

Gradient-Adjusted Underdamped Langevin Dynamics for Sampling*

Xinzhe Zuo[†], Stanley Osher[†], and Wuchen Li[‡]

Abstract. Sampling from a target distribution is a fundamental problem with wide-ranging applications in scientific computing and machine learning. Traditional Markov chain Monte Carlo (MCMC) algorithms, such as the unadjusted Langevin algorithm (ULA), derived from the overdamped Langevin dynamics, have been extensively studied. From an optimization perspective, the Kolmogorov forward equation of the overdamped Langevin dynamics can be treated as the gradient flow of the relative entropy in the space of probability densities embedded with Wasserstein-2 metrics. Several efforts have also been devoted to including momentum-based methods, such as underdamped Langevin dynamics, for faster convergence of sampling algorithms. Recent advances in optimization have demonstrated the effectiveness of primal-dual damping and Hessian-driven damping dynamics in achieving faster convergence when solving optimization problems. Motivated by these developments, we introduce a class of stochastic differential equations (SDEs) called gradient-adjusted underdamped Langevin dynamics (GAUL), which add stochastic perturbations in primal-dual damping dynamics and Hessian-driven damping dynamics from optimization. We prove that GAUL admits the correct invariant distribution, whose marginal is the target distribution. The proposed method outperforms overdamped and underdamped Langevin dynamics regarding convergence speed in the total variation distance for Gaussian target distributions. Moreover, using the Euler–Maruyama discretization, we show that the mixing time toward a biased target distribution only depends on the square root of the condition number of the target covariance matrix. In addition, we propose another discretization scheme based on the splitting method, which yields a smaller first-order asymptotic bias than the Euler–Maruyama scheme when sampling a Gaussian distribution. Numerical experiments for non-Gaussian target distributions, such as Bayesian regression problems and Bayesian neural networks, further illustrate the advantages of our approach over classical methods based on overdamped or underdamped Langevin dynamics. We also compare with the randomized Hamiltonian Monte Carlo method, showing that it achieves competitive performance.

Key words. Hessian-driven damping dynamics, primal-dual damping dynamics, Nesterov’s method, Langevin dynamics, optimal convergence rate

MSC codes. 37M25, 65C05, 82C31

DOI. 10.1137/24M1702015

1. Introduction. Sampling from a target distribution is a long-standing quest and has numerous applications in scientific computing, including Bayesian statistical inference

*Received by the editors October 14, 2024; accepted for publication (in revised form) July 22, 2025; published electronically October 10, 2025.

<https://doi.org/10.1137/24M1702015>

Funding: The first author is partially supported by AFOSR YIP award FA9550-23-1-0087. The first and second authors are partially funded by STROBE NSF STC DMR 1548924, AFOSR MURI FA9550-18-502, and ONR N00014-20-1-2787. The third author is partially supported by the AFOSR YIP award FA9550-23-1-0087, NSF DMS-2245097, and NSF RTG: 2038080.

[†]Department of Mathematics, University of California, Los Angeles, CA 90095 USA (zxz@math.ucla.edu, sjo@math.ucla.edu).

[‡]Department of Mathematics, University of South Carolina, Columbia, SC 29208 USA (wuchen@mailbox.sc.edu).

[55, 63, 52, 37], Bayesian inverse problems [67, 41, 29, 35], as well as Bayesian neural networks [77, 2, 73, 42, 54, 60]. In this direction, various algorithms have been developed to sample a target distribution $\pi \propto \exp(-f)$ for a given function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, where π is only known up to a normalization constant. In this area, a simple and popular algorithm is the unadjusted Langevin algorithm (ULA):

$$(1.1) \quad \mathbf{x}_{k+1} = \mathbf{x}_k - h\nabla f(\mathbf{x}_k) + \sqrt{2h}\mathbf{z}_k,$$

where $\mathbf{x}_k \in \mathbb{R}^d$, k is the iteration number, f is assumed to be a differentiable function, $h > 0$ is a step size, and \mathbf{z}_k is a d -dimensional random variable with independently and identically distributed (i.i.d) entries following standard Gaussian distributions. The ULA algorithm (1.1) comes from the forward Euler discretization of a stochastic differential equation (SDE) known as overdamped Langevin dynamics:

$$(1.2) \quad d\mathbf{x}_t = -\nabla f(\mathbf{x}_t)dt + \sqrt{2d}\mathbf{B}_t,$$

where $\mathbf{x}_t \in \mathbb{R}^d$ and \mathbf{B}_t is a standard d -dimensional Brownian motion. Under some mild conditions on f , it has been shown that the SDE (2.15) has a unique strong solution $\{\mathbf{x}_t, t \geq 0\}$ that is a Markov process [64, 58]. Moreover, the distribution of \mathbf{x}_t converges to the invariant distribution $\pi \propto \exp(-f)$ as $t \rightarrow \infty$. The asymptotic convergence guarantees of (1.1) have been established decades ago [71, 36, 57, 70]. In more recent years, nonasymptotic behaviors of (1.1) have also been explored by several works [25, 26, 32, 27, 21, 75].

An important result by [43] states that the Kolmogorov forward equation of Langevin dynamics corresponds to the gradient flow of the relative entropy functional in the space of probability density functions with the Wasserstein-2 metric. This observation serves as a bridge between the sampling community and the optimization community by studying optimization problems in Wasserstein-2 space. In the field of optimization, Nesterov's accelerated gradient [62] is a first order algorithm for finding the minimum of a convex/strongly convex objective function f . The intuition is that Nesterov's method incorporates momentum into the updates. It is much faster than the traditional gradient descent method, in the sense that the convergence speed for convex functions is $\mathcal{O}(\frac{1}{k^2})$ where k is the number of iterations compared to $\mathcal{O}(\frac{1}{k})$ for gradient descent. The convergence speed of Nesterov's method for L -smooth, m -strongly convex functions is $\mathcal{O}(\exp(-k/\sqrt{\kappa}))$, where $\kappa = L/m$ is the condition number of f compared to $\mathcal{O}(\exp(-k/\kappa))$ for gradient descent. By taking the step size to 0, one obtains a second-order ODE for Nesterov's method called the Nesterov's accelerated gradient flow or Nesterov's ODE [68, 5]. In recent years, one extends the gradient flow of the relative entropy into Nesterov's accelerated gradient flow [68], which is explored in [76, 69, 53] from different perspectives. For the optimization in Wasserstein-2 space perspective, [76, 69, 19] study a class of accelerated dynamics with depending on the score function, i.e., the gradient of logarithm of density function. This results in the approximation of a nonlinear partial differential equation, known as the damped Euler equation [16]. In this case, the optimal choices of parameters for sampling a target distribution share similarities with the classical Nesterov's accelerated gradient flow. On the other hand, from a stochastic dynamics perspective, a line of research has been devoted to study the accelerated version of Langevin dynamics, known as the underdamped Langevin dynamics [15, 22, 53, 78]. As explained later in section 2.2,

the underdamped Langevin dynamics consists of a deterministic component and a stochastic component. The deterministic component exactly corresponds to the Nesterov's accelerated gradient flow. The marginal of invariant distribution in x -axis satisfies the target distribution. However, the optimal choice of parameters in underdamped Langevin dynamics might not directly follow the classical Nesterov's method [22].

Recently, [79] proposed to use the primal-dual hybrid gradient (PDHG) method [18, 74] to solve unconstrained optimization problems. The original PDHG method is designed for optimization problem with linear constraints. [79] formulated the optimality condition $\nabla f(\mathbf{x}) = 0$ of a strongly convex function f into the solution of a saddle point problem

$$\inf_{\mathbf{x} \in \mathbb{R}^d} \sup_{\mathbf{p} \in \mathbb{R}^d} \langle \nabla f(\mathbf{x}), \mathbf{p} \rangle - \frac{\gamma}{2} \|\mathbf{p}\|^2,$$

where $\gamma > 0$ is a selected regularization parameter. They proceed by using the PDHG algorithm with appropriate preconditioners to solve the above saddle point problem. By taking the limit as the step size goes to zero, their algorithm yields a continuous-time flow, which is a second-order ordinary differential equation (ODE) called the primal-dual damping (PDD) dynamics. In particular, the PDD dynamic contains Nesterov's ODE [68]. In other words, Nesterov's ODE is a special case of PDD dynamics. The PDD dynamics also shares similarities with the Hessian-driven damping dynamics that has been studied in recent years [5, 3, 4]. The main difference between the PDD dynamics and the Nesterov's ODE is a second-order term $\nabla^2 f(\mathbf{x})\dot{\mathbf{x}}$ that appears in the former. This term is also presented in the Hessian driven damping dynamics. It has been observed that the PDD dynamics and the Hessian driven damping dynamics yield faster convergence toward the global minimum than the traditional gradient flow and Nesterov's ODE. Therefore, it is natural to extend the PDD dynamics and Hessian driven damping dynamics to SDEs for sampling a target distribution.

In this paper, we take inspirations from [79, 3] to design a system of SDE called gradient-adjusted underdamped Langevin dynamics (GAUL) that resembles the primal-dual damping dynamics and the Hessian driven damping dynamics. Consider

$$(1.3) \quad \begin{pmatrix} d\mathbf{x}_t \\ d\mathbf{p}_t \end{pmatrix} = \begin{pmatrix} -a\mathbf{C}\nabla f(\mathbf{x}_t)dt + \mathbf{C}\mathbf{p}_t dt \\ -\nabla f(\mathbf{x}_t)dt - \gamma\mathbf{p}_t dt \end{pmatrix} + \sqrt{\begin{pmatrix} 2a\mathbf{C} & \mathbf{I} - \mathbf{C} \\ \mathbf{I} - \mathbf{C} & 2\gamma\mathbf{I} \end{pmatrix}} \begin{pmatrix} d\mathbf{B}_t^{(1)} \\ d\mathbf{B}_t^{(2)} \end{pmatrix}$$

for some constants $a, \gamma > 0$, whose detailed choices will be explained later. \mathbf{C} is a preconditioner such that the diffusion matrix in front of the Brownian motion term is well-defined and positive semidefinite. And $\mathbf{B}_t^{(i)}$ is a standard Brownian motion in \mathbb{R}^d for $i = 1, 2$. The superscript on \mathbf{B}_t indicates that $\mathbf{B}_t^{(1)}$ and $\mathbf{B}_t^{(2)}$ are independent. We show that the invariant distribution GAUL (1.3) is the desired target distribution of the form $\frac{1}{Z} \exp(-f(\mathbf{x}) - \|\mathbf{p}\|^2/2)$. Noticeably, the \mathbf{x} -marginal distribution is the target distribution π . Additionally, we demonstrate that for a quadratic function f , GAUL achieves the exponential convergence and outperforms both overdamped and underdamped Langevin dynamics. A series of numerical examples are provided to demonstrate the advantage of the proposed method.

To illustrate the main idea, we summarize main theoretical results into the following informal theorem.

Theorem 1.1 (Informal). Suppose that $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is given by $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \Lambda \mathbf{x}$ with a symmetric positive definite matrix $\Lambda \in \mathbb{R}^{d \times d}$ with eigenvalues $s_1 \geq s_2 \geq \dots \geq s_d > 0$. Let $\kappa = s_1/s_d$ be the condition number of matrix Λ . And let $\mathbf{C} = \mathbf{I}$.

- (1) Denote by $\rho_x(\mathbf{x}, t)$ the law of \mathbf{x}_t driven by (1.3), and $\pi(\mathbf{x}) \propto \exp(-f(\mathbf{x}))$ the target distribution. Let $a > 0$, $\gamma = as_d + 2\sqrt{s_d}$. Then it takes at most $t = \mathcal{O}(\log(d/\delta)/(as_d + 2\sqrt{s_d}))$ for the total variation distance between $\rho_x(\mathbf{x}, t)$ and $\pi(\mathbf{x})$ to decrease to δ .
- (2) Denote by $\tilde{\rho}_x(\mathbf{x}, k)$ the law of \mathbf{x} after k iterations of the Euler–Maruyama discretization of (1.3). Suppose $\sqrt{s_1} - \sqrt{s_d} \geq 2$, $a = 1$, $\gamma = s_d + 2\sqrt{s_d}$ and consider the Euler–Maruyama discretization of (1.3) with step size $h = 1/5s_1$. Then it takes at most $N = \mathcal{O}(\log(d/\delta)/(\kappa^{-1} + (\kappa s_1)^{-1/2}))$ iterations for the total variation distance between $\tilde{\rho}_x(\mathbf{x}, k)$ and $\tilde{\pi}(\mathbf{x})$ to decrease to δ , where $\tilde{\pi}(\mathbf{x})$ is a biased target distribution given by (SM1.24).
- (3) When taking $a = \frac{2}{\sqrt{s_1} - \sqrt{s_d}}$, $\gamma = as_d + 2\sqrt{s_d}$ and $h = \frac{1}{2(as_1 + \gamma)}$, we can improve the number of iterations in (2) to $N = \mathcal{O}(\sqrt{\kappa} \log(d/\delta))$.

The detailed version of Theorem 1.1 is given in Theorem 3.9, Theorem 3.16, and Theorem 3.17. It is worth noting that GAUL (1.3) reduces to underdamped Langevin dynamics when $a = 0$ and $\mathbf{C} = \mathbf{I}$. Our theorem implies that in the Gaussian case, GAUL converges to the target measure faster than underdamped Langevin dynamics. In particular, we demonstrate that the Euler–Maruyama discretization admits a mixing time proportional to the square root of the condition number of covariance matrix. While this work primarily focuses on Gaussian distributions, our numerical experiments also explore non-log-concave target distributions in Bayesian linear regressions and Bayesian neural networks, which demonstrate potential advantages of GAUL over overdamped and underdamped Langevin dynamics. Extending these results to more general distributions and discretization schemes is an important future research direction. The choice of preconditioner \mathbf{C} is tricky as one needs to guarantee that the diffusion matrix in (1.3) is positive semidefinite. Therefore, we mainly focus on the case when $\mathbf{C} = \mathbf{I}$. We address our results for $\mathbf{C} \neq \mathbf{I}$ in Remark 3.10 and Remark 3.20. For $\mathbf{C} = \mathbf{I}$, [51] also explored dynamics (1.3), which they called Hessian-free high-resolution (HFHR) dynamics. For this closely related work, we provide some comparisons later in Remark 2.4.

This paper is organized as follows. In section 2, we review the connection between optimization methods and sampling dynamics, which leads to the construction of our proposed SDE called gradient-adjusted underdamped Langevin dynamics (GAUL). Our main results are presented in section 3, where we prove the exponential convergence of GAUL to the target distribution when the target measure follows a Gaussian distribution. We also study the Euler–Maruyama discretization of GAUL and prove its linear convergence to a biased target distribution. In addition, we propose another splitting-based discretization scheme *BAGOGAB* and show that it admits a smaller first-order asymptotic bias than the Euler–Maruyama discretization when sampling the Gaussian target distribution. Lastly, in section 4, we present several numerical examples to compare GAUL with overdamped, underdamped Langevin dynamics, as well as the randomized Hamiltonian Monte Carlo method.

2. Preliminaries. In this section, we briefly review the relation among Euclidean gradient flows, overdamped Langevin dynamics, and Wasserstein gradient flows. We then draw the connection between the underdamped Langevin dynamics and Nesterov’s ODEs. We next

review primal-dual damping (PDD) flows [79] and Hessian driven damping dynamics. Finally, we introduce a new SDE called gradient-adjusted underdamped Langevin dynamics (GAUL) for sampling, which resembles the PDD flow and the Hessian-driven damping dynamics with designed stochastic perturbations in terms of Brownian motions.

2.1. Gradient descent, unadjusted Langevin algorithms, and optimal transport gradient flows. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable convex function with L -Lipschitz gradient. The classical gradient descent algorithm for finding the global minimum of $f(\mathbf{x})$ is an iterative algorithm that reads

$$(2.1) \quad \mathbf{x}_{k+1} = \mathbf{x}_k - h\nabla f(\mathbf{x}_k),$$

where $h > 0$ is the step size. When f is convex and the step size is not too large, this algorithm converges at a rate of $\mathcal{O}(k^{-1})$. When f is m -strongly convex, the same algorithm can be shown to converge at a rate of $\mathcal{O}((1 - m/L)^k)$ if the step size is chosen appropriately. The gradient descent algorithm (2.1) can be understood as the forward Euler time discretization of the gradient flow

$$(2.2) \quad \dot{\mathbf{x}}(t) = -\nabla f(\mathbf{x}(t)),$$

where $\mathbf{x}(t)$ describes a trajectory in \mathbb{R}^d that travels in the direction of the steepest descent. Similar convergence results can be obtained for the gradient flow (2.2). When f is convex, the gradient flow (2.2) converges at a rate of $\mathcal{O}(t^{-1})$. When f is assumed to be m -strongly convex, the gradient flow (2.2) converges at a rate of $\mathcal{O}(\exp(-mt))$.

