume that there exist  $\mathbf{x}_0$  and a constant  $\alpha > 0$  such that  $\int \exp(\alpha \|\mathbf{x} - \mathbf{x}_0\|_2^2) d\nu(\mathbf{x}) < \infty$ . Then for any probability measure  $\mu$  on  $\mathbb{R}^d$ , it satisfies

$$\mathcal{W}_2(\mu, \nu) \leq C_\nu \left( \sqrt{D_{\text{KL}}(\mu|\nu)} + (D_{\text{KL}}(\mu|\nu)/2)^{1/4} \right),$$

where  $C_\nu$  is defined as

$$C_\nu = \inf_{\mathbf{x}_0 \in \mathbb{R}^d, \alpha > 0} \sqrt{\frac{1}{\alpha} \left( \frac{3}{2} + \log \int \exp(\alpha \|\mathbf{x} - \mathbf{x}_0\|_2^2) d\nu(\mathbf{x}) \right)}.$$

*Proof of Lemma 4.4.* Let  $P_K, Q_K$  denote the probability measures for GLD iterate  $\mathbf{X}_K$  and SGLD iterate  $\mathbf{Y}_K$  respectively. Applying Lemma C.8 to probability measures  $P_K$  and  $Q_K$  yields

$$|\mathbb{E}[F_n(\mathbf{Y}_K)] - \mathbb{E}[F_n(\mathbf{X}_K)]| \leq (C_1\sqrt{\Gamma} + C_2)\mathcal{W}_2(Q_K, P_K), \quad (\text{C.6})$$

where  $C_1, C_2 > 0$  are absolute constants and  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$  is the upper bound for both  $\mathbb{E}[\|\mathbf{X}_k\|_2^2]$  and  $\mathbb{E}[\|\mathbf{Y}_k\|_2^2]$  according to Lemma C.6. We further bound the  $\mathcal{W}_2$  distance via the KL-divergence by Lemma C.9 as follows

$$\mathcal{W}_2(Q_K, P_K) \leq \Lambda(\sqrt{D_{\text{KL}}(Q_K||P_K)} + \sqrt[4]{D_{\text{KL}}(Q_K||P_K)}), \quad (\text{C.7})$$

where  $\Lambda = \sqrt{3/2 + \log \mathbb{E}_{P_K}[\exp(\|\mathbf{X}_K\|_2^2)]}$ . Applying Lemma C.7 we obtain  $\Lambda = \sqrt{(6 + 2\Gamma)K\eta}$ . Therefore, we only need to bound the KL-divergence between density functions  $P_K$  and  $Q_K$ . To this end, we introduce a continuous-time Markov process  $\{\mathbf{D}(t)\}_{t \geq 0}$  to bridge the gap between diffusion  $\{\mathbf{X}(t)\}_{t \geq 0}$  and its numerical approximation  $\{\mathbf{X}_k\}_{k=0,1,\dots,K}$ . Define

$$d\mathbf{D}(t) = b(\mathbf{D}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t), \quad (\text{C.8})$$

where  $b(\mathbf{D}(t)) = -\sum_{k=0}^{\infty} \nabla F(\mathbf{X}(\eta k))\mathbb{1}\{t \in [\eta k, \eta(k+1))\}$ . Integrating (C.8) on interval  $[\eta k, \eta(k+1))$  yields

$$\mathbf{D}(\eta(k+1)) = \mathbf{D}(\eta k) - \eta \nabla F(\mathbf{D}(\eta k)) + \sqrt{2\eta\beta^{-1}} \cdot \boldsymbol{\epsilon}_k,$$

where  $\boldsymbol{\epsilon}_k \sim N(\mathbf{0}, \mathbf{I}_{d \times d})$ . This implies that the distribution of random vector  $(\mathbf{X}_1, \dots, \mathbf{X}_K)$  is equivalent to that of  $(\mathbf{D}(\eta), \dots, \mathbf{D}(\eta K))$ . Similarly, for  $\mathbf{Y}_k$  we define

$$d\widetilde{\mathbf{M}}(t) = c(\widetilde{\mathbf{M}}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t),$$

where the drift coefficient is defined as  $c(\widetilde{\mathbf{M}}(t)) = -\sum_{k=0}^{\infty} g_k(\widetilde{\mathbf{M}}(\eta k))\mathbb{1}\{t \in [\eta k, \eta(k+1))\}$  and  $g_k(\mathbf{x}) = 1/B \sum_{i \in I_k} \nabla f_i(\mathbf{x})$  is a mini-batch of the full gradient with  $I_k$  being a random subset of  $\{1, 2, \dots, n\}$  of size  $B$ . Now we have that the distribution of random vector  $(\mathbf{Y}_1, \dots, \mathbf{Y}_K)$  is equivalent to that of  $(\widetilde{\mathbf{M}}(\eta), \dots, \widetilde{\mathbf{M}}(\eta K))$ . However, the process  $\widetilde{\mathbf{M}}(t)$  is not Markov due to the randomness of the stochastic gradient  $g_k$ . Therefore, we define the following Markov process which has the same one-time marginals as

$$d\mathbf{M}(t) = h(\mathbf{M}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t), \quad (\text{C.9})$$

where  $h(\cdot) = -\mathbb{E}[g_k(\widetilde{\mathbf{M}}(\eta k))\mathbb{1}\{t \in [\eta k, \eta(k+1))\}|\widetilde{\mathbf{M}}(t) = \cdot]$  is the conditional expectation of the left end point of the interval which  $\widetilde{\mathbf{M}}(t)$  lies in. Let  $\mathbb{P}_t$  denote the distribution of  $\mathbf{D}(t)$  and  $\mathbb{Q}_t$  denote the distribution of  $\mathbf{M}(t)$ . By (C.8) and (C.9), the Radon-Nikodym derivative of  $\mathbb{P}_t$  with respect to  $\mathbb{Q}_t$  is given by the following Girsanov formula [38]

$$\begin{aligned} \frac{d\mathbb{P}_t}{d\mathbb{Q}_t}(\mathbf{M}) &= \exp \left\{ \sqrt{\frac{\beta}{2}} \int_0^t (h(\mathbf{M}(s)) - b(\mathbf{M}(s)))^\top (d\mathbf{M}(s) - h(\mathbf{M}(s))ds) \right. \\ &\quad \left. - \frac{\beta}{4} \int_0^t \|h(\mathbf{M}(s)) - b(\mathbf{M}(s))\|_2^2 ds \right\}. \end{aligned}$$

Since Markov processes  $\{\mathbf{D}(t)\}_{t \geq 0}$  and  $\{\mathbf{M}(t)\}_{t \geq 0}$  are constructed based on Markov chains  $\mathbf{X}_k$  and  $\mathbf{Y}_k$ , by data-processing inequality the K-L divergence between  $P_K$  and  $Q_K$  can be bounded by

$$\begin{aligned} D_{KL}(Q_K||P_K) &\leq D_{KL}(\mathbb{Q}_{\eta K}||\mathbb{P}_{\eta K}) \\ &= -\mathbb{E} \left[ \log \left( \frac{d\mathbb{P}_{\eta K}}{d\mathbb{Q}_{\eta K}}(\mathbf{M}) \right) \right] \\ &= \frac{\beta}{4} \int_0^{\eta K} \mathbb{E}[\|h(\mathbf{M}(r)) - b(\mathbf{M}(r))\|_2^2] dr, \end{aligned} \quad (\text{C.10})$$

where in the last equality we used the fact that  $d\mathbf{B}(t)$  follows Gaussian distribution independently for any  $t \geq 0$ . By definition, we know that both  $h(\mathbf{M}(r))$  and  $b(\mathbf{M}(r))$  are step functions when

