

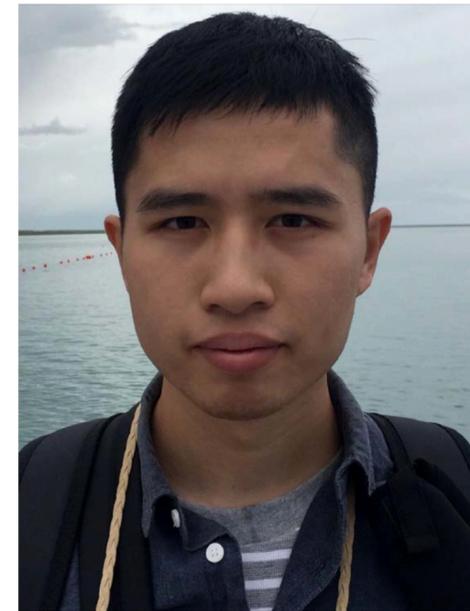
# Global Convergence of Langevin Dynamics Based Algorithms for Nonconvex Optimization



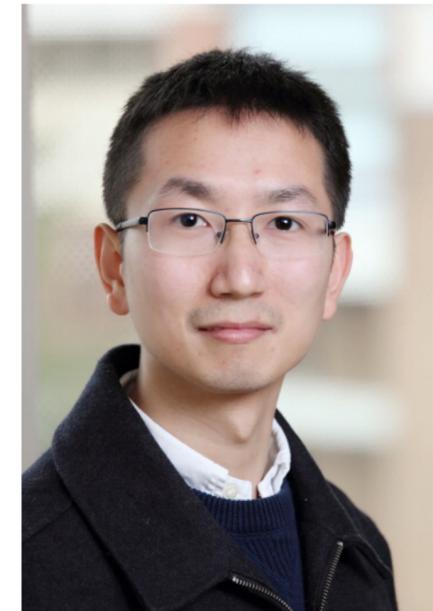
**Pan Xu\***  
UCLA



**Jinghui Chen\***  
UVa



**Difan Zou**  
UCLA



**Quanquan Gu**  
UCLA

**UCLA**



# Langevin Dynamics

## Langevin Dynamics:

$$d\mathbf{X}(t) = \underbrace{-\nabla F_n(\mathbf{X}(t))dt}_{\text{drift term}} + \underbrace{\sqrt{2\beta^{-1}}d\mathbf{B}(t)}_{\text{diffusion term}},$$

- $\beta$ : inverse temperature parameter
- $B(t)$ : standard Brownian motion

**Asymptotic property** (Roberts & Tweedie, 1996): converges to a stationary distribution

$$\pi(d\mathbf{x}) \propto \exp(-\beta F_n(\mathbf{x}))$$

Implication: The stationary distribution concentrates on the global minima.

# Gradient Langevin Dynamics

## Langevin Dynamics:

$$d\mathbf{X}(t) = -\nabla F_n(\mathbf{X}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t),$$

## Gradient Langevin Dynamics (GLD, aka. Langevin Monte Carlo):

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

- $\eta$  is the step size
- $\boldsymbol{\epsilon}_k$  is a standard Gaussian random vector

