# Constrained Exploration via Reflected Replica Exchange Stochastic Gradient Langevin Dynamics

Haoyang Zheng[1], Hengrong Du[2], Qi Feng[3], Wei Deng[4], Guang Lin[1]

[1] Purdue University  [2] Vanderbilt University  [3] Florida State University  [4] Morgan Stanley

## Problem Formulation

How **expensive** is it to generate samples which will converge to a target probability density $\pi$:

$$\mathrm{d}\pi(x_1, x_2) = \frac{1}{Z} \underbrace{e^{-\frac{U(x_1)}{\tau_1} - \frac{U(x_2)}{\tau_2}}}_{p(x_1, x_2)} \mathrm{d}x_1 \mathrm{d}x_2,$$

where $Z$ is a normalizing constant:

$$Z = \int_{\Omega \times \Omega} p(x_1, x_2) \mathrm{d}x_1 \mathrm{d}x_2.$$

## Motivations

- reSGLD may **over-explore** if high-temp chain delves too deeply into the distribution tails.
- It deteriorates the model's stability and lead to **poor predictions**.
- We proposed reflected reSGLD, which utilizes **reflection** steps within a bounded domain for **constrained non-convex** exploration.
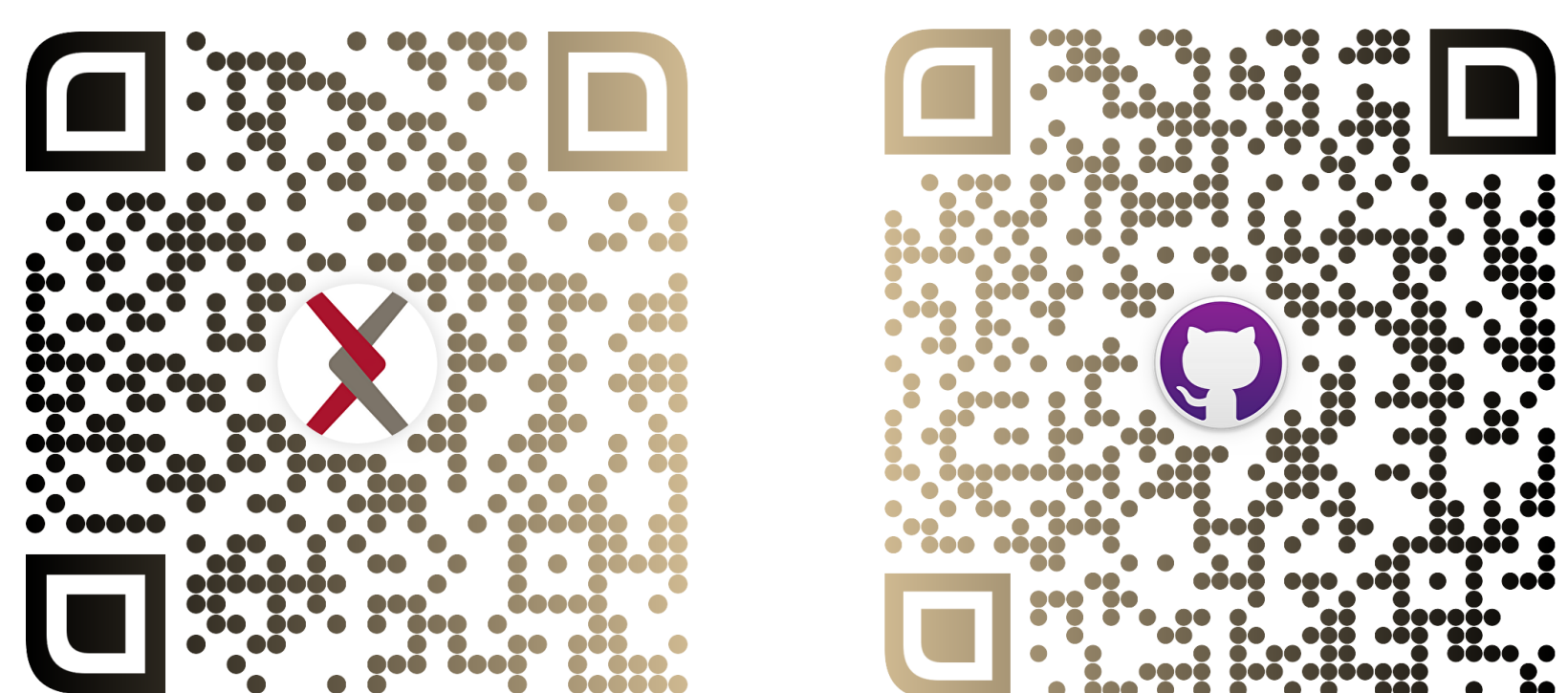
## Contributions

**Theory:**

- We prove the proposed method outperforms the naïve reSGLD.
- Reducing the domain diameter enhances mixing rates with a **quadratic** behavior.