While the goal of optimization is to find the global minimum of f , the goal of sampling algorithm is to sample from a distribution of the form $\frac{1}{Z_1} \exp(-f(\mathbf{x}))$, where the normalization constant $Z_1 > 0$ is assumed to be finite, i.e.,

$$Z_1 = \int_{\mathbb{R}^d} e^{-f(\mathbf{x})} d\mathbf{x} < +\infty.$$

The classical unadjusted Langevin algorithm (ULA) given in (1.1) is a simple modification to the gradient descent method. Recall that ULA is given by

$$(2.3) \quad \mathbf{x}_{k+1} = \mathbf{x}_k - h\nabla f(\mathbf{x}_k) + \sqrt{2h}\mathbf{z}_k,$$

where \mathbf{z}_k is a d -dimensional standard Gaussian random variable and h is the step size. We obtain (2.3) from (2.1) by adding a Gaussian noise term \mathbf{z}_k scaled by $\sqrt{2h}$. Similar to how (2.1) can be viewed as the Euler discretization of (2.2), ULA (2.3) represents the forward Euler discretization of the overdamped Langevin dynamics:

$$(2.4) \quad d\mathbf{x}_t = -\nabla f(\mathbf{x}_t)dt + \sqrt{2}d\mathbf{B}_t,$$

where \mathbf{B}_t is a standard d -dimensional Brownian motion. Denote by $\rho(\mathbf{x}, t)$ the probability density function for \mathbf{x}_t . Then the Kolmogorov forward equation (also known as the Fokker-Planck equation) of the overdamped Langevin dynamics (2.4) is given as

$$(2.5) \quad \frac{\partial \rho}{\partial t} = \nabla \cdot (\rho \nabla f) + \Delta \rho.$$

Clearly, $\pi(\mathbf{x}) = \frac{1}{Z_1} \exp(-f(\mathbf{x}))$ is a stationary solution of the Fokker–Planck equation (2.5). In other words, note that $\nabla\pi = -\pi\nabla f$, then

$$0 = \partial_t \pi = \nabla \cdot (\pi \nabla f) + \Delta \pi = \nabla \cdot ((\pi \nabla f + \nabla \pi)).$$

In the literature, one can also study the gradient drift Fokker–Planck equation (2.5) from a gradient flow point of view. This means that (2.5) is a gradient flow in the probability space embedded with a Wasserstein-2 metric. We review some facts on a formal manner; see rigorous treatment in [1].

Define the probability space on \mathbb{R}^d with finite second-order moment:

$$\mathcal{P}(\mathbb{R}^d) = \left\{ \rho(\cdot) \in C^\infty : \int_{\mathbb{R}^d} \rho(\mathbf{x}) d\mathbf{x} = 1, \int_{\mathbb{R}^d} |\mathbf{x}|^2 \rho(\mathbf{x}) d\mathbf{x} < \infty, \rho(\cdot) \geq 0 \right\}.$$

We note that $\mathcal{P}(\mathbb{R}^d)$ can be equipped with the L_2 –Wasserstein metric g_W at each $\rho \in \mathcal{P}(\mathbb{R}^d)$ to form a Riemannian manifold $(\mathcal{P}(\mathbb{R}^d), g_W)$. Let $\mathcal{F} : \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}$ be an energy functional on $\mathcal{P}(\mathbb{R}^d)$. To be more precise, denote the Wasserstein gradient operator of functional $\mathcal{F}(\rho)$ at the density function $\rho \in \mathcal{P}(\mathbb{R}^d)$, such that

$$\text{grad}_W \mathcal{F}(\rho) := -\nabla \cdot \left(\rho \nabla \frac{\delta}{\delta \rho} \mathcal{F}(\rho) \right),$$

where $\frac{\delta}{\delta \rho}$ is the L_2 –first variation with respect to ρ . This yields that the gradient descent flow in the Wasserstein-2 space satisfies

$$\frac{\partial \rho}{\partial t} = -\text{grad}_W \mathcal{F}(\rho) = \nabla \cdot \left(\rho \nabla \frac{\delta}{\delta \rho} \mathcal{F}(\rho) \right).$$

The above PDE is also named the *Wasserstein gradient descent flow*, in short Wasserstein gradient flows, which depend on the choices of the energy functionals $\mathcal{F}(\rho)$.

An important example observed by [43] is as follows. Consider the relative entropy functional, also named Kullback–Leibler(KL) divergence

$$\mathcal{F}(\rho) := \text{D}_{\text{KL}}(\rho || \pi) = \int_{\mathbb{R}^d} \rho(\mathbf{x}) \log \left(\frac{\rho(\mathbf{x})}{\pi(\mathbf{x})} \right) d\mathbf{x}.$$

One can show that the Fokker–Planck equation (2.5) is the gradient flow of the relative entropy in $(\mathcal{P}(\mathbb{R}^d), g_W)$. Upon recognizing $\frac{\delta}{\delta \rho} \text{D}_{\text{KL}}(\rho || \pi) = \log \left(\frac{\rho}{\pi} \right) + 1$, we obtain that (2.5) can be expressed as

$$\begin{aligned} \frac{\partial \rho}{\partial t} &= -\text{grad}_W \text{D}_{\text{KL}}(\rho || \pi) = \nabla \cdot \left(\rho \nabla \log \left(\frac{\rho}{\pi} \right) \right) \\ (2.6) \quad &= \nabla \cdot (\rho \nabla \log \rho) - \nabla \cdot (\rho \nabla \log \pi) \\ &= \Delta \rho + \nabla \cdot (\rho \nabla f), \end{aligned}$$

where we use facts that $\rho \nabla \log \rho = \nabla \rho$ and $\nabla \log \pi = -\nabla f$.

We note that the gradient of the logarithm of the density function, i.e., $\nabla \log \rho$, is often called the score function. The analysis of score functions are essential in understanding the convergence behavior of the Fokker–Planck equation (2.5) toward its invariant distribution; see related analytical studies in [34].

2.2. Nesterov's ODEs and underdamped Langevin dynamics. Consider the problem of minimizing $f : \mathbb{R}^d \rightarrow \mathbb{R}$ for some convex function f with L -Lipschitz gradient. [62] proposed the following iterations:

$$(2.7a) \quad \mathbf{x}_{k+1} = \mathbf{p}_k - h \nabla f(\mathbf{p}_k),$$

$$(2.7b) \quad \mathbf{p}_{k+1} = \mathbf{x}_{k+1} + \gamma_k (\mathbf{x}_{k+1} - \mathbf{x}_k),$$

where $\gamma_k = (k-1)/(k-2)$. [62] showed that the above method converges at a rate of $\mathcal{O}(k^{-2})$ instead of $\mathcal{O}(k^{-1})$ which is the convergence rate of the classical gradient descent method. If f is further assumed to be m -strongly convex, then taking $h = 1/L$ and $\gamma_k = \frac{1-\sqrt{\kappa}}{1+\sqrt{\kappa}}$, where $\kappa = L/m$, yields a convergence rate of $\mathcal{O}(\exp(-k/\sqrt{\kappa}))$. This is also considerably faster than gradient descent, which is $\mathcal{O}((1-\kappa^{-1})^k)$. [68] showed that the continuous-time limit of Nesterov's accelerated gradient method [62] satisfies a second order ODE:

$$(2.8) \quad \ddot{\mathbf{x}} + \gamma_t \dot{\mathbf{x}} + \nabla f(\mathbf{x}) = 0.$$

If f is a convex function, then $\gamma_t = 3/t$; if f is a m -strongly convex function, then $\gamma_t = \gamma = 2\sqrt{m}$. As observed in [56], (2.8) can be formulated as a damped Hamiltonian system:

$$(2.9) \quad \begin{pmatrix} \dot{\mathbf{x}} \\ \dot{\mathbf{p}} \end{pmatrix} = \begin{pmatrix} 0 & \mathbf{I} \\ -\mathbf{I} & 0 \end{pmatrix} \begin{pmatrix} \nabla_x H(\mathbf{x}, \mathbf{p}) \\ \nabla_p H(\mathbf{x}, \mathbf{p}) \end{pmatrix} + \begin{pmatrix} 0 & \mathbf{I} \\ -\mathbf{I} & -\gamma_t \mathbf{I} \end{pmatrix} \begin{pmatrix} \nabla_x H(\mathbf{x}, \mathbf{p}) \\ \nabla_p H(\mathbf{x}, \mathbf{p}) \end{pmatrix},$$

where the Hamiltonian function is defined as $H(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \|\mathbf{p}\|^2/2$, $\mathbf{p} \in \mathbb{R}^d$. On the other hand, the underdamped Langevin dynamics for sampling $\Pi(\mathbf{x}, \mathbf{p}) \propto \exp(-f(\mathbf{x}) - \|\mathbf{p}\|^2/2)$ is given by the system of SDE:

$$\begin{aligned} d\mathbf{x}_t &= \mathbf{p}_t dt, \\ d\mathbf{p}_t &= -\nabla f(\mathbf{x}_t) dt - \gamma_t \mathbf{p}_t dt + \sqrt{2\gamma_t} d\mathbf{B}_t, \end{aligned}$$

where γ_t is some damping parameter, and \mathbf{B}_t is a d -dimensional standard Brownian motion. This can be reformulated as

$$(2.10) \quad \begin{pmatrix} d\mathbf{x}_t \\ d\mathbf{p}_t \end{pmatrix} = \begin{pmatrix} 0 & \mathbf{I} \\ -\mathbf{I} & -\gamma_t \mathbf{I} \end{pmatrix} \begin{pmatrix} \nabla_x H(\mathbf{x}, \mathbf{p}) \\ \nabla_p H(\mathbf{x}, \mathbf{p}) \end{pmatrix} dt + \begin{pmatrix} 0 & 0 \\ 0 & \sqrt{2\gamma_t} \mathbf{I} \end{pmatrix} d\mathbf{B}_t,$$

where \mathbf{B}_t is a $2d$ -dimensional standard Brownian motion. Observe that by adding a suitable Brownian motion term (the last term on the right hand side of (2.10)) to (2.9), Nesterov's accelerated gradient method for convex optimization becomes an algorithm for sampling $\Pi(\mathbf{x}, \mathbf{p}) = \frac{1}{Z} \exp(-f(\mathbf{x}) - \|\mathbf{p}\|^2/2)$, where $Z := \int_{\mathbb{R}^{2d}} \exp(-f(\mathbf{x}) - \|\mathbf{p}\|^2/2) d\mathbf{x} d\mathbf{p} < +\infty$ is a normalization constant. Moreover, the \mathbf{x} -marginal of $\Pi(\mathbf{x}, \mathbf{p})$ is simply $\pi(\mathbf{x}) = \frac{1}{Z_1} \exp(-f(\mathbf{x}))$ up to a normalizing constant $Z_1 := \int_{\mathbb{R}^{2d}} \exp(-f(\mathbf{x}) - \|\mathbf{p}\|^2/2) d\mathbf{x} d\mathbf{p} < +\infty$. Therefore, (2.10) can be used to sample distributions of the form $\exp(-f(\mathbf{x}))/Z_1$. We postpone the proofs in terms of Fokker-Planck equations and their invariant distributions in Propositions 2.1 and 2.2.

2.3. Primal-dual damping dynamics and Hessian driven damping dynamics. Recently, [79] proposed to solve an unconstrained strongly convex optimization problem using the PDHG method by considering the saddle point problem

$$\inf_{\mathbf{x} \in \mathbb{R}^d} \sup_{\mathbf{p} \in \mathbb{R}^d} \langle \nabla f(\mathbf{x}), \mathbf{p} \rangle - \frac{\gamma}{2} \|\mathbf{p}\|^2,$$

where γ is a damping parameter, and $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is m -strongly convex. Note that the saddle point $(\mathbf{x}^*, \mathbf{p}^*)$ for the above inf-sup problem satisfies $\nabla f(\mathbf{x}^*) = \mathbf{p}^* = 0$. Then the primal-dual damping (PDD) algorithm [79] admits the following iterations:

$$\begin{aligned}\mathbf{p}_{k+1} &= \frac{1}{1 + \tau_1 \gamma} \mathbf{p}_k + \frac{\tau_1}{1 + \tau_1 \gamma} \nabla f(\mathbf{x}_k), \\ \tilde{\mathbf{p}}_{k+1} &= \mathbf{p}_{k+1} + \omega(\mathbf{p}_{k+1} - \mathbf{p}_k), \\ \mathbf{x}_{k+1} &= \mathbf{x}_k - \tau_2 \mathbf{C}(\mathbf{x}_k) \tilde{\mathbf{p}}_{k+1},\end{aligned}$$

where $\tau_1, \tau_2 > 0$ are dual and primal step sizes, $\omega > 0$ is an extrapolation parameter, and $\mathbf{C} \in \mathbb{R}^{d \times d}$ is a preconditioning positive definite matrix that could depend on \mathbf{x}_k and t . The continuous-time limit of the PDD algorithm can be obtained by letting $\tau_1, \tau_2 \rightarrow 0$ while keeping $\tau_1 \omega \rightarrow a$ for some $a > 0$. This yields a second-order ODE called the PDD flow:

$$(2.11) \quad \ddot{\mathbf{x}} + \left(\gamma + a\mathbf{C}\nabla^2 f(\mathbf{x}) - \dot{\mathbf{C}}\mathbf{C}^{-1} \right) \dot{\mathbf{x}} + \mathbf{C}\nabla f(\mathbf{x}) = 0.$$

In the case when \mathbf{C} is constant, (2.11) reads

$$(2.12) \quad \ddot{\mathbf{x}} + \left(\gamma + a\mathbf{C}\nabla^2 f(\mathbf{x}) \right) \dot{\mathbf{x}} + \mathbf{C}\nabla f(\mathbf{x}) = 0.$$

And when $\mathbf{C} = \mathbf{I}$, the PDD flow simplifies to

$$(2.13) \quad \ddot{\mathbf{x}} + \gamma \dot{\mathbf{x}} + a\nabla^2 f(\mathbf{x}) \dot{\mathbf{x}} + \nabla f(\mathbf{x}) = 0.$$

This corresponds to the Hessian driven damping dynamic [3] when $\gamma = 2\sqrt{m}$. The terminology ‘Hessian driven damping’ comes from the Hessian term $\nabla^2 f(\mathbf{x}) \dot{\mathbf{x}}$ in (2.13), which is controlled by a constant $a \geq 0$. When $a = 0$, (2.13) reduces to Nesterov’s ODE (2.8). As in dynamics (2.9), we can express (2.11) as

$$(2.14) \quad \begin{pmatrix} \dot{\mathbf{x}} \\ \dot{\mathbf{p}} \end{pmatrix} = \begin{pmatrix} -a\mathbf{C} & \mathbf{C} \\ (\gamma a - 1)\mathbf{I} & -\gamma\mathbf{I} \end{pmatrix} \begin{pmatrix} \nabla_x H(\mathbf{x}, \mathbf{p}) \\ \nabla_p H(\mathbf{x}, \mathbf{p}) \end{pmatrix},$$

where as before the Hamiltonian function is $H(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \|\mathbf{p}\|^2/2$. Note that one of the key differences between (2.9) and (2.14) is that the top left block of the preconditioner matrix is nonzero in (2.14), which gives rise to the Hessian damping term $\nabla^2 f(\mathbf{x}) \dot{\mathbf{x}}$. Throughout this paper, we focus on the dynamical system (2.14).