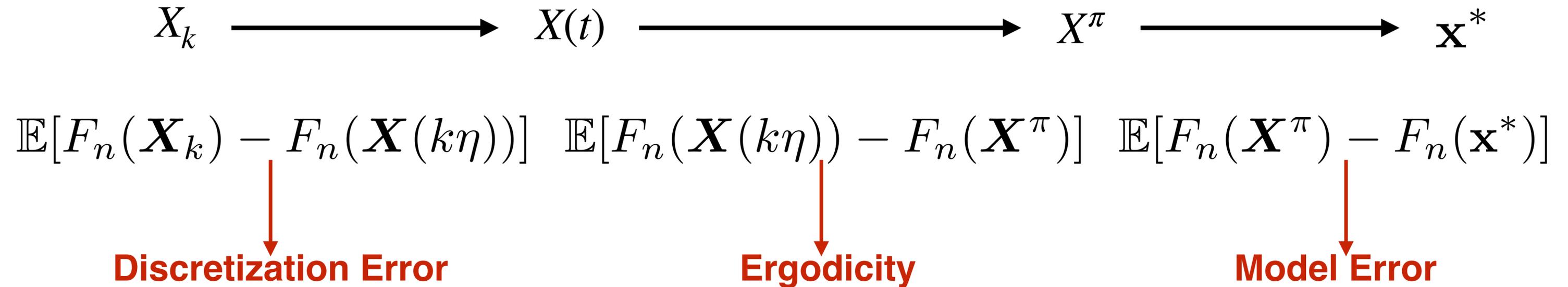
## Goal: bound the Optimization Error

$$\mathbb{E}[F_n(\mathbf{X}_k) - F_n(\mathbf{x}^*)] \quad \mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} F_n(\mathbf{x})$$

# Decomposition of Optimization Error

**Goal:** bound the **Optimization Error**  $\mathbb{E}[F_n(\mathbf{X}_k) - F_n(\mathbf{x}^*)]$

**Decomposition:** (Raginsky et al., 2017)

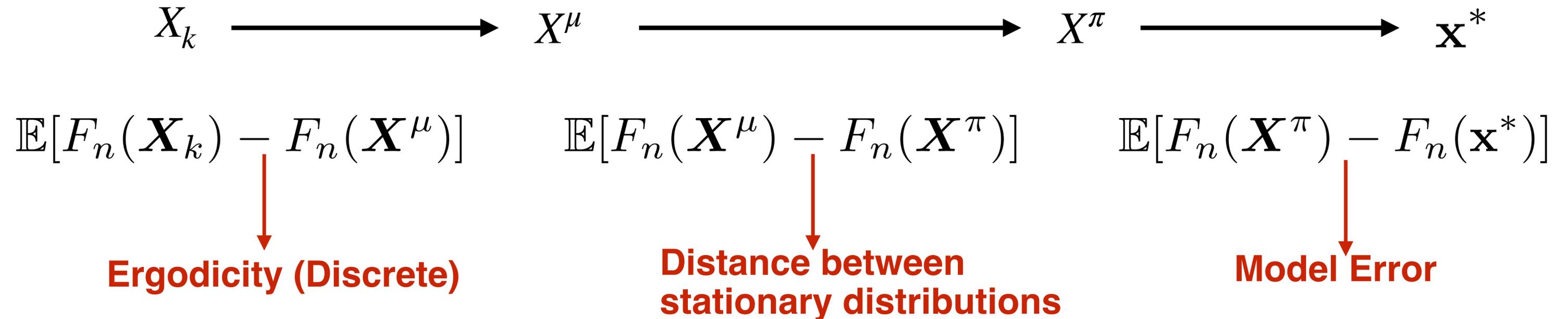


**Iteration complexity:**  $\mathbb{E}[F_n(\mathbf{X}_n) - F_n(\mathbf{x}^*)] \leq \epsilon + O\left(\frac{d}{\beta} \log \frac{\beta}{d}\right)$   
 $k = \tilde{O}\left(\frac{1}{\epsilon^4 \lambda^{*5}} \log^5 \frac{1}{\epsilon}\right)$   
 $\downarrow$   
**Model Error**

# Novel Decomposition for Faster Rates

**Goal:** bound the **Optimization Error**  $\mathbb{E}[F_n(\mathbf{X}_k) - F_n(\mathbf{x}^*)]$

**Decomposition (this paper):**



**Iteration complexity:**  $\mathbb{E}[F_n(\mathbf{X}_n) - F_n(\mathbf{x}^*)] \leq \epsilon + O\left(\frac{d}{\beta} \log \frac{\beta}{d}\right)$

$k = \tilde{O}\left(\frac{1}{\epsilon^4 \lambda^{*5}} \log^5 \frac{1}{\epsilon}\right) \longrightarrow k = \tilde{O}\left(\frac{1}{\epsilon \lambda^*} \log \frac{1}{\epsilon}\right)$

