

The School of Mathematics



THE UNIVERSITY
of EDINBURGH

Constrained Langevin Algorithms for Score-Based Diffusion

by

Tadgh Kelly

Dissertation Presented for the Degree of
MSc in Computational Applied Mathematics

August 2023

Supervised by
Prof. Benedict Leimkuhler

Abstract

Score-based models have been seen to perform remarkably well at generating data in a range of domains. Score-based generative modelling (SBGM) consists of a noising stochastic differential equation (SDE) which is used to transform the datapoints to a known distribution, and a denoising SDE which transforms back from the known distribution to the data distribution. Constrained generative modelling has applications in many fields where the data is supported on a manifold rather than on Euclidean space. Most work on constrained score-based diffusion has focused on Riemannian manifolds using the Riemannian metric, and has used overdamped Langevin dynamics to noise and denoise the data. Here we examine the use of g-BAOAB, an integrator for constrained underdamped Langevin dynamics, and in doing so we focus on general manifolds defined by holonomic constraints, by considering the constraint manifold as an embedded manifold in Euclidean space. The approach is demonstrated on a problem of human pose generation, as well as sampling from a distribution on a sphere, and is compared to a continuous normalizing “Moser flow” model. **I'm not sure exactly how I should phrase my results, since they aren't very quantifiable.**

Acknowledgments

Here come your acknowledgments ...

Own Work Declaration

Here comes your own work declaration

Contents

1	Introduction	1
1.1	Motivation	2
1.2	Outline of thesis	3
2	Background	4
2.1	Hamiltonian Systems	4
2.1.1	Symplectic Integrators	4
2.1.2	Splitting methods	5
2.2	Constrained Dynamics	5
2.2.1	Example: The Simple Pendulum	6
2.2.2	Constraint-preserving integrators	7
2.3	RATTLE and SHAKE	7
2.3.1	Iterative methods	8
2.4	Probability measures on manifolds	8
3	Langevin Dynamics	9
3.1	Splitting methods for Langevin dynamics	9
3.2	Constrained Langevin dynamics	10
3.2.1	Geodesic Integration	11
3.2.2	g-BAOAB	12
3.3	Example	12
4	Score based diffusion	14
4.1	Score approximation	14
4.1.1	Fisher divergence on Riemannian manifolds	15
4.2	Perturbing data	16
4.3	Time Reversal	17
4.3.1	Time reversal of constrained Langevin dynamics	17
5	Connection with ODEs	17
5.1	Normalizing flows	18
5.1.1	Kullback-Leibler divergence	18
5.1.2	Continuous normalizing flows	19
5.2	Riemanninan normalizing flows	19
5.3	Score-based flow	20
5.3.1	Score-based flow on Riemannian manifolds	21
5.3.2	Score-based flow for constrained (underdamped) Langevin dynamics with zero force.	21
5.4	Moser Flows	22
5.4.1	Training the model	23
6	Experiments	24
6.1	3D Human Poses Estimation	24
6.2	3D Human Poses Generation	26
6.2.1	Adding noise	26
6.2.2	Score Matching	27
6.3	Results	29
7	Comparison with Moser flow	29
7.1	Moser flow model	29
7.2	Score-based model	30
7.3	Comparison / Results	31

8 Related work	32
8.1 Stochastic interpolants	32
8.2 Pushforward of Euclidean normalizing flows	33
9 Limitations	33
10 Conclusion	33
Appendices	33
A Computational cost	33
A.1 Score-based generative modelling	33
A.2 Moser flows	33
B Activation functions	34
C Assumptions	34
C.1 Theorem ??	34
C.2 Theorems 5.2 and 5.3	34