$r \in [\eta k, \eta(k+1))$  for any  $k$ . This observation directly yields

$$\begin{aligned} \int_0^{\eta K} \mathbb{E}[\|h(\mathbf{M}(r)) - b(\mathbf{M}(r))\|_2^2] dr &\leq \sum_{k=0}^{K-1} \int_{\eta k}^{\eta(k+1)} \mathbb{E}[\|g_k(\widetilde{\mathbf{M}}(\eta k)) - \nabla F_n(\widetilde{\mathbf{M}}(\eta k))\|_2^2] dr \\ &= \eta \sum_{k=0}^{K-1} \mathbb{E}[\|g_k(\mathbf{Y}_k) - \nabla F_n(\mathbf{Y}_k)\|_2^2], \end{aligned}$$

where the first inequality is due to Jensen's inequality and the convexity of function  $\|\cdot\|^2$ , and the last equality is due to the equivalence in distribution. By Lemmas C.5 and C.6, we further have

$$\int_0^{\eta K} \mathbb{E}[\|h(\mathbf{M}(r)) - b(\mathbf{M}(r))\|_2^2] dr \leq \frac{4\eta K(n-B)(M\Gamma + G)^2}{B(n-1)}. \quad (\text{C.11})$$

Submitting (C.10) and (C.11) into (C.7), we have

$$\begin{aligned} \mathcal{W}_2(Q_K, P_K) &\leq \Lambda \left( \sqrt{\frac{\beta\eta K(n-B)(M\Gamma + G)^2}{B(n-1)}} + \sqrt[4]{\frac{\beta\eta K(n-B)(M\Gamma + G)^2}{B(n-1)}} \right) \\ &\leq \Lambda \sqrt{\frac{\beta\eta K\sqrt{n-B}(M\Gamma + G)^2}{\sqrt{B(n-1)}}}. \end{aligned} \quad (\text{C.12})$$

Combining (C.6) with (C.12), we obtain the expected function value gap between SGLD and GLD:

$$|\mathbb{E}[F(\mathbf{Y}_k)] - \mathbb{E}[F(\mathbf{X}_k)]| \leq C_1 \Gamma \sqrt{K\eta} \left[ \frac{\beta\eta K\sqrt{n-B}(M\sqrt{\Gamma} + G)^2}{\sqrt{B(n-1)}} \right]^{1/2},$$

where we adopt the fact that  $K\eta > 1$  and assume that  $C_1 \geq C_2$ . □

#### C.4 Proof of Lemma 4.5

Similar to the proof of Lemma 4.4, to bound the difference between  $F_n(\mathbf{X}_K)$  and  $F_n(\mathbf{Z}_K)$ , we need the following lemmas.

**Lemma C.10.** Under Assumptions 3.1 and 3.2, for each iteration  $k = sL + \ell$  in Algorithm 3, it holds that

$$\mathbb{E}\|\widetilde{\nabla}_k - \nabla F_n(\mathbf{Z}_k)\|_2^2 \leq \frac{M^2(n-B)}{B(n-1)} \mathbb{E}\|\mathbf{Z}_k - \widetilde{\mathbf{Z}}^{(s)}\|_2^2,$$

where  $\widetilde{\nabla}_k = 1/B \sum_{i_k \in I_k} (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\widetilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\widetilde{\mathbf{Z}}^{(s)})$  and  $B = |I_k|$  is the mini-batch size.

*Proof of Lemma 4.5.* Denote  $Q_K^Z$  as the probability density functions for  $\mathbf{Z}_K$ . For the simplicity of notation, we omit the index  $Z$  in the remaining part of this proof when no confusion arises. Similar as in the proof of Lemma 4.4, we first apply Lemma C.8 to probability measures  $P_K$  for  $\mathbf{X}_K$  and  $Q_K^Z$  for  $\mathbf{Z}_K$ , and obtain the following upper bound of function value gap

$$|\mathbb{E}[F_n(\mathbf{Z}_K)] - \mathbb{E}[F_n(\mathbf{X}_K)]| \leq (C_1\sqrt{\Gamma} + C_2)\mathcal{W}_2(Q_K^Z, P_K), \quad (\text{C.13})$$

where  $C_1, C_2 > 0$  are absolute constants and  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$  is the upper bound for both  $\mathbb{E}[\|\mathbf{X}_k\|_2^2]$  and  $\mathbb{E}[\|\mathbf{Z}_k\|_2^2]$  according to Lemma C.6. Further by Lemma C.9, the  $\mathcal{W}_2$  distance can be bounded by

$$\mathcal{W}_2(Q_K^Z, P_K) \leq \Lambda(\sqrt{D_{\text{KL}}(Q_K^Z||P_K)} + \sqrt[4]{D_{\text{KL}}(Q_K^Z||P_K)}), \quad (\text{C.14})$$

where  $\Lambda = \sqrt{3/2 + \log \mathbb{E}_{P_K}[e^{\|\mathbf{X}_K\|_2^2}]}$ . Applying Lemma C.7 we obtain  $\Lambda = \sqrt{(6 + 2\Gamma)K\eta}$ . Therefore, we need to bound the KL-divergence between density functions  $P_K$  and  $Q_K^Z$ . Similar to the proof of Lemma 4.4, we define a continuous-time Markov process associated with  $\mathbf{Z}_k$  as follows

$$d\widetilde{\mathbf{N}}(t) = p(\widetilde{\mathbf{N}}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t),$$

where  $p(\tilde{\mathbf{N}}(t)) = -\sum_{k=0}^{\infty} \tilde{\nabla}_k \mathbb{1}\{t \in [\eta k, \eta(k+1))\}$  and  $\tilde{\nabla}_k$  is the semi-stochastic gradient at  $k$ -th iteration of SVRG-LD. We have that the distribution of random vector  $(\mathbf{Z}_1, \dots, \mathbf{Z}_K)$  is equivalent to that of  $(\tilde{\mathbf{N}}(\eta), \dots, \tilde{\mathbf{N}}(\eta K))$ . However,  $\tilde{\mathbf{N}}(t)$  is not Markov due to the randomness of  $\tilde{\nabla}_k$ . We define the following Markov process which has the same one-time marginals as  $\tilde{\mathbf{N}}(t)$

$$d\mathbf{N}(t) = q(\mathbf{N}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t), \quad (\text{C.15})$$

where  $q(\cdot) = -\mathbb{E}[\tilde{\nabla}_k \mathbb{1}\{t \in [\eta k, \eta(k+1))\} | p(\tilde{\mathbf{N}}(t)) = \cdot]$ . Let  $\mathbb{Q}_t^Z$  denote the distribution of  $\mathbf{N}(t)$ . By (C.8) and (C.15), the Radon-Nikodym derivative of  $\mathbb{P}_t$  with respect to  $\mathbb{Q}_t^Z$  is given by the Girsanov formula [38]

$$\begin{aligned} \frac{d\mathbb{P}_t}{d\mathbb{Q}_t^Z}(\mathbf{N}) &= \exp \left\{ \sqrt{\frac{\beta}{2}} \int_0^t (q(\mathbf{N}(r)) - b(\mathbf{N}(r)))^\top (d\mathbf{N}(r) - h(\mathbf{N}(r))dr) \right. \\ &\quad \left. - \frac{\beta}{4} \int_0^t \|q(\mathbf{N}(r)) - b(\mathbf{N}(r))\|_2^2 dr \right\}. \end{aligned}$$