**Model Error**

# Global Convergence of Variants of GLD

## Stochastic Gradient Langevin Dynamics (SGLD):

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

- unbiased stochastic gradient  $\mathbb{E}[\nabla G(X)|X] = \nabla F_n(X)$

## Stochastic Variance Reduced Gradient Langevin Dynamics (SVRG-LD):

$$\mathbf{Z}_{k+1} = \mathbf{Z}_k - \eta \tilde{\nabla}_k + \sqrt{2\eta\beta^{-1}} \cdot \boldsymbol{\epsilon}_k,$$

- semi-stochastic gradient  $\tilde{\nabla}_k = \nabla G_k(\mathbf{Z}_k) - \nabla G_k(\tilde{\mathbf{Z}}^{(s)}) + \nabla F_n(\tilde{\mathbf{Z}}^{(s)})$
- $\tilde{\mathbf{Z}}^{(s)}$  is a snapshot of  $\mathbf{Z}_k$ , updated after every  $m$  iterations.

# Thanks!

Poster session:  
10:45 AM -- 12:45 PM  
@Room 210 & 230 AB #46



## Global Convergence of Langevin Dynamics Based Algorithms for Nonconvex Optimization

Pan Xu\*<sup>†</sup> and Jinghui Chen\*<sup>‡</sup> and Difan Zou<sup>†</sup> and Quanquan Gu<sup>†</sup>

<sup>†</sup>University of California, Los Angeles

<sup>‡</sup>University of Virginia



### Problem Setup and Background

#### Optimization problem:

$$\min_{\mathbf{x}} F(\mathbf{x}) := 1/n \sum_{i=1}^n f_i(\mathbf{x}),$$

Assumptions on the function:

- $f_i$  is  $M$ -smooth:  $\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|_2 \leq M\|\mathbf{x} - \mathbf{y}\|_2 \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ .
- $F$  is  $(m, b)$ -dissipative:  $\langle \nabla F(\mathbf{x}), \mathbf{x} \rangle \geq m\|\mathbf{x}\|_2^2 - b, \forall \mathbf{x} \in \mathbb{R}^d$ .

Note that  $F$  is nonconvex.

- Langevin Dynamics:** stochastic differential equation
$$d\mathbf{X}(t) = -\nabla F(\mathbf{X}(t))dt + \sqrt{2\beta^{-1}}d\mathbf{B}(t),$$

where the parameters are

- $\beta > 0$  is called the inverse temperature parameter.
- $\mathbf{B}(t)$  is a standard Brownian motion in  $\mathbb{R}^d$ .

- Asymptotic property:** the distribution of stochastic process  $\mathbf{X}(t)$  converges to the following stationary distribution

$$\pi \propto \exp(-\beta F(\mathbf{x})).$$

- $\pi$  concentrates on the **global minimizer** of  $F$ .
- Discretize it to obtain optimization algorithm.

### Langevin Dynamics Based Algorithms

#### Gradient Langevin Dynamics (GLD)

$$\mathbf{X}_{k+1} = \mathbf{X}_k - \eta \nabla F(\mathbf{X}_k) + \sqrt{2\eta/\beta} \epsilon_k$$

- $\epsilon_k$ : an additive standard Gaussian noise
- $\eta$ : step size
- Converges fast  $\odot$ ; computation is high when  $n$  is large  $\ominus$

#### Stochastic Gradient Langevin Dynamics (SGLD)

$$\mathbf{Y}_{k+1} = \mathbf{Y}_k - \eta/B \sum_{i \in I_k} \nabla f_i(\mathbf{Y}_k) + \sqrt{2\eta/\beta} \epsilon_k$$

- $\nabla f_i(\mathbf{Y}_k)$ : unbiased stochastic gradient, i.e.,  $\mathbb{E}[\nabla f_i(\mathbf{x})] = \nabla F(\mathbf{x})$
- $I_k$ : a subset of  $\{1, \dots, n\}$  with  $|I_k| = B$
- Reduces the per iteration gradient complexity  $\odot$ ; converges slowly  $\ominus$