2.4. Gradient-adjusted underdamped Langevin dynamics. We design a sampling dynamics that resembles the PDD flow and the Hessian driven damping with stochastic perturbations by Brownian motions. Our goal is still to sample a distribution proportional to $\exp(-f(\mathbf{x}))$ for some $f: \mathbb{R}^d \rightarrow \mathbb{R}$. Let $H(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \|\mathbf{p}\|^2/2$. And denote by $\mathbf{X} = (\mathbf{x}, \mathbf{p}) \in \mathbb{R}^{2d}$. We consider the following SDE.

$$(2.15) \quad d\mathbf{X}_t = -\mathbf{Q}\nabla H(\mathbf{X}_t)dt + \sqrt{2\text{sym}(\mathbf{Q})}d\mathbf{B}_t,$$

where $\mathbf{Q} \in \mathbb{R}^{2d \times 2d}$ is of the form

$$(2.16) \quad \mathbf{Q} = \begin{pmatrix} a\mathbf{C} & -\mathbf{C} \\ \mathbf{I} & \gamma\mathbf{I} \end{pmatrix},$$

for some constant $a, \gamma \in \mathbb{R}$, and symmetric positive definite $\mathbf{C} \in \mathbb{R}^{d \times d}$. $\nabla H(\mathbf{X}_t) = (\nabla_x H(\mathbf{X}_t), \nabla_p H(\mathbf{X}_t))^T$. And $\text{sym}(\mathbf{Q}) = \frac{1}{2}(\mathbf{Q} + \mathbf{Q}^T)$ is the symmetrization of \mathbf{Q} . We assume that $\text{sym}(\mathbf{Q})$ is positive semidefinite. Throughout this paper, we will limit our discussion to $a, \gamma \geq 0$. \mathbf{B}_t is a $2d$ -dimensional standard Brownian motion. Observe that when $a = 0$, (2.15) reduces to underdamped Langevin dynamics (2.10). When $a > 0$, (2.15) has an additional gradient term $a\mathbf{C}\nabla f(\mathbf{x}_t)$ in the $d\mathbf{x}_t$ equation. Thus, we call (2.15) gradient-adjusted underdamped Langevin dynamics. Let us examine the probability density function $\rho(\mathbf{X}, t)$ of the diffusion governed by (2.15). This is described by the following Fokker–Planck equation:

$$(2.17) \quad \frac{\partial \rho}{\partial t} = \nabla \cdot (\mathbf{Q}\nabla H\rho) + \sum_{i,j=1}^{2d} \frac{\partial^2}{\partial X_i \partial X_j} (Q_{ij}\rho).$$

We assume that f is differentiable and ∇f is a smooth Lipschitz vector field. This ensures that the Fokker–Planck equation (2.17) has a smooth solution when $t > 0$ for a given initial condition, such that $\rho(\mathbf{X}, 0) \geq 0$ and $\int_{\mathbb{R}^{2d}} \rho(\mathbf{X}, 0) d\mathbf{X} = 1$.

Denote by $\Pi(\mathbf{X}) = \frac{1}{Z} e^{-H(\mathbf{X})}$, where Z is a normalization constant such that $\Pi(\mathbf{X})$ integrates to one on \mathbb{R}^{2d} . We show that $\Pi(\mathbf{X})$ is the invariant distribution of (2.17). First, we have the following decomposition for (2.17).

Proposition 2.1 ([34] Proposition 1). *The Fokker–Planck equation (2.17) can be decomposed as*

$$(2.18) \quad \frac{\partial \rho}{\partial t} = \nabla \cdot \left(\rho \text{sym}(\mathbf{Q}) \nabla \log \frac{\rho}{\Pi} \right) + \nabla \cdot (\rho \Gamma),$$

where

$$(2.19) \quad \begin{aligned} \Gamma(\mathbf{X}) &:= \text{sym}(\mathbf{Q}) \nabla \log(\Pi(\mathbf{X})) + \mathbf{Q} \nabla H(\mathbf{X}) \\ &= \frac{1}{2}(\mathbf{Q} - \mathbf{Q}^T) \nabla H(\mathbf{X}). \end{aligned}$$

In particular, the following equality holds:

$$\nabla \cdot (\Pi(\mathbf{X}) \Gamma(\mathbf{X})) = 0.$$

The proof is presented in section SM3. Observe that the first term on the right-hand side of (2.18) is a Kullback–Leibler (KL) divergence functional that appears in a Fokker–Planck equation associated with the overdamped Langevin dynamics (2.5). The second term is due to the fact that the drift term $-\mathbf{Q}\nabla H$ in (2.15) is a nongradient vector field. One can also decompose (2.15) into reversible and nonreversible parts of SDEs [6].

$$d\mathbf{X}_t = -(\mathbf{Q}_1 \nabla H(\mathbf{X}_t) + \mathbf{Q}_2 \nabla H(\mathbf{X}_t)) dt + \sqrt{2\mathbf{Q}_1} d\mathbf{B}_t,$$

where

$$\mathbf{Q}_1 = \begin{pmatrix} a\mathbf{C} & \frac{1}{2}(\mathbf{I} - \mathbf{C}) \\ \frac{1}{2}(\mathbf{I} - \mathbf{C}) & \gamma\mathbf{I} \end{pmatrix}, \quad \mathbf{Q}_2 = \begin{pmatrix} 0 & \frac{1}{2}(-\mathbf{I} - \mathbf{C}) \\ \frac{1}{2}(\mathbf{I} + \mathbf{C}) & 0 \end{pmatrix}.$$

Proposition 2.2. $\Pi(\mathbf{X})$ is an invariant distribution for (2.17).

The proof is based on a straightforward calculation: When $\rho = \Pi$, we have $\nabla \cdot (\rho\Gamma) = 0$, and therefore $\frac{\partial \rho}{\partial t} = 0$. For completeness, we have included this calculation in section SM3. This shows that $\Pi(\mathbf{X})$ is indeed the invariant distribution of (2.17). Like the underdamped Langevin dynamics, the \mathbf{x} -marginal of the invariant distribution is $\exp(-f(\mathbf{x}))$ up to some normalization constant. Therefore, (2.15) can be used for sampling $\frac{1}{Z_1} \exp(-f(\mathbf{x}))$ by first jointly sampling $\mathbf{X} = (\mathbf{x}, \mathbf{p})$ and then taking out the \mathbf{x} -marginal.

Remark 2.3. GAUL can also be viewed as a preconditioned overdamped Langevin dynamics on the space of $(\mathbf{x}, \mathbf{p}) \in \mathbb{R}^{2d}$. Designing optimal preconditioning matrix and optimal diffusion matrix have been studied in literature; see [17, 7, 38, 47, 39, 20, 50, 49]. In particular, [50] considered the necessary condition on the optimal diffusion coefficient by studying the spectral gap of the generator associated with the SDE, which requires the solution to an optimization subproblem. While the problem considered by [50] is more general, our diffusion matrix (2.16) is much simpler and does not require solving an optimization problem. Another closely related work is [49], which considered preconditioning of the form $\mathbf{Q} = \mathbf{I} + \mathbf{J}$. Here \mathbf{I} is the identity matrix and \mathbf{J} is skew-symmetric, i.e., $\mathbf{J} = -\mathbf{J}^T$. [49] studied the optimal \mathbf{J} when the potential f is a quadratic function, which is also the focus of this work.

Remark 2.4. In [51], the authors also studied (1.3) with $\mathbf{C} = \mathbf{I}$ which they called Hessian-free high-resolution (HFHR) dynamics. They considered potential functions f that are L -smooth and m -strongly convex. They proved a convergence rate of $\frac{\sqrt{m}}{2\sqrt{\kappa}}$ in continuous time in terms of Wasserstein-2 distance between the target and sample measure. [51] used the randomized midpoint method [65] combined with Strang splitting as their discretization and showed an iteration complexity of $\tilde{\mathcal{O}}(\sqrt{d}/\varepsilon)$. Specifically, [51] showed that for a two-dimensional Gaussian target measure, under the optimal choice of parameter (damping parameter γ and step size h) for underdamped Langevin dynamics with Euler–Maruyama discretization, the convergence rate is $\mathcal{O}((1 - \kappa^{-1})^k)$. This rate is recovered in Corollary 3.18. On the other hand, [51] showed that under their choice of parameter for HRHF, the convergence rate is $\mathcal{O}((1 - 2\kappa^{-1})^k)$, which is a slight improvement compared with underdamped Langevin dynamics. In this work, we performed a detailed eigenvalue analysis of GAUL on Gaussian target measure. We showed that under our choice of parameters (γ, a, h) , the convergence rate toward the biased target measure is $\mathcal{O}((1 - c\sqrt{\kappa})^k)$ for some constant c .

3. Analysis of GAUL on quadratic potential functions. In this section, we establish the convergence rate for the proposed SDE (2.17) toward the target distribution following a Gaussian distribution.

3.1. Problem setup. In this subsection, we present the main problem addressed in this paper. We are interested in sampling from a distribution whose probability density function is proportional to $\exp(-f(\mathbf{x}))$ for $f: \mathbb{R}^d \rightarrow \mathbb{R}$. In this paper, we focus on a concrete example in which the potential function f is quadratic and thus the target distribution is a Gaussian distribution. Let

$$(3.1) \quad f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \Sigma_*^{-1} \mathbf{x},$$

where $\mathbf{x} \in \mathbb{R}^d$ and $\Sigma_* \succ 0$ is a symmetric positive definite matrix in $\mathbb{R}^{d \times d}$. Define

$$(3.2) \quad \tilde{\Sigma} = \begin{pmatrix} \Sigma_* & 0 \\ 0 & \mathbf{I} \end{pmatrix}.$$

As in the previous section, denote by $\mathbf{X} = (\mathbf{x}, \mathbf{p}) \in \mathbb{R}^{2d}$. And $H(\mathbf{X}) = f(\mathbf{x}) + \|\mathbf{p}\|^2/2$. Then, we can write

$$(3.3) \quad H(\mathbf{X}) = \frac{1}{2} \mathbf{X}^T \begin{pmatrix} \Sigma_*^{-1} & 0 \\ 0 & \mathbf{I} \end{pmatrix} \mathbf{X} := \frac{1}{2} \mathbf{X}^T \tilde{\Sigma}^{-1} \mathbf{X}.$$

Define the target density $\pi: \mathbb{R}^{2d} \rightarrow \mathbb{R}$ to be

$$(3.4) \quad \Pi(\mathbf{X}) = \frac{1}{Z} \exp(-H(\mathbf{X})),$$

where $H(\mathbf{X})$ is given by (3.3) and $Z = \int_{\mathbb{R}^{2d}} \exp(-H(\mathbf{X})) d\mathbf{X}$ is a normalization constant such that $\Pi(\mathbf{X})$ integrates to one on \mathbb{R}^{2d} . We also define the \mathbf{x} -marginal target density to be

$$(3.5) \quad \pi(\mathbf{x}) = \frac{1}{Z_1} \exp(-f(\mathbf{x})),$$

where $f(\mathbf{x})$ is given by (3.1) and $Z_1 = \int_{\mathbb{R}^d} \exp(-f(\mathbf{x})) d\mathbf{x}$ is a normalization constant.

Remark 3.1. Note that for any symmetric positive definite Σ_* , we have that $\Sigma_*^{-1} = \mathbf{P}\Lambda\mathbf{P}^T$ for some orthogonal matrix \mathbf{P} and diagonal matrix $\Lambda = \text{diag}(s_1, \dots, s_d)$ with $s_1 \geq \dots \geq s_d > 0$. By a change of variable $\mathbf{y} = \mathbf{P}^T \mathbf{x}$, one can rewrite $f(\mathbf{x})$ in terms of \mathbf{y} , such that

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \Sigma_*^{-1} \mathbf{x} = \frac{1}{2} \mathbf{x}^T \mathbf{P}\Lambda\mathbf{P}^T \mathbf{x} = \frac{1}{2} \mathbf{y}^T \Lambda \mathbf{y}.$$

For simplicity of notation, we assume that $\mathbf{P} = \mathbf{I}$ and $\Sigma_*^{-1} = \Lambda$ is a diagonal matrix. We denote by $\kappa = s_1/s_d$ the condition number of f . We will also assume that $s_1 > 1 > s_d$ throughout this paper. Furthermore, to simplify our analysis, we consider $\mathbf{C} = \text{diag}(c_1, \dots, c_d)$.

3.2. Continuous time analysis. In this subsection, we study the convergence of GAUL. In particular, we analyze the convergence of the Fokker–Planck equation (2.17) to the target density (3.4), (3.5) by directly studying an ODE system of the covariance of the distribution.

Proposition 3.2. *Let \mathbf{X}_t be the solution of (2.15) where $H(\mathbf{X})$ is given by (3.3), and $\mathbf{X}_0 \sim \mathcal{N}(0, \mathbf{I}_{2d \times 2d})$. Then $\mathbf{X}_t \sim \mathcal{N}(0, \Sigma(t))$ where the covariance $\Sigma(t)$ satisfies the following matrix ODE:*

$$(3.6) \quad \dot{\Sigma}(t) = 2 \text{sym}(\mathbf{Q}(\mathbf{I} - \tilde{\Sigma}^{-1} \Sigma(t))).$$

Moreover, (3.6) is well-defined and has a solution for all $t \geq 0$, such that $\Sigma(t)$ is symmetric semipositive definite.

The proof is postponed in section SM3. We denote by $\Sigma_{ij}(t) \in \mathbb{R}^{d \times d}$ the block components of $\Sigma(t) \in \mathbb{R}^{2d \times 2d}$:

$$\Sigma(t) = \begin{pmatrix} \Sigma_{11}(t) & \Sigma_{12}(t) \\ \Sigma_{12}^T(t) & \Sigma_{22}(t) \end{pmatrix}.$$

Then we can write (3.6) in terms of the block components.

Corollary 3.3. *The componentwise covariance matrix $\Sigma_{ij}(t)$ satisfies the following ODE system:*

$$(3.7a) \quad \dot{\Sigma}_{11} = -2a(\text{sym}(\mathbf{C}\Sigma_*^{-1}\Sigma_{11}) - \mathbf{C}) + 2\text{sym}(\mathbf{C}\Sigma_{12}),$$

$$(3.7b) \quad \dot{\Sigma}_{22} = -2\text{sym}(\Sigma_*^{-1}\Sigma_{12}) - 2\gamma(\Sigma_{22} - \mathbf{I}),$$

$$(3.7c) \quad \dot{\Sigma}_{12} = -a\mathbf{C}\Sigma_*^{-1}\Sigma_{12} - (\mathbf{C} - \mathbf{C}\Sigma_{22}) + (\mathbf{I} - \Sigma_{11}\Sigma_*^{-1}) - \gamma\Sigma_{12}.$$

Moreover, with initial conditions $\Sigma_{11}(0) = \Sigma_{22}(0) = \mathbf{I}$ and $\Sigma_{12}(0) = 0$, the stationary states of $\Sigma_{11}(t)$, $\Sigma_{22}(t)$ and $\Sigma_{12}(t)$ are given by Σ_* , \mathbf{I} and 0 respectively.

From now on, we consider $\mathbf{C} = \mathbf{I}$ in our analysis. We address our results for $\mathbf{C} \neq \mathbf{I}$ in Remark 3.10 and Remark 3.20. Note that when $\mathbf{C} = \mathbf{I}$, we have $\mathbf{Q} = \text{sym}(\mathbf{Q})$ is always positive semidefinite for $a, \gamma \geq 0$. Our next theorem makes sure that the stationary state of (3.6) is actually unique and characterizes the convergence speed of the covariance matrix toward its stationary state.