Since Markov processes  $\{\mathbf{D}(t)\}_{t \geq 0}$  and  $\{\mathbf{N}(t)\}_{t \geq 0}$  are constructed based on  $\mathbf{X}_k$  and  $\mathbf{Z}_k$ , by data-processing inequality the K-L divergence between  $P_K$  and  $Q_K^Z$  in (C.14) can be bounded by

$$\begin{aligned} D_{\text{KL}}(Q_K^Z \| P_K) &\leq D_{\text{KL}}(\mathbb{Q}_{\eta K}^Z \| \mathbb{P}_{\eta K}) \\ &= -\mathbb{E} \left[ \log \left( \frac{d\mathbb{P}_{\eta K}}{d\mathbb{Q}_{\eta K}^Z}(\mathbf{N}) \right) \right] \\ &= \frac{\beta}{4} \int_0^{\eta K} \mathbb{E} [\|q(\mathbf{N}(r)) - b(\mathbf{N}(r))\|_2^2] dr. \end{aligned} \quad (\text{C.16})$$

where in the last equality we used the fact that  $d\mathbf{B}(t)$  follows Gaussian distribution independently for any  $t \geq 0$ . By definition, we know that both  $q(\mathbf{N}(r))$  and  $b(\mathbf{N}(r))$  are step functions when  $r \in [\eta k, \eta(k+1))$  for any  $k$ . This observation directly yields

$$\begin{aligned} \int_0^{\eta K} \mathbb{E} [\|q(\mathbf{N}(r)) - b(\mathbf{N}(r))\|_2^2] dr &\leq \sum_{k=0}^{K-1} \int_{\eta k}^{\eta(k+1)} \mathbb{E} [\|\tilde{\nabla}_k(\tilde{\mathbf{N}}(\eta k)) - \nabla F_n(\tilde{\mathbf{N}}(\eta k))\|_2^2] dr \\ &= \eta \sum_{k=0}^{K-1} \mathbb{E} [\|\tilde{\nabla}_k(\mathbf{Z}_k) - \nabla F_n(\mathbf{Z}_k)\|_2^2], \end{aligned}$$

where the first inequality is due to Jensen's inequality and the convexity of function  $\|\cdot\|_2^2$ , and the last equality is due to the equivalence in distribution. Combine the above results we obtain

$$\begin{aligned} D_{\text{KL}}(Q_K^Z \| P_K) &\leq \frac{\beta\eta}{4} \sum_{k=0}^{K-1} \mathbb{E} [\|\tilde{\nabla}_k - \nabla F_n(\mathbf{Z}_k)\|_2^2] \\ &\leq \frac{\beta\eta}{4} \sum_{s=0}^{K/L} \sum_{\ell=0}^{L-1} \mathbb{E} [\|\tilde{\nabla}_{sL+\ell} - \nabla F_n(\mathbf{Z}_{sL+\ell})\|_2^2], \end{aligned} \quad (\text{C.17})$$

where the second inequality follows the fact that  $k = sL + \ell \leq (s+1)L$  for some  $\ell = 0, 1, \dots, L-1$ . Applying Lemma C.10, the inner summation in (C.17) yields

$$\sum_{\ell=0}^{L-1} \mathbb{E} [\|\tilde{\nabla}_{sL+\ell} - \nabla F_n(\mathbf{Z}_{sL+\ell})\|_2^2] \leq \sum_{\ell=0}^{L-1} \frac{M^2(n-B)}{B(n-1)} \mathbb{E} \|\mathbf{Z}_{sL+\ell} - \tilde{\mathbf{Z}}^{(s)}\|_2^2. \quad (\text{C.18})$$

Note that we have

$$\begin{aligned} &\mathbb{E} \|\mathbf{Z}_{sL+\ell} - \tilde{\mathbf{Z}}^{(s)}\|_2^2 \\ &= \mathbb{E} \left\| \sum_{u=0}^{\ell-1} \eta (\nabla f_{i_{sL+u}}(\mathbf{Z}_{sL+u}) - \nabla f_{i_{sL+u}}(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)})) - \sum_{u=0}^{\ell-1} \sqrt{\frac{2\eta}{\beta}} \epsilon_{sL+u} \right\|_2^2 \\ &\leq \ell \sum_{u=0}^{\ell-1} \mathbb{E} [2\eta^2 \|\nabla f_{i_{sL+u}}(\mathbf{Z}_{sL+u}) - \nabla f_{i_{sL+u}}(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)})\|_2^2] + \sum_{u=0}^{\ell-1} \frac{4\eta d}{\beta} \\ &\leq 4\ell\eta \left( 9\ell\eta(M^2\Gamma^2 + G^2) + \frac{d}{\beta} \right), \end{aligned} \quad (\text{C.19})$$



where the first inequality holds due to the triangle inequality for the first summation term, the second one follows from Lemma D.1 and Lemma C.6. Submit (C.19) back into (C.18) we have

$$\begin{aligned} \sum_{\ell=0}^{L-1} \mathbb{E}[\|\tilde{\nabla}_{sL+\ell} - \nabla F_n(\mathbf{Z}_{sL+\ell})\|_2^2] &\leq \frac{4\eta M^2(n-B)}{B(n-1)} \sum_{\ell=0}^{L-1} \left( 9\ell^2\eta(M^2\Gamma^2 + G^2) + \frac{\ell d}{\beta} \right) \\ &\leq \frac{4\eta M^2(n-B)}{B(n-1)} \left( 3L^3\eta(M^2\Gamma + G^2) + \frac{dL^2}{2\beta} \right), \end{aligned} \quad (\text{C.20})$$

Since (C.20) does not depend on the outer loop index  $i$ , submitting it into (C.17) yields

$$\frac{\beta\eta}{4} \sum_{k=0}^{K-1} \mathbb{E}[\|\tilde{\nabla}_k - \nabla F_n(\mathbf{Z}_k)\|_2^2] \leq \frac{\eta^2 K L M^2(n-B)(3L\eta\beta(M^2\Gamma + G^2) + d/2)}{B(n-1)}. \quad (\text{C.21})$$

Combining (C.13), (C.14) (C.17) and (C.21), we obtain

$$|\mathbb{E}[F_n(\mathbf{Z}_K)] - \mathbb{E}[F_n(\mathbf{X}_K)]| \leq C_1\Gamma\sqrt{K\eta} \left[ \frac{\eta^2 K L M^2(n-B)(3L\eta\beta(M^2\Gamma + G^2) + d/2)}{B(n-1)} \right]^{1/4}.$$

where we use the fact that  $K\eta > 1$ ,  $\eta < 1$  and assume that  $C_1 \geq C_2$ .  $\square$

## D Proof of Auxiliary Lemmas

In this section, we prove additional lemmas used in Appendix C.

### D.1 Proof of Lemma C.2

*Proof.* Applying Itô's Lemma yields

$$dV(\mathbf{X}(t)) = -2\langle \mathbf{X}(t), \nabla F_n(\mathbf{X}(t)) \rangle dt + \frac{2d}{\beta} dt + 2\sqrt{\frac{2}{\beta}} \langle \mathbf{X}(t), d\mathbf{B}(t) \rangle. \quad (\text{D.1})$$

Multiplying  $e^{2mt}$  to both sides of the above equation, where  $m > 0$  is the dissipative constant, we obtain

$$\begin{aligned} 2me^{2mt}V(\mathbf{X}(t))dt + e^{2mt}dV(\mathbf{X}(t)) &= 2me^{2mt}V(\mathbf{X}(t))dt - 2e^{2mt}\langle \mathbf{X}(t), \nabla F_n(\mathbf{X}(t)) \rangle dt \\ &\quad + \frac{2d}{\beta}e^{2mt}dt + \sqrt{\frac{8}{\beta}}e^{2mt}\langle \mathbf{X}(t), d\mathbf{B}(t) \rangle. \end{aligned}$$