#### Stochastic Variance Reduced Gradient Langevin Dynamics (SVRG-LD)

$$\tilde{\nabla}_k = 1/B \sum_{i_k \in I_k} (\nabla f_{i_k}(\mathbf{Z}_k) - \nabla f_{i_k}(\tilde{\mathbf{Z}}^{(s)})) + \tilde{\mathbf{W}}$$

$$\mathbf{Z}_{k+1} = \mathbf{Z}_k - \eta \tilde{\nabla}_k + \sqrt{2\eta/\beta} \epsilon_k$$

- $\tilde{\mathbf{Z}}^{(s)}$  is a snapshot of  $\mathbf{Z}_k$  every  $L$  iterations.
- $\tilde{\mathbf{W}} = \nabla F(\tilde{\mathbf{Z}}^{(s)})$  is the full gradient at  $\tilde{\mathbf{Z}}^{(s)}$ .
- Multiple-epoch algorithm (each epoch has  $L$  iterations).
- Reduces the per iteration gradient complexity  $\odot$  and converges faster than SGLD  $\ominus$

### Theoretical Results

- GLD:** Under smoothness and dissipative assumptions, assume  $\eta \lesssim \epsilon$ , GLD achieves  $\mathbb{E}[F(\mathbf{X}_K)] - \mathbb{E}[F(\mathbf{x}^*)] \leq \epsilon + O(d/\beta)$ .

$$\text{Iteration complexity: } K = O(d\epsilon^{-1}\lambda^{-1} \cdot \log(1/\epsilon))$$

- $\mathbf{x}^*$  = argmin  $F(x)$  is the global minimizer.
- $O(d/\beta)$  is the model error of Langevin dynamics.
- $\mathbf{X}_K$  is called an *almost minimizer* of  $F$ .
- $\lambda = O(e^{-d})$  is the spectral gap of Markov process  $\mathbf{X}_k$ .

- SGLD:** Under the same conditions, if  $\eta \lesssim \epsilon$ ,  $B \gtrsim d^6/(\lambda\epsilon)^4 \log^4(1/\epsilon)$ , SGLD achieves  $\mathbb{E}[F(\mathbf{Y}_K)] - \mathbb{E}[F(\mathbf{x}^*)] \leq \epsilon + O(d/\beta)$

$$\text{Iteration complexity: } K = O(d\epsilon^{-1}\lambda^{-1} \cdot \log(1/\epsilon))$$

$B$  is the mini-batch size chosen in SGLD.

- SVRG-LD:** Under the same conditions, if we choose  $\eta \lesssim \epsilon$ , SVRG-LD achieves  $\mathbb{E}[F(\mathbf{Z}_K)] - \mathbb{E}[F(\mathbf{x}^*)] \leq \epsilon$ .

$$\text{Iteration complexity: } K = O(Ld^6B^{-1}\lambda^{-4}\epsilon^{-4} \cdot \log^4(1/\epsilon) + 1/\epsilon)$$

- Comparison of **gradient complexity** with state-of-the-art:

Table: Gradient complexities to converge to the almost minimizer.

	GLD	SGLD	SVRG-LD
[Raginsky et al., (2017)]	$\tilde{O}(\frac{n}{\epsilon}) \cdot e^{O(d)}$	$\tilde{O}(\frac{1}{\epsilon}) \cdot e^{O(d)}$	N/A
This paper	$\tilde{O}(\frac{n}{\epsilon}) \cdot e^{O(d)}$	$\tilde{O}(\frac{1}{\epsilon}) \cdot e^{O(d)}$	$\tilde{O}(\frac{\sqrt{n}}{\epsilon^{3/2}}) \cdot e^{O(d)}$

Choose  $B = \sqrt{n}\epsilon^{-3/2}$  and  $L = \sqrt{n}\epsilon^{3/2}$  for SVRG-LD.

### Decomposition of Optimization Error

- Goal:** bound the optimization error  $\mathbb{E}[F(\mathbf{X}_k)] - F(\mathbf{x}^*)$

#### Decomposition:

$$\mathbb{E}[F(\mathbf{X}_k)] - F(\mathbf{x}^*) = \underbrace{\mathbb{E}[F(\mathbf{X}_k)] - F(\mathbf{X}^\mu)}_{I_1} + \underbrace{\mathbb{E}[F(\mathbf{X}^\mu)] - F(\mathbf{X}^\pi)}_{I_2} + \underbrace{\mathbb{E}[F(\mathbf{X}^\pi)] - F(\mathbf{x}^*)}_{I_3}$$

- $\mu$ : the stationary distribution of **discrete-time process**  $\mathbf{X}_k$
- $\pi$ : the stationary distribution of **continuous-time process**  $\mathbf{X}(t)$

$I_1$  Geometric ergodicity of GLD

$I_2$  Distance between two stationary distributions

$I_3$  Gap between Langevin diffusion and global minimum

#### Comparison with existing decomposition approach

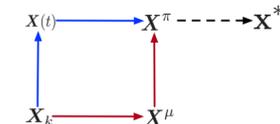


Figure: Blue arrow: decomposition scheme in [Raginsky et al., (2017)]; Red arrow: decomposition scheme in this paper.

- Bypass the discretization error between  $\mathbf{X}_k$  and  $\mathbf{X}(t)$ .
- Directly analyze the convergence to stationarity of  $\mathbf{X}_k$

### Proof Road Map

#### Lemma 1 (Bounding $I_1$ )

Under smoothness and dissipative assumptions, GLD has a unique invariant measure  $\mu$  on  $\mathbb{R}^d$ . It holds that

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

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

- $\mu$  is the stationary distribution of **discrete-time process**  $\mathbf{X}_k$

#### Lemma 2 (Bounding $I_2$ )

Under the same conditions, the invariant measures  $\mu$  and  $\pi$  satisfy

$$|\mathbb{E}[F(\mathbf{X}^\mu)] - \mathbb{E}[F(\mathbf{X}^\pi)]| \leq C_\psi\eta/\beta,$$

$C_\psi > 0$  is a constant depending on the generator of Langevin diffusion.

#### Lemma 3 (Bounding $I_3$ )

Under the same conditions, the error  $I_3$  can be bounded by

$$\mathbb{E}[F(\mathbf{X}^\pi)] - F(\mathbf{x}^*) \leq \frac{d}{2\beta} \log\left(\frac{eM(m\beta/d + 1)}{m}\right).$$

- Combining Lemmas 1, 2 & 3 yields the results for GLD.

### Proof for SGLD & SVRG-LD

- Decomposition of the optimization error of SGLD

$$\mathbb{E}[F(\mathbf{Y}_k)] - F(\mathbf{x}^*) = \mathbb{E}[F(\mathbf{Y}_k) - F(\mathbf{X}_k)] + \mathbb{E}[F(\mathbf{X}_k)] - F(\mathbf{x}^*)$$

- Lemma 4 (The distance between SGLD and GLD)**

Under smoothness and dissipative assumptions, the outputs of SGLD ( $\mathbf{Y}_k$ ) and GLD ( $\mathbf{X}_k$ ) satisfy

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

where  $C_1$  is an absolute constant and  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$ .

- Combining results for GLD and Lemma 4 yields the results for SGLD.

- Decomposition of the optimization error of SVRG-LD

$$\mathbb{E}[F(\mathbf{Z}_k)] - F(\mathbf{x}^*) = \mathbb{E}[F(\mathbf{Z}_k) - F(\mathbf{X}_k)] + \mathbb{E}[F(\mathbf{X}_k)] - F(\mathbf{x}^*)$$

#### Lemma 5 (The distance between SVRG-LD and GLD)

Under the same conditions, the outputs of SVRG-LD ( $\mathbf{Z}_k$ ) and GLD ( $\mathbf{X}_k$ ) satisfy

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

where  $C_1$  is an absolute constant,  $\Gamma = 2(1 + 1/m)(b + 2G^2 + d/\beta)$  and  $L$  is length of each epoch.

- Combining previous results and Lemma 5 completes the proof.