Theorem 3.4. *Let \mathbf{X}_t be the solution of (2.15) where $H(\mathbf{X})$ is given by (3.3), and $\mathbf{X}_0 \sim \mathcal{N}(0, \mathbf{I}_{2d \times 2d})$. Then $\Sigma(t)$ converges to the unique stationary state $\tilde{\Sigma}$ given in (3.2). The optimal choice of γ is given by $\gamma^* = as_d + 2\sqrt{s_d}$ under which we have $\|\Sigma_{11}(t) - \Sigma_*\|_F = \mathcal{O}(te^{-(2as_d + 2\sqrt{s_d})t})$ and $\|\Sigma_{22}(t) - \mathbf{I}\|_F = \mathcal{O}(te^{-(2as_d + 2\sqrt{s_d})t})$ for $t \geq 1$.*

Proof. As mentioned in Remark 3.1, we consider $\Sigma_*^{-1} = \Lambda$. By our assumption on \mathbf{X}_0 , (3.7) implies that $\Sigma_{11}(t)$, $\Sigma_{22}(t)$, and $\Sigma_{12}(t)$ will be diagonal matrices for all $t > 0$. This simplifies the ODE system (3.7). After some manipulation, we obtain

$$(3.8) \quad \begin{pmatrix} \dot{\Sigma}_{11} \\ \dot{\Sigma}_{22} \\ \dot{\Sigma}_{22} \end{pmatrix} = \underbrace{\begin{pmatrix} -2a\mathbf{C}\Sigma_*^{-1} & -2\gamma\mathbf{C}\Sigma_* & -\mathbf{C}\Sigma_* \\ 0 & 0 & \mathbf{I} \\ 2\Sigma_*^{-2} & 2(-1 - a\gamma)\mathbf{C}\Sigma_*^{-1} - 2\gamma^2\mathbf{I} & -3\gamma\mathbf{I} - a\mathbf{C}\Sigma_*^{-1} \end{pmatrix}}_{\mathcal{D}} \begin{pmatrix} \Sigma_{11} \\ \Sigma_{22} \\ \dot{\Sigma}_{22} \end{pmatrix} + \mathbf{T},$$

where

$$\mathbf{T} = \begin{pmatrix} 2a\mathbf{C} + 2\gamma\mathbf{C}\Sigma_* \\ 0 \\ 2a\gamma\Sigma_*^{-1}\mathbf{C} + 2\gamma^2\mathbf{I} + 2\Sigma_*^{-1}\mathbf{C} - 2\Sigma_*^{-1} \end{pmatrix},$$

and $\mathbf{C} = \mathbf{I}$. We have already seen in Corollary 3.3 that the stationary state of $\Sigma(t)$ is $\tilde{\Sigma}$ given in (3.2). To show uniqueness, we compute the eigenvalues of \mathcal{D} :

$$\begin{aligned} \lambda_0^{(i)} &= -as_i - \gamma, \\ \lambda_1^{(i)} &= -as_i - \gamma - \sqrt{\gamma^2 - 2a\gamma s_i + s_i(-4 + a^2 s_i)}, \\ \lambda_2^{(i)} &= -as_i - \gamma + \sqrt{\gamma^2 - 2a\gamma s_i + s_i(-4 + a^2 s_i)}, \end{aligned}$$

where s_i 's are the diagonal elements of Λ for $i = 1, \dots, d$. It is clear that 0 is not an eigenvalue of \mathcal{D} . Therefore, $\tilde{\Sigma}$ is the unique stationary state for $\Sigma(t)$. The convergence speed of (3.8) is essentially controlled by the largest real part of the eigenvalues of \mathcal{D} . Note that for all i ,

$$\Re(\lambda_2^{(i)}) \geq \Re(\lambda_0^{(i)}) \geq \Re(\lambda_1^{(i)}),$$

where $\Re(z)$ denotes the real part of $z \in \mathbb{C}$. Therefore, to characterize the convergence speed of (3.8), it suffices to control $\max_i \Re(\lambda_2^{(i)})$. By Lemma SM1.7, we know that for any given $a \geq 0$, the optimal choice of γ is

$$\gamma^* = \arg \min_{\gamma > 0} \max_i \Re(\lambda_2^{(i)}) = as_d + 2\sqrt{s_d}.$$

With $\gamma = \gamma^*$, we get that

$$\max_{i,j} \Re(\lambda_j^{(i)}) \leq \max_i \Re(\lambda_2^{(i)}) \leq -2as_d - 2\sqrt{s_d}.$$

This leads to

$$(3.9) \quad \left\| \begin{pmatrix} \Sigma_{11}(t) - \Sigma_* \\ \Sigma_{22}(t) - \mathbf{I} \\ \dot{\Sigma}_{22}(t) \end{pmatrix} \right\|_{\mathbb{F}} \leq C_1 t e^{-(2as_d + 2\sqrt{s_d})t},$$

which is valid for $t \geq 1$. The constant C_1 depends on d, s_1, s_d^{-1} at most polynomially according to Lemma SM1.8. Note that the extra t dependence comes from the repeated eigenvalue $\lambda_0^{(d)} = \lambda_1^{(d)} = \lambda_2^{(d)}$ when $\gamma = \gamma^*$. By a triangle inequality, we get

$$\|\Sigma_{11} - \Sigma_*\|_{\mathbb{F}} \leq \left\| \begin{pmatrix} \Sigma_{11}(t) - \Sigma_* \\ \Sigma_{22}(t) - \mathbf{I} \\ \dot{\Sigma}_{22}(t) \end{pmatrix} \right\|_{\mathbb{F}} \leq C_1 t e^{-(2as_d + 2\sqrt{s_d})t}.$$

And similarly,

$$\|\Sigma_{22} - \mathbf{I}\|_{\mathbb{F}} \leq C_1 t e^{-(2as_d + 2\sqrt{s_d})t}. \quad \blacksquare$$

Remark 3.5. The choice $a = 0$ corresponds to underdamped Langevin dynamics (UL). Taking $a > 0$ gives an extra factor of $e^{-2as_d t}$ in terms of convergence.

Definition 3.6 (Mixing time). *The total variation between two probability measures \mathcal{P} and \mathcal{Q} over a measurable space $(\mathbb{R}^d, \mathcal{F})$ is*

$$\text{TV}(\mathcal{P}, \mathcal{Q}) = \sup_{A \in \mathcal{F}} |\mathcal{P}(A) - \mathcal{Q}(A)|.$$

Let \mathcal{T}_p be an operator on the space of probability distributions. Assume that $\mathcal{T}_p^k(\nu_0) \rightarrow \nu$ as $k \rightarrow \infty$ for some initial distribution ν_0 and stationary distribution ν . The discrete δ -mixing time ($\delta \in (0, 1)$) is given by

$$t_{\text{mix}}^{\text{dis}}(\delta; \nu_0, \nu) = \min\{k \mid \text{TV}(\mathcal{T}_p^k(\nu_0), \nu) \leq \delta\}.$$

Similarly, if $\mathcal{T}_p(t; \cdot)$ is an operator for each $t \geq 0$ with $\mathcal{T}_p(0; \cdot) = \text{id}(\cdot)$, assume that $\mathcal{T}_p(t; \nu_0) \rightarrow \nu$ as $t \rightarrow \infty$. The continuous δ -mixing time ($\delta \in (0, 1)$) is given by

$$t_{\text{mix}}^{\text{cont}}(\delta; \nu_0, \nu) = \min\{t \mid \text{TV}(\mathcal{T}_p(t; \nu_0), \nu) \leq \delta\}.$$

Theorem 3.7 ([30]). Let $\mu \in \mathbb{R}^d$, Σ_1, Σ_2 be two positive definite covariance matrices, and $\lambda_1, \dots, \lambda_d$ denote the eigenvalues of $\Sigma_1^{-1}\Sigma_2 - \mathbf{I}$. Then the total variation satisfies

$$\mathrm{TV}(\mathcal{N}(\mu, \Sigma_1), \mathcal{N}(\mu, \Sigma_2)) \leq \frac{3}{2} \min \left\{ 1, \sqrt{\sum_{i=1}^d \lambda_i^2} \right\}.$$

A straightforward corollary follows from Schur decomposition theorem.

Corollary 3.8. We have

$$\mathrm{TV}(\mathcal{N}(\mu, \Sigma_1), \mathcal{N}(\mu, \Sigma_2)) \leq \frac{3}{2} \min \{1, \|\Sigma_1^{-1}\Sigma_2 - \mathbf{I}\|_{\mathrm{F}}\}.$$

Using Theorem 3.4 and Corollary 3.8, we obtain the following mixing time theorem when the potential function f is quadratic.

Theorem 3.9 (Continuous mixing time). Consider the same setting as in Theorem 3.4. Consider $0 < \delta \ll 1$. Then

$$t_{\mathrm{mix}}^{\mathrm{cont}}(\delta; \nu_0, \pi) \leq \frac{\mathcal{O}(\log(d) + \log(\kappa)) + \log(1/\delta)}{as_d + 2\sqrt{s_d}}.$$

Here ν_0 is the distribution of \mathbf{x} , which is $\mathcal{N}(0, \mathbf{I}_{d \times d})$. π is the target density in the \mathbf{x} variable given in (3.5).

Proof. We shall use Corollary 3.8 with

$$\Sigma_1 = \Sigma_*, \quad \Sigma_2 = \Sigma_{11}(t).$$

We have

$$\begin{aligned} \|\Sigma_1^{-1}\Sigma_2 - \mathbf{I}\|_{\mathrm{F}} &= \|\Sigma_*^{-1}(\Sigma_{11}(t) - \Sigma_*)\|_{\mathrm{F}} \\ &\leq C_1 t e^{-(2as_d + 2\sqrt{s_d})t} s_1. \end{aligned}$$

By a direct computation, we get

$$t_{\mathrm{mix}}^{\mathrm{cont}}(\delta; \nu_0, \pi) \leq \frac{\log(\tilde{C}_1/\delta)}{as_d + 2\sqrt{s_d}},$$

where $\tilde{C}_1 = \frac{3}{2}C_1 s_1$. By Lemma SM1.8, we have that

$$t_{\mathrm{mix}}^{\mathrm{cont}}(\delta; \nu_0, \pi) \leq \frac{\mathcal{O}(\log(d\kappa)) + \log(1/\delta)}{as_d + 2\sqrt{s_d}}. \quad \blacksquare$$

Remark 3.10. When $\mathbf{C} = \mathrm{diag}(c_1, \dots, c_d)$ and $\mathrm{sym}(\mathbf{Q}) \succeq 0$ in (2.16), our proof can be easily adapted to show similar results in Theorem 3.9:

$$t_{\mathrm{mix}}^{\mathrm{cont}}(\delta; \nu_0, \pi) \leq \frac{\mathcal{O}(\log(d) + \log(\hat{\kappa})) + \log(1/\delta)}{a\hat{s}_d + 2\sqrt{\hat{s}_d}},$$

where \hat{s}_i is the i -th largest eigenvalue of matrix $\mathbf{C}\Sigma_*^{-1}$. And $\hat{\kappa} = \hat{s}_1/\hat{s}_d$. In other words, the matrix \mathbf{C} can be viewed as a preconditioner for the target covariance matrix in the sampling problem.

3.3. Discrete time analysis. To implement (2.15), we need to consider its time discretization. As discretization is not the focus of this paper, we will only analyze the simplest discretization using the Euler–Maruyama method in section A.

Let us first make a few observations regarding the discretization in section A. After a straightforward computation, we obtain the following update rule.

Proposition 3.11. *The Euler–Maruyama discretization of (2.15) given in section A with step size h can be written in the following form:*

$$(3.10) \quad \begin{pmatrix} \mathbf{x}_{n+1} \\ \mathbf{p}_{n+1} \end{pmatrix} = \mathbf{A} \begin{pmatrix} \mathbf{x}_n \\ \mathbf{p}_n \end{pmatrix} + \mathbf{L}z,$$

where

$$(3.11) \quad \mathbf{A} = \mathbf{I}_{2d \times 2d} - h \underbrace{\begin{pmatrix} a\Lambda & -\mathbf{I}_{d \times d} \\ \Lambda & \gamma \mathbf{I}_{d \times d} \end{pmatrix}}_{\mathbf{G}}, \quad \mathbf{L} = \begin{pmatrix} \sqrt{2ah}\mathbf{I} & 0 \\ 0 & \sqrt{2\gamma h}\mathbf{I} \end{pmatrix}.$$

And \mathbf{z} is a $2d$ -dimensional Brownian motion, i.e., $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}_{2d \times 2d})$.

Using (3.10), we can derive the evolution of the mean and covariance at each time step. As before, let us denote by $\mathbf{X}_n = (\mathbf{x}_n, \mathbf{p}_n)$.

Corollary 3.12. *Suppose that $\mathbb{E}(\mathbf{x}_0) = \mathbb{E}(\mathbf{p}_0) = 0$. Then*

$$(3.12) \quad \text{cov}(\mathbf{X}_{n+1}, \mathbf{X}_{n+1}) = \mathbf{A} \text{cov}(\mathbf{X}_n, \mathbf{X}_n) \mathbf{A}^T + \mathbf{L} \mathbf{L}^T.$$

Proof. From (3.10), it is clear that $\mathbb{E}(\mathbf{x}_n) = \mathbb{E}(\mathbf{p}_n) = 0$ for all $n \geq 0$. We calculate

$$\begin{aligned} \text{cov}(\mathbf{X}_{n+1}, \mathbf{X}_{n+1}) &= \mathbb{E}(\mathbf{A} \mathbf{X}_n \mathbf{X}_n^T \mathbf{A}^T + \mathbf{A} \mathbf{X}_n \mathbf{z}^T \mathbf{L}^T + \mathbf{L} \mathbf{z} \mathbf{X}_n^T \mathbf{A}^T + \mathbf{L} \mathbf{z} \mathbf{z}^T \mathbf{L}^T) \\ &= \mathbf{A} \text{cov}(\mathbf{X}_n, \mathbf{X}_n) \mathbf{A}^T + \mathbf{L} \mathbf{L}^T. \end{aligned}$$

Corollary 3.13. *Denote by \mathbf{Y}^* a solution to the fixed point equation $\mathbf{Y} = \mathbf{A} \mathbf{Y} \mathbf{A}^T + \mathbf{L} \mathbf{L}^T$. And let $\mathbf{Y}_n = \text{cov}(\mathbf{X}_n, \mathbf{X}_n) - \mathbf{Y}^*$. Then*

$$\mathbf{Y}_{n+1} = \mathbf{A} \mathbf{Y}_n \mathbf{A}^T.$$

Theorem 3.14. *Suppose $a \geq \frac{2}{\sqrt{s_1} - \sqrt{s_d}}$ and the step size h satisfies $0 < h < 1/(as_1 + \gamma)$ and $\gamma = \gamma^* = as_d + 2\sqrt{s_d}$. Then there exists a unique \mathbf{Y}^* satisfying*

$$\mathbf{Y}^* = \mathbf{A} \mathbf{Y}^* \mathbf{A}^T + \mathbf{L} \mathbf{L}^T.$$

Moreover, the iteration $\mathbf{Y}_{k+1} = \mathbf{A} \mathbf{Y}_k \mathbf{A}^T + \mathbf{L} \mathbf{L}^T$ converges to \mathbf{Y}^* linearly: $\|\mathbf{Y}_k - \mathbf{Y}^*\|_{\text{F}} \leq \tilde{C} h^2 k^2 (1 - \frac{h}{2}(as_d + \sqrt{s_d}))^{2k-2}$, where the constant $\tilde{C} = d^2 \cdot \mathcal{O}(\text{poly}(\kappa))$.

Proof. Existence: we directly compute this stationary point in Lemma SM1.17. Uniqueness: by Lemma SM1.14 and Corollary SM1.10, we see that \mathbf{Y}^* is unique. The convergence rate is proved in Lemma SM1.14 and Theorem SM1.16. ■

Corollary 3.15. *Let \mathbf{Y}_{11}^* be the upper left $d \times d$ block of \mathbf{Y}^* . Then \mathbf{Y}_{11}^* is a biased estimate of the true covariance Σ_* . In particular, denote by $\tilde{\pi} = \mathcal{N}(0, \mathbf{Y}_{11}^*)$. We have*

$$\text{TV}(\pi, \tilde{\pi}) = h\mathcal{O}(\sqrt{d\kappa}) + \mathcal{O}(h^2).$$

Proof. By Corollary SM1.18 and Corollary 3.8, we finish the proof. ■

Theorem 3.16 (Discrete mixing time). *Suppose $\sqrt{s_1} - \sqrt{s_d} \geq 2$. We take $a = 1$, $\gamma = \gamma^* = s_d + 2\sqrt{s_d}$, $h = 1/5s_1$. If we use the Euler–Maruyama scheme for (2.15), then for $0 < \delta \ll 1$,*

$$(3.13) \quad t_{\text{mix}}^{\text{dis}}(\delta; \nu_0, \tilde{\pi}) = \mathcal{O}\left(\frac{\log(\kappa) + \log(1/\delta) + \log(d)}{\frac{1}{\kappa} + \frac{1}{\sqrt{\kappa s_1}}}\right).$$

Here ν_0 is the distribution of \mathbf{x} , which is $\mathcal{N}(0, \mathbf{I}_{d \times d})$. $\tilde{\pi}$ is the target density in the \mathbf{x} variable which is a zero mean Gaussian distribution with a variance given by (SM1.24).

Proof. Note that from our previous notation, we have that

$$\text{cov}(\mathbf{x}_k, \mathbf{x}_k) = \begin{pmatrix} \mathbf{I}_{d \times d} & 0 \end{pmatrix} \text{cov}(\mathbf{X}_k, \mathbf{X}_k) \begin{pmatrix} \mathbf{I}_{d \times d} \\ 0 \end{pmatrix} =: \tilde{\mathbf{Y}}_k.$$

Moreover, let us define

$$\tilde{\mathbf{Y}}^* = \begin{pmatrix} \mathbf{I}_{d \times d} & 0 \end{pmatrix} \mathbf{Y}^* \begin{pmatrix} \mathbf{I}_{d \times d} \\ 0 \end{pmatrix}$$

to be the limiting covariance in the \mathbf{x} variable for the discretization (\mathbf{Y}^* is defined in Theorem 3.14). Clearly, we have that

$$(3.14) \quad \|\tilde{\mathbf{Y}}_k - \tilde{\mathbf{Y}}^*\|_{\text{F}} \leq \|\mathbf{Y}_k - \mathbf{Y}^*\|_{\text{F}} \leq \tilde{C}h^2k^2 \left(1 - \frac{h}{2}(as_d + \sqrt{s_d})\right)^{2k-2}.$$

Using Corollary 3.8, we compute

$$\begin{aligned} \|(\tilde{\mathbf{Y}}^*)^{-1}\tilde{\mathbf{Y}}_k - \mathbf{I}\|_{\text{F}} &= \|(\tilde{\mathbf{Y}}^*)^{-1}(\tilde{\mathbf{Y}}_k - \tilde{\mathbf{Y}}^*)\|_{\text{F}} \\ &\leq \|(\tilde{\mathbf{Y}}^*)^{-1}\|_{\text{F}} \|\tilde{\mathbf{Y}}_k - \tilde{\mathbf{Y}}^*\|_{\text{F}}. \end{aligned}$$

By Lemma SM1.17, $\tilde{\mathbf{Y}}^*$ is a diagonal matrix. Therefore $(\tilde{\mathbf{Y}}^*)^{-1}$ is also a diagonal matrix. Moreover, from (SM1.24), we see that $\|(\tilde{\mathbf{Y}}^*)^{-1}\|_{\text{F}} \leq \sqrt{d}\mathcal{O}(\text{poly}(\kappa))$. Therefore, we obtain

$$\begin{aligned} \|(\tilde{\mathbf{Y}}^*)^{-1}\tilde{\mathbf{Y}}_k - \mathbf{I}\|_{\text{F}} &\leq d^{5/2} \cdot \mathcal{O}(\text{poly}(\kappa))h^2k^2 \left(1 - \frac{h}{2}(s_d + \sqrt{s_d})\right)^{2k-2} \\ &\leq d^{5/2} \cdot \mathcal{O}(\text{poly}(\kappa))h^2k^2 e^{-(k-1)h(s_d + \sqrt{s_d})}, \end{aligned}$$

where we used $1 - x \leq e^{-x}$ for $x \in \mathbb{R}$ to get the second inequality. Letting $h = 1/5s_1$ and taking logarithm on both hand sides, we conclude that

$$t_{\text{mix}}^{\text{dis}}(\delta; \nu_0, \tilde{\pi}) \leq \frac{\mathcal{O}(\log(d)) + \mathcal{O}(\log(\kappa)) + \log(1/\delta)}{\frac{1}{10}\left(\frac{1}{\kappa} + \frac{1}{\sqrt{\kappa s_1}}\right)}. \quad \blacksquare$$

The above theorem gives a convergence rate when $a = 1$. At this stage, a natural question arises: given the optimal γ that we proved in Theorem 3.4, is there a better choice of a ? In Theorem 3.14, we require that $a \geq \frac{2}{\sqrt{s_1 - \sqrt{s_d}}}$ and $0 < h < \frac{1}{as_1 + \gamma}$ to guarantee the existence and uniqueness of the stationary point of the covariance equation (3.12). Since the proof of Theorem 3.16 relies on (3.14), let us plug in $h = \frac{1}{\beta(as_1 + \gamma)}$ into (3.14) for some $\beta > 1$. The dominating factor in (3.14) is the exponential decay term $(1 - \frac{h}{2}(as_d + \sqrt{s_d}))^{2k-2}$. We can optimize this term over a . We have

$$\frac{d}{da} \left(1 - \frac{as_d + \sqrt{s_d}}{2\beta(as_1 + as_d + 2\sqrt{s_d})} \right) > 0.$$

This suggests that we take $a = \frac{2}{\sqrt{s_1 - \sqrt{s_d}}}$. Indeed, such choice of a gives a much better convergence rate in terms of the condition number κ as shown in the next theorem.

Theorem 3.17 (A better choice of a). *The denominator of the mixing time given in Theorem 3.16 can be improved to $\kappa^{-1/2}$ by choosing $a = \frac{2}{\sqrt{s_1 - \sqrt{s_d}}}$, $\gamma = as_d + 2\sqrt{s_d}$, and $h = \frac{1}{2(as_1 + \gamma)}$. To be more precise, we have*

$$(3.15) \quad t_{\text{mix}}^{\text{dis}}(\delta; \nu_0, \tilde{\pi}) = \mathcal{O} \left(\frac{\log(\kappa) + \log(1/\delta) + \log(d)}{\frac{1}{\sqrt{\kappa}}} \right).$$

Proof. The proof will be very similar to that of Theorem 3.16. We start with (3.14). And we can explicitly calculate

$$\begin{aligned} 1 - \frac{h}{2}(as_d + \sqrt{s_d}) &= 1 - \frac{as_d + \sqrt{s_d}}{4(as_1 + as_d + 2\sqrt{s_d})} \\ &= 1 - \frac{2s_d + \sqrt{s_d}(\sqrt{s_1} - \sqrt{s_d})}{8(s_1 + s_d + \sqrt{s_d}(\sqrt{s_1} - \sqrt{s_d}))} \\ &= 1 - \frac{\sqrt{s_1 s_d} + s_d}{8(s_1 + \sqrt{s_1 s_d})} \\ &\leq 1 - \frac{1}{16\sqrt{\kappa}}. \end{aligned}$$

The rest of the proof is the same as the proof of Theorem 3.16, and we will suppress it for brevity. ■

The following corollary follows from Lemma SM1.15 and the proof of Theorem 3.16.

Corollary 3.18 (Underdamped Langevin mixing time). *Suppose $a = 0$, $\gamma = 2\sqrt{s_d}$, $h = \sqrt{s_d}/s_1$. If we use the Euler–Maruyama scheme for (2.15), then for $0 < \delta \ll 1$,*

$$(3.16) \quad t_{\text{mix}}^{\text{dis}}(\delta; \nu_0, \tilde{\pi}) = \mathcal{O} \left(\frac{\log(\kappa) + \log(1/\delta) + \log(d)}{\frac{1}{\kappa}} \right),$$

ν_0 is the distribution of \mathbf{x} , which is $\mathcal{N}(0, \mathbf{I}_{d \times d})$. $\tilde{\pi}$ is the target density in the \mathbf{x} variable which is a zero mean Gaussian with variance given by (SM1.24) with $a = 0$.

Remark 3.19. $a = 0$ in (2.15) corresponds to the underdamped Langevin dynamics. In this case, we show in Lemma SM1.15 that to guarantee convergence (to a biased target) the step size restriction on h is more strict than when $a = 1$. In particular, when $a = 0$ it follows from Lemma SM1.15 that the choice $h = 1/5s_1$ does not guarantee convergence if $s_d < 10^{-2}$. Comparing (3.15) and (3.16), we see that the mixing time for GAUL beats that of underdamped Langevin dynamics under the Euler–Maruyama discretization. We are aware that this does not imply the same result will hold when comparing the mixing time toward the true target distribution $\pi(\mathbf{x})$ given in (3.5), due to the presence of bias in the Euler–Maruyama scheme.

Remark 3.20. When $\mathbf{C} = \text{diag}(c_1, \dots, c_d)$ and $\text{sym}(\mathbf{Q}) \succeq 0$ in (2.16), we also have a similar mixing time described in Theorem 3.17, which is $\mathcal{O}(\sqrt{\hat{\kappa}}(\log(\hat{\kappa}) + \log(1/\delta) + \log(d)))$ when $a = \frac{2}{\sqrt{\hat{s}_1} - \sqrt{\hat{s}_d}}$, $\gamma = a\hat{s}_d + 2\sqrt{\hat{s}_d}$, and $h = \frac{1}{2(a\hat{s}_1 + \gamma)}$. The notation \hat{s}_i and $\hat{\kappa}$ are defined in Remark 3.10.

We introduce another discretization scheme inspired by the \mathcal{BAO} splitting [46, 45, 13, 48] for underdamped Langevin dynamics, which is able to achieve an $\mathcal{O}(\sqrt{\kappa})$ convergence [48]. We name our method $\mathcal{BAGOGAB}$ splitting, details of which are postponed to section SM2. We have the following corollary regarding the first order discretization error of $\mathcal{BAGOGAB}$ splitting when sampling the same Gaussian distribution $\mathcal{N}(0, \Sigma_*)$.

Corollary 3.21. With the same choice of parameter $a = \frac{2}{\sqrt{s_1} - \sqrt{s_d}}$ and $\gamma = \gamma^* = as_d + 2\sqrt{s_d}$ as in the Euler–Maruyama discretization, the covariance in the (\mathbf{x}, \mathbf{p}) variables satisfies a similar equation as in (3.12) (see Corollary SM2.2 for details). Abusing notation, denote by \mathbf{Y}^* the stationary point to the covariance equation in Corollary SM2.2 and denote by \mathbf{Y}_{11}^* the upper $d \times d$ block of \mathbf{Y}^* (i.e., the covariance in the \mathbf{x} variable). Then

$$\text{TV}(\pi, \mathcal{N}(0, \mathbf{Y}_{11}^*)) = h\mathcal{O}\left(\frac{1}{\sqrt{s_1}}\text{tr}(\Sigma_*^{-1})\right) + \mathcal{O}(h^2).$$

Proof. By Corollary 3.8 and Corollary SM2.3, we finish the proof. ■

Comparing with Corollary 3.15, we see that the $\mathcal{BAGOGAB}$ splitting achieves a smaller bias than the Euler–Maruyama scheme when sampling a Gaussian distribution. This is also demonstrated in our numerical experiments. To analyze the convergence speed of $\mathcal{BAGOGAB}$, one would use the same eigenvalue analysis as we did for the Euler–Maruyama scheme. The exact eigenvalues for $\mathcal{BAGOGAB}$ splitting is much more complicated than the Euler–Maruyama scheme. However, it follows that the eigenvalues of the Euler–Maruyama discretization are exactly the first order Taylor expansion of the eigenvalues of $\mathcal{BAGOGAB}$ splitting (see Corollary SM2.5). Therefore, we use the same choice of γ , a , and h for $\mathcal{BAGOGAB}$ in our numerical experiments to verify its acceleration.

Remark 3.22. When the target potential f is not a quadratic function, it is more technical in proving the convergence speed. A common technique to prove convergence in the Wasserstein-2 distance is by a coupling argument (see [22, 28]). [15] proved L_2 convergence under a Poincarè-type inequality using Bochner’s formula. In the L_1 distance and KL divergence, [34] designs the convergence analysis toward these problems. We leave the convergence analysis of general f with optimal choices of preconditioned matrices \mathbf{Q} for future work.

4. Numerical experiment. In this section, we implement several numerical examples to compare the proposed SDE with the overdamped (labeled ‘OL’) and underdamped (labeled ‘UL’) Langevin dynamics. Recall that ‘OL’ corresponds to the choice $a = 1, \gamma = 0$, and ‘UL’ corresponds to $a = 0$ in (2.15). We also include the *BACAB* [45, 46] splitting method for UL, which is labeled as ‘UL_splitting’. We set $\mathbf{C} = \mathbf{I}$. We compare two discretizations of GAUL: using the Euler–Maruyama scheme (GAUL_EM) and using the splitting method (GAUL_Splitting) detailed in section SM2. In the meantime, we also compare with randomized Hamiltonian Monte Carlo (RHMC) [61, 9, 11]. We use the `dynamic_hmc` module from Blackjax to implement RHMC. This module allows users to define an integration step function that generates the next pseudo or quasi-random number of integration steps in the sequence. In our experiment, this integration step function is set to return a geometric random variable with success probability 0.03, clipped at [5, 100].

4.1. Gaussian examples.

4.1.1. One dimension. We begin with a simple example, a one dimensional Gaussian distribution with zero mean. In Fig. 1, we consider two cases where the variances are given by 0.01 and 100, respectively. We first sample $M = 10^5$ particles from $\mathcal{N}(0, \mathbf{I}_{2 \times 2})$ (although our experiment is in one dimension, we need both \mathbf{x} and \mathbf{p} variables). When measuring the convergence speed, we use KL divergence in Gaussian distributions to measure the change of covariances. Note that we will only measure the KL divergence in the x variable since we are primarily interested in sampling distribution of the form $\frac{1}{Z}e^{-f(x)}$. In this experiment, we can make use of the fact that the sample distribution and the target distribution are both Gaussians. And the KL divergence between two centered Gaussians has a closed form expression:

$$(4.1) \quad D_{\text{KL}}(\Sigma(t), \tilde{\Sigma}) = \frac{1}{2}(\text{tr}(\Sigma(t)\tilde{\Sigma}^{-1}) - \log \det(\Sigma(t)\tilde{\Sigma}^{-1}) - d).$$

In this one dimensional example, we study two cases where $\tilde{\Sigma} = 0.01$ or 100. $\Sigma(t)$ can be approximated by the unbiased sample variance. For $\tilde{\Sigma} = 0.01$, we choose time step size $h = 10^{-4}$, total number of steps $N = 400$, $\gamma_{ul} = 2\tilde{\Sigma}^{-1/2} = 20$, and $\gamma_{pdd} = 2\tilde{\Sigma}^{-1/2} + \tilde{\Sigma}^{-1} = 120$. For $\tilde{\Sigma} = 100$, we choose the time step size $h = 10^{-2}$, total number of steps $N = 600$,

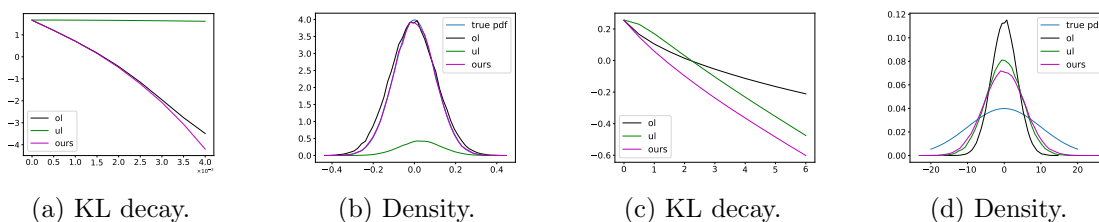


Figure 1. Convergence and density comparisons of three methods. (a) and (c): KL divergence between the sample and the target distribution, which is a one-dimensional Gaussian with zero mean and variance 0.01 (a), 100 (c). ‘ol’ represents overdamped Langevin dynamics; ‘ul’ represents underdamped Langevin dynamics. x -axis represents time and y -axis is in \log_{10} scale. (b) and (d): density comparison at the end of the experiment between the three methods and the true density.

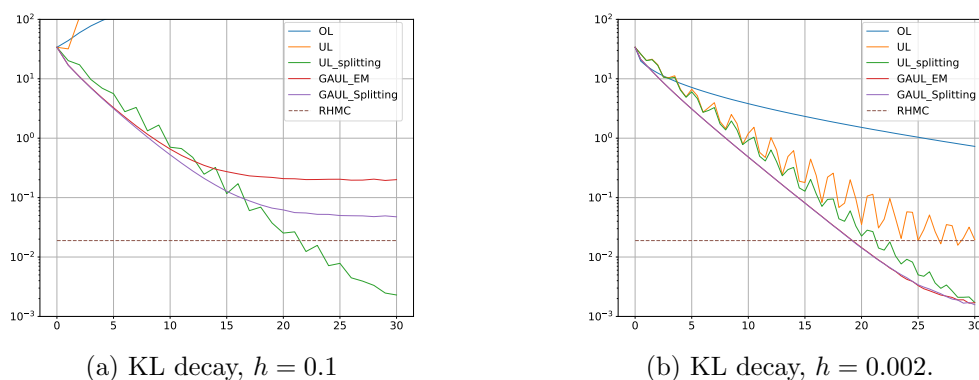


Figure 2. Convergence plots of different methods in terms of KL divergence between the sample and target distribution. (a): $h = 0.1$. (b): $h = 0.002$. The x -axis represents time. When h is large, it is evident that the GAUL_Splitting yields better bias than the GAUL_EM method. With a smaller h , both GAUL_Splitting and GAUL_EM are able to achieve a smaller error than RHMC. Note that UL_splitting achieves the best performance for both small and large h as it has zero asymptotic bias [46].

$\gamma_{ul} = 2\tilde{\Sigma}^{-1/2} = 0.2$, and $\gamma_{pdd} = 2\tilde{\Sigma}^{-1/2} + \tilde{\Sigma}^{-1} = 0.21$. In Fig. 1, we observe that our proposed method outperforms both overdamped and underdamped Langevin dynamics in both cases.

4.1.2. 20 dimensions. Let the target distribution be a 20-dimensional Gaussian with zero mean and covariance given by \mathbf{QDQ}^T , where \mathbf{D} a diagonal matrix with diagonal entries given by $0.05 + 5i$ for $i = 0, \dots, 19$, and \mathbf{Q} is an orthogonal matrix. The last entry of \mathbf{D} has the largest variance, which is $\sigma_{\max}^2 = 95.05$. Therefore, we choose $a = \frac{2}{\sigma_{\min}^{-1/2} - \sigma_{\max}^{-1/2}}$, $\gamma_{ul} = 2\sigma_{\max}^{-1}$, and $\gamma_{pdd} = 2\sigma_{\max}^{-1} + a\sigma_{\max}^{-2}$. In this example, we use $M = 10^5$ particles and (1) time step size $h = 10^{-1}$ and run for 300 steps; (2) time step size $h = 2 \times 10^{-3}$ and run for 15000 steps. In both experiments, we run four independent chains of RHMC, each of which produces 25000 particles. The KL divergence can still be computed using (4.1). All results are presented in Fig. 2.

4.2. Mixture of Gaussian.

4.2.1. Strongly log-concave. Consider the problem of sampling from a mixture of Gaussian distributions $\mathcal{N}(\alpha, \mathbf{I})$ and $\mathcal{N}(-\alpha, \mathbf{I})$, whose density satisfies:

$$p(\mathbf{x}) = \frac{1}{2(2\pi)^{d/2}} \left(e^{-\|\mathbf{x}-\alpha\|_2^2/2} + e^{-\|\mathbf{x}+\alpha\|_2^2/2} \right).$$

The corresponding potential is given as

$$(4.2) \quad f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \alpha\|_2^2 - \log \left(1 + e^{-2\mathbf{x}^\top \alpha} \right),$$

$$(4.3) \quad \nabla f(\mathbf{x}) = \mathbf{x} - \alpha + 2\alpha(1 + e^{2\mathbf{x}^\top \alpha})^{-1}.$$

Following [33, 26], we set $\alpha = (1/2, 1/2)$ and $d = 2$. This choice of parameters yields strong convexity parameter $m = 1/2$ and Lipschitz constant $L = 1$. We choose $a = \frac{2}{\sqrt{L} - \sqrt{m}}$, $\gamma_{ul} =$

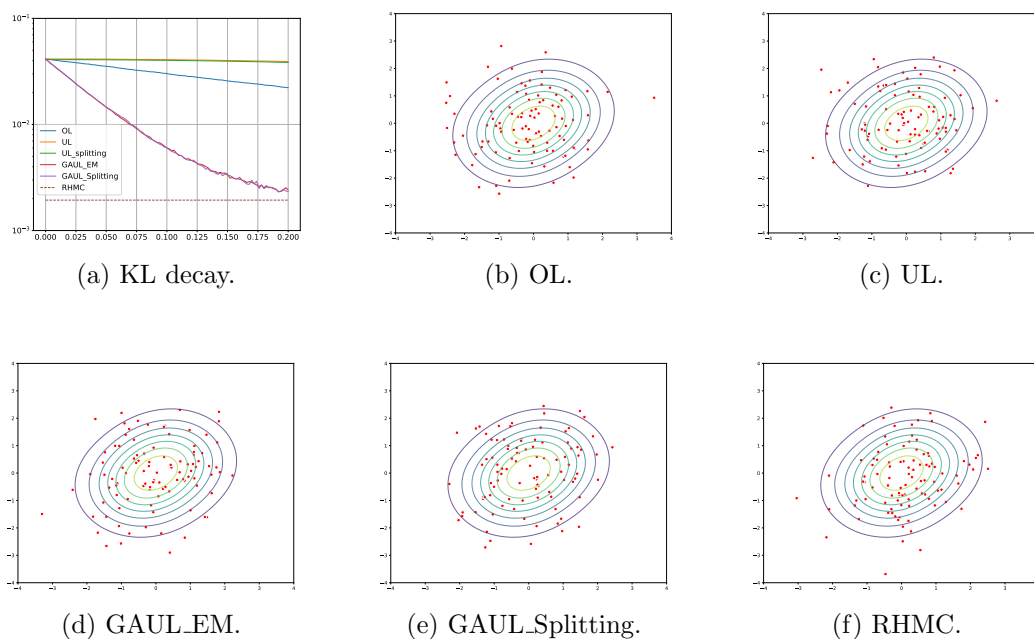


Figure 3. Convergence and scatter plots. (a): KL divergence between the sample and target distribution, which is a mixture of two unit variance Gaussians located at $(1/2, 1/2)$ and $(-1/2, -1/2)$. x -axis represents time and y -axis is in \log_{10} scale. (b)–(f): scatter plot of different methods. Contour of the true density is also provided for comparison.

$2m^{1/2}$, and $\gamma_{pdd} = 2m^{1/2} + am$. Initially particles are sampled from $\mathcal{N}(0, \mathbf{I})$. We use time step $h = 2 \times 10^{-4}$ and run for 1000 steps for OL, UL and GAUL. We use 5×10^5 particles and $n^2 = 2500$ bins to approximate the KL divergence between the sample points and the target distribution (see Remark 4.1). We run 250 independent chains of RHMC, each of which produces 2000 particles. The results are shown in Fig. 3.

Remark 4.1. To compute the KL divergence between sample points and a non-Gaussian target distribution in two dimension, we first get the 2d histogram of the samples points using n^2 bins (n in each dimension). We then use this 2d histogram as an approximation of the empirical distribution of the samples. Similarly, we can get a discretized target distribution by evaluating the target distribution at the center of each bins. Finally, we can compute the discrete KL divergence using n^2 values from the histogram and the discretized target distribution.

4.2.2. Non log-concave. We also consider the same example as in section 4.2.1 with $\alpha = (3, 3)$. As the distance between the two Gaussians increases, the target density is no longer log-concave. We use time step size $h = 10^{-3}$ and run for 2000 steps. We use $a = 1$, $\gamma_{ul} = \sqrt{2}$, and $\gamma_{pdd} = \sqrt{2} + 1/2$. We use 5×10^5 particles and $n^2 = 2500$ bins to evaluate the KL divergence. We run 250 independent chains of RHMC, each of which produces 2000 particles. The results are demonstrated in Fig. 4.

4.3. Quadratic cosine. Consider a potential function given by a quadratic function and a cosine term:

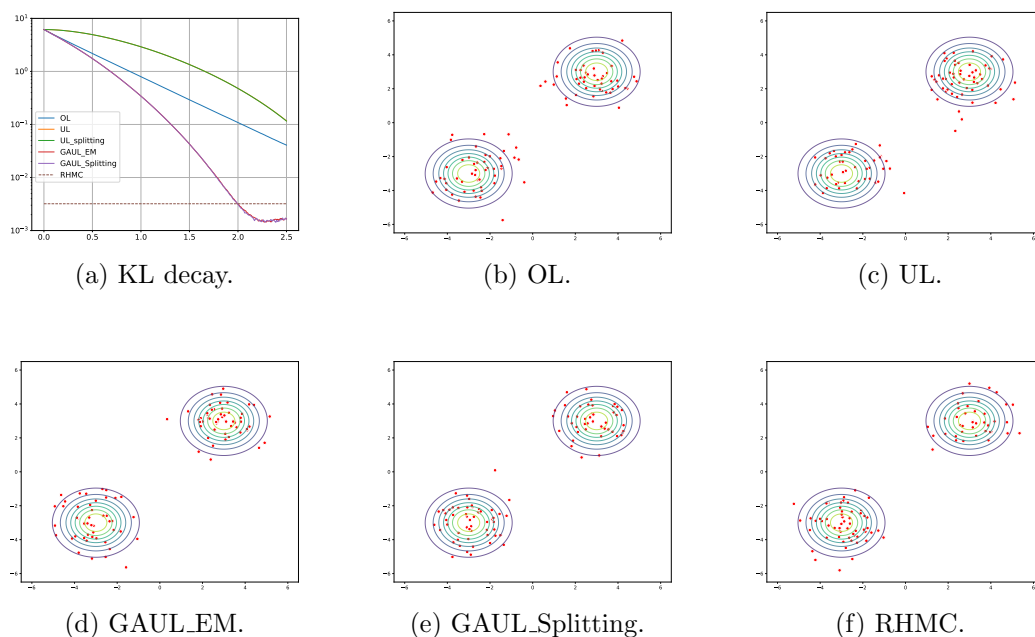


Figure 4. Convergence and scatter plots for mixture of Gaussians centered at $(3, 3)$ and $(-3, -3)$.

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T B^{-1} \mathbf{x} - \cos(\mathbf{c}^T \mathbf{x}),$$

where $B = \mathbf{P} \text{diag}(1, 25) \mathbf{P}^T$ for an orthogonal matrix \mathbf{P} and $\mathbf{c} = \sqrt{0.95} (1, 1)^T$. Here \mathbf{P} is generated by using `torch.linalg.qr(torch.randn(d))` in Pytorch, where $d = 2$ is the dimension. We set $a = 1$, $\gamma_{ul} = 2m^{1/2}$, and $\gamma_{pdd} = 2m^{1/2} + m$, where we choose $m = 1/25$. We use time step size $h = 10^{-2}$ and run for 4000 steps. We use 10^6 particles and $n^2 = 2500$ bins to evaluate the KL divergence. We run 100 independent chains of RHMC, each of which produces 10000 particles. The results are demonstrated in Fig. 5.

4.4. Bimodal. We consider a two-dimensional bimodal distribution studied in [76] whose target density has the following form:

$$p(\mathbf{x}) \propto \exp(-2(\|\mathbf{x}\| - 3)^2) \left[\exp(-2(x_1 - 3)^2) + \exp(-2(x_1 + 3)^2) \right].$$

The corresponding potential function is given by

$$f(\mathbf{x}) = 2(\|\mathbf{x}\| - 3)^2 - \log \left[\exp(-2(x_1 - 3)^2) + \exp(-2(x_1 + 3)^2) \right].$$

The gradient is

$$\begin{aligned} \nabla f(\mathbf{x}) = & \frac{4(x_1 - 3) \exp(-2(x_1 - 3)^2) + 4(x_1 + 3) \exp(-2(x_1 + 3)^2)}{\exp(-2(x_1 - 3)^2) + \exp(-2(x_1 + 3)^2)} \mathbf{e}_1 \\ & + 4 \frac{(\|\mathbf{x}\| - 3) \mathbf{x}}{\|\mathbf{x}\|}, \end{aligned}$$

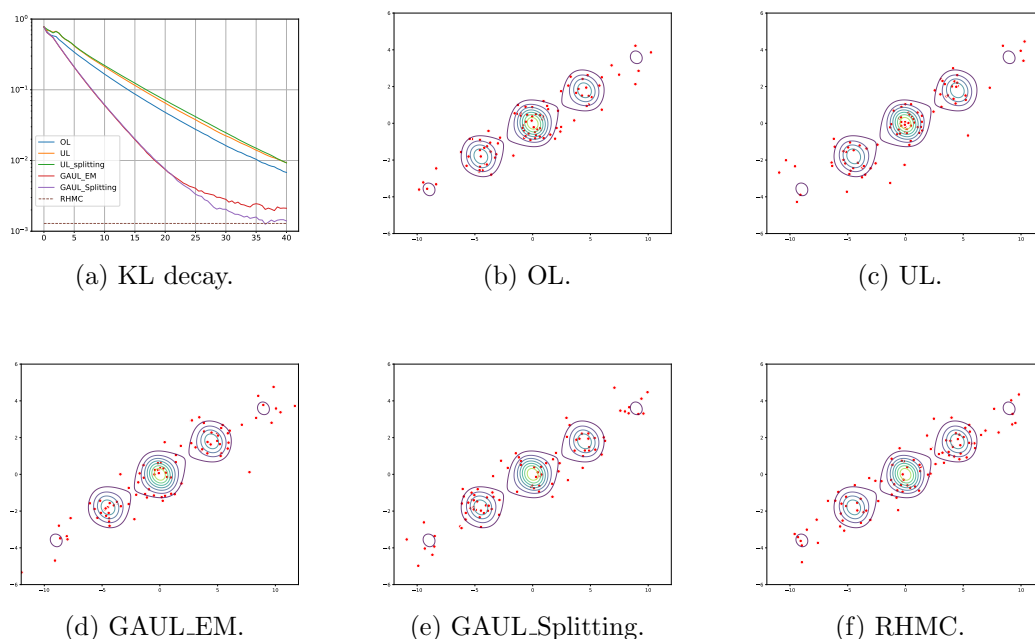


Figure 5. Convergence and scatter plots for the quadratic cosine example.

where $e_1 = (1, 0)^T$ is the first standard coordinate vector. We set $\gamma_{ul} = 2m^{1/2}$ and $\gamma_{pdd} = 2m^{1/2} + m$, where we choose $m = 1/2$. We use time step size $h = 2 \times 10^{-3}$ and run for 800 iterations. We use 10^6 particles and $n^2 = 2500$ bins to evaluate the KL divergence. We run 250 independent chains of RHMC, each of which produces 2000 particles. The results are shown in Fig. 6.

4.5. Bayesian logistic regression. We consider the Bayesian logistic regression problem studied in [33, 26, 72]. We give a brief description of the problem. Suppose we are given a feature matrix $X \in \mathbb{R}^{n \times d}$ with rows $x_i \in \mathbb{R}^d$. Correspondingly we are given $Y \in \{0, 1\}^n$ the binary response vector for each of the covariates in our feature matrix. The logistic model for the probability of $y_i = 1$ given $x_i \in \mathbb{R}^d$ and a parameter $\theta \in \mathbb{R}^d$ is

$$(4.4) \quad \mathbb{P}(y_i = 1 | x_i, \theta) = \frac{\exp(\theta^T x_i)}{1 + \exp(\theta^T x_i)}.$$

Suppose we impose a prior distribution on the parameter $\theta \sim \mathcal{N}(0, \Sigma_X)$, where $\Sigma_X = \frac{1}{n} X^T X$ is the sample covariance of X . Then the posterior distribution for θ can be calculated by

$$p(\theta | X, Y) \propto \exp \left[Y^T X \theta - \sum_{i=1}^n \log(1 + \exp(\theta^T x_i)) - \frac{\alpha}{2} \theta^T \Sigma_X \theta \right],$$

where $\alpha > 0$ is a regularization parameter. The potential function is

$$f(\theta) = -Y^T X \theta + \sum_{i=1}^n \log(1 + \exp(\theta^T x_i)) + \frac{\alpha}{2} \theta^T \Sigma_X \theta.$$

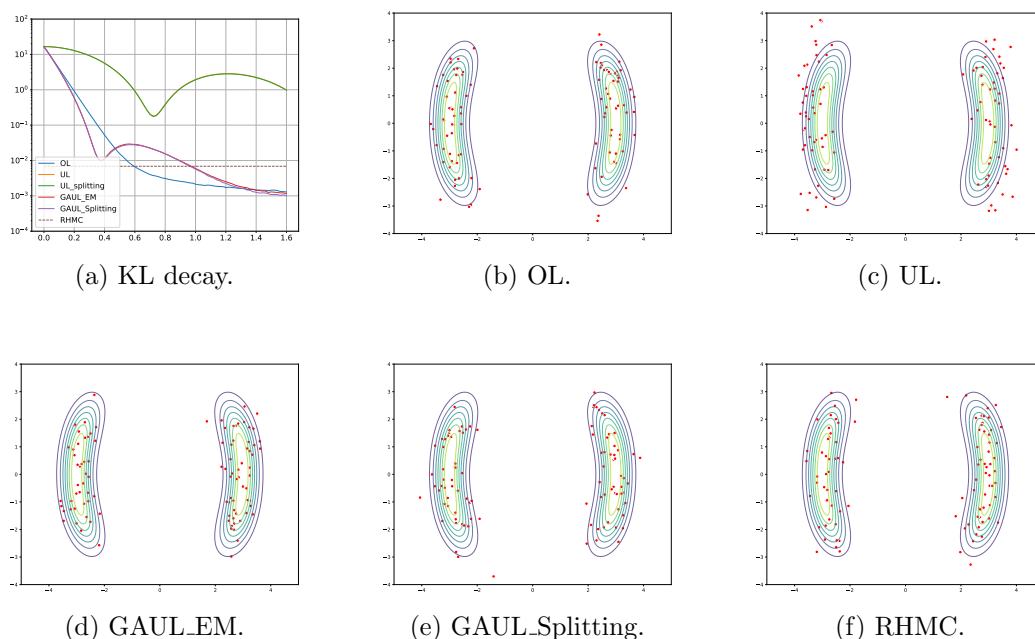


Figure 6. Convergence and scatter plots for bimodal distribution.