**Experiments:**

- We introduces the novel use of the method in dynamic system identification.
- Extensive testing of r2SGLD against multi-modal distribution simulation and large-scale deep learning tasks.

## Methodology

### Reflected Replica Exchange Langevin Diffusion:

The system dynamics are described by the following SEDs:

$$\mathrm{d}\boldsymbol{\beta}_t^{(1)} = -\nabla U(\boldsymbol{\beta}_t^{(1)})\mathrm{d}t + \sqrt{2\tau_1}\ \mathrm{d}W_t^{(1)} + \nu(\boldsymbol{\beta}_t^{(1)})L^{(1)}(\mathrm{d}t),$$

$$\mathrm{d}\boldsymbol{\beta}_t^{(2)} = -\nabla U(\boldsymbol{\beta}_t^{(2)})\mathrm{d}t + \sqrt{2\tau_2}\ \mathrm{d}W_t^{(2)} + \nu(\boldsymbol{\beta}_t^{(2)})L^{(2)}(\mathrm{d}t),$$

**Notes:** $W_t$ is Wiener process; $\nu(\boldsymbol{\beta}_t)$ is inner unit vector; $L$ is independent local time.

### The Swap Function:

$$S(\boldsymbol{\beta}_t^{(1)}, \boldsymbol{\beta}_t^{(2)}) := e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right)\left(U(\boldsymbol{\beta}_t^{(1)}) - U(\boldsymbol{\beta}_t^{(2)})\right)}$$

### The r2SGLD Algorithm:

**Input** Initial parameters $\widetilde{\boldsymbol{\beta}}_1^{(1)}, \widetilde{\boldsymbol{\beta}}_1^{(2)}$; temperatures $\tau_1, \tau_2$; learn rate $\eta$.

  **for** $k = 1, 2, \cdots, K$ **do**

    **Sampling Step**

$$\widetilde{\boldsymbol{\beta}}_{k+1}^{(i)} = \mathcal{R}\left(\widetilde{\boldsymbol{\beta}}_k^{(i)} - \eta\nabla\widetilde{U}(\widetilde{\boldsymbol{\beta}}_k^{(i)}) + \sqrt{2\eta\tau_1}\boldsymbol{\xi}_k^{(i)}\right), i = 1, 2.$$

    **Swapping Step**

    Generate a uniform random number $u \in [0, 1]$.

    Compute $\widetilde{S}(\widetilde{\boldsymbol{\beta}}_k^{(1)}, \widetilde{\boldsymbol{\beta}}_k^{(2)}) = e^{\left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right)\left(\widetilde{U}(\widetilde{\boldsymbol{\beta}}_k^{(1)}) - \widetilde{U}(\widetilde{\boldsymbol{\beta}}_k^{(2)}) - \left(\frac{1}{\tau_1} - \frac{1}{\tau_2}\right)\frac{\widetilde{\sigma}^2}{C}\right)}.$

    **if** $u < \widetilde{S}$ **then**

      Swap $\widetilde{\boldsymbol{\beta}}_{k+1}^{(1)}$ and $\widetilde{\boldsymbol{\beta}}_{k+1}^{(2)}$.

**Output** Parameters $\{\widetilde{\boldsymbol{\beta}}_k^{(1)}\}_{k=1}^{K+1}$.

## Theoretical Analysis

### (a) Assumptions

*Assumption* A1. $\Omega$ is a compact domain with a boundary $\partial\Omega$ whose *second fundamental form* is bounded below by some constant $\kappa \leq 0$.

*Assumption* A2. The function $U \in C^2(\Omega)$. Since $\Omega$ is compact, there exists an $L > 0$ such that for all $x, y \in \Omega$,

$$\|\nabla U(x) - \nabla U(y)\| \leq L\|x - y\|.$$

### (c) Discretization Analysis

**Theorem 3.11** (Discretization error). *Assume that the domain $\Omega$ is convex, and Assumptions A1, A2 hold true, then*

$$\mathcal{W}_1(\mu_T, \widetilde{\mu}_T) \leq \widetilde{\mathcal{O}}\left(\eta^{1/4} + \sqrt{\max_k \mathbb{E}[\|\phi_k\|^2]} + \sqrt{\eta^{-1/2}\max_k\sqrt{\mathbb{E}[\|\psi_k\|^2]}}\right).$$

*where $\widetilde{\mu}_T$ denotes the distribution of $\widetilde{\boldsymbol{\beta}}_T^\eta$, which is the continuous-time interpolation for r2SGLD, $\phi_k := \nabla\widetilde{U} - \nabla U$ is the noise in the stochastic gradient, and $\psi_k := \widetilde{S} - S$ is the noise in the stochastic swapping rate.*

### (b) Continuous-time Analysis

**Theorem 3.4.** *Given any initial measure $\mu_0$ for which $\mathrm{d}\mu_0/\mathrm{d}\pi$ satisfies (5), the $\chi^2$-divergence to the invariant distribution $\pi$ decays exponentially according to:*

$$\chi^2(\mu_t\|\pi) \leq \chi^2(\mu_0\|\pi)\exp\left(-2t(1 + \eta_S)C_\mathrm{P}^{-1}\right), \quad (11)$$

*where $\eta_S := \inf\limits_{t>0} \dfrac{\mathcal{E}_S\left(\frac{\mathrm{d}\mu_t}{\mathrm{d}\pi}\right)}{\mathcal{E}\left(\frac{\mathrm{d}\mu_t}{\mathrm{d}\pi}\right)} - 1$ is the acceleration effect in $\chi^2$-divergence.*

**Theorem 3.8.** *Given any initial measure $\mu_0$ for which $\mathrm{d}\mu_0/\mathrm{d}\pi \geq 0$ and satisfying (5), the 2-Wasserstein distance between $\mu_t$ and $\pi$ satisfies the following accelerated exponential decay estimate:*

$$\mathcal{W}_2(\mu_t, \pi) \leq \sqrt{2C_\mathrm{LS}D(\mu_0\|\pi)}\exp\left(-t(1 + \delta_S)C_\mathrm{LS}^{-1}\right), \quad (15)$$

*where $\delta_S := \inf\limits_{t>0} \dfrac{\mathcal{E}_S\left(\sqrt{\frac{\mathrm{d}\mu_t}{\mathrm{d}\pi}}\right)}{\mathcal{E}\left(\sqrt{\frac{\mathrm{d}\mu_t}{\mathrm{d}\pi}}\right)} - 1$ is the acceleration effect in $\mathcal{W}_2$-distance.*

## Experiments

### Dynamic System Identification



### Constrained Multi-modal Simulations



Truth    R-SGLD

R-cycSGLD    r2SGLD



Kullback–Leibler Divergence

### Image Classification

| METHODS | METRICS (RESNET20) | | |
|---|---|---|---|
| | ACC (%) ↑ | NLL ↓ | BRIER (‰) ↓ |
| SGDM | $72.13 \pm 0.60$ | $9667 \pm 108$ | $2.78 \pm 0.05$ |
| SGHMC | $72.47 \pm 0.45$ | $9543 \pm 157$ | $2.75 \pm 0.05$ |
| cycSGHMC | $73.49 \pm 0.17$ | $8913 \pm 76$ | $2.65 \pm 0.02$ |
| reSGHMC | $75.01 \pm 0.14$ | $8552 \pm 69$ | $2.50 \pm 0.01$ |
| R-SGDM | $72.43 \pm 0.35$ | $9626 \pm 94$ | $2.75 \pm 0.03$ |
| R-SGHMC | $72.85 \pm 0.51$ | $9501 \pm 167$ | $2.73 \pm 0.05$ |
| R-cycSGHMC | $73.77 \pm 0.22$ | $8953 \pm 52$ | $2.62 \pm 0.02$ |
| **r2SGHMC** | $\mathbf{75.38 \pm 0.17}$ | $\mathbf{8489 \pm 66}$ | $\mathbf{2.46 \pm 0.02}$ |

| METHODS | METRICS (RESNET56) | | |
|---|---|---|---|
| | ACC (%) ↑ | NLL ↓ | BRIER (‰) ↓ |
| SGDM | $74.40 \pm 0.71$ | $9724 \pm 169$ | $3.59 \pm 0.23$ |
| SGHMC | $74.22 \pm 0.66$ | $9723 \pm 214$ | $3.23 \pm 0.21$ |
| cycSGHMC | $77.98 \pm 0.61$ | $8303 \pm 161$ | $3.19 \pm 0.20$ |
| reSGHMC | $78.87 \pm 0.44$ | $7406 \pm 130$ | $2.94 \pm 0.06$ |
| R-SGDM | $74.70 \pm 0.68$ | $9507 \pm 106$ | $3.53 \pm 0.18$ |
| R-SGHMC | $75.10 \pm 0.55$ | $9232 \pm 158$ | $3.36 \pm 0.23$ |
| R-cycSGHMC | $78.41 \pm 0.67$ | $7711 \pm 144$ | $3.12 \pm 0.11$ |
| **r2SGHMC** | $\mathbf{79.39 \pm 0.30}$ | $\mathbf{7155 \pm 91}$ | $\mathbf{2.89 \pm 0.02}$ |



Mixing Rates vs. Domain Diameters

### Takeaway

- The proposed r2SGLD algorithm performs the best.
- A smaller domain diameter in multi-modal simulation can improve the mixing rate.
- Large initial learning rate in CIFAR100 facilitates exploration.