We integrate the above equation from time 0 to  $t$  and have

$$\begin{aligned} V(\mathbf{X}(t)) &= e^{-2mt}V(\mathbf{X}_0) + 2m \int_0^t e^{2m(s-t)}V(\mathbf{X}(s))ds - 2 \int_0^t e^{2m(s-t)}\langle \mathbf{X}(s), \nabla F_n(\mathbf{X}(s)) \rangle ds \\ &\quad + \frac{2d}{\beta} \int_0^t e^{2m(s-t)}ds + 2\sqrt{\frac{2}{\beta}} \int_0^t e^{2m(s-t)}\langle \mathbf{X}(s), d\mathbf{B}(s) \rangle. \end{aligned} \quad (\text{D.2})$$

Note that by Assumption 3.2, we have

$$\begin{aligned} -2 \int_0^t e^{2m(s-t)}\langle \mathbf{X}(s), \nabla F_n(\mathbf{X}(s)) \rangle ds &\leq -2 \int_0^t e^{2m(s-t)}(m\|\mathbf{X}(s)\|_2^2 - b)ds \\ &= -2m \int_0^t e^{2m(s-t)}V(\mathbf{X}(s))ds + \frac{b+m}{m}(1 - e^{-2mt}). \end{aligned} \quad (\text{D.3})$$

Combining (D.2) and (D.3), and taking expectation over  $\mathbf{X}(t)$  with initial point  $\mathbf{x}$ , we get

$$\begin{aligned} \mathbb{E}^\mathbf{x}[V(\mathbf{X}(t))] &\leq e^{-2mt}V(\mathbf{x}) + \frac{b+m}{m}(1 - e^{-2mt}) + \frac{d}{m\beta}(1 - e^{-2mt}) \\ &= e^{-2mt}V(\mathbf{x}) + \frac{b+m+d/\beta}{m}(1 - e^{-2mt}), \end{aligned}$$

where we employed the fact that  $d\mathbf{B}(s)$  follows Gaussian distribution with zero mean and is independent with  $\mathbf{X}(s)$ .  $\square$

## D.2 Proof of Lemma C.3

Here we provide a sketch of proof to refine the parameters in the results of Mattingly et al. [40]. For detailed proof, we refer interested readers to Theorem 7.3 in Mattingly et al. [40].

*Proof.* Denote  $\kappa = 2M(b + m + d)/m$  according to Lemma C.2 where  $b, m$  are the dissipative parameters. We also define  $\phi = \rho^d$  with some constant  $0 < \rho < 1$ . Let  $\{\mathbf{X}_{l\tau}\}_{l=0,1,\dots}$  be a sub-sampled chain from  $\{\mathbf{X}_k\}_{k=0,1,\dots}$  at sample rate  $\tau > 0$ . By the proof of Theorem 2.5 in Mattingly et al. [40], we obtain the following result

$$|\mathbb{E}[g(\mathbf{X}_{l\tau})] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq \kappa[\bar{V} + 1](1 - \phi)^{\alpha l\tau} + \sqrt{2}V(\mathbf{x}_0)\delta^{l\tau}\kappa^{\alpha l\tau/2}\frac{1}{\sqrt{\phi}}, \quad (\text{D.4})$$

where  $\mathbf{X}^\mu$  follows the invariant distribution of Markov process  $\{\mathbf{X}_k\}_{k=0,1,\dots}$ ,  $\bar{V} = 2 \sup_{\mathbf{x} \in \mathcal{C}} V(\mathbf{x})$  is a bounded constant,  $\delta \in (e^{-2m\eta}, 1)$  is a constant, and  $\alpha \in (0, 1)$  is chosen small enough such that  $\delta\kappa^{\alpha/2} \leq 1$ . In particular, we choose  $\alpha \in (0, 1)$  such that  $\delta\kappa^{\alpha/2} \leq (1 - \phi)^\alpha$ , which yields

$$\alpha \leq \frac{\log(1/\delta)}{\log(\sqrt{\kappa}/(1 - \phi))} \leq \frac{\log(1/\delta)}{\log(\sqrt{\kappa})},$$

where the last inequality is due to  $1 - \phi < 1$ . Submitting the choice of  $\alpha$  into (D.4) we have

$$\begin{aligned} |\mathbb{E}[g(\mathbf{X}_{l\tau})] - \mathbb{E}[g(\mathbf{X}^\mu)]| &\leq \frac{2\sqrt{2}\kappa}{\sqrt{\phi}}[\bar{V} + 1]V(\mathbf{x}_0)(1 - \phi)^{l\tau \log(1/\delta)/\log(\sqrt{\kappa})} \\ &= \frac{2\sqrt{2}\kappa}{\sqrt{\phi}}[\bar{V} + 1]V(\mathbf{x}_0)e^{l\tau \log(r)}, \end{aligned} \quad (\text{D.5})$$

where  $r = (1 - \phi)^{\log(1/\delta)/\log(\sqrt{\kappa})}$  is defined as the contraction parameter. Note that by Taylor's expansion we have

$$\log r = \log(1 - (1 - r)) = -(1 - r) - \frac{(1 - r)^2}{2} - \frac{(1 - r)^3}{3} - \dots \leq -(1 - r), \quad (\text{D.6})$$

when  $|1 - r| \leq 1$ . By definition  $r = (1 - \phi)^{\log(1/\delta)/\log(\sqrt{\kappa})}$  and  $\phi = \rho^d$  where  $\rho \in (0, 1)$  is a constant. Since it is more interesting to deal with the situation where dimension parameter  $d$  is large enough and not negligible, we can always assume that  $|\phi| = \rho^d$  is sufficiently small such that for any  $0 < \zeta < 1$

$$(1 - \phi)^\zeta = 1 - \zeta\phi + \zeta(\zeta - 1)/2\phi^2 + \dots + \binom{\zeta}{n}(-\phi)^n + \dots \leq 1 - \zeta\phi \quad (\text{D.7})$$

by Taylor's expansion. Submitting (D.6) and (D.7) into (D.5) yields

$$|\mathbb{E}[g(\mathbf{X}_{l\tau})] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq \frac{2\sqrt{2}\kappa}{\sqrt{\phi}}[\bar{V} + 1]V(\mathbf{x}_0) \exp\left(-\frac{2ml\tau\eta\rho^d}{\log(\kappa)}\right), \quad (\text{D.8})$$

where we chose  $\delta = e^{-m\eta}$ . Next we need to prove that the unsampled chain is also exponential ergodic. Let  $k = l\tau + j$  with  $j = 0, 1, \dots, \tau - 1$ . We immediately get

$$|\mathbb{E}[g(\mathbf{X}_{l\tau+j})] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq \frac{2\sqrt{2}\kappa}{\sqrt{\phi}}[\bar{V} + 1]\mathbb{E}[V(\mathbf{X}_j)] \exp\left(-\frac{2ml\tau\eta\rho^d}{\log(\kappa)}\right).$$

Since the GLD approximation (2.1) of Langevin is ergodic when sampled at rate  $\tau = 1$ , we have  $k = l\tau = l$  and  $j = 0$ . Note that by Lemma A.2 in Mattingly et al. [40], we have  $\mathcal{C} = \{\mathbf{x} : V(\mathbf{x}) \leq \kappa/e^{-m\eta}\}$ , which implies that  $\bar{V} = \kappa e^{m\eta}$ . Thus we obtain

$$|\mathbb{E}[g(\mathbf{X}_k)] - \mathbb{E}[g(\mathbf{X}^\mu)]| \leq C\kappa\rho^{-d/2}(\kappa e^{m\eta} + 1) \exp\left(-\frac{2mk\eta\rho^d}{\log(\kappa)}\right),$$

where we used the fact that  $\mathbf{x}_0 = \mathbf{0}$  and  $C > 0$  is an absolute constant.  $\square$

### D.3 Proof of Lemma C.4

To prove Lemma C.4, we choose the test function in Poisson equation (1.3) as  $g = F_n$ . Given the Poisson equation, suppose we choose  $g$  as  $F_n$ , the distance between the time average of the GLD process and the expectation of  $F_n$  over the Gibbs measure can be expressed by

$$\frac{1}{K} \sum_{k=1}^K F_n(\mathbf{X}_k) - \bar{F} = \frac{1}{K} \sum_{k=1}^K \mathcal{L}\psi(\mathbf{X}_k). \quad (\text{D.9})$$

Note that by [41, 51], we know the Poisson equation (1.3) defined by the generator of Langevin dynamics has a unique solution  $\psi$  under Assumptions 3.1 and 3.2. According to Theorem 3.2 in [23], the  $p$ -th order derivatives of  $\psi$  can be bounded by some polynomial growth function with sophisticated coefficients ( $p = 0, 1, 2$ ). To simplify the presentation, we hence follow the convention in the line of literature [12, 51] and assume that  $\mathbb{E}[\|\nabla^p \psi(\mathbf{X}_k)\|]$  can be further upper bounded by a constant  $C_\psi$  for all  $\{\mathbf{X}_k\}_{k \geq 0}$  and  $p = (0, 1, 2)$ , which is determined by the Langevin diffusion and its Poisson equation. In fact, Erdogdu et al. [23] showed that the upper bound of derivatives (up to fourth order) of  $\psi$  only requires the dissipative and smooth assumptions. We refer interested readers to [23] for more details on deriving the  $C_\psi$  for Langevin diffusion. We show that the case  $p = 0$  can be easily verified as follows. By Assumption 3.1, using a similar argument as in the proof of Lemma 4.1, we bound  $F_n(\mathbf{x})$  by a quadratic function  $V(\mathbf{x})$

$$F_n(\mathbf{x}) \leq \frac{M}{2} V(\mathbf{x}) = \frac{M}{2} (C_0 + \|\mathbf{x}\|_2^2).$$

Applying Assumption 3.2 and Theorem 13 in Vollmer et al. [51] we have

$$|\psi(\mathbf{x})| \leq C_1(1 + \|\mathbf{x}\|_2^2) \leq C_2 V(\mathbf{x}). \quad (\text{D.10})$$

Note that by Assumptions 3.1 and 3.2 we can verify that a quadratic  $V(\mathbf{x})$  and  $p^* = 2$  satisfy Assumption 12 in [51] and therefore we obtain that for all  $p \leq p^*$ , we have

$$\sup_k \mathbb{E} V^p(\mathbf{X}_k) \leq \infty. \quad (\text{D.11})$$

Combining (D.10) and (D.11) we show that  $\psi(\mathbf{X}_k)$  is bounded in expectation.

*Proof.* For the simplicity of notation, we first assume that  $\beta = 1$  and then show the result for arbitrary  $\beta$  by a scaling technique. Note that for the continuous-time Markov process  $\{\mathbf{D}(t)\}_{t \geq 0}$  defined in (C.8), the distribution of random vector  $(\mathbf{X}_1, \dots, \mathbf{X}_K)$  is equivalent to that of  $(\mathbf{D}(\eta), \dots, \mathbf{D}(\eta K))$ . Let  $\psi$  be the solution of Poisson equation  $\mathcal{L}\psi = g - \int g(\mathbf{x})\pi(d\mathbf{x})$ . Since we have  $\mathbb{E}[\psi(\mathbf{X}_k)|\mathbf{X}_0 = \mathbf{x}] = \mathbb{E}[\psi(\mathbf{D}(\eta k))|\mathbf{D}_0 = \mathbf{x}]$ . We denote  $\mathbb{E}[\psi(\mathbf{D}(\eta k))|\mathbf{D}_0 = \mathbf{x}]$  by  $\mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta k))]$ . By applying (A.2), we compute the Taylor expansion of  $\mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta k))]$  at  $\mathbf{D}(\eta(k-1))$ :

$$\mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta k))] = \mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta(k-1)))] + \eta \mathbb{E}^\mathbf{x}[\mathcal{L}\psi(\mathbf{D}(\eta(k-1)))] + O(\eta^2).$$

Note that the remainder also depends on the second order derivative of the Poisson equation and are bounded by constant  $C_\psi$ . Take average over  $k = 1, \dots, K$  and rearrange the equation we have

$$\frac{1}{\eta K} (\mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta K))] - \psi(\mathbf{x})) + O(\eta) = \frac{1}{K} \sum_{k=1}^K \mathbb{E}^\mathbf{x}[\mathcal{L}\psi(\mathbf{D}(\eta(k-1)))]. \quad (\text{D.12})$$

Submit the Poisson equation (D.9) into the above equation (D.12) we have

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}^\mathbf{x}[F_n(\mathbf{X}_k)] - \bar{F} &= \frac{1}{K} \sum_{k=1}^K \mathbb{E}^\mathbf{x}[\mathcal{L}\psi(\mathbf{X}_{k-1})] = \frac{1}{K} \sum_{k=1}^K \mathbb{E}^\mathbf{x}[\mathcal{L}\psi(\mathbf{D}(\eta(k-1)))] \\ &= \frac{1}{\eta K} (\mathbb{E}^\mathbf{x}[\psi(\mathbf{D}(\eta K))] - \psi(\mathbf{x})) + O(\eta) \\ &= \frac{1}{\eta K} (\mathbb{E}^\mathbf{x}[\psi(\mathbf{X}_K)] - \psi(\mathbf{x})) + O(\eta), \end{aligned}$$

where the second and the fourth equation hold due to the fact that the distribution of  $\{\mathbf{X}_k\}$  is the same as the distribution of  $\{\mathbf{D}(\eta k)\}$ . We have assumed that  $\psi(\mathbf{X}_k)$  and its first and second order

derivatives are bounded by constant  $C_\psi$  in expectation over the randomness of  $\mathbf{X}_k$ . Therefore, we are able to obtain the following conclusion

$$\left| \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}^\mathbf{x}[F_n(\mathbf{X}_k)] - \bar{F} \right| \leq C_\psi \left( \frac{1}{\eta K} + \eta \right).$$

This completes the proof for the case  $\beta = 1$ . In order to apply our analysis to the case where  $\beta$  can take any arbitrary constant value, we conduct the same scaling argument as in (C.2).

$$\left| \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}^\mathbf{x}[F_n(\mathbf{X}_k)] - \bar{F} \right| \leq C_\psi \left( \frac{1}{\eta' K} + \eta' \right) = C_\psi \left( \frac{\beta}{\eta K} + \frac{\eta}{\beta} \right).$$

This completes the proof.  $\square$

#### D.4 Proof of Lemma C.5

We first lay down the following lemma on the bounds of gradient of  $f_i$ .

**Lemma D.1.** For any  $\mathbf{x} \in \mathbb{R}^d$ , it holds that

$$\|\nabla f_i(\mathbf{x})\|_2 \leq M\|\mathbf{x}\|_2 + G$$

for constant  $G = \max_{i=1, \dots, n} \{\|\nabla f_i(\mathbf{x}^*)\|_2\} + bM/m$ .

*Proof of Lemma C.5.* Let  $\mathbf{u}_i(\mathbf{x}) = \nabla F(\mathbf{x}) - \nabla f_i(\mathbf{x})$ , consider

$$\begin{aligned} \mathbb{E} \left\| \frac{1}{B} \sum_{i \in I_k} \mathbf{u}_i(\mathbf{x}) \right\|_2^2 &= \frac{1}{B^2} \mathbb{E} \sum_{i \neq i' \in I_k} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{B-1}{Bn(n-1)} \sum_{i \neq i'} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{B-1}{Bn(n-1)} \sum_{i, i'} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) - \frac{B-1}{B(n-1)} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{n-B}{B(n-1)} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2, \end{aligned} \tag{D.13}$$

where the last equality is due to the fact that  $1/n \sum_{i=1}^n \mathbf{u}_i(\mathbf{x}) = 0$ . By Lemma D.1 we have  $\|\nabla f_i(\mathbf{x})\|_2 \leq M\|\mathbf{x}\|_2 + G$ , therefore we have  $\|\nabla F(\mathbf{x})\|_2 \leq M\|\mathbf{x}\|_2 + G$  and consequently,  $\|\mathbf{u}_i(\mathbf{x})\|_2 \leq 2(M\|\mathbf{x}\|_2 + G)$ . Thus (D.13) can be further bounded as:

$$\mathbb{E} \left\| \frac{1}{B} \sum_{i \in I_k} \mathbf{u}_i(\mathbf{x}) \right\|_2^2 \leq \frac{n-B}{B(n-1)} 4(M\|\mathbf{x}\|_2 + G)^2.$$

This completes the proof.  $\square$

#### D.5 Proof of Lemma C.6

In this section, we provide the proof of  $L^2$  bound of GLD and SVRG-LD iterates  $\mathbf{X}_k$  and  $\mathbf{Z}_k$ . Note that a similar result of SGLD has been proved in Raginsky et al. [45] and thus we omit the corresponding proof for the simplicity of presentation.

*Proof of Lemma C.6. Part I:* We first prove the the upper bound for GLD. By the definition in (2.1), we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{X}_{k+1}\|_2^2] &= \mathbb{E}[\|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2] + \sqrt{\frac{8\eta}{\beta}} \mathbb{E}[\langle \mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k), \boldsymbol{\epsilon}_k \rangle] + \frac{2\eta}{\beta} \mathbb{E}[\|\boldsymbol{\epsilon}_k\|_2^2] \\ &= \mathbb{E}[\|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2] + \frac{2\eta d}{\beta}, \end{aligned}$$

where the second equality follows from that  $\epsilon_k$  is independent on  $\mathbf{X}_k$ . Now we bound the first term

$$\begin{aligned}\mathbb{E}[\|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2] &= \mathbb{E}[\|\mathbf{X}_k\|_2^2] - 2\eta \mathbb{E}[\langle \mathbf{X}_k, \nabla F_n(\mathbf{X}_k) \rangle] + \eta^2 \mathbb{E}[\|\nabla F_n(\mathbf{X}_k)\|_2^2] \\ &\leq \mathbb{E}[\|\mathbf{X}_k\|_2^2] + 2\eta(b - m\mathbb{E}[\|\mathbf{X}_k\|_2^2]) + 2\eta^2(M^2\mathbb{E}[\|\mathbf{X}_k\|_2^2] + G^2) \\ &= (1 - 2\eta m + 2\eta^2 M^2)\mathbb{E}[\|\mathbf{X}_k\|_2^2] + 2\eta b + 2\eta^2 G^2,\end{aligned}$$

where the inequality follows from Assumption 3.2, Lemma D.1 and triangle inequality. Substitute the above bound back and we will have

$$\mathbb{E}[\|\mathbf{X}_{k+1}\|_2^2] \leq (1 - 2\eta m + 2\eta^2 M^2)\mathbb{E}[\|\mathbf{X}_k\|_2^2] + 2\eta b + 2\eta^2 G^2 + \frac{2\eta d}{\beta}. \quad (\text{D.14})$$

For sufficient small  $\eta$  that satisfies  $\eta \leq \min\{1, m/(2M^2)\}$ , there are only two cases we need to take into account:

If  $1 - 2\eta m + 2\eta^2 M^2 \leq 0$ , then from (D.14) we have

$$\mathbb{E}[\|\mathbf{X}_{k+1}\|_2^2] \leq 2\eta b + 2\eta^2 G^2 + \frac{2\eta d}{\beta} \leq \|\mathbf{X}_0\|_2^2 + 2\left(b + G^2 + \frac{d}{\beta}\right). \quad (\text{D.15})$$

If  $0 < 1 - 2\eta m + 2\eta^2 M^2 \leq 1$ , then iterate (D.14) and we have

$$\mathbb{E}[\|\mathbf{X}_k\|_2^2] \leq (1 - 2\eta m + 2\eta^2 M^2)^k \|\mathbf{X}_0\|_2^2 + \frac{\eta b + \eta^2 G^2 + \frac{\eta d}{\beta}}{\eta m - \eta^2 M^2} \leq \|\mathbf{X}_0\|_2^2 + \frac{2}{m} \left(b + G^2 + \frac{d}{\beta}\right). \quad (\text{D.16})$$

Combine (D.15) and (D.16) and we have

$$\mathbb{E}[\|\mathbf{X}_k\|_2^2] \leq \|\mathbf{X}_0\|_2^2 + \left(2 + \frac{2}{m}\right) \left(b + G^2 + \frac{d}{\beta}\right) = 2\left(1 + \frac{1}{m}\right) \left(b + G^2 + \frac{d}{\beta}\right),$$

where the equation holds by choosing  $\mathbf{X}_0 = \mathbf{0}$ .

**Part II:** Now we prove the  $L^2$  bound for SVRG-LD, i.e.,  $\mathbb{E}[\|\mathbf{Z}_k\|_2^2]$ , by mathematical induction. Since  $\tilde{\nabla}_k = 1/B \sum_{i_k \in I_k} (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)}))$ , we have

$$\begin{aligned}\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] &= \mathbb{E}[\|\mathbf{Z}_k - \eta \tilde{\nabla}_k\|_2^2] + \sqrt{\frac{8\eta}{\beta}} \mathbb{E}[\langle \mathbf{Z}_k - \eta \tilde{\nabla}_k, \epsilon_k \rangle] + \frac{2\eta}{\beta} \mathbb{E}[\|\epsilon_k\|_2^2] \\ &= \mathbb{E}[\|\mathbf{Z}_k - \eta \tilde{\nabla}_k\|_2^2] + \frac{2\eta d}{\beta},\end{aligned} \quad (\text{D.17})$$

where the second equality follows from the fact that  $\epsilon_k$  is independent of  $\mathbf{Z}_k$  and standard Gaussian. We prove it by induction. First, consider the case when  $k = 1$ . Since we choose the initial point at  $\mathbf{Z}_0 = \mathbf{0}$ , we immediately have

$$\begin{aligned}\mathbb{E}[\|\mathbf{Z}_1\|_2^2] &= \mathbb{E}[\|\mathbf{Z}_0 - \eta \tilde{\nabla}_0\|_2^2] + \sqrt{\frac{8\eta}{\beta}} \mathbb{E}[\langle \mathbf{Z}_0 - \eta \tilde{\nabla}_0, \epsilon_0 \rangle] + \frac{2\eta}{\beta} \mathbb{E}[\|\epsilon_0\|_2^2] \\ &= \eta^2 \mathbb{E}[\|\nabla F_n(\mathbf{Z}_0)\|_2^2] + \frac{2\eta d}{\beta} \\ &\leq \eta^2 G^2 + \frac{2\eta d}{\beta},\end{aligned}$$

where the second equality holds due to the fact that  $\tilde{\nabla}_0 = \nabla F_n(\mathbf{Z}_0)$  and the inequality follows from Lemma D.1. For sufficiently small  $\eta$  we can see that the conclusion of Lemma C.6 holds for  $\mathbb{E}[\|\mathbf{Z}_1\|_2^2]$ , i.e.,  $\mathbb{E}[\|\mathbf{Z}_1\|_2^2] \leq \Gamma$ , where  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$ . Now assume that the conclusion holds for all iteration from 1 to  $k$ , then for the  $(k+1)$ -th iteration, by (D.17) we have,

$$\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] = \mathbb{E}[\|\mathbf{Z}_k - \eta \tilde{\nabla}_k\|_2^2] + \frac{2\eta d}{\beta}, \quad (\text{D.18})$$

For the first term on the R.H.S of (D.18) we have

$$\begin{aligned}\mathbb{E}[\|\mathbf{Z}_k - \eta \tilde{\nabla}_k\|_2^2] &= \mathbb{E}[\|\mathbf{Z}_k - \eta \nabla F_n(\mathbf{Z}_k)\|_2^2] + 2\eta \mathbb{E}[\langle \mathbf{Z}_k - \eta \nabla F_n(\mathbf{Z}_k), \nabla F_n(\mathbf{Z}_k) - \tilde{\nabla}_k \rangle] \\ &\quad + \eta^2 \mathbb{E}[\|\nabla F_n(\mathbf{Z}_k) - \tilde{\nabla}_k\|_2^2] \\ &= \underbrace{\mathbb{E}[\|\mathbf{Z}_k - \eta \nabla F_n(\mathbf{Z}_k)\|_2^2]}_{T_1} + \underbrace{\eta^2 \mathbb{E}[\|\nabla F_n(\mathbf{Z}_k) - \tilde{\nabla}_k\|_2^2]}_{T_2},\end{aligned} \quad (\text{D.19})$$

where the second equality holds due to the fact that  $\mathbb{E}[\tilde{\nabla}_k] = \nabla F_n(\mathbf{Z}_k)$ . For term  $T_1$ , we can further bound it by

$$\begin{aligned}\mathbb{E}[\|\mathbf{Z}_k - \eta \nabla F_n(\mathbf{Z}_k)\|_2^2] &= \mathbb{E}[\|\mathbf{Z}_k\|_2^2] - 2\eta \mathbb{E}[\langle \mathbf{Z}_k, \nabla F_n(\mathbf{Z}_k) \rangle] + \eta^2 \mathbb{E}[\|\nabla F_n(\mathbf{Z}_k)\|_2^2] \\ &\leq \mathbb{E}[\|\mathbf{Z}_k\|_2^2] + 2\eta(b - m\mathbb{E}[\|\mathbf{Z}_k\|_2^2]) + 2\eta^2(M^2\mathbb{E}[\|\mathbf{Z}_k\|_2^2] + G^2) \\ &= (1 - 2\eta m + 2\eta^2 M^2)\mathbb{E}[\|\mathbf{Z}_k\|_2^2] + 2\eta b + 2\eta^2 G^2,\end{aligned}$$

where the inequality follows from Lemma D.1 and triangle inequality. For term  $T_2$ , by Lemma C.10 we have

$$\mathbb{E}\|\nabla F_n(\mathbf{Z}_k) - \tilde{\nabla}_k\|_2^2 \leq \frac{M^2(n-B)}{B(n-1)}\mathbb{E}\|\mathbf{Z}_k - \tilde{\mathbf{Z}}^{(s)}\|_2^2 \leq \frac{2M^2(n-B)}{B(n-1)}\left(\mathbb{E}\|\mathbf{Z}_k\|_2^2 + \mathbb{E}\|\tilde{\mathbf{Z}}^{(s)}\|_2^2\right).$$

Submit the above bound back into (D.17) we have

$$\begin{aligned}\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] &\leq \left(1 - 2\eta m + 2\eta^2 M^2\left(1 + \frac{n-B}{B(n-1)}\right)\right)\mathbb{E}[\|\mathbf{Z}_k\|_2^2] \\ &\quad + \frac{2\eta^2 M^2(n-B)}{B(n-1)}\mathbb{E}\|\tilde{\mathbf{Z}}^{(s)}\|_2^2 + 2\eta b + 2\eta^2 G^2 + \frac{2\eta d}{\beta}.\end{aligned}\quad (\text{D.20})$$

Note that by assumption we have  $\mathbb{E}\|\mathbf{Z}_j\|_2^2 \leq \Gamma$  for all  $j = 1, \dots, k$  where  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$ , thus (D.20) can be further bounded as:

$$\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] \leq \underbrace{\left(1 - 2\eta m + 2\eta^2 M^2\left(1 + \frac{2(n-B)}{B(n-1)}\right)\right)}_{C_\lambda} \Gamma + 2\eta b + 2\eta^2 G^2 + \frac{2\eta d}{\beta}.\quad (\text{D.21})$$

For sufficient small  $\eta$  that satisfies

$$\eta \leq \min\left(1, \frac{m}{2M^2(1 + 2(n-B)/(B(n-1)))}\right),$$

there are only two cases we need to take into account:

If  $C_\lambda \leq 0$ , then from (D.21) we have

$$\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] \leq 2\eta b + 2\eta^2 G^2 + \frac{2\eta d}{\beta} \leq 2\left(b + G^2 + \frac{d}{\beta}\right).\quad (\text{D.22})$$

If  $0 < C_\lambda \leq 1$ , then iterate (D.21) and we have

$$\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] \leq C_\lambda^{k+1}\|\mathbf{Z}_0\|_2^2 + \frac{\eta b + \eta^2 G^2 + \frac{\eta d}{\beta}}{\eta m - \eta^2 M^2\left(1 + \frac{2(n-B)}{B(n-1)}\right)} \leq \frac{2}{m}\left(b + G^2 + \frac{d}{\beta}\right).\quad (\text{D.23})$$

Combining (D.22) and (D.23), we have

$$\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2] \leq 2\left(1 + \frac{1}{m}\right)\left(b + 2G^2 + \frac{d}{\beta}\right).$$

Thus we show that when  $\mathbb{E}[\|\mathbf{Z}_j\|_2^2], j = 1, \dots, k$  are bounded,  $\mathbb{E}[\|\mathbf{Z}_{k+1}\|_2^2]$  is also bounded. By mathematical induction we complete the proof.  $\square$

## D.6 Proof of Lemma C.7

*Proof.* We have the following equation according to the update of GLD in (2.1),

$$\begin{aligned}\mathbb{E}[\exp(\|\mathbf{X}_{k+1}\|_2^2)] &= \mathbb{E}\exp\left(\left\|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k) + \sqrt{\frac{2\eta}{\beta}}\boldsymbol{\epsilon}_k\right\|_2^2\right) \\ &= \mathbb{E}\exp\left(\left\|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\right\|_2^2 + \sqrt{\frac{8\eta}{\beta}}\langle \mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k), \boldsymbol{\epsilon}_k \rangle + \frac{2\eta}{\beta}\|\boldsymbol{\epsilon}_k\|_2^2\right).\end{aligned}\quad (\text{D.24})$$

Let  $H(\mathbf{x}) = \exp(\|\mathbf{x}\|_2^2)$ , we have  $\mathbb{E}[H(\mathbf{X}_{k+1})] = \mathbb{E}_{\mathbf{X}_k}[\mathbb{E}[H(\mathbf{X}_{k+1})|\mathbf{X}_k]]$ . Thus we can first compute the conditional expectation on the R.H.S of (D.24) given  $\mathbf{X}_k$ , then compute the expectation with respect to  $\mathbf{X}_k$ . Note that  $\epsilon_k$  follows standard multivariate normal distribution, i.e.,  $\epsilon_k \sim N(\mathbf{0}, \mathbf{I}_{d \times d})$ . Then it can be shown that

$$\begin{aligned} & \mathbb{E} \left[ \exp \left( \sqrt{\frac{8\eta}{\beta}} \langle \mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k), \epsilon_k \rangle + \frac{2\eta}{\beta} \|\epsilon_k\|_2^2 \right) \middle| \mathbf{X}_k \right] \\ &= \frac{1}{(1 - 4\eta/\beta)^{d/2}} \exp \left( \frac{4\eta}{\beta - 4\eta} \|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2 \right) \end{aligned}$$

holds as long as  $\beta > 4\eta$ . Plugging the above equation into (D.24), we have

$$\mathbb{E}[H(\mathbf{X}_{k+1})] = \frac{1}{(1 - 4\eta/\beta)^{d/2}} \mathbb{E}_{\mathbf{X}_k} \left[ \exp \left( \frac{\beta}{\beta - 4\eta} \|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2 \right) \right]. \quad (\text{D.25})$$

Note that by Assumption 3.2 and Lemma D.1 we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}_k} \exp \left( \frac{\beta}{\beta - 4\eta} \|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2 \right) \\ &= \mathbb{E}_{\mathbf{X}_k} \exp \left( \frac{\beta}{\beta - 4\eta} (\|\mathbf{X}_k\|_2^2 - 2\eta \langle \mathbf{X}_k, \nabla F_n(\mathbf{X}_k) \rangle + \eta^2 \|\nabla F_n(\mathbf{X}_k)\|_2^2) \right) \\ &\leq \mathbb{E}_{\mathbf{X}_k} \exp \left( \frac{\beta}{\beta - 4\eta} (\|\mathbf{X}_k\|_2^2 - 2\eta(m\|\mathbf{X}_k\|_2^2 - b) + 2\eta^2(M^2\|\mathbf{X}_k\|_2^2 + G^2)) \right) \\ &= \mathbb{E}_{\mathbf{X}_k} \exp \left( \frac{\beta}{\beta - 4\eta} ((1 - 2\eta m + 2\eta^2 M^2)\|\mathbf{X}_k\|_2^2 + 2b\eta + 2\eta^2 G^2) \right). \end{aligned}$$

Consider sufficiently small  $\eta$  such that  $\eta < m/M^2$ . Then for  $\beta$  satisfying  $\beta \geq \max\{2/(m - M^2\eta), 4\eta\}$ , we have  $\beta(1 - 2\eta m + 2\eta^2 M^2)/(\beta - 4\eta) \leq 1$ . Therefore, the above expectation can be upper bounded by

$$\mathbb{E}_{\mathbf{X}_k} \exp \left( \frac{\beta}{\beta - 4\eta} \|\mathbf{X}_k - \eta \nabla F_n(\mathbf{X}_k)\|_2^2 \right) \leq \exp \left( \frac{2\eta\beta(b + \eta G^2)}{\beta - 4\eta} \right) \mathbb{E}[H(\mathbf{X}_k)].$$

Substituting the above inequality into (D.25), it follows that

$$\begin{aligned} \mathbb{E}[H(\mathbf{X}_{k+1})] &\leq \frac{1}{(1 - 4\eta/\beta)^{d/2}} \exp \left( \frac{2\eta\beta(b + \eta G^2)}{\beta - 4\eta} \right) \mathbb{E}[H(\mathbf{X}_k)] \\ &\leq \exp \left( \frac{2\eta(\beta b + \eta\beta G^2 + d)}{\beta - 4\eta} \right) \mathbb{E}[H(\mathbf{X}_k)], \end{aligned}$$

where we used the fact that  $\log(1/(1 - x)) \leq x/(1 - x)$  for  $0 < x < 1$  and that

$$\log \left( \frac{1}{(1 - 4\eta/\beta)^{d/2}} \right) = \frac{d}{2} \log \left( \frac{1}{1 - 4\eta/\beta} \right) \leq \frac{2d\eta/\beta}{1 - 4\eta/\beta} = \frac{2\eta d}{\beta - 4\eta}.$$

Then we are able to show by induction that

$$\mathbb{E}[H(\mathbf{X}_k)] \leq \exp \left( \frac{2k\eta(\beta b + \eta\beta G^2 + d)}{\beta - 4\eta} \right) \mathbb{E}[H(\|\mathbf{X}_0\|_2)],$$

which immediately implies that

$$\log \mathbb{E}[\exp(\|\mathbf{X}_k\|_2^2)] \leq \|\mathbf{X}_0\|_2^2 + \frac{2\beta(b + G^2) + 2d}{\beta - 4\eta} k\eta,$$

where we assume that  $\eta \leq 1$  and  $\beta > 4\eta$ .

□

## D.7 Proof of Lemma C.10

*Proof.* Since by Algorithm 3 we have  $\tilde{\nabla}_k = (1/B) \sum_{i_k \in I_k} (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)}))$ , therefore,

$$\mathbb{E}[\|\tilde{\nabla}_k - \nabla F_n(\mathbf{Z}_k)\|_2^2] = \mathbb{E}\left\|\frac{1}{B} \sum_{i_k \in I_k} (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)}) - \nabla F_n(\mathbf{Z}_k))\right\|_2^2.$$

Let  $\mathbf{u}_i = \nabla F_n(\mathbf{Z}_k) - \nabla F_n(\tilde{\mathbf{Z}}^{(s)}) - (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}^{(s)}))$ .

$$\begin{aligned} \mathbb{E}\left\|\frac{1}{B} \sum_{i \in I_k} \mathbf{u}_i(\mathbf{x})\right\|_2^2 &= \frac{1}{B^2} \mathbb{E} \sum_{i \neq i' \in I_k} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{B-1}{Bn(n-1)} \sum_{i \neq i'} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{B-1}{Bn(n-1)} \sum_{i, i'} \mathbf{u}_i(\mathbf{x})^\top \mathbf{u}_{i'}(\mathbf{x}) - \frac{B-1}{B(n-1)} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 + \frac{1}{B} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2 \\ &= \frac{n-B}{B(n-1)} \mathbb{E} \|\mathbf{u}_i(\mathbf{x})\|_2^2, \end{aligned} \tag{D.26}$$

where the last equality is due to the fact that  $1/n \sum_{i=1}^n \mathbf{u}_i(\mathbf{x}) = 0$ . Therefore, we have

$$\begin{aligned} \mathbb{E}[\|\tilde{\nabla}_k - \nabla F_n(\mathbf{Z}_k)\|_2^2] &\leq \frac{n-B}{B(n-1)} \mathbb{E} \|\mathbf{u}_i\|_2^2 \\ &= \frac{n-B}{B(n-1)} \mathbb{E} \|\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}) - \mathbb{E}[\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}})]\|_2^2 \\ &\leq \frac{n-B}{B(n-1)} \mathbb{E} \|\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}})\|_2^2 \\ &\leq \frac{M^2(n-B)}{B(n-1)} \mathbb{E} \|\mathbf{Z}_k - \tilde{\mathbf{Z}}\|_2^2, \end{aligned} \tag{D.27}$$

where the second inequality holds due to the fact that  $\mathbb{E}[\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\|_2^2] \leq \mathbb{E}[\|\mathbf{x}\|_2^2]$  and the last inequality follows from Assumption 3.1. This completes the proof.  $\square$

## E Proof of Auxiliary Lemmas in Appendix D

### E.1 Proof of Lemma D.1

*Proof.* By Assumption 3.2 we obtain

$$\langle \mathbf{x}^*, \nabla F_n(\mathbf{x}^*) \rangle \geq m \|\mathbf{x}^*\|_2^2 - b.$$

Note that  $\mathbf{x}^*$  is the minimizer for  $F_n$ , which implies that  $\nabla F_n(\mathbf{x}^*) = \mathbf{0}$  and threfore  $\|\mathbf{x}^*\|_2 \leq b/m$ . By Assumption 3.1 we further have

$$\|\nabla f_i(\mathbf{x})\|_2 \leq \|\nabla f_i(\mathbf{x}^*)\|_2 + M \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \|\nabla f_i(\mathbf{x}^*)\|_2 + \frac{bM}{m} + M \|\mathbf{x}\|_2.$$

The proof is completed by setting  $G = \max_{i=1, \dots, n} \{\|\nabla f_i(\mathbf{x}^*)\|_2\} + bM/m$ .  $\square$