List of Tables

List of Figures

1	Facial mages generated by implicit and explicit generative models.	1
2	Facial images generated by a diffusion model [10].	1
3	Moser flow and score-based generative models trained on earth sciences data gathered by Mathieu and Nickel [8]. Blue and red dots represent training and test datapoints, respectively, while the learned density is green-blue, with blue indicating larger values.	2
4	Forward and symplectic Euler integrator for the simple pharmonic oscillator with mass $m = 1$	4
5	The simple pendulum.	6
6	Density compared to g-BAOAB for the function $U(x, y, z) = (x^2 - 1)^2 + (y^2 - 1)^2 + (z^2 - 1)^2$	13
7	Convergence of ϕ , θ , ϕ^2 and θ^2 for different stepsizes h . (note plot with math h)	13
8	10 point rolling average of error per iteration for different values of K_r	13
9	Convergence in expected potential $U(\mathbf{q})$ for different values of γ	14
10	First, second and third eigenpose.	25
11	Mean joint-to-joint error per frame and negative log density per iteration	25
12	Target pose, constrained optimization algorithm pose and g-BAOAB mean pose after 2,000 g-BAOAB iterations.	26
13	Initial pose being perturbed using the g-BAOAB algorithm, with a step size of 0.01.	27
14	Monte Carlo integration of $(1 - \nabla \hat{s})$	28
15	Score-based diffusion neural network architecture and training loss [progression] .	29
16	Samples from the data distribution	30
17	Samples from a score-based diffusion model using g-BAOAB	30
18	Training data for the Moser flow and SBG models (Note that the colours have been normalized differently than in Figure 6 in order to match the scale in ??) .	31
19	Moser flow neural network architecture	31
20	Result of a Moser flow model fit and a score-based generative model fit to data generated by g-BAOAB waiting for the score-based model to train	31

1 Introduction

Generative modelling aims to train a generative model to sample from the same distribution as the data. There are many applications of generative modelling, including image syntheses and inpainting [5, 35] and speech and music synthesis [26], but there are also applications in many other fields, since sampling from an unknown distribution is a common problem.

Early generative modelling techniques could be grouped into two categories; likelihood-based models and implicit generative models. Likelihood-based models, as demonstrated in Figure 1b such as Variational Auto-Encoders (VAEs) learn the probability density function using maximum likelihood [9]. Likelihood-based models require a restricted model to ensure a tractable normalizing constant, which commonly means that VAEs are based on fitting Gaussian-based models. Implicit generative models generate data without explicitly learning the density function, and are usually in the form of General Adversarial Networks (GANs), in which one neural network generates the data and another called the “discriminator” aims to discriminate between true and artificial data. However training for GANs is famously unstable [37].

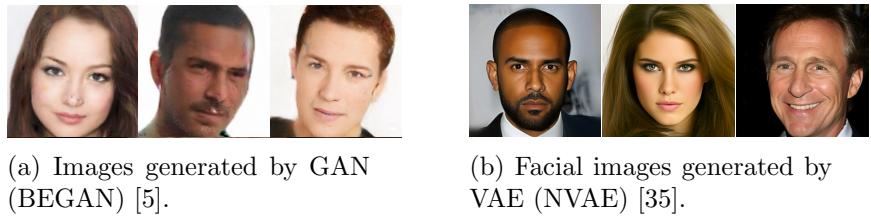


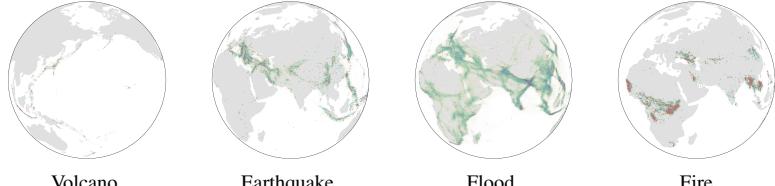
Figure 1: Facial images generated by implicit and explicit generative models.

In recent years, diffusion models have been shown to be very effective, and have become an active area of research, outperforming GANs at a number of tasks [?]. Diffusion models are a family of generative model which perturb data by injecting noise, and then learn to reverse the noising process in order to generate samples from the original data distribution [37]. Score-based diffusion is based on learning the score function of the noising process, the gradient of the log of the marginal probability density function, which eliminates the need for the normalizing constant of the density function. Another important fact about score-based models is that they allow a model for the score to be trained without explicit knowledge of the true score function, by minimizing the Fisher divergence, which can be done by sampling from the data distribution [12, 32]. Thus the only information we need to train a score-based model is available from the data samples. While score matching can be performed without adding noise to the data, the motivation for the noising process comes from the difficulty in training a score-based model in areas of low density, where there may be very few data points. Adding noise aids the training process in these regions by preventing overfitting and counteracting biases that may be introduced in these regions [33].

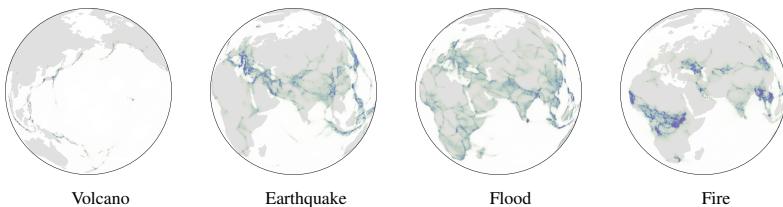


Figure 2: Facial images generated by a diffusion model [10].

Constrained data arises in many areas, including protein modelling and image recognition as well as many other fields [29, 23]. Some that have been examined in the context of diffusion models are geoscience, in which there are many events which occur on the surface of the earth, as shown in Figure 3, and human pose generation, in which the length of the bones are constrained to be constant [9]. In settings like these, the data often lie on Riemannian manifolds rather than in Euclidean space. This motivates the use of constrained diffusion algorithms, in order to ensure that the samples lie on the manifold in question, which wouldn't be the case with standard diffusion models.



(a) Results of a Moser flow model [27].



(b) Results of a score-based model [8].

Figure 3: Moser flow and score-based generative models trained on earth sciences data gathered by Mathieu and Nickel [8]. Blue and red dots represent training and test datapoints, respectively, while the learned density is green-blue, with blue indicating larger values.

There are two main methods of generative modelling which have been used for modelling on a constraint manifold; score-based diffusion and Continuous Normalizing Flows (CNFs). Recent research into score-based diffusion on Riemannian manifolds has shown promising results, with results from Euclidean models often having direct analogs on Riemannian manifold [8]. In terms of CNFs on manifolds there are generally two different ways of generating a normalizing flow on the manifold. For manifolds which are topologically equivalent to Euclidean space, there are methods which involve the pushforward of a Euclidean flow along an invertible map $f : \mathbb{R}^n \rightarrow \mathcal{M}$, normally taken to be the exponential map, however, such a map is not necessarily available unless \mathcal{M} is topologically equivalent to Euclidean space. Alternatively, Riemannian geometry, specifically the Riemannian divergence, can be used to define flows on the manifold directly [24]. This method is more general in that it can be used to construct CNFs on manifolds which are not topologically equivalent to Euclidean space. A third option for embedded submanifolds of Euclidean space is to focus on the embedding from Euclidean space to define a CNF on the manifold, as touched upon by Rozen et al [27]. This method will allow us to work with holonomic constraints of the same form as those used in constrained Langevin dynamics, which will allow a comparison between the two methods.

1.1 Motivation

While overdamped Langevin dynamics have often been used in SBD, underdamped Langevin dynamics has been used less frequently, and in particular in the case of constrained SBD, has not been studied. Using underdamped Langevin dynamics may be preferable to the overdamped equations, since it has been observed to converge more quickly to the invariant distribution [7]. This motivates the study of the underdamped Langevin equations in the context of constrained

SBD.

1.2 Outline of thesis

This report is structured as follows; in Section 2 we introduce some background concepts related to constrained Hamiltonian systems. In Section 3, we introduce Langevin dynamics with constraints and the g-BAOAB integrator, and study its properties on a toy problem. In Section 4 we introduce score-based diffusions on Euclidean space and generalized to manifolds, and the time-reversal of the constrained underdamped Langevin dynamics. In Section ?? we relate the SBD model to CNFs, and discusses the limitations of underdamped Langevin dynamics in this area. We also introduce Moser flows, a type of CNF which is particularly well suited to certain Riemannian manifolds. Section ?? contains experiments in which the Moser flow model is compared to the SBD model using g-BAOAB.

2 Background

2.1 Hamiltonian Systems

A Hamiltonian system is a system of differential equations characterized by its Hamiltonian function $H : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\frac{d}{dt}\mathbf{q} = \nabla_{\mathbf{p}}H(\mathbf{q}, \mathbf{p}), \quad (2.1)$$

$$\frac{d}{dt}\mathbf{p} = -\nabla_{\mathbf{q}}H(\mathbf{q}). \quad (2.2)$$

Then when the Hamiltonian is the total kinetic plus potential energy, $H(\mathbf{q}, \mathbf{p}) = \frac{\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}}{2} + U(\mathbf{q}, \mathbf{p})$, we have

$$\frac{d}{dt}\mathbf{q} = \mathbf{M}^{-1}\mathbf{p}, \quad (2.3)$$

$$\frac{d}{dt}\mathbf{p} = -\nabla_{\mathbf{q}}U(\mathbf{p}, \mathbf{q}). \quad (2.4)$$

In this paper we will omit the mass matrix by setting $\mathbf{M} = \mathbf{I} \in \mathbb{R}^{n \times n}$, since this is often done in data science settings when we aren't dealing with actual physical systems.

2.1.1 Symplectic Integrators

The evolution of a Hamiltonian system is symplectic, which means that it has an area form $\omega = dp \wedge dq$ in the phase space, that measures areas of phase space,

$$A = \int_S \omega, \quad (2.5)$$

for a given region S . Symplectic integrators are designed to maintain the symplectic structure of the phase space, which gives them the desirable property of having bounded energy fluctuations [?]. A numerical method is called symplectic if the area form is preserved exactly:

$$d\mathbf{q}_{n+1} \wedge d\mathbf{p}_{n+1} = d\mathbf{q}_n \wedge d\mathbf{p}_n. \quad (2.6)$$

While we can also construct integrators designed specifically to conserve energy, their performance is often poorer than a simpler symplectic method [18]. As shown in Figure 4 by comparing the symplectic Euler method to the forward Euler method for the simple harmonic oscillator, the energy fluctuates for both methods but the fluctuations remain bounded by the symplectic method.

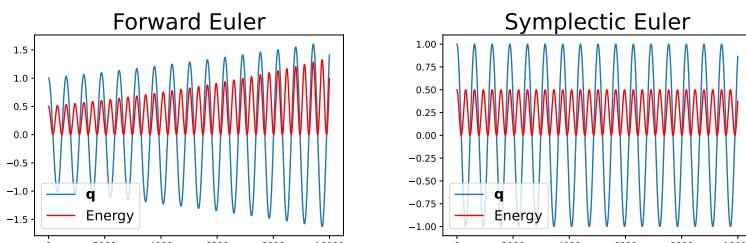


Figure 4: Forward and symplectic Euler integrator for the simple pharmonic oscillator with mass $m = 1$.

2.1.2 Splitting methods

Splitting methods provide a way to create integrators for Hamiltonian systems and in particular, a way to create symplectic integrators. Suppose we can split the Hamiltonian H into a sum of Hamiltonians,

$$H(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^k H_i(\mathbf{q}, \mathbf{p}), \quad (2.7)$$

with each of the Hamiltonians H_i giving an explicitly solvable system. Since explicitly solving the equations gives a symplectic flow map, and compositions of symplectic maps are themselves symplectic maps, this gives a symplectic integrator. In Hamiltonian systems, this often involves splitting the Hamiltonian into kinetic and potential energy. First we introduce the flow map. Consider an initial value problem:

$$\dot{\mathbf{q}} = \mathbf{f}(\mathbf{q}), \quad \mathbf{q}(0) = \mathbf{q}_0. \quad (2.8)$$

Then the *flow map* map \mathcal{F}_t maps points in the phase space according to the time-evolution of the problem:

$$\mathcal{F}_t(\mathbf{q}_0) = \mathbf{q}_t. \quad (2.9)$$

If $H(\mathbf{p}, \mathbf{q}) = H_1(\mathbf{p}, \mathbf{q}) + H_2(\mathbf{p}, \mathbf{q})$, and H_1 and H_2 have flow maps \mathcal{F}^1 and \mathcal{F}^2 , then it can be shown that the map given by their composition,

$$\mathcal{F}^{1+2} = \mathcal{F}^1 \circ \mathcal{F}^2, \quad (2.10)$$

is a second order approximation of the true flow \mathcal{F} [16].

2.2 Constrained Dynamics

We will begin by introducing mechanical systems with holonomic constraints. For a constrained system, a particle is subject to two types of force: applied forces which are defined using the potential energy function U , and constraint forces which are the forces which force the particle to remain on the constraint surface [18]. D'Alembert's principle states that the constraining force acts along the normal direction to the constraint surface, therefore if \mathbf{F}_g is the constraint force we can write $\mathbf{F}_g = \lambda \nabla_{\mathbf{q}}$ for some scalar λ . λ is unknown, and it can be solved for by requiring that the particle remains on the constraint surface.

Consider algebraic constraints of the form

$$g_i(\mathbf{q}) = 0, \quad i = 1, \dots, m, \quad (2.11)$$

where $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is smooth for all $i \in [1, \dots, m]$. This mechanical system is holonomic and constant in time. If the Hamiltonian is again given by $H(\mathbf{q}, \mathbf{p}) = \frac{\mathbf{p}^T \mathbf{p}}{2} + U(\mathbf{q}, \mathbf{p})$, then we have

$$\frac{d}{dt} \mathbf{q} = \mathbf{p} \quad (2.12a)$$

$$\frac{d}{dt} \mathbf{p} = \mathbf{F} - \sum_{i=1}^m \lambda_i g_i(q) \quad (2.12b)$$

$$g_i(\mathbf{q}) = 0, \quad i = 1 \dots m, \quad (2.12c)$$

where we have omitted the mass matrix and let $\mathbf{F} = -\nabla_{\mathbf{q}} U$. We can rewrite this system as:

$$\frac{d}{dt}\mathbf{q} = \mathbf{p} \quad (2.13a)$$

$$\frac{d}{dt}\mathbf{p} = \mathbf{F}(\mathbf{q}) - \mathbf{G}^T(\mathbf{q})\boldsymbol{\Lambda} \quad (2.13b)$$

where $\mathbf{g}(\mathbf{q})$ is the vector of constraints $\mathbf{g}(\mathbf{q}) = [g_1(\mathbf{q}) \ g_2(\mathbf{q}) \ \dots \ g_m(\mathbf{q})]^T$, $\boldsymbol{\Lambda} = [\lambda_1, \dots, \lambda_m]^T$ and $\mathbf{G}(\mathbf{q}) \in \mathbb{R}^{m \times n}$ represents the Jacobian matrix of the constraint vector evaluated at \mathbf{q} .

We will assume that $\{\nabla_{\mathbf{q}} g_i(\mathbf{q})\}_{i=1,\dots,m}$ forms a linearly independent set for each $\mathbf{q} \in \mathbb{R}^n$, or equivalently that $\mathbf{G}(\mathbf{q})$ has full rank for each $\mathbf{q} \in \mathbb{R}^n$. This lets us consider the constraint surface as a smooth algebraic manifold \mathcal{M} , which implies that it can be equipped with a Riemannian metric. While there are many results which make use of Riemannian geometry, for a general set of holonomic constraints the Riemannian metric will not be known, so we will avoid the use of Riemannian metrics and local coordinates, and instead focus on the manifold as a submanifold of Euclidean space.

2.2.1 Example: The Simple Pendulum

Consider the simple pendulum shown below, where the bob has mass 1 and the length of the rod is 1. We let $g_1(\mathbf{q})$ be the constraint so that $g_1(\mathbf{q}) = \|\mathbf{q}\|_2^2 - 1$. The tension then acts along the normal of the constraint, so we can write $\vec{T} = \lambda \nabla_{\mathbf{q}} g_1$, which gives the following equations:

$$\frac{d}{dt}\mathbf{q} = \mathbf{p} \quad (2.14a)$$

$$\frac{d}{dt}\mathbf{p} = \mathbf{F} - \lambda \nabla_{\mathbf{q}} g_1 \quad (2.14b)$$

$$g_1(\mathbf{q}) = 0. \quad (2.14c)$$

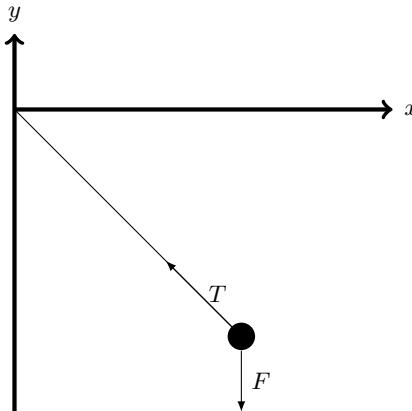


Figure 5: The simple pendulum.

Using the fact that $\frac{d}{dt}g_1 = 0$ gives:

$$\frac{d}{dt}g_1 = \frac{d}{d\mathbf{q}}g_1 \frac{d}{dt}\mathbf{q} + \frac{d}{d\mathbf{p}}g_1 \frac{d}{dt}\mathbf{p} = 0 \quad (2.15)$$

$$\implies \ddot{g}_1 = 2\dot{\mathbf{q}}^T \dot{\mathbf{q}} + 2\mathbf{q}^T \ddot{\mathbf{q}} = 0 \quad (2.16)$$

$$\implies \mathbf{q}^T (\mathbf{F} - 2\lambda\mathbf{q}) + \dot{\mathbf{q}}^T \dot{\mathbf{q}} = 0 \quad (2.17)$$

$$\implies 2\lambda\mathbf{q}^T \mathbf{q} = (\mathbf{q}^T \mathbf{F} + \dot{\mathbf{q}}^T \dot{\mathbf{q}}) \quad (2.18)$$

$$\implies \lambda = \frac{\mathbf{q}^T \mathbf{F} + \dot{\mathbf{q}}^T \dot{\mathbf{q}}}{2\mathbf{q}^T \mathbf{q}} \quad (2.19)$$

In this case this gives

$$\lambda = \frac{1}{2} \left(\frac{[x \ y]^T \mathbf{F} + \dot{x}^2 + \dot{y}^2}{x^2 + y^2} \right),$$

where we have replaced the vectors \mathbf{q} and \mathbf{p} with their x and y components. In the general case, we have a function for Λ , given by:

$$\Lambda(\mathbf{q}, \mathbf{p}) = -(\mathbf{G}\mathbf{G}^T)^{-1}[\mathbf{G}\nabla_q U(\mathbf{q}) + g_{\mathbf{q}\mathbf{q}} < \mathbf{p}, \mathbf{p} >], \quad (2.20)$$

where \mathbf{G} is used as shorthand for $\mathbf{G}(\mathbf{q})$, $\mathbf{G}\mathbf{G}^T$ is non-singular since \mathbf{G} has full rank and $g_{\mathbf{q}\mathbf{q}} < \mathbf{p}, \mathbf{p} >$ represents a vector with m elements, with the i th element given by

$$(g_i)_{\mathbf{q}\mathbf{q}} < \mathbf{p}, \mathbf{p} > = \sum_{k=1}^n \sum_{j=1}^n \frac{d^2 g_i}{d\mathbf{q}_k d\mathbf{q}_j} \mathbf{p}_k \mathbf{p}_j. \quad (2.21)$$

However, applying a numerical integration scheme based on a direct discretisation of Equations 2.14 and 2.20 leads to issues, in that the solutions won't satisfy the constraint exactly, and errors could propagate forward leading to solutions which lie far from the constraint surface. This motivates another class of integrators designed to deal with this problem. **My solution is off by a sign from the one in the general case...**

2.2.2 Constraint-preserving integrators

This section introduces a number of examples of constraint-preserving integrators, which aim to explicitly and exactly preserve the constraints at each step. Constraint-preserving integrators generally use iterative methods to find solutions that satisfy the constraints. Certain numerical integration methods can be modified so that solutions exactly preserve the constraints.

2.3 RATTLE and SHAKE

Two popular methods for constrained integration of ODEs are RATTLE and SHAKE. Although SHAKE is not technically symplectic since the velocities aren't tangent to the constraint manifold [20], it was later improved by RATTLE, which is symplectic. SHAKE was proposed by Ryckaert et al., and based on a standard leapfrog discretisation of the ODE [28].

$$\mathbf{q}^{n+1} = \mathbf{q}^n + \Delta t \mathbf{p}^{n+1/2} \quad (2.22a)$$

$$\mathbf{p}^{n+1/2} = \mathbf{p}^{n-1/2} - \Delta t \nabla_{\mathbf{q}} U(\mathbf{q}^n) - \Delta t \mathbf{G}(\mathbf{q}^n)^T \lambda^n \quad (2.22b)$$

$$\mathbf{g}(\mathbf{q}^{n+1}) = \mathbf{0} \quad (2.22c)$$

$$\mathbf{p}^n = \frac{1}{2}(\mathbf{p}^{n+1/2} + \mathbf{p}^{n-1/2}). \quad (2.22d)$$

RATTLE was proposed later by Andersen to correct SHAKE by also enforcing the tangency constraint, the requirement that the velocity should be perpendicular to the normal of the constraint manifold [3].

$$\mathbf{q}^{n+1} = \mathbf{q}^n + \Delta t \mathbf{p}^{n+1/2} \quad (2.23a)$$

$$\mathbf{p}^{n+1/2} = \mathbf{p}^n - \frac{\Delta t}{2} \nabla_{\mathbf{q}} U(\mathbf{q}^n) - \frac{\Delta t}{2} \mathbf{G}(\mathbf{q}^n)^T \lambda_r^n \quad (2.23b)$$

$$g_i(\mathbf{q}^{n+1}) = 0, \quad i = 1 \dots m. \quad (2.23c)$$

$$\mathbf{p}^{n+1} = \mathbf{p}^{n+1/2} - \frac{\Delta t}{2} \nabla_{\mathbf{q}} U(\mathbf{q}^{n+1}) - \frac{\Delta t}{2} \mathbf{G}(\mathbf{q}^{n+1})^T \lambda_v^{n+1} \quad (2.23d)$$

$$0 = \mathbf{G}(\mathbf{q}^{n+1}) \mathbf{p}^{n+1} \quad (2.23e)$$

Where λ_r and λ_v are chