



**HAL**  
open science

# Generalized and hybrid Metropolis-Hastings overdamped Langevin algorithms

Romain Poncet

► **To cite this version:**

Romain Poncet. Generalized and hybrid Metropolis-Hastings overdamped Langevin algorithms. 2017. hal-01440499

**HAL Id: hal-01440499**

**<https://hal.science/hal-01440499>**

Preprint submitted on 19 Jan 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# GENERALIZED AND HYBRID METROPOLIS-HASTINGS OVERDAMPED LANGEVIN ALGORITHMS

ROMAIN PONCET<sup>1</sup>

<sup>1</sup> CMAP, Ecole Polytechnique, CNRS, Université Paris-Saclay, 91128 Palaiseau, France;  
romain.poncet@cmap.polytechnique.fr

**Abstract.** It has been shown that the nonreversible overdamped Langevin dynamics enjoy better convergence properties in terms of spectral gap and asymptotic variance than the reversible one ([12, 13, 16, 25, 20, 21, 8]). In this article we propose a variance reduction method for the Metropolis-Hastings Adjusted Langevin Algorithm (MALA) that makes use of the good behaviour of these nonreversible dynamics. It consists in constructing a nonreversible Markov chain (with respect to the target invariant measure) by using a Generalized Metropolis-Hastings adjustment on a lifted state space. We present two variations of this method and we discuss the importance of a well-chosen proposal distribution in terms of average rejection probability. We conclude with numerical experimentations to compare our algorithms with the MALA, and show variance reduction of several orders of magnitude in some favourable toy cases.

## 1. INTRODUCTION

This article proposes a new class of MCMC algorithms whose objective is to compute expectations

$$\pi(f) := \mathbb{E}_\pi(f) = \int_{\mathbb{R}^d} f(x)\pi(dx), \quad (1.1)$$

for a given observable  $f$ , with respect to a probability measure  $\pi(dx)$  absolutely continuous, with respect to the Lebesgue measure, with density  $\pi(x) = e^{-U(x)}$ . We suppose, as it is the case in many practical situations, that  $\pi$  is only known up to a multiplicative constant.

Many techniques have been developed to solve this problem. Deterministic quadratures can be very efficient at low dimension. Yet, in the high dimensional case, these methods tend to become inefficient, and MCMC methods can be used instead. The basic idea is to construct an ergodic Markov chain with respect to  $\pi$ , and to approximate  $\pi(f)$  by the time average of this Markov chain. There are infinitely many ways to construct such a discrete time process. The general idea is to use an approximate time discretization of a time-continuous process known to be ergodic with respect to  $\pi$ . Generally, we cannot expect this discrete time process to be ergodic with respect to  $\pi$ . Thus one can use a Metropolis-Hastings acceptance-rejection step that ensures the detailed balance, and thus makes the chain reversible and ergodic with respect to  $\pi$ . In the case of the Euler-Maruyama discretization of the overdamped Langevin dynamics,

$$dX_t = \nabla \log \pi(X) + \sqrt{2}dW_t, \quad (1.2)$$

with  $(W_t)_{t \geq 0}$  a standard Brownian motion in  $\mathbb{R}^d$ , the method is called Metropolis Adjusted Langevin Algorithm (MALA, [22]).

---

*Key words and phrases.* Non-reversible diffusions, Langevin samplers, Markov Chain Monte Carlo, Metropolis-Hastings, variance reduction, lifting method, MCMC.

Yet, it has been noticed in several contexts that departing from the reversibility can improve the performances of MCMC methods. This article aims to propose a generalization of the standard MALA that can be able construct nonreversible Markov chains that can outperform classical MALA.

**1.1. Nonreversible dynamics.** On the continuous time setting, analysis have been carried out to compare the convergence properties of some time-continuous dynamics that are ergodic with respect to  $\pi$  [12, 13, 16, 25, 20, 21, 8], based on two kinds of optimality criterion: the speed of convergence toward equilibrium, measured in terms of spectral gap in  $L^2(\pi)$ , and asymptotic variance for the time averages. Obviously from a computational point of view an increase of the spectral gap enables to reduce the burn-in, and a reduction of the asymptotic variance leads to a decrease of the computational complexity of the corresponding MCMC method. These analysis compare, for different vector fields  $\gamma$ , the overdamped Langevin dynamics given by,

$$dX_t = \nabla \log \pi(X_t)dt + \gamma(X_t)dt + \sqrt{2}dW_t. \quad (1.3)$$

Under the condition of non explosion and that the vector field  $\gamma$  is taken such that  $\nabla \cdot (\gamma\pi) = 0$ , this dynamic is ergodic with respect to  $\pi$ . Such vector fields can be constructed easily: for any skew-symmetric matrix  $J$ , the vector field  $\gamma$  defined by  $\gamma(x) = J\nabla \log \pi(x)$  satisfies this divergence-free equation. Moreover, under this hypothesis, the following equation,

$$dX_t = \gamma(X_t)dt, \quad (1.4)$$

conserves the energy  $U$ , which justifies the Hamiltonian denomination of the term  $\gamma$ . Moreover, this dynamic is time reversible if and only if  $\gamma = 0$ , and in this case, the detailed balance is satisfied (that is to say that the generator of the diffusion (1.3) is self-adjoint in  $L^2(\pi)$ ). It is well known that among all vector fields  $\gamma$  satisfying the non explosion condition and such that  $\nabla \cdot (\gamma\pi) = 0$ , the dynamics given by (1.2) in the reversible case ( $\gamma = 0$ ) has the worse rate of convergence in terms of spectral gap in  $L^2(\pi)$  [12, 13, 16]. Recent work has been done to construct divergence free (with respect to  $\pi$ ) perturbations of the drift that achieve optimal convergence properties in the Gaussian case [16, 25]. Recent works also show that breaking the non reversibility with such divergence free perturbations on the drift also leads to improvement on the asymptotic variance. It is shown in [20] that the asymptotic variance decreases under the addition of the irreversible drift. Moreover, it has been shown in [21] that the asymptotic variance is monotonically decreasing with respect to the growth of the drift, and the limiting behavior for infinitely strong drifts is characterized. More recently, in [8] the authors investigate the dependence of the asymptotic variance on the strength of the nonreversible perturbation.

On the discrete time setting, classical methods that depart from reversible sampling consist in hybrid (Hamiltonian) MCMC [7, 17] and generalized hybrid MCMC methods [15]. In the former method, the drift direction is chosen isotropically at each time step and long time Hamiltonian integration is then carried out in this direction. The latter can be seen as a generalization of the former that brings some inertia in the direction of the Hamiltonian dynamics. Another class of nonreversible samplers is composed by lifting methods. They are designed in the discrete state space case, to construct a Markov chain that satisfies some skew detailed balance [6, 5, 11, 24]. They consist in increasing the state space to take into account a privileged drift direction that is explored more efficiently. More recently, Bierkens proposed an extension of the classical Metropolis-Hastings algorithm to generate unbiased nonreversible Markov chain [2]. This is achieved by modifying the acceptance probability to depart from detailed balance. Eventually, a recent and quite different approach has been proposed in [3] in the big data settings to circumvent the poor scalability of standard MCMC methods. The authors construct a continuous time piecewise deterministic Markov process. It is a constant velocity model where the velocity direction switches at random times with a rate depending on the target distribution.

**1.2. Outline.** We propose in this article a bias-free algorithm similar to MALA that aims to exploit the asymptotic variance reduction of the nonreversible time-continuous process. The idea is to construct a Markov chain with invariant measure  $\pi$ , by discretizing an equation of the form (1.3), instead of equation (1.2), enhanced with an acceptance-rejection step. The main difficulty consists in unbiasing the unadjusted chain. Indeed, it is not worth considering the use of a standard Metropolis-Hastings acceptance probability since it is designed to impose detailed balance with respect to the target distribution  $\pi$ , and thus to define a reversible Markov chain with respect to  $\pi$ . It would lead to a poor average acceptance ratio. An elegant way would be to construct an adequate acceptance probability with respect to this proposal, to ensure a high average acceptance ratio. In the setting of [2], it would consist in finding a good vorticity kernel. Yet, we are not able to exactly do this. Instead, we propose a class of lifted algorithms that rely on these unadjusted chains. More precisely the first algorithm is a generalized Metropolis-Hastings algorithm in an enhanced state space, and the second one can be seen as the analogous of the generalized hybrid Monte Carlo method for the overdamped Langevin equation.

In section 2 we present the first algorithm (generalized MALA). We discuss in 2.1 how its performances are closely related to the choice of the transition kernel of the unadjusted chain. In section 2.2 we prove geometric convergence of the Markov chain constructed with this algorithm under some hypotheses, that ensure the existence of a central limit theorem. Then, in section 3 we propose a modification of this algorithm (the generalized hybrid MALA). In section 4 we present numerical comparisons of these algorithms with respect to classical MALA, which is followed by concluding remarks.

## 2. GENERALIZED MALA

In this section, we construct a nonreversible Markov chain, ergodic with respect to a target distribution  $\pi$  known up to a normalizing constant. The algorithm is similar to MALA in the sense that it constructs a Markov chain from the discretization of an overdamped Langevin dynamic, augmented with an acceptance-rejection step that makes it ergodic with respect to  $\pi$ . The difference is that we construct a Markov chain on the discretization of a nonreversible Langevin equation to try to benefit from the smaller asymptotic variance of this kind of Markov processes, than the reversible ones. The main issue is then to choose a right acceptance probability that preserves the good ergodic properties of the underlying Markov process.

To state our algorithm, we slightly modify Equation (1.3). For  $\xi \in \mathbb{R}$ , we consider the diffusions,

$$dX_t = \nabla \log \pi(X_t)dt + \xi \gamma(X_t)dt + \sqrt{2}dW_t, \tag{2.1}$$

with divergence-free condition  $\nabla \cdot (\gamma\pi) = 0$ . This way,  $\xi$  specifies the direction and the intensity of the nonreversibility. We denote now by  $Q^\xi$  a proposal kernel that correspond to some discretization of the diffusions (2.1) with parameter  $\xi$ . We propose the following algorithm that we call Generalized MALA (GMALA),

**Algorithm 2.1** (Generalized MALA). *Let  $h > 0$ ,  $(x_0, \xi_0) \in \mathbb{R}^d \times \mathbb{R}$  be an initial point and an initial direction. Iterate on  $n \geq 0$ .*

- (1) *Sample  $y^{n+1}$  according to  $Q^{\xi^n}(x^n, dy)$ .*
- (2) *Accept the move with probability*

$$A^{\xi^n}(x^n, y^{n+1}) = 1 \wedge \frac{\pi(y^{n+1})Q^{-\xi^n}(y^{n+1}, x^n)}{\pi(x^n)Q^{\xi^n}(x^n, y^{n+1})}. \tag{2.2}$$

*and set  $(x^{n+1}, \xi^{n+1}) = (y^{n+1}, \xi^n)$ ; otherwise set  $(x^{n+1}, \xi^{n+1}) = (x^n, -\xi^n)$ .*

The important part of this algorithm is that the direction  $\xi^n$  of the Hamiltonian exploration must be inverted at each rejection to ensure its unbiasedness. A good choice of  $Q^\xi$  is given in the next section.

Unbiasedness is obvious since this algorithm is actually built as a Generalized Metropolis-Hastings algorithm on the increased state space  $E = \mathbb{R}^d \times \{-\xi_0, \xi_0\}$ . To simplify notations, we denote by  $x_\xi$  all element  $(x, \xi) \in E$ . We set  $S$  the involutive transformation defined for all element  $x_\xi \in E$  by  $S(x_\xi) = x_{-\xi}$ . We extend the definition of  $\pi$  on  $E$  by  $\pi(x_\xi) = \frac{\pi(x)}{2}$  for  $x_\xi \in E$ . Obviously  $\pi$  is unchanged by  $S$ . Then, the algorithm constructs a Markov chain with transition kernel density  $P$  given by,

$$P(x_\xi, y_\eta) = Q^\xi(x, y)A^\xi(x, y)\mathbb{1}_\xi(\eta) + \delta_{S(x_\xi)}(y_\eta) \left( 1 - \int_{\mathbb{R}^d} Q^\xi(x, z)A^\xi(x, z)dz \right),$$

where  $\mathbb{1}_{x_i}$  denotes the characteristic function and  $\delta_{S(x_\xi)}$  a Dirac delta function, that satisfies the following skew detailed balance,

$$\forall x_\xi, y_\eta \in E, \quad \pi(x_\xi)P(x_\xi, y_\eta) = \pi(y_\eta)P(S(y_\eta), S(x_\xi)). \quad (2.3)$$

Heuristically, we hope that the discretization of the time-continuous process (1.3) specified by  $Q^\xi$  benefits from the same good behavior in terms of asymptotic variance reduction. Since the Markov chain is constructed as parts of the discretization of the time-continuous dynamics between the rejections, then the closer to one is the average acceptance probability, the longer are these parts, and the more we can hope to benefit from this good behavior. Thus, to give a hint about the relevance of this algorithm, we compute in section 2.1 these average acceptance probabilities, and we show that they are of the same order of those for MALA for some well-chosen discretization. Moreover, we numerically show in section 4 that it can outperform MALA by several order of magnitude in terms of asymptotic variance.

Yet, before showing these results, we present a heuristic justification about this algorithm. In the reversible case, the time-continuous equation satisfies the detailed balance with respect to  $\pi$ ,

$$\forall x, y \in \mathbb{R}^d, \forall h > 0, \quad \pi(x)P_h(x, y) = \pi(y)P_h(y, x),$$

where  $P_h(x, y)$  denotes the density of the transition kernel to go from  $x$  to  $y$  after a time  $h$  for the dynamics (1.2). Moreover, MALA imposes that the Markov chain also satisfies the detailed balance with respect to  $\pi$ . In the nonreversible case, the time continuous process (2.1) satisfies a skew detailed balance as stated by the following lemma,

**Lemma 2.1.** *For all  $x, y \in \mathbb{R}^d$ , for all  $\xi \in \mathbb{R}$  and for all  $h > 0$ , the following relation holds,*

$$\pi(x)P_h^\xi(x, y) = \pi(y)P_h^{-\xi}(y, x), \quad (2.4)$$

where  $P_h^\xi(x, dy)$  is the transition probability measure of the  $h$ -skeleton of the process  $(X_t^\xi)_{t \geq 0}$ , solution of Equation (2.1).

The analogous of the reversible case would be to construct a  $\pi$ -invariant Markov chain that satisfies the same kind of skew detailed balance. Because the skew detailed balance (2.4) can be seen as a detailed balance on the enhanced state space  $E$  up to the transformation  $S$ ,

$$\forall x, y \in \mathbb{R}^d, \forall \xi \in \mathbb{R}, \forall h > 0 \quad \pi(x)P_h(x_\xi, y_\xi) = \pi(y)P_h(S(y_\xi), S(x_\xi)), \quad (2.5)$$

where  $P_h(x_\xi, y_\xi) = P_h^\xi(x, y)$ , the Generalized Metropolis-Hastings method given by algorithm 2.1 is actually the classical way to construct such a Markov chain. Nevertheless, the main difference between the time-continuous and the discrete time dynamics is the fact that the latter requires some direction switching of the nonreversible component of the dynamics. As stated previously, the idea of lifting the state space to construct a nonreversible chain has been used in the discrete state space setting. Yet, the idea is quite different. With classical lifting methods, the goal is to switch between several directions with well-chosen probabilities, to quickly explore the state space in all of these directions. In our case, we do not aim to switch between several nonreversible directions. Yet, we are forced to do so at each

rejection, and we have no choice but to reverse the current nonreversible directions. That is to say that we control neither the probability of switching nor the direction.

*Proof of Lemma 2.1.* We denote by  $\mathcal{L}^\xi$  the generator of diffusion (2.1). One can show that,

$$\mathcal{L}^\xi = \mathcal{S} + \xi \mathcal{A},$$

where  $\mathcal{S} = \nabla \log \pi(x) \nabla \cdot$ , and  $\mathcal{A} = \gamma \cdot \nabla$ . Then, one can show that  $\mathcal{S}$  and  $\mathcal{A}$  are respectively the symmetric and the antisymmetric parts of  $\mathcal{L}^1$ , with respect to  $L^2(\pi)$ . Moreover, one can show that for all  $\xi \in \mathbb{R}$ ,  $\mathcal{L}^\xi$  and  $\mathcal{L}^{-\xi}$  are adjoint from one another. This amounts to show that  $D((\mathcal{L}^{-\xi})^*) \subset D(\mathcal{L}^\xi)$ , which follows from the injectivity of the operator  $(\mathcal{L}^{-\xi})^* - iI$ , and the surjectivity of the operator  $\mathcal{L}^\xi - iI$ . It follows that the semigroups  $e^{\mathcal{L}^\xi t}$  and  $e^{\mathcal{L}^{-\xi} t}$  are adjoint from one another (Corollary 10.6 [18]). Then, for all bounded functions  $f$  and  $g$ ,

$$\begin{aligned} \int_{(\mathbb{R}^d)^2} f(x)g(y)P_h^\xi(x, dy)\pi(dx) &= \int_{\mathbb{R}^d} f(x)e^{\mathcal{L}^\xi h}g(x)\pi(dx), \\ &= \int_{\mathbb{R}^d} e^{(\mathcal{L}^\xi)^* h}f(x)g(x)\pi(dx), \\ &= \int_{\mathbb{R}^d} e^{\mathcal{L}^{-\xi} h}f(x)g(x)\pi(dx), \\ &= \int_{(\mathbb{R}^d)^2} f(y)g(x)P_h^{-\xi}(x, dy)\pi(dx). \end{aligned}$$

□

**2.1. Choice of the proposition kernel.** The method presented above can be tuned with the choice of the proposition kernel  $Q^\xi$ . It should be chosen such that it approximates the law of the transition probability of the  $h$ -skeleton of the process  $(X_t^\xi)_{t \geq 0}$  solution of equation (2.1). The basic idea is to define  $Q^\xi$  as the density of the law of an approximate discretization of equation (2.1). Doing so, we can hope that the skew-detailed balance for the unadjusted chain would be almost satisfied (since it is satisfied for the time-continuous dynamics), that is to say we can hope to benefit from an acceptance ratio close to one. Moreover, we recall that the choice of the proposal kernel is the only degree of freedom that enables to tune the rate of the direction switching. Formally, the closer the proposal kernel is from the transition kernel of the  $h$ -skeleton of the continuous process, the higher the acceptance ratio is (and thus the less frequent is the direction switching).

The basic idea would be to propose according to the Maruyama-Euler approximation of equation (2.1), as it is done for MALA. We denote by  $Q_1^\xi$  this kernel. That is to say  $Q_1^\xi(x, dy)$  is given by the law of  $y$ , solution of

$$y = x - h\nabla U(x) - h\xi\gamma(x) + \sqrt{2h}\chi, \quad (2.6)$$

where  $\chi$  a standard normal deviate. Then,  $Q_1^\xi$  is given by

$$Q_1^\xi(x, dy) = \frac{1}{(4\pi h)^{d/2}} \exp\left(-\frac{1}{4h}\|y - (x - h\nabla U(x) - h\xi\gamma(x))\|^2\right) dy. \quad (2.7)$$

Sadly, even though the simplicity of this proposal is appealing, this proposition kernel leads to an average rejection rate of order  $h$  when  $\xi \neq 0$  and  $\gamma$  is non-linear as stated in Proposition 2. It is significantly worse than MALA that enjoys an average rejection rate of order  $h^{3/2}$  (see [4]). As shown later by the numerical simulations, this bad rejection rate is not a pure theoretical problem: it forbids to use large discretization steps  $h$ , and thus the method only generates highly correlated samples.

To overcome this problem, we propose to implicit the resolution of the nonreversible term in equation (2.1) with a centered point discretization. More precisely, we propose a move  $y_\xi$  from  $x_\xi$ , such that  $y$  would be distributed according to the solution of

$$\Phi_x^{h\xi}(y) = x - h\nabla U(x) + \sqrt{2h}\chi, \quad (2.8)$$

where  $\chi$  a standard normal deviate, and  $\Phi_x^{h\xi}$  is the function defined by,

$$\Phi_x^{h\xi} : \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad y \mapsto y + h\xi\gamma \left( \frac{x+y}{2} \right). \quad (2.9)$$

Existence of such a proposal kernel, denoted by  $Q_2^\xi$ , is ensured under the hypothesis that  $U$  is twice differentiable with uniformly bounded second derivative.

**Assumption 1.** *We suppose that  $U$  is twice differentiable with uniformly bounded second derivative.*

**Proposition 1.** *Under Assumption 1, there exists  $h_0 > 0$ , such that for all  $h < h_0$ , for all  $x_\xi \in E$ , for all  $\chi \in \mathbb{R}^d$ , there exists a unique  $y \in \mathbb{R}^d$ , solution of equation (2.8), and if  $\chi$  denotes a standard normal deviate, the function  $Y_\xi$  defined almost surely as the solution of equation (2.8) is a well-defined random variable, with law  $Q_2^\xi(x, dy)$ , given by,*

$$Q_2^\xi(x, dy) = \frac{1}{(4\pi h)^{d/2}} |\text{Jac } \phi_x^{\xi h}(y)| \exp \left( \frac{1}{4h} \|\phi_x^{h\xi}(y) - (x - h\nabla \log \pi(x))\|^2 \right) dy. \quad (2.10)$$

Moreover

$$x - y = \sqrt{2h}\chi + O \left( h \|\nabla U(x)\| + h^{3/2} \|\chi\| \right),$$

and there exists  $C > 0$ , independent of  $x$ ,  $\chi$ , and  $h$ , such that,

$$\|\nabla U(y)\| \leq (1 + Ch) \|\nabla U(x)\| + C\sqrt{2h} \|\chi\|.$$

*Proof of Proposition 1.* The first point is a corollary of Lemma 2.2 below. The second point can be proven by making use of the fact that  $\nabla U$  is Lipschitz.  $\square$

The following lemma is used to prove Proposition 1. We can note that the  $h_0$  given by this lemma depends on the Lipschitz coefficient of  $\nabla U$ .

**Lemma 2.2.** *Under Assumption 1, there exists  $h_0 > 0$ , such that for all  $h < h_0$ , for all  $x_\xi \in E$ , the function  $\Phi_x^{h\xi} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , defined by (2.9) is a  $C^1$ -diffeomorphism.*

*Proof of Lemma 2.2.* The proof that  $\Phi_x^{h\xi}$  is bijective relies on Picard fixed point theorem. For any fixed  $z \in \mathbb{R}^d$ , we define  $\Psi$  on  $\mathbb{R}^d$  by  $\Psi(y) = z - h\xi J\nabla U \left( \frac{x+y}{2} \right)$ . Then, since  $\nabla U$  is supposed to be Lipschitz, then for small enough  $h$  this application is a contraction mapping from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ . The fact that  $\Phi_x^{h\xi}$  is everywhere differentiable is clear, and the fact that its inverse is everywhere differentiable as well comes from the fact that the determinant of the Jacobian of  $\Phi_x^{h\xi}$  is strictly positive for sufficiently small  $h$ .  $\square$

Sampling from the this proposal kernel can be done by a fixed point method when no analytical solution is available since the proof of Lemma 2.2 uses a Picard fixed point argument. This transition kernel leads to a rejection rate of order  $h^{3/2}$  (see Proposition 2), which is an improvement from  $Q_1^\xi$ . The main advantage of this kernel is that the computation of  $|\text{Jac } \phi_x^{\xi h}(y)|$  can be avoided in the case where  $\gamma$  is defined by  $\gamma(x) = J\nabla \log \pi(x)$ , with  $J$  a skew symmetric matrix. Indeed, only the ratio  $|\text{Jac } \phi_x^{\xi h}(y)| / |\text{Jac } \phi_y^{-\xi h}(x)|$  is required to compute the acceptance probability  $A(x, y)$ , and this ratio is equal to 1 if  $\gamma$  is of gradient type, as stated by the following lemma.

**Lemma 2.3.** For all  $x, y \in \mathbb{R}^d$ , and for all  $h > 0$ ,

$$|\text{Jac } \phi_x^{\xi h}(y)| = |\text{Jac } \phi_y^{-\xi h}(x)|$$

*Proof of lemma 2.3.* This result uses the more general fact that for any skew-symmetric matrix  $A$  and any symmetric matrix  $S$ , the matrices  $Id + AS$  and  $Id - AS$  have same determinant. This statement is equivalent to say that  $\chi_{AS} = \chi_{-AS}$ , where we denote by  $\chi_M$  the characteristic polynomial of any square matrix  $M$ . This last statement is true since for any square matrices  $A$  and  $S$  (of the same size)  $\chi_{AS} = \chi_{SA}$ . Then using the transposition,  $\chi_{AS} = \chi_{A^t S^t}$ . Eventually using the fact that  $A^t = -A$  and  $S^t = S$ , it comes  $\chi_{AS} = \chi_{-AS}$ . The result follows from the fact that the matrix  $\text{Jac } \phi_x^{\xi h}(y)$  is of the form  $Id + \xi AS$ .  $\square$

We denote by  $\alpha_{h,\xi}^1$  and  $\alpha_{h,\xi}^2$  respectively the acceptance probability for proposal kernels  $Q_1^\xi$  and  $Q_2^\xi$  (defined respectively by (2.7) and (2.10)). The following proposition provides an upper bound on the moments of the rejection probability.

**Proposition 2.** Suppose that  $U$  is three times differentiable with bounded second and third derivatives. Then for all  $l \geq 1$ , there exists  $C(l) > 0$  and  $h_0 > 0$  such that for all positive  $h < h_0$ , and for all  $x \in \mathbb{R}^d$ ,

$$\begin{aligned} \mathbb{E} \left[ (1 - \alpha_{h,\xi}^1(x, Y_{h,\xi}^1))^{2l} \right] &\leq C(l)(1 + \|\nabla U(x)\|^{4l})h^{2l} \\ \mathbb{E} \left[ (1 - \alpha_{h,\xi}^2(x, Y_{h,\xi}^2))^{2l} \right] &\leq C(l)(1 + \|\nabla U(x)\|^{4l})h^{3l} \end{aligned}$$

*Proof of Proposition 2.* To deal with both results at once, we define for all  $x \in \mathbb{R}^d$ , for all  $\xi \in \{-1, 1\}$ , and for all  $\theta \in \{0, 1\}$ , the random variable  $Y_{x_\xi, \theta}$  that satisfies the following implicit equation almost surely,

$$Y_{x_\xi, \theta} = x - h\nabla U(x) - h\xi J \left( \theta \nabla U \left( \frac{x + Y_{x_\xi, \theta}}{2} \right) + (1 - \theta)\nabla U(x) \right) + \sqrt{2h}\chi,$$

where  $\chi$  is a standard normal deviate in  $\mathbb{R}^d$ . Well-posedness of  $Y_{x_\xi, \theta}$  is given by Proposition 1. We set  $R_\theta^\xi$  the associated proposal kernel. We get  $R_0^\xi = Q_1^\xi$  and  $R_1^\xi = Q_2^\xi$ . We define the Metropolis-Hastings ratio  $r(x_\xi, y_\xi)$  for proposing  $y_\xi$  from  $x_\xi$  with kernel  $R_\theta^\xi$  by,

$$r(x_\xi, y_\xi) = \frac{\pi(y_\xi)R_\theta(y_\xi, x_\xi)}{\pi(x_\xi)R_\theta(x_\xi, y_\xi)},$$

where we set  $\chi \in \mathbb{R}^d$  such that,

$$y = x - h\nabla U(x) - h\xi J \left( \theta \nabla U \left( \frac{x + y}{2} \right) + (1 - \theta)\nabla U(x) \right) + \sqrt{2h}\chi,$$



Then, a straightforward computation gives,

$$\begin{aligned}
\log(r(x_\xi, y_\xi)) &= U(x) - U(y) + \langle y - x, \nabla U(x) \rangle \\
&+ \frac{1}{2} \langle y - x, \nabla U(y) - \nabla U(x) \rangle \\
&- \frac{\xi}{2} (1 - \theta) \langle y - x, J(\nabla U(y) - \nabla U(x)) \rangle \\
&+ \frac{h}{4} \left( \|\nabla U(x)\|^2 - \|\nabla U(y)\|^2 \right) \\
&+ \frac{h}{2} \xi \theta \langle J \nabla U \left( \frac{x+y}{2} \right), \nabla U(x) + \nabla U(y) \rangle \\
&+ \frac{h}{2} \xi^2 \theta (1 - \theta) \langle J \nabla U \left( \frac{x+y}{2} \right), J(\nabla U(x) - \nabla U(y)) \rangle \\
&+ \frac{h}{4} (1 - \theta)^2 \xi^2 \left( \|J \nabla U(x)\|^2 - \|J \nabla U(y)\|^2 \right).
\end{aligned}$$

A Taylor expansion of these terms, making use of Proposition 1 and noticing that  $\theta(1 - \theta) = 0$  for  $\theta \in \{0, 1\}$  leads to

$$\begin{aligned}
\log(r(x_\xi, y_\xi)) &= -\xi h(1 - \theta) \langle \chi, J D^2 U(x) \cdot \chi \rangle \\
&+ O(h^{3/2} (\|\nabla U(x)\|^2 + \|\chi\|^2 + \|\chi\|^3)) \\
&= O(h(1 - \theta) \|\chi\|^2 + h^{3/2} (\|\nabla U(x)\|^2 + \|\chi\|^2 + \|\chi\|^3)).
\end{aligned}$$

Moreover,

$$\begin{aligned}
(1 - 1 \wedge r(x_\xi, y_\xi))^{2l} &= O(\log(r(x_\xi, y_\xi))^{2l}) \\
&= O(h^{2l} (1 - \theta)^{2l} \|\chi\|^{4l} + h^{3l} (\|\nabla U(x)\|^{4l} + \|\chi\|^{4l} + \|\chi\|^{6l})),
\end{aligned}$$

and thus there exists  $C(l) > 0$  such that,

$$\mathbb{E} \left[ (1 - 1 \wedge r(x_\xi, Y_{x_\xi, \theta}))^{2l} \right] \leq C(l) (1 + \|\nabla U(x)\|^{4l}) h^{3l} + C(l) h^{2l} (1 - \theta)^{2l}$$

□

The previous method is quite efficient since no computation of the Hessian of  $\log \pi$  is required. Nevertheless, global Lipschitzness of  $\nabla U$  is required to justify the method. We propose a last kernel, denoted by  $Q_3^\xi$ , that does not require this hypothesis but still require  $\gamma$  to be of gradient type and the computation of the Hessian of  $U$ . More precisely,  $Q_3^\xi(x, dy)$  is the law of  $y$ , solution of

$$\left( Id + \frac{h\xi}{2} J \text{Hess}(U)(x) \right) (y - x) = -h(Id + \xi J) \nabla U(x) + \sqrt{2h} \chi, \quad (2.11)$$

where  $\chi$  a standard normal deviate. Then,  $Q_3^\xi$  is given by

$$Q_3^\xi(x, dy) = \frac{1}{(4\pi h)^{d/2}} \det M^\xi(x) \exp \left( \frac{1}{4h} \|M^\xi(x)(y - x) + h(Id + \xi J) \nabla U(x)\|^2 \right) dy, \quad (2.12)$$

with  $M^\xi(x) = Id + \frac{h\xi}{2} J \text{Hess} U(x)$ . This transition kernel offers also a rejection rate  $a_{h,\xi}^3$  of order  $h^{3/2}$ .

**Proposition 3.** *Suppose that  $U$  is three times differentiable with bounded second and third derivatives. Then for all  $l \geq 1$ , there exists  $C(l) > 0$  and  $h_0 > 0$  such that for all positive  $h < h_0$ , and for all  $x \in \mathbb{R}^d$ ,*

$$\mathbb{E} \left[ (1 - \alpha_{h,\xi}^3(x, Y_{h,\xi}^2))^{2l} \right] \leq C(l) (1 + \|\nabla U(x)\|^{4l}) h^{3l}$$

*Proof of Proposition 3.* The proof involves the same arguments as in Proposition 2 and is left to the reader.  $\square$

Even though Kernel  $Q_2^\xi$  enables to circumvent the bad acceptance ratio of the explicit proposal given by  $Q_1^\xi$ , it raises some difficulties. First, it requires the potential  $U$  to be globally Lipschitz, which is quite restrictive. To use the GMALA method in the non globally Lipschitz case, one can resort to importance sampling (to get back to a globally Lipschitz setting), or use a globally Lipschitz truncation of the potential to build the proposition kernel. The same idea is used in MALTA [22]. This last trick might lead to high average rejection ratio in the truncated regions of the potential. Moreover, the computation of the proposed move requires the resolution of a nonlinearly implicit equation, that can be solved by a fixed point method, which is more costly than the Euler-Maruyama discretization used by MALA. Specific methods of preconditioning should be considered to accelerate the convergence of the fixed point.

It would be interesting to propose a proposition kernel based on an explicit scheme, that would achieve a better acceptance ratio than  $Q_1^\xi$ , with a lower computational cost than  $Q_2^\xi$ . The question of decreasing the Metropolis-Hastings rejection rate has been recently studied in [9]. The authors propose a proposition kernel constructed from an explicit scheme, which is a correction (at order  $h^{3/2}$ ) of the standard Euler-Maruyama proposal. Nevertheless, in our case this approach would not enable us to construct a proposition kernel based on an explicit scheme with a high average acceptance ratio. Work has also been done in this direction with the Metropolis-Hastings algorithm with delayed rejection, proposed in [23], and more recently [1]. Yet, it is unclear whether this approach could be use efficiently in our case.

**2.2. Convergence of GMALA.** In this section, we only treat the case in which GMALA is used with  $Q_2^\xi$ . To obtain a central limit theorem for the Markov chain built with GMALA, it is convenient to prove the geometric convergence in total variation norm toward the target measure  $\pi$ . Because we require  $\nabla U$  to be globally Lipschitz to ensure well-posedness of GMALA with transition kernel  $Q_2^\xi$ , MALA would be likely to benefit from geometric convergence in this case. Indeed, it is well-known that MALA can be geometrically ergodic when the tails of  $\pi$  are heavier than Gaussian, under suitable hypotheses (see Theorem 4.1 [22]). For strictly lighter tails, we cannot hope for geometric convergence (see Theorem 4.2 [22]). This part is devoted to showing that under some hypotheses, GMALA can benefit from geometric convergence.

Classically for MALA, the convergence follows on from the aperiodicity and the irreducibility of the chain, since by construction the chain is positive with invariant measure  $\pi$ . Under these conditions, the geometric convergence can be proven by exhibiting a suitable Foster-Lyapunov function  $V$  such that,

$$\lim_{|x| \rightarrow \infty} \frac{PV(x)}{V(x)} < 1. \quad (2.13)$$

In our case, the situation is slightly different. The Markov chain built with GMALA is still aperiodic and phi-irreducible. This is a simple consequence of the surjectivity of the application  $\Phi_x^{h\xi}$  in Lemma 2.2. To be able to prove a drift condition such as (2.13), it is usually required to be able to show that the acceptance probability for a proposed move starting from  $x_\xi$  does not vanish in expectation for large  $x$ , which is not always true in our setting. More precisely, we can only say that the maximum of the two average acceptance probabilities for the proposed moves starting from  $x_\xi$  and  $x_{-\xi}$  does not vanishes for large  $x$ . Then one strategy could be to choose a Foster-Lyapunov function  $V$  that decreases in expectation when  $\xi$  is changed to  $-\xi$  after a rejection. An other strategy, that we present in the following and that seems to be more natural and more easily generalizable to potential  $U$  that does not satisfy the specific hypotheses we use in this section. We show that the odd and the even subsequences of the Markov chain

converge geometrically quickly to the target measure  $\pi$  by showing a drift condition on  $P^2$ : the transition probability kernel of the two-steps Markov chain defined by

$$P^2(x, dy) = \int_E P(x, dz)P(z, dy), \quad \forall x \in E,$$

under the following assumption.

**Assumption 2.** *We suppose that  $U$  is three times differentiable, and that,*

- (1) *eigenvalues of  $D^2U(x)$  are uniformly upper bounded and lower bounded away from 0, for  $x$  outside of a ball centered in 0,*
- (2) *the product  $\|D^3U(x)\| \|\nabla U(x)\|$  is bounded for  $x$  outside of a ball centered in 0.*

**Proposition 4.** *Suppose assumption 2. There exists  $s > 0$  and  $h_0 > 0$  such that for all  $h \leq h_0$ , for  $\xi \in \{-1, 1\}$ ,*

$$\lim_{\|x\| \rightarrow +\infty} \frac{P^2 V_s(x_\xi)}{V_s(x_\xi)} = 0, \quad (2.14)$$

where the Foster-Lyapunov function  $V$  is defined by  $V_s(x_\xi) = \exp(sU(x))$ .

**Remark 2.1.** *In order to prove the drift condition (2.13), one can use the Foster-Lyapunov  $\tilde{V}$  defined by,*

$$\tilde{V}_s(x_\xi) = \exp\left(sU(x) + s \frac{\xi h^2 \langle \nabla U(x), D^2U(x) J \nabla U(x) \rangle}{|\xi h^2 \langle \nabla U(x), D^2U(x) J \nabla U(x) \rangle|}\right),$$

for  $s$  small enough. Yet, this proof is left to the reader, but uses the arguments developed in the following.

Assumption 2 is not meant to be sharp. We are not striving for optimality here, but instead we aim to propose a simple criterion which is likely to be satisfied in smooth cases.

The first step to prove equation (2.14), is to show that the proposed move decreases the Lyapunov function  $V_s$  in expectation.

**Lemma 2.4.** *Under assumption 1, there exists  $s_0 > 0$  such that for all  $s < s_0$ , there exists  $C_1, C_2 > 0$  such that for all  $h < h_0$  given by Proposition 1, and for all  $x \in \mathbb{R}^d$  large enough (depending on  $h$ ),*

$$\mathbb{E}[V_s(Y_{h,\xi})] \leq V_s(x) e^{-s \frac{3h(1-C_1h)}{4} \|\nabla U(x)\|^2 + sC_2},$$

where  $Y_{h,\xi}$  is defined in Proposition 1. Thus, for  $h$  small enough,

$$\mathbb{E}[V_s(Y_{h,\xi})] \leq V_s(x) e^{-s \frac{h}{2} \|\nabla U(x)\|^2 + sC_2},$$

and for  $x$  large enough,

$$\mathbb{E}[V_s(Y_{h,\xi})] \leq V_s(x) e^{-s \frac{h}{4} \|\nabla U(x)\|^2}.$$

*Proof of Lemma 2.4.* For all  $y \in \mathbb{R}^d$ , we set  $\chi \in \mathbb{R}^d$  such that,

$$y = x - h \nabla U(x) - h \xi J \left( \theta \nabla U \left( \frac{x+y}{2} \right) + (1-\theta) \nabla U(x) \right) + \sqrt{2h} \chi,$$

A Taylor expansion of  $y$  in  $h$  yields,

$$\begin{aligned} U(y) &= U(x) + \nabla U(x) \cdot (y-x) + O(\|y-x\|^2) \\ &= U(x) - h \|\nabla U(x)\|^2 + \sqrt{2h} \nabla U(x) \cdot \chi + h \xi \nabla U(x) \cdot J \nabla U \left( \frac{x+y}{2} \right) + O(\|x-y\|^2). \end{aligned}$$

Noticing that

$$h\xi\nabla U(x) \cdot J\nabla U\left(\frac{x+y}{2}\right) = h\xi\nabla U(x) \cdot J(\nabla U\left(\frac{x+y}{2}\right) - \nabla U(x)),$$

the Cauchy-Schwarz and the triangular inequality yield,

$$\left| h\xi\nabla U(x) \cdot J\nabla U\left(\frac{x+y}{2}\right) \right| \leq Ch \|\nabla U(x)\| \|x-y\|.$$

Thus, by Young inequality and Proposition 1, for  $h \leq 1$ ,

$$U(y) \leq U(x) - \frac{3h(1-Ch)}{4} \|\nabla U(x)\|^2 + O(\|x\|^2),$$

The conclusion holds for small enough  $s$  such that  $e^{sO(\|G\|^2)}$  is integrable, where  $G$  is a standard normal deviate in  $\mathbb{R}^d$ .  $\square$

Classically ([22]), proofs of geometric convergence of MALA require to show that the average acceptance ratio of the proposed move from  $x$ , does not vanishes when  $x$  is large. Namely that there exists  $\varepsilon > 0$  such that

$$I(x) = \{y : \alpha(x, y) \leq 0\}$$

asymptotically has  $q$ -measure 0, where  $\alpha$  is the acceptance probability. Our case is slightly different since it is possible that the average acceptance ratio of the proposed move from  $x_\xi \in E$  vanishes when  $\|x\| \rightarrow +\infty$ , which leads to a rejection and to the switching  $\xi \leftarrow -\xi$  with probability close to one. Yet the average acceptance probability of the next proposed move from  $x_{-\xi}$  is close to one. This behavior is described by the following Lemma.

**Lemma 2.5.** *Under Assumption 2, there exists  $h_0 > 0$  such that for all  $h < h_0$  and for all  $\varepsilon > 0$ , there exists  $C(h, \varepsilon) > 0$  such that for all  $x_\xi \in E$  such that  $\|x\| \geq C(h, \varepsilon)$ ,*

$$\xi \langle \nabla U(x), D^2U(x) \cdot J\nabla U(x) \rangle \geq 0 \implies \mathbb{P}(\alpha_{h,\xi}^2(x_\xi, y_\xi) = 1) \geq 1 - \varepsilon,$$

where  $y$  is defined by Equation (2.8).

*Proof of Lemma 2.5.* We recall that  $\alpha_{h,\xi}^2$  is defined for all  $x_\xi, y_\xi \in E$  by  $\alpha_{h,\xi}^2(x_\xi, y_\xi) = 1 \wedge e^{r(x_\xi, y_\xi)}$ , with,

$$\begin{aligned} r(x_\xi, y_\xi) &= U(x) - U(y) + \langle \sqrt{2h}\chi, \nabla U(x) \rangle + \frac{1}{2} \langle \sqrt{2h}\chi, \nabla U(y) - \nabla U(x) \rangle \\ &\quad - h \|\nabla U(x)\|^2 - h \langle \nabla U(y) - \nabla U(x), \nabla U(x) \rangle - \frac{h}{4} \|\nabla U(x) - \nabla U(y)\|^2. \end{aligned}$$

The proof is done by computing a Taylor expansion in  $h$  of this quantity to evaluate its sign in the asymptotic case  $\|x\| \rightarrow +\infty$ . We denote by  $\cdot$  the matrix vector product, and by  $:$  and  $\dot{\cdot}$  respectively the

double and triple dot products.

$$\begin{aligned}
& U(x) - U(y) + \langle \sqrt{2h}\chi, \nabla U(x) \rangle - h \|\nabla U(x)\|^2 \\
&= h\xi \nabla U(x) \cdot J \left( \nabla U \left( \frac{x+y}{2} \right) - \nabla U(x) \right) - \frac{1}{2} D^2 U(x) : (y-x)^2 + O(\|D^3 U(x) : (y-x)^3\|) \\
&= \frac{h}{2} \xi \nabla U(x) \cdot J (D^2 U(x) \cdot (y-x) + O(D^3 U(x) : (y-x)^2)) \\
&\quad - \frac{1}{2} D^2 U(x) : (y-x)^2 + O(\|D^3 U(x) : (y-x)^3\|) \\
& \frac{1}{2} \langle \sqrt{2h}\chi, \nabla U(y) - \nabla U(x) \rangle = \frac{1}{2} \langle \sqrt{2h}\chi, D^2 U(x) \cdot (y-x) + O(D^3 U(x) : (y-x)^2) \rangle, \\
&\quad - h \langle \nabla U(y) - \nabla U(x), \nabla U(x) \rangle = -h \langle D^2 U(x) \cdot (y-x) + O(D^3 U(x) : (y-x)^2), \nabla U(x) \rangle \\
&\quad - \frac{h}{4} \|\nabla U(x) - \nabla U(y)\|^2 = hO(\|x-y\|^2).
\end{aligned}$$

Collecting these terms, and using Proposition 1 and Assumption 2, it comes for  $h \leq 1$

$$r(x_\xi, y_\xi) = -\frac{h}{2} \langle \nabla U(x), D^2 U(x) \cdot (y-x) \rangle + O(h^3 \|\nabla U(x)\|^2) + O(\|\chi\|^2 + \|\chi\|^3),$$

where the big  $O$  are independent of  $h$ . and eventually,

$$\begin{aligned}
r(x_\xi, y_\xi) &= \frac{h^2}{2} (\langle \nabla U(x), D^2 U(x) \cdot \nabla U(x) \rangle + \xi \langle \nabla U(x), D^2 U(x) \cdot J \nabla U(x) \rangle) \\
&\quad + O(h^3 \|\nabla U(x)\|^2) + O(\|\chi\|^2 + \|\chi\|^3),
\end{aligned}$$

Then, using the definite positivity of  $D^2 U(x)$  and the fact that  $\|\nabla U(x)\|$  is coercive, for all  $h$  small enough and for all  $\varepsilon > 0$ , there exists  $C(h, \varepsilon) > 0$  such that for all  $x \in \mathbb{R}^d$  such that  $\|x\| \geq C(h, \varepsilon)$ ,

$$\xi \langle \nabla U(x), D^2 U(x) \cdot J \nabla U(x) \rangle \geq 0 \implies \mathbb{P}(r(x_\xi, y_\xi) \geq 0) \geq 1 - \varepsilon.$$

□

Proposition 4 can now be obtained as a consequence of Lemmas 2.5 and 2.4.

*Proof of proposition 4.* Because of the acceptance-rejection step, we get for all  $x_\xi \in E$ ,

$$\begin{aligned}
P^2 V(x) &= \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) \int_{\mathbb{R}^d} q(y_\xi, z_\xi) \alpha_{h,\xi}^2(y_\xi, z_\xi) V(z_\xi) dz dy \\
&\quad + \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) V(y_{-\xi}) \left( 1 - \int_{\mathbb{R}^d} q(y_\xi, z_\xi) \alpha_{h,\xi}^2(y_\xi, z_\xi) dz \right) dy \\
&\quad + \left( 1 - \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) dy \right) \left( \int_{\mathbb{R}^d} q(x_{-\xi}, y_{-\xi}) \alpha_{h,\xi}^2(x_{-\xi}, y_{-\xi}) V(y_{-\xi}) dy \right) \\
&\quad + V(x_{-\xi}) \left( 1 - \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) dy \right) \left( 1 - \int_{\mathbb{R}^d} q(x_{-\xi}, y_{-\xi}) \alpha_{h,\xi}^2(x_{-\xi}, y_{-\xi}) dy \right).
\end{aligned}$$

Simply using  $\alpha_{h,\xi}^2 \leq 1$  leads to,

$$\begin{aligned}
 P^2V(x) &\leq \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \int_{\mathbb{R}^d} q(y_\xi, z_\xi) V(z_\xi) dz dy \\
 &\quad + \int_{\mathbb{R}^d} q(x_\xi, y_\xi) V(y_{-\xi}) \left( 1 - \int_{\mathbb{R}^d} q(y_\xi, z_\xi) \alpha_{h,\xi}^2(y_\xi, z_\xi) dz \right) dy \\
 &\quad + \left( 1 - \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) dy \right) \left( \int_{\mathbb{R}^d} q(x_{-\xi}, y_{-\xi}) V(y_{-\xi}) dy \right) \\
 &\quad + V(x_{-\xi}) \left( 1 - \int_{\mathbb{R}^d} q(x_\xi, y_\xi) \alpha_{h,\xi}^2(x_\xi, y_\xi) dy \right) \left( 1 - \int_{\mathbb{R}^d} q(x_{-\xi}, y_{-\xi}) \alpha_{h,\xi}^2(x_{-\xi}, y_{-\xi}) dy \right).
 \end{aligned}$$

And eventually Lemmas 2.5 and 2.4 gives,

$$\begin{aligned}
 \frac{P^2V(x)}{V(x)} &\leq e^{-s\frac{h}{2}\|\nabla U(x)\|^2 + 2sC_2} + e^{-s\frac{h}{2}\|\nabla U(x)\|^2 + sC_2} \\
 &\quad + (1 - \mathbb{E}[\alpha_{h,\xi}^2(x_\xi, Y_{h,\xi})]) e^{-s\frac{h}{2}\|\nabla U(x)\|^2 + sC_2} \\
 &\quad + (1 - \mathbb{E}[\alpha_{h,\xi}^2(x_\xi, Y_{h,\xi})]) (1 - \mathbb{E}[\alpha_{h,\xi}^2(x_{-\xi}, Y_{h,-\xi})]).
 \end{aligned}$$

And thus,

$$\lim_{\|x\| \rightarrow +\infty} \frac{P^2V(x)}{V(x)} = 0$$

□

### 3. GENERALIZED HYBRID MALA

We propose in this part a second algorithm. It is based on a splitting method by solving by turns the reversible and the non reversible parts of equation (2.1) with measure preserving schemes. The idea is then to integrate the pure reversible equation

$$dx = -\nabla U(x)dt + dW_t. \tag{3.1}$$

by using MALA, and to integrate the pure Hamiltonian equation

$$dx = -\xi J \nabla U(x)dt. \tag{3.2}$$

by an hybrid Monte-Carlo method based on a suitable Hamiltonian integrator. This method presents some theoretical advantages on Generalized MALA. Before presenting them, we begin with explaining the algorithm.

We denote by  $\Psi_t^\xi$  the flow of the Hamiltonian equation (3.2) on the time interval  $[0, t]$ , and by  $\Phi_h^\xi$  a numerical integrator for (3.2) on a time step  $h$ . We precise later necessary conditions on this integrator that ensure unbiasedness of the algorithm.

**Algorithm 3.1 (Generalized Hybrid MALA).** *Let  $x_0$  be an initial point. Set  $\xi = \pm 1$ . Let  $h > 0$ . Iterate on  $n \geq 0$ ,*

- (1) *Integration of the reversible part (3.1):*  
 MALA is used to sample  $x_{n+1/2}$  from  $x_n$ , with time-step  $h$ .
- (2) *Integration of the non reversible part (3.2):*
  - (a) *Compute  $\tilde{x}_{n+1} = \Phi_h^\xi(x_{n+1/2})$ .*

(b) Set  $x_{n+1} = \tilde{x}_{n+1}$  with probability

$$\beta_{h,\xi}(x_{n+1/2}) = \min(1, \exp(U(x_{n+1/2}) - U(\tilde{x}_{n+1}))) \min\left(1, \exp(U(x_{n+1/2}) - U(\Phi_h^\xi(x_{n+1/2})))\right).$$

Otherwise set  $x_{n+1} = x_{n+1/2}$  and  $\xi \leftarrow -\xi$ .

Similarly to the Hybrid Monte-Carlo algorithm, the first step enables to explore the state space across the iso-potential lines, whereas the second step enables to explore it along the iso-potential lines.

To ensure unbiasedness of the second step, the integrator  $\Phi_h^\xi$  must satisfy the following properties,

$$\Phi_h^\xi = (\Phi_h^{-\xi})^{-1}, \quad (3.3)$$

$$\det \text{Jac } \Phi_h^\xi = 1. \quad (3.4)$$

These properties are classical for hybrid Monte-Carlo methods (see Chapter 2. [17]).

**Lemma 3.1.** *Under conditions (3.3) and (3.4), the second step leaves the measure  $\pi$  invariant.*

*Proof of Lemma 3.1.* The proof can be found in [17]. It consists in seeing this step as a generalized Metropolis-Hastings step, with proposal  $Q(x_\xi, y_\eta) = \delta_{\Phi_h^\xi(x)} \delta_\xi(\eta)$ , and symmetric operator  $S(x_\xi) = x_{-\xi}$ . Then, it is enough to show that the Metropolis-Hastings ratio  $r$  defined by the following Radon-Nikodym derivative

$$r(x_\xi, y_\eta) = \frac{Q(S(y_\eta), S(dx_\xi))\pi(dy_\eta)}{Q(x_\xi, dy_\eta)\pi(dx_\xi)},$$

is equal to  $\exp(\log \pi(y) - \log \pi(x))$ . □

For example, for all  $x_\xi \in E$  the centered point integrator  $\Phi_h^\xi(x)$  defined by the solution  $y$  of the following equation

$$y = x - h\xi J\nabla \log \pi \left( \frac{x+y}{2} \right), \quad (3.5)$$

satisfies these properties, under some assumptions that ensure well-posedness.

**Lemma 3.2.** *Under Assumption 1, there exists  $h_0 > 0$  such that for all positive  $h < h_0$ , there exists a unique solution to equation 3.5, and the integrator  $\Phi_h^\xi$  is well-defined and satisfies equations (3.3) and (3.4).*

*Proof of Lemma 3.2.* The proof follows on from Lemma 2.2. □

The main benefit of this algorithm with respect to GMALA, is the better average acceptance ratio of the Hybrid step with the centered point integrator, which is of order  $O(h^3)$  instead of  $O(h^{3/2})$  for GMALA, which may enable to reduce the rate of switching directions.

**Lemma 3.3.** *Suppose Assumption 2. Then, for all  $l \geq 1$ , there exists  $h_0 > 0$  such that for all positive  $h < h_0$ , and for all  $x \in \mathbb{R}^d$ ,*

$$(1 - \beta_{h,\xi}(x))^{2l} \leq C \|\nabla U(x)\|^{2l} h^{6l}.$$

*Proof of Lemma 3.3.* For  $x \in \mathbb{R}^d$ , we set  $y = \Phi_h^\xi(x)$ . A Taylor expansion of  $U(y)$  yields,

$$|U(y) - U(x)| = O(h^3 \|D^3 U(x)\| \|\nabla U(x)\|^3).$$

□

The centered point integrator is an example of integrator for the Hamiltonian dynamics, that can be used as soon as the potential  $U$  is globally Lipschitz. Actually, similarly to GMALA in the non globally Lipschitz case, one can construct an approximate integrator  $\Phi_h^\xi$  by using a truncation of  $\nabla U$  to ensure its global Lipschitzness. Yet, this trick may lead to high rejection rates in the areas where  $\nabla U$  is truncated. Again, other strategies can be used like importance sampling or even a change of variable in the integrand (1.1). The GHMALA algorithm is especially interesting when we are able to integrate efficiently the Hamiltonian dynamics. We show later on the numerical experimentations two examples of specific integrators that enable to improve significantly the performance of GHMALA with respect to GMALA with proposal kernel  $Q_2^\xi$ . We can recall for example, the case of separated Hamiltonian dynamics that can be integrated with explicit volume preserving schemes which are time-reversible and symmetric ([19]).

The authors propose in [8] a similar splitted scheme where the hybrid step is replaced by fourth-order Runge-Kutta method. Even though this choice leads to a biased estimator, it enabled to get rid of the centered point integrator, that may be more costly than the Runge-Kutta method in terms of computation time.

#### 4. NUMERICAL EXPERIMENTATIONS

**4.1. Anisotropic distribution.** In this section, we test our nonreversible MALA algorithm by computing an observable  $f$  with respect to a two dimensional anisotropic distribution. More precisely we want to estimate  $\mathbb{E}[f(X)]$  with  $f((x_1, x_2)) = x_1^2 1_{x_1 > 15}$ , and  $X \sim \pi(x)$  with  $\pi(x) \propto e^{-V(x)}$ , where

$$V((x_1, x_2)) = \frac{x_1^2}{\sqrt{1 + 50x_1^2}} + x_2^2.$$

Such a distribution is more stretched out in the  $x_1$  direction rather than the  $x_2$ . We expect our lifted algorithm to be more favorable than MALA when the anisotropy is strong. Indeed, MALA tends to perform a slower exploration of the state space in the  $x_1$  direction with respect to  $x_2$  the direction. The lifted algorithm is supposed to correct this problem since the Hamiltonian dynamics should lead to a fast exploration of the iso-potential lines, which are stretched in the direction  $x_1$ . We choose as the skew-symmetric matrix  $J$  defined by,

$$J = \alpha \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \quad (4.1)$$

with  $\alpha$  a real parameter.

To compare MALA with GMALA and GHMALA, different optimal time steps parameters should be used. There is a tradeoff between achieving high average acceptance ratios (obtained with small time steps) and small correlation between the successive samples (obtained with large time-steps). More precisely MALA usually gives its best results with a large  $h$  that ensures a significant average rejection ratio. On the contrary, GMALA and GHMALA require much smaller time steps to avoid the regress that happens with the direction switching at each rejection. Thus, MALA should be used with a higher time step than the two others. Moreover, as GMALA with proposal kernel  $Q_1^\xi$  leads to a worse average acceptance ratio than  $Q_2^\xi$ , the former requires a smaller time-step than the latter. Figure 1 shows the average acceptance ratio with respect to the time step  $h$  for GMALA with proposal kernels  $Q_1^\xi$  and  $Q_2^\xi$ . We can verify that the average rejection ratio of GMALA with  $Q_1^\xi$  is indeed much worse than MALA and is of order  $h$ . In practice, this limitation leads to very correlated successive samples, and thus bad asymptotic variances for the estimators based on such Markov chain. In the following, we always consider GMALA with proposal kernel  $Q_2^\xi$  instead of  $Q_1^\xi$ .



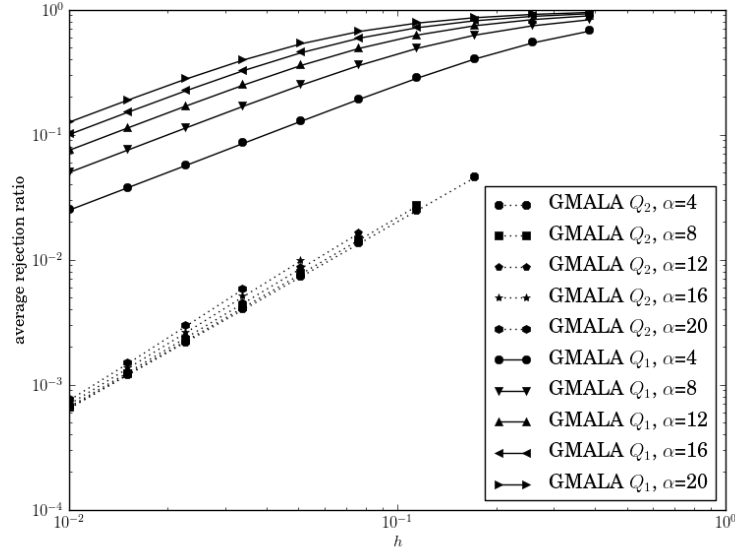


FIGURE 1. Comparison of average rejection ratio for GMALA with proposal kernels  $Q_1^\xi$  and  $Q_2^\xi$ .

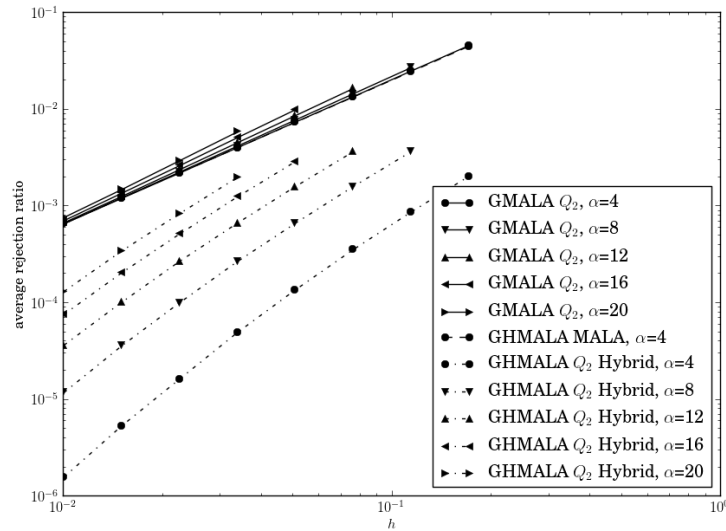


FIGURE 2. Comparison of average rejection ratio for MALA, GMALA (with proposal kernel  $Q_2^\xi$ ) and GHMALA.

Figure 2 shows the acceptance ratio with respect to the time step  $h$  for those three algorithms. GHMALA is composed by two steps: the first one consists of MALA and the second one consists of a Hybrid

iteration. Both steps are composed by an acceptance step. The average acceptance rate of the first step does not depend on  $\alpha$ , and is denoted by *GHMALA MALA* in the legend. Moreover, it is the same as a plain MALA. The average acceptance rate of the second step is denoted by *GHMALA Hybrid*. We also observe that the rejection ratio for GMALA with  $Q_2^\xi$  is close to MALA's and scales as  $h^{3/2}$ . This is due to the fact that the non reversibility contributes at order  $h^2$  in the expression of the average rejection ratio. We can also verify the upper bound on the rejection probability of the hybrid step, given by Lemma 3.3.

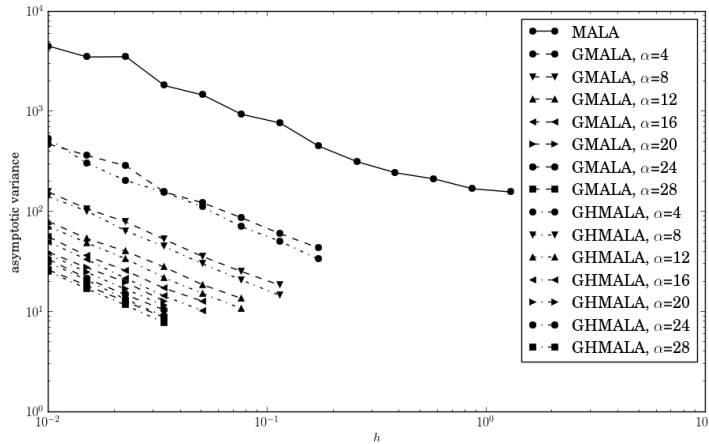


FIGURE 3. Variance comparison of MALA, GMALA ( $Q_2^\xi$ ) and GHMALA on the anisotropic distribution

We compare now the asymptotic variance of the estimators build with MALA, GMALA and GHMALA. To do so, we compute 1000 independent estimators composed by the time average of  $10^5$  samples. We plot in Figure 3 the empirical relative variance of these estimators, We observe that GMALA performs much better than MALA by a factor 20. We should precise that a reduction in the variance does not necessarily mean a reduction in computing time since one iteration of GMALA is more costly than MALA since it requires a fixed point iteration.

The question of choosing  $\alpha$  is quite natural. In the case of the time-continuous process, it is known that the decrease in asymptotic variance is monotonic with the intensity  $\alpha$ , which suggests to use the algorithms with large  $\alpha$  [20, 8, 14]. Yet, well-posedness of the proposal kernel  $Q_2^\xi$  requires the product  $\alpha h$  to be small enough to ensure convergence of the fixed point. From a computational point of view, the cost of the method is proportional to the number of Picard iterations, which scales like  $-\log \rho$ , where  $\rho$  denotes the contraction ratio (that scales as  $\alpha h$ ). Then, the choice of parameters  $\alpha$  and  $h$  should take into account these two effects.

**4.2. Warped Gaussian distribution.** This example deals with a non quadratic potential. Both GMALA and GHMALA can be adapted to this case. The simple idea is to build a proposal kernel with a truncation of  $\nabla U$ , to make it globally Lipschitz. This is the same idea as the one used for MALTA ([22]). Moreover, it is also possible to choose a more efficient integrator than the centered point integrator defined by Equation (3.5).

To illustrate these two methods, we test now our algorithms in the case of a two-dimensional warped Gaussian distribution. This toy case has been introduced in [10] and used as a benchmark for variance

reduction methods based on nonreversible Langevin samplers in [8]. More precisely, we aim to estimate  $\mathbb{E}[f(X)]$  with the observable  $f$  and  $X$  distributed with  $\pi$ , defined by,

$$f(x) = x_1^2 + x_2^2, \quad \pi(x) \propto e^{-V(x)}, \quad \text{with} \quad V(x) = \frac{x_1^2}{100} + \left(x_2 + \frac{x_1^2}{20} - 5\right)^2.$$

We define the skew symmetric matrix  $J$  by (4.1). It appeared in [8] that the nonreversible Langevin dynamics enables to reduce the asymptotic variance by several orders of magnitude.

We propose to implement GMALA using a truncated drift to make it globally Lipschitz to ensure the well-posedness of the method, and to implement GHMALA using a specific integrator of the Hamiltonian dynamics. We show that GHMALA performs better than GMALA and MALA in this case. More precisely, the integrator is defined as the centered point integrator, after the symplectic change of variable  $\psi$  defined for all  $(x_1, x_2) \in \mathbb{R}^2$  by

$$\psi(x_1, x_2) = \left(x_1, x_2 + \frac{x_1^2}{20} - 5\right).$$

Typically, this integrator enables to solve this dynamics for larger time steps than GMALA with proposal kernel  $Q_2^\epsilon$ , and thus reduces the asymptotic variance. Figure 4 displays the asymptotic variance for the estimators build with these algorithms. It is computed as the empirical variance of 2000 independent estimators constructed as the time average of  $10^5$  iterations of the algorithms. We can observe that

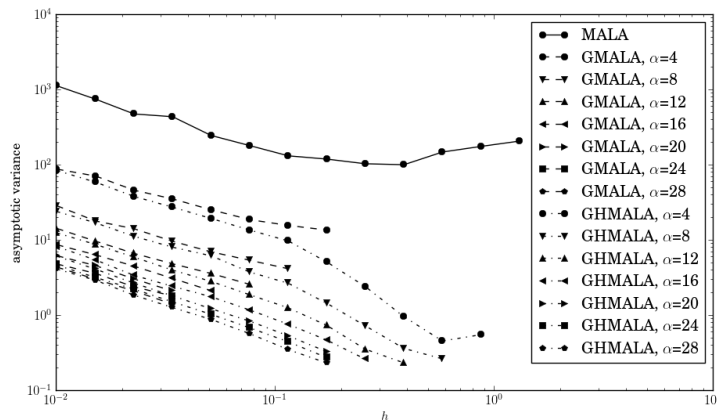


FIGURE 4. Variance comparison of MALA, GMALA and GHMALA on the warped Gaussian distribution

GMALA and GHMALA performs similarly for small time steps  $h$ . Yet, for larger time-steps, it is not possible to define the proposal kernel for GMALA, whereas it is still the case for GHMALA. Eventually, we achieve a variance reduction of about a factor 500 with GHMALA and 60 with GMALA, compared with classical MALA.

**4.3. Quartic Gaussian distribution.** This toy case aims to present a particular case where GHMALA can be used without implicit integrator, and in a non globally Lipschitz case, which may enable to reduce the computational cost by avoiding the fixed point iteration. Again, we aim to estimate  $\mathbb{E}[f(X)]$  where

$X$  is distributed with  $\pi$ , and where,

$$f(x) = x_1^2 + x_2^2, \quad \pi(x) \propto e^{-V(x)}, \quad \text{with} \quad V(x) = \frac{x_1^2}{100} + x_2^4.$$

We define the skew-symmetric matrix  $J$  by (4.1). In this case, the Hamiltonian dynamics defines then a separable system and volume-preserving explicit methods can be used (see [19]). More precisely, we define  $\Phi_h^\xi$  for all  $x = (x_1, x_2) \in \mathbb{R}^2$  by  $\Phi_h^\xi(x) = (y_1, y_2)$ , where,

$$\begin{aligned} y_1^{1/2} &= x_1 - \frac{h}{2} \alpha \xi \frac{\partial V}{\partial x_1}(x) \\ y_2 &= x_2 + h \alpha \xi \frac{\partial V}{\partial x_2}((y_1^{1/2}, x_2)) \\ y_1 &= y_1^{1/2} - \frac{h}{2} \alpha \xi \frac{\partial V}{\partial x_1}((y_1^{1/2}, y_2)) \end{aligned}$$

Thus, the non Lipschitz nonlinearity of  $\nabla V$  is not an obstacle to the well-posedness of this integrator.

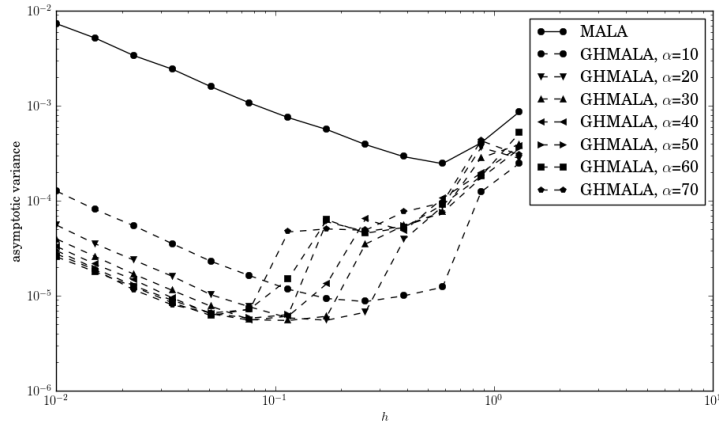


FIGURE 5. Variance comparison of MALA and GHMALA on the quartic Gaussian distribution

We plot in figure 5 the asymptotic variance for the time average estimator build with GHMALA in the same way as for the warped Gaussian distribution. We can observe that the decrease in variance between MALA and GHMALA for small  $h$  is around 280. Nevertheless, for larger  $h$ , the explicit integration is not accurate enough, which leads to an increase in the asymptotic variance for larger time steps. Eventually, the smallest asymptotic variance of the time average estimator of GHMALA is around 50 times lower than the smallest asymptotic variance for MALA.

### CONCLUSION

We presented a class of unbiased algorithm that enables to benefit from the variance reduction of the nonreversible Langevin equations (1.3) with respect to the reversible dynamics (1.2). More precisely, we presented two variations of these algorithms. The first one (GMALA) can be viewed as a lifting method, and more specifically as a generalized Metropolis Hastings methods on a lifted state space. The second one (GHMALA), similar to the first one, can be viewed as a Generalized Hybrid Monte-Carlo method.

Numerical experimentations show that variance reductions (compared with classical MALA) of several orders of magnitude can be achieved for potentials concentrated on a lower dimensional submanifold. We also expect these algorithms to perform better in the case of entropic barriers. The main difficulty is, in the case of GMALA, to use a proposal that enables to achieve a sufficiently high average acceptance ratio (to compete with MALA). For example this can be done by using a mid-point discretization. Even though this scheme is implicit, the computation of the Metropolis-Hastings acceptance probability does not require the computation of the Hessian of  $\log \pi$ . In the case of GHMALA, numerical experimentations show that the choice of a suitable Hamiltonian integrator may lead to large improvements and computational cost reduction compared with the mid-point method.

**Acknowledgement :** This work was supported by a public grant as part of the Investissement d’avenir project, reference ANR-11-LABX-0056-LMH, LabEx LMH.

#### REFERENCES

1. Marco Banterle, Clara Grazian, Anthony Lee, and Christian P Robert, *Accelerating metropolis-hastings algorithms by delayed acceptance*, arXiv preprint arXiv:1503.00996 (2015).
2. Joris Bierkens, *Non-reversible metropolis-hastings*, Statistics and Computing (2015), 1–16.
3. Joris Bierkens, Paul Fearnhead, and Gareth Roberts, *The zig-zag process and super-efficient sampling for bayesian analysis of big data*, arXiv preprint arXiv:1607.03188 (2016).
4. Nawaf Bou-Rabee and Eric Vanden-Eijnden, *Pathwise accuracy and ergodicity of metropolized integrators for SDEs*, Comm. Pure Appl. Math. **63** (2010), no. 5, 655–696. MR 2583309
5. Fang Chen, László Lovász, and Igor Pak, *Lifting markov chains to speed up mixing*, Proceedings of the Thirty-first Annual ACM Symposium on Theory of Computing (New York, NY, USA), STOC ’99, ACM, 1999, pp. 275–281.
6. Persi Diaconis, Susan Holmes, and Radford M. Neal, *Analysis of a nonreversible markov chain sampler*, Ann. Appl. Probab. **10** (2000), no. 3, 726–752.
7. Simon Duane, A.D. Kennedy, Brian J. Pendleton, and Duncan Roweth, *Hybrid monte carlo*, Physics Letters B **195** (1987), no. 2, 216 – 222.
8. A. B. Duncan, T. Lelièvre, and G. A. Pavliotis, *Variance reduction using nonreversible langevin samplers*, Journal of Statistical Physics **163** (2016), no. 3, 457–491.
9. Max Fathi and Gabriel Stoltz, *Improving dynamical properties of stabilized discretizations of overdamped langevin dynamics*, arXiv preprint arXiv:1505.04905 (2015).
10. Heikki Haario, Eero Saksman, and Johanna Tamminen, *An adaptive metropolis algorithm*, Bernoulli **7** (2001), no. 2, 223–242.
11. K Hukushima and Y Sakai, *An irreversible markov-chain monte carlo method with skew detailed balance conditions*, Journal of Physics: Conference Series **473** (2013), no. 1, 012012.
12. Chii-Ruey Hwang, Shu-Yin Hwang-Ma, and Shuenn-Jyi Sheu, *Accelerating gaussian diffusions*, Ann. Appl. Probab. **3** (1993), no. 3, 897–913.
13. ———, *Accelerating diffusions*, Ann. Appl. Probab. **15** (2005), no. 2, 1433–1444.
14. Chii-Ruey Hwang, Raoul Normand, and Sheng-Jih Wu, *Variance reduction for diffusions*, Stochastic Processes and their Applications **125** (2015), no. 9, 3522 – 3540.
15. A.D. Kennedy and Brian Pendleton, *Cost of the generalised hybrid monte carlo algorithm for free field theory*, Nuclear Physics B **607** (2001), no. 3, 456 – 510.
16. T Lelièvre, F Nier, and GA Pavliotis, *Optimal non-reversible linear drift for the convergence to equilibrium of a diffusion*, Journal of Statistical Physics **152** (2013), 237–274.
17. Tony Lelièvre, Mathias Rousset, and Gabriel Stoltz, *Free energy computations : a mathematical perspective*, Imperial College Press, London, Hackensack (N.J.), Singapore, 2010.
18. A. Pazy, *Semigroups of linear operators and applications to partial differential equations*, Applied Mathematical Sciences, vol. 44, Springer-Verlag, New York, 1983. MR 710486
19. M-Z Qin and W-J Zhu, *Volume-preserving schemes and numerical experiments*, Computers & Mathematics with Applications **26** (1993), no. 4, 33–42.

20. Luc Rey-Bellet and Konstantinos Spiliopoulos, *Irreversible langevin samplers and variance reduction: a large deviations approach*, *Nonlinearity* **28** (2015), no. 7, 2081.
21. Luc Rey-Bellet and Konstantinos Spiliopoulos, *Variance reduction for irreversible langevin samplers and diffusion on graphs*, *Electron. Commun. Probab.* **20** (2015), 16 pp.
22. Gareth O. Roberts and Richard L. Tweedie, *Exponential convergence of langevin distributions and their discrete approximations*, *Bernoulli* **2** (1996), no. 4, 341–363.
23. Luke Tierney and Antonietta Mira, *Some adaptive monte carlo methods for bayesian inference*, *Statistics in Medicine* **18** (1999), no. 17-18, 2507–2515.
24. Marija Vucelja, *Lifting—a nonreversible markov chain monte carlo algorithm*, arXiv preprint arXiv:1412.8762 (2014).
25. Sheng-Jih Wu, Chii-Ruey Hwang, and Moody T. Chu, *Attaining the optimal gaussian diffusion acceleration*, *Journal of Statistical Physics* **155** (2014), no. 3, 571–590.