Its gradient is

$$\nabla f(\theta) = -X^T Y + \sum_{i=1}^n \frac{x_i}{1 + \exp(-\theta^T x_i)} + \alpha \Sigma_X \theta.$$

As shown in [33], the Hessian of f is upper bounded by $L = (0.25n + \alpha)\lambda_{\max}(\Sigma_X)$ and lower bounded by $m = \alpha\lambda_{\min}$. To generate X and Y , we set $x_{i,j}$ to be independent Rademacher random variables for each i and j . And each y_i is generated according to (4.4) with $\theta = \theta^* = (1, 1)^T$. We set $\alpha = 0.5$, $d = 2$, $n = 50$, $\gamma_{ul} = 2m^{1/2}$, and $\gamma_{pdd} = 2m^{1/2} + m$. To sample the posterior distribution, we use time step size $h = 10^{-2}$ and run for 40 iterations. The initial distribution of particles is $\mathcal{N}(0, L^{-1}\mathbf{I})$. As for evaluation metric, we directly evaluate the KL divergence between the sampled posterior and the true posterior. We use 10^6 particles and $n^2 = 2500$ bins to evaluate the KL divergence as before. We run 10000 independent chains of RHMC, each of which produces 100 particles. This is different from the choice by [33] and [72], where [33] compared the samples with θ^* . [72] compared samples with the true minimizer of $f(\theta)$, i.e., the maximum a posteriori (MAP) estimate in the Bayesian optimization literature. We believe that directly measuring the KL divergence gives a better understanding of how ‘close’ our samples are to the true posterior distribution. The results are presented in Fig. 7.

4.6. Bayesian neural network. In this section, we compare GAUL with overdamped, underdamped Langevin dynamics, and randomized Hamiltonian Monte Carlo in training Bayesian neural network. We test a one-hidden-layer fully connected neural network with 50 hidden neurons and ReLU activation function on the UCI concrete dataset. We use $h = 10^{-3}$, $a = 0.1$, and $\gamma = 0.5$. For each method, we sample $M = 20$ particles (each particle corresponds to a neural network) and take the average output as the final output. Fig. 8a and

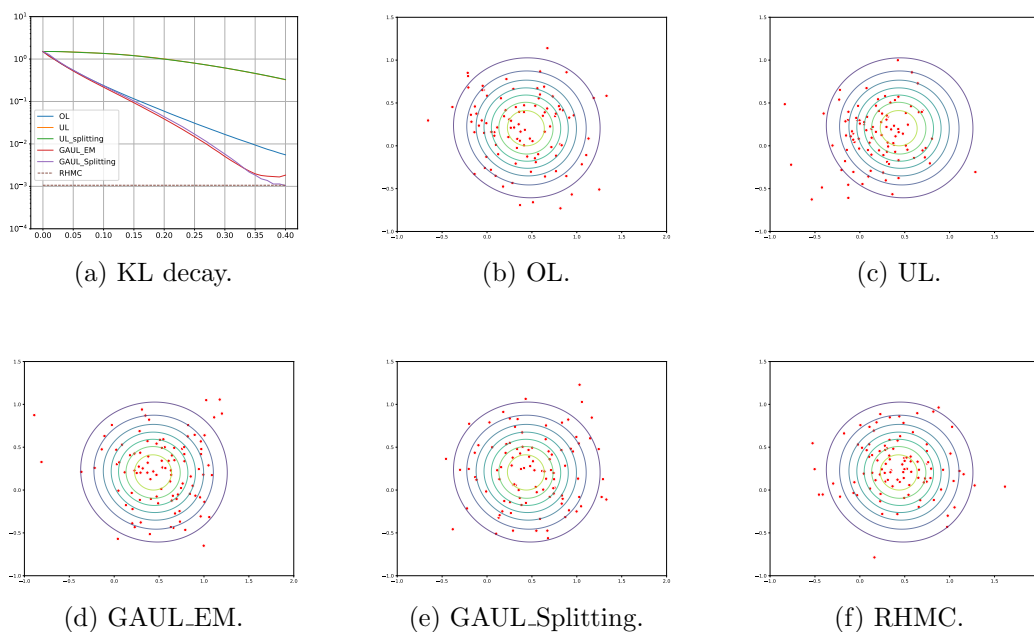


Figure 7. Convergence and scatter plots for Bayesian logistic regression.

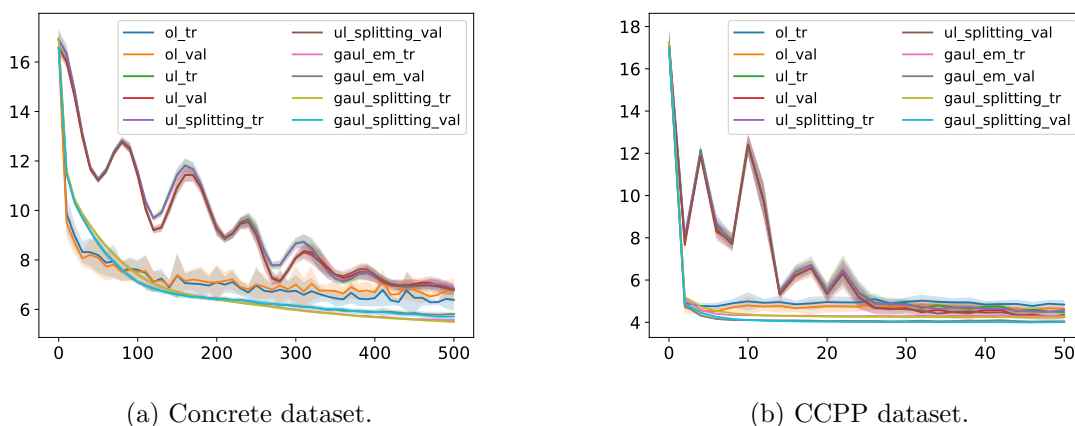


Figure 8. Convergence comparison. x -axis represents number of epochs. y -axis represents rmse averaged over 10 experiments.

Table 1 show the rMSE averaged over 10 experiments. We see that ‘ul’ can achieve smaller training and validation error than ‘ol’. However, ‘ul’ also exhibits a slow start and an oscillatory behavior at the beginning of training as is commonly seen in acceleration methods in optimization. GAUL can get rid of the oscillation and achieve a even smaller training and validation error as is demonstrated in Table 1. We have also tested out the three methods using the Combined Cycle Power Plant (CCPP) dataset. We choose the same parameter as the concrete experiment. The results are presented in Fig. 8b and Table 1. Due to the nature of the dynamic_hmc in Blackjax, it is difficult to plot the training and validation curve over time. The final result for RHMC is included in Table 1.

Table 1

Training and validation RMSE for Bayesian neural network (transposed view).

	concrete tr err	concrete val err	ccpp tr err	ccpp val err
ol	6.43 ± 0.28	6.67 ± 0.22	4.84 ± 0.22	4.63 ± 0.25
ul	6.19 ± 0.16	6.26 ± 0.11	4.48 ± 0.11	4.25 ± 0.11
ul_splitting	6.24 ± 0.10	6.24 ± 0.14	4.59 ± 0.16	4.36 ± 0.19
gaul_em	5.57 ± 0.09	5.81 ± 0.09	4.28 ± 0.03	4.04 ± 0.04
gaul_splitting	5.49 ± 0.03	5.71 ± 0.11	4.25 ± 0.02	4.01 ± 0.03
rhmc	5.67 ± 0.17	5.94 ± 0.21	4.49 ± 0.11	4.26 ± 0.11

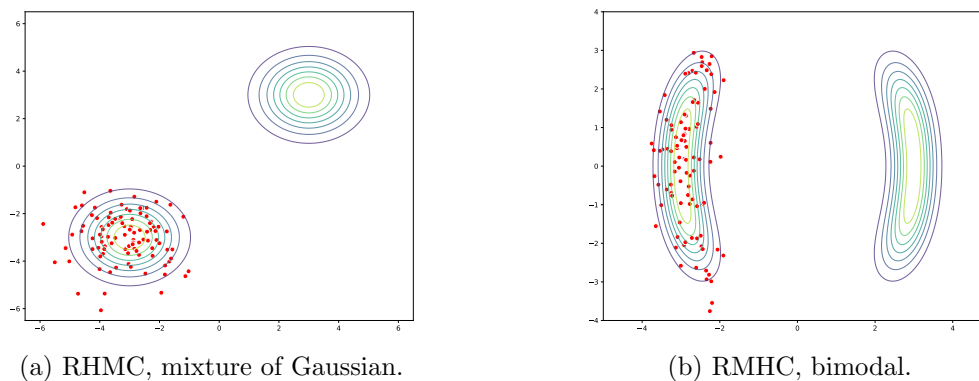


Figure 9. Mode collapse for RHMC. Both figures show the last 100 particles in a single chain.

4.7. Summary of numerical experiments. We have compared overdamped Langevin, underdamped Langevin, GAUL (Euler–Maruyama scheme and *BAGOCAB* splitting), and the randomized Hamiltonian Monte Carlo method over several log-concave and non log-concave examples. In all experiments, GAUL converges faster than both OL and UL. In addition, GAUL_Splitting usually induces a smaller bias than GAUL_EM. In terms of final error, GAUL performs comparably to RHMC, if not better. Thanks to ergodicity, RHMC is able to use a single particle to generate a multiple instances of the target distribution, making it efficient for sampling large amount of particles. However, for multimodal distributions (e.g., section 4.2.2 and section 4.2.1), if we only look at the last 100 particles from a single chain, we observe a mode collapse phenomenon shown in Fig. 9. This is not the case for other methods considered in our experiments. Therefore, we also need to run multiple independent chains for RHMC to avoid mode collapse.

5. Conclusions. In this work, we introduced gradient-adjusted underdamped Langevin dynamics (GAUL) inspired by primal-dual damping dynamics and Hessian-driven damping dynamics. We demonstrated that GAUL admits the correct invariant target distribution $\pi \propto \exp(-f)$ under appropriate conditions and achieves exponential convergence for quadratic functions, outperforming both the overdamped and underdamped Langevin dynamics in terms of convergence speed. Our numerical experiments further illustrate the practical advantages of GAUL, showcasing faster convergence and more efficient sampling compared to classical methods, such as overdamped and underdamped Langevin dynamics.

We also note a connection between the primal-dual damping dynamics and GAUL. A key challenge in the primal-dual damping algorithm is the design of preconditioner matrices, which can accelerate the algorithm's convergence compared to the gradient descent method. In the context of solving a linear problem where f is a quadratic function and the diffusion constant is zero, [79] demonstrates that the convergence rate depends on the square root of the smallest eigenvalue. In this paper, we extend the study from a sampling perspective, where f is also a quadratic function, but the diffusion is nonzero. Toward a Gaussian target distribution, GAUL converges to a biased target distribution with the mixing time depending on $\sqrt{\kappa}$. This is in contrast with overdamped and underdamped Langevin sampling algorithms.

Several possible future directions are worth exploring. First, can we show that GAUL converges faster than overdamped and underdamped Langevin dynamics for more general potentials, which is beyond the current study of Gaussian distributions? One common assumption is that the potential f is strongly log-concave [14, 23, 24, 25, 31, 33, 40, 44, 51]. Recently, [15] proved that for a class of distributions that satisfy a Poincaré-type inequality, underdamped Langevin dynamics converges in L_2 with rate $\exp(-\sqrt{mt})$, where m is the Poincaré constant. Then it is interesting to study for the same class of distributions, whether GAUL could converge at an even faster rate. Another direction is to study the convergence of GAUL under different metrics. From a more practical perspective, designing new time discretization schemes [65, 22, 59, 72, 51] for implementing GAUL is also an important direction. We proved that using the Euler–Maruyama discretization, GAUL will converge to a biased target distribution, which is not surprising since ULA is also biased. We also proved that $BAGOCAB$ splitting could achieve a smaller first-order bias than Euler–Maruyama. This is also verified in our numerical experiments. Therefore, another promising direction could be to combine GAUL with MCMC methods [8, 33, 12, 66, 10], such as Metropolis–Hastings algorithms, to design a hybrid method with accept/reject options so that the algorithm converges to the correct target distribution in the discrete-time update. Finally, choosing the preconditioner \mathbf{C} to accelerate convergence is an important topic. The difficulty of picking \mathbf{C} arises from the positive semidefinite constraint on $\text{sym}(\mathbf{Q})$ in (2.16), which we should explore in future work.

Appendix A. Euler–Maruyama discretization. The Euler–Maruyama scheme of (2.15) with step size h and $\mathbf{C} = \mathbf{I}$ reads

$$(A.1a) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - a\nabla f(\mathbf{x}_t)h + \mathbf{p}_t h + \sqrt{2ah}\mathbf{z}^{(1)},$$

$$(A.1b) \quad \mathbf{p}_{t+1} = \mathbf{p}_t - \nabla f(\mathbf{x}_t)h - \gamma\mathbf{p}_t h + \sqrt{2\gamma h}\mathbf{z}^{(2)}.$$

$\mathbf{z}^{(i)}$ is a standard Gaussian random variable for $i = 1, 2$.

REFERENCES

- [1] L. AMBROSIO, N. GIGLI, AND G. SAVARÉ, *Gradient Flows: In Metric Spaces and in the Space of Probability Measures*, Springer Science & Business Media, 2008.
- [2] C. ANDRIEU, N. DE FREITAS, A. DOUCET, AND M. I. JORDAN, *An introduction to MCMC for machine learning*, Machine Learn., 50 (2003), pp. 5–43, <https://doi.org/10.1023/A:1020281327116>.
- [3] H. ATTOUCH, Z. CHBANI, J. FADILI, AND H. RIAHI, *First-order optimization algorithms via inertial systems with Hessian driven damping*, Math. Program., 193 (2022), pp. 113–155.

- [4] H. ATTOUCH, Z. CHBANI, J. FADILI, AND H. RIAHI, *Convergence of iterates for first-order optimization algorithms with inertia and Hessian driven damping*, *Optimization*, 72 (2023), pp. 1199–1238.
- [5] H. ATTOUCH, Z. CHBANI, AND H. RIAHI, *Fast proximal methods via time scaling of damped inertial dynamics*, *SIAM J. Optim.*, 29 (2019), pp. 2227–2256, <https://doi.org/10.1137/18M1230207>.
- [6] A. BARP, S. TAKAO, M. BETANCOURT, A. ARNAUDON, AND M. GIROLAMI, *A Unifying and Canonical Description of Measure-Preserving Diffusions*, preprint, [arXiv:2105.02845](https://arxiv.org/abs/2105.02845), 2021.
- [7] C. H. BENNETT, *Mass tensor molecular dynamics*, *J. Comput. Phys.*, 19 (1975), pp. 267–279.
- [8] J. BESAG, *Comments on “Representations of knowledge in complex systems” by U. Grenander and MI Müller*, *J. Roy. Statist. Soc. Ser. B*, 56 (1994), 4.
- [9] M. BETANCOURT, *A Conceptual Introduction to Hamiltonian Monte Carlo*, preprint, [arXiv:1701.02434](https://arxiv.org/abs/1701.02434), 2017.
- [10] N. BOU-RABEE AND S. OBERDÖRSTER, *Mixing of Metropolis-adjusted Markov chains via couplings: The high acceptance regime*, *Electron. J. Probab.*, 29 (2024), pp. 1–27, <https://doi.org/10.1214/24-EJP1150>.
- [11] N. BOU-RABEE AND J. M. SANZ-SERNA, *Randomized Hamiltonian Monte Carlo*, *Ann. Appl. Probab.*, 27 (2017), pp. 2159–2194, <https://doi.org/10.1214/16-AAP1255> (accessed 2025-04-21).
- [12] N. BOU-RABEE AND E. VANDEN-ELJNDEN, *Pathwise accuracy and ergodicity of Metropolized integrators for SDEs*, *Comm. Pure Appl. Math.*, 63 (2010), pp. 655–696, <https://doi.org/10.1002/cpa.20306>.
- [13] G. BUSSI AND M. PARRINELLO, *Accurate sampling using Langevin dynamics*, *Phys. Rev. E—Statistical, Nonlinear, and Soft Matter Physics*, 75 (2007), 056707, <https://doi.org/10.1103/PhysRevE.75.056707>.
- [14] Y. CAO, J. LU, AND L. WANG, *Complexity of Randomized Algorithms for Underdamped Langevin Dynamics*, preprint, [arXiv:2003.09906](https://arxiv.org/abs/2003.09906), 2020.
- [15] Y. CAO, J. LU, AND L. WANG, *On explicit L_2 -convergence rate estimate for underdamped Langevin dynamics*, *Arch. Ration. Mech. Anal.*, 247 (2023), 90.
- [16] J. A. CARRILLO, Y.-P. CHOI, AND O. TSE, *Convergence to equilibrium in Wasserstein distance for damped Euler equations with interaction forces*, *Commun. Math. Phys.*, 365 (2019), pp. 329–361, <https://doi.org/10.1007/s00220-018-3276-8>.
- [17] F. CASAS, J. M. SANZ-SERNA, AND L. SHAW, *Split Hamiltonian Monte Carlo revisited*, *Stat. Comput.*, 32 (2022), 86, <https://doi.org/10.1007/s11222-022-10149-4>.
- [18] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications to imaging*, *J. Math. Imaging Vis.*, 40 (2011), pp. 120–145, <https://doi.org/10.1007/s10851-010-0251-1>.
- [19] S. CHEN, Q. LI, O. TSE, AND S. J. WRIGHT, *Accelerating Optimization over the Space of Probability Measures*, preprint, [arXiv:2310.04006](https://arxiv.org/abs/2310.04006), 2023.
- [20] Y. CHEN, D. Z. HUANG, J. HUANG, S. REICH, AND A. M. STUART, *Gradient Flows for Sampling: Mean-field Models, Gaussian Approximations and Affine Invariance*, preprint, [arXiv:2302.11024](https://arxiv.org/abs/2302.11024), 2023.
- [21] X. CHENG AND P. BARTLETT, *Convergence of Langevin MCMC in KL-divergence*, in *Algorithmic Learning Theory*, PMLR, 2018, pp. 186–211.
- [22] X. CHENG, N. S. CHATTERJI, P. L. BARTLETT, AND M. I. JORDAN, *Underdamped Langevin MCMC: A non-asymptotic analysis*, in *Conference on Learning Theory*, PMLR, 2018, pp. 300–323.
- [23] S. CHEWI, P. R. GERBER, C. LU, T. LE GOUIC, AND P. RIGOLLET, *The query complexity of sampling from strongly log-concave distributions in one dimension*, in *Conference on Learning Theory*, PMLR, 2022, pp. 2041–2059.
- [24] S. CHEWI, C. LU, K. AHN, X. CHENG, T. LE GOUIC, AND P. RIGOLLET, *Optimal dimension dependence of the Metropolis-adjusted Langevin algorithm*, in *Conference on Learning Theory*, PMLR, 2021, pp. 1260–1300.
- [25] A. DALALYAN, *Further and stronger analogy between sampling and optimization: Langevin Monte Carlo and gradient descent*, in *Conference on Learning Theory*, PMLR, 2017, pp. 678–689.
- [26] A. S. DALALYAN, *Theoretical guarantees for approximate sampling from smooth and log-concave densities*, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 79 (2017), pp. 651–676, <https://doi.org/10.1111/rssb.12183>.

- [27] A. S. DALALYAN AND A. KARAGULYAN, *User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient*, *Stochastic Process. Appl.*, 129 (2019), pp. 5278–5311, <https://doi.org/10.1016/j.spa.2019.02.016>.
- [28] A. S. DALALYAN AND L. RIOU-DURAND, *On sampling from a log-concave density using kinetic Langevin diffusions*, *Bernoulli*, 26 (2020), pp. 1956–1988, <https://doi.org/10.3150/19-BEJ1178>.
- [29] M. DASHTI AND A. M. STUART, *The Bayesian Approach to Inverse Problems*, preprint, [arXiv:1302.6989](https://arxiv.org/abs/1302.6989), 2013.
- [30] L. DEVROYE, A. MEHRABIAN, AND T. REDDAD, *The Total Variation Distance between High-Dimensional Gaussians with the Same Mean*, preprint, [arXiv:1810.08693](https://arxiv.org/abs/1810.08693), 2018.
- [31] A. DURMUS, S. MAJEWSKI, AND B. MIASOJEDOW, *Analysis of Langevin Monte Carlo via convex optimization*, *J. Mach. Learn. Res.*, 20 (2019), pp. 1–46.
- [32] A. DURMUS AND E. MOULINES, *Nonasymptotic convergence analysis for the unadjusted Langevin algorithm*, *Ann. Appl. Probab.*, 27 (2017), pp. 1551–1587, <https://doi.org/10.1214/16-AAP1238>.
- [33] R. DWIVEDI, Y. CHEN, M. J. WAINWRIGHT, AND B. YU, *Log-concave sampling: Metropolis-Hastings algorithms are fast*, *J. Mach. Learn. Res.*, 20 (2019), pp. 1–42.
- [34] Q. FENG, X. ZUO, AND W. LI, *Fisher Information Dissipation for Time Inhomogeneous Stochastic Differential Equations*, preprint, [arXiv:2402.01036](https://arxiv.org/abs/2402.01036), 2024.
- [35] A. GARBUNO-INIGO, F. HOFFMANN, W. LI, AND A. M. STUART, *Interacting Langevin diffusions: Gradient structure and ensemble Kalman sampler*, *SIAM J. Appl. Dyn. Syst.*, 19 (2020), pp. 412–441, <https://doi.org/10.1137/19M1251655>.
- [36] S. B. GELFAND AND S. K. MITTER, *Simulated annealing type algorithms for multivariate optimization*, *Algorithmica*, 6 (1991), pp. 419–436, <https://doi.org/10.1007/BF01759052>.
- [37] A. GELMAN, J. B. CARLIN, H. S. STERN, AND D. B. RUBIN, *Bayesian Data Analysis*, Chapman and Hall/CRC, 1995, <https://doi.org/10.1201/9780429258411>.
- [38] M. GIROLAMI AND B. CALDERHEAD, *Riemann manifold Langevin and Hamiltonian Monte Carlo methods*, *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 73 (2011), pp. 123–214, <https://doi.org/10.1111/j.1467-9868.2010.00765.x>.
- [39] J. GOODMAN AND J. WEARE, *Ensemble samplers with affine invariance*, *Comm. Appl. Math. Comput. Sci.*, 5 (2010), pp. 65–80, <https://doi.org/10.2140/camcos.2010.5.65>.
- [40] Y. HE, K. BALASUBRAMANIAN, AND M. A. ERDOGDU, *On the ergodicity, bias and asymptotic normality of randomized midpoint sampling method*, *Adv. Neural Inf. Process. Syst.*, 33 (2020), pp. 7366–7376.
- [41] J. IDIER, *Bayesian Approach to Inverse Problems*, John Wiley & Sons, 2013.
- [42] P. IZMAILOV, S. VIKRAM, M. D. HOFFMAN, AND A. G. WILSON, *What are Bayesian neural network posteriors really like?*, in *International Conference on Machine Learning*, PMLR, 2021, pp. 4629–4640.
- [43] R. JORDAN, D. KINDERLEHRER, AND F. OTTO, *The variational formulation of the Fokker–Planck equation*, *SIAM J. Math. Anal.*, 29 (1998), pp. 1–17, <https://doi.org/10.1137/S0036141096303359>.
- [44] Y. T. LEE, R. SHEN, AND K. TIAN, *Logsmooth gradient concentration and tighter runtimes for Metropolized Hamiltonian Monte Carlo*, in *Conference on Learning Theory*, PMLR, 2020, pp. 2565–2597.
- [45] B. LEIMKUHLE AND C. MATTHEWS, *Rational construction of stochastic numerical methods for molecular sampling*, *Appl. Math. Res. eXpress*, 2013, (2013), pp. 34–56.
- [46] B. LEIMKUHLE AND C. MATTHEWS, *Robust and efficient configurational molecular sampling via Langevin dynamics*, *J. Chem. Phys.*, 138 (2013), 174102, <https://doi.org/10.1063/1.4802990>.
- [47] B. LEIMKUHLE, C. MATTHEWS, AND J. WEARE, *Ensemble preconditioning for Markov chain Monte Carlo simulation*, *Stat. Comput.*, 28 (2018), pp. 277–290, <https://doi.org/10.1007/s11222-017-9730-1>.
- [48] B. J. LEIMKUHLE, D. PAULIN, AND P. A. WHALLEY, *Contraction and convergence rates for discretized kinetic Langevin dynamics*, *SIAM J. Numer. Anal.*, 62 (2024), pp. 1226–1258, <https://doi.org/10.1137/23M1556289>.
- [49] T. LELIÈVRE, F. NIER, AND G. A. PAVLIOTIS, *Optimal non-reversible linear drift for the convergence to equilibrium of a diffusion*, *J. Stat. Phys.*, 152 (2013), pp. 237–274, <https://doi.org/10.1007/s10955-013-0769-x>.
- [50] T. LELIÈVRE, G. A. PAVLIOTIS, G. ROBIN, R. SANTET, AND G. STOLTZ, *Optimizing the Diffusion of Overdamped Langevin Dynamics*, preprint, [arXiv:2404.12087](https://arxiv.org/abs/2404.12087), 2024.

- [51] R. LI, H. ZHA, AND M. TAO, *Hessian-free high-resolution Nesterov acceleration for sampling*, in International Conference on Machine Learning, PMLR, 2022, pp. 13125–13162.
- [52] J. S. LIU, *Monte Carlo Strategies in Scientific Computing*, Springer Ser. Statist. 10, Springer, 2001.
- [53] Y.-A. MA, N. S. CHATTERJI, X. CHENG, N. FLAMMARION, P. L. BARTLETT, AND M. I. JORDAN, *Is there an analog of Nesterov acceleration for gradient-based MCMC?*, *Bernoulli*, 27 (2021), pp. 1942–1992, <https://doi.org/10.3150/20-BEJ1297>.
- [54] D. J. MACKAY, *Bayesian neural networks and density networks*, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, 354 (1995), pp. 73–80, [https://doi.org/10.1016/0168-9002\(94\)00931-7](https://doi.org/10.1016/0168-9002(94)00931-7).
- [55] D. J. MACKAY, *Information Theory, Inference and Learning Algorithms*, Cambridge University Press, 2003.
- [56] C. J. MADDISON, D. PAULIN, Y. W. TEH, B. O'DONOGHUE, AND A. DOUCET, *Hamiltonian Descent Methods*, preprint, [arXiv:1809.05042](https://arxiv.org/abs/1809.05042), 2018.
- [57] J. C. MATTINGLY, A. M. STUART, AND D. J. HIGHAM, *Ergodicity for SDEs and approximations: Locally Lipschitz vector fields and degenerate noise*, *Stochastic Process. Appl.*, 101 (2002), pp. 185–232, [https://doi.org/10.1016/S0304-4149\(02\)00150-3](https://doi.org/10.1016/S0304-4149(02)00150-3).
- [58] S. P. MEYN AND R. L. TWEEDIE, *Markov Chains and Stochastic Stability*, Springer Science & Business Media, 2012.
- [59] W. MOU, Y.-A. MA, M. J. WAINWRIGHT, P. L. BARTLETT, AND M. I. JORDAN, *High-Order Langevin Diffusion Yields an Accelerated MCMC Algorithm*, preprint, [arXiv:1908.10859](https://arxiv.org/abs/1908.10859), 2019.
- [60] R. M. NEAL, *Bayesian Learning for Neural Networks*, Lect. Notes Statist. 118, Springer Science & Business Media, 2012.
- [61] R. M. NEAL, ET AL., *MCMC using Hamiltonian dynamics*, in Handbook of Markov Chain Monte Carlo, Chapman & Hall/CRC, 2011, pp. 113–162.
- [62] Y. E. NESTEROV, *A method of solving a convex programming problem with convergence rate $\mathcal{O}(\frac{1}{k^2})$* , in Doklady Akademii Nauk, Vol. 269, Russian Academy of Sciences, 1983, pp. 543–547.
- [63] C. P. ROBERT, G. CASELLA, AND G. CASELLA, *Monte Carlo Statistical Methods*, Vol. 2, Springer, 1999.
- [64] G. O. ROBERTS AND R. L. TWEEDIE, *Exponential convergence of Langevin distributions and their discrete approximations*, *Bernoulli*, 2 (1996), pp. 341–363, <https://doi.org/10.2307/3318418>.
- [65] R. SHEN AND Y. T. LEE, *The randomized midpoint method for log-concave sampling*, *Adv. Neural Inf. Process. Syst.*, 32 (2019).
- [66] G. STOLTZ, M. ROUSSET, ET AL., *Free Energy Computations: A Mathematical Perspective*, World Scientific, 2010.
- [67] A. M. STUART, *Inverse problems: A Bayesian perspective*, *Acta Numer.*, 19 (2010), pp. 451–559.
- [68] W. SU, S. BOYD, AND E. J. CANDÈS, *A differential equation for modeling Nesterov's accelerated gradient method: Theory and insights*, *J. Mach. Learn. Res.*, 17 (2016), pp. 1–43.
- [69] A. TAGHVAEI AND P. MEHTA, *Accelerated flow for probability distributions*, in International Conference on Machine Learning, PMLR, 2019, pp. 6076–6085.
- [70] D. TALAY, *Stochastic Hamiltonian systems: Exponential convergence to the invariant measure, and discretization by the implicit Euler scheme*, *Markov Process. Related Fields*, 8 (2002), pp. 163–198.
- [71] D. TALAY AND L. TUBARO, *Expansion of the global error for numerical schemes solving stochastic differential equations*, *Stoch. Anal. Appl.*, 8 (1990), pp. 483–509, <https://doi.org/10.1080/07362999008809220>.
- [72] H. Y. TAN, S. OSHER, AND W. LI, *Noise-Free Sampling Algorithms via Regularized Wasserstein Proximality*, preprint, [arXiv:2308.14945](https://arxiv.org/abs/2308.14945), 2023.
- [73] Y. W. TEH, A. THIÉRY, AND S. J. VOLLMER, *Consistency and fluctuations for stochastic gradient Langevin dynamics*, *J. Mach. Learn. Res.*, 17 (2016), 7.
- [74] T. VALKONEN, *A primal-dual hybrid gradient method for nonlinear operators with applications to MRI*, *Inverse Problems*, 30 (2014), 055012, <https://doi.org/10.1088/0266-5611/30/5/055012>.
- [75] S. VEMPALA AND A. WIBISONO, *Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices*, in Geometric Aspects of Functional Analysis, Lecture Notes in Mathematics 2327, R. Eldan, B. Klartag, A. Litvak, and E. Milman, eds., Springer, Cham, 2023, https://doi.org/10.1007/978-3-031-26300-2_15.

- [76] Y. WANG AND W. LI, *Accelerated information gradient flow*, J. Sci. Comput., 90 (2022), pp. 1–47, <https://doi.org/10.1007/s10915-021-01709-3>.
- [77] M. WELLING AND Y. W. TEH, *Bayesian learning via stochastic gradient Langevin dynamics*, in Proceedings of the 28th International Conference on Machine Learning (ICML-11), Citeseer, 2011, pp. 681–688.
- [78] S. ZHANG, S. CHEWI, M. LI, K. BALASUBRAMANIAN, AND M. A. ERDOGDU, *Improved discretization analysis for underdamped Langevin Monte Carlo*, in The Thirty Sixth Annual Conference on Learning Theory, PMLR, 2023, pp. 36–71.
- [79] X. ZUO, S. OSHER, AND W. LI, *Primal-dual damping algorithms for optimization*, Ann. Math. Sci. Appl., 9 (2024), pp. 467–504, <https://doi.org/10.4310/AMSA.2024.v9.n2.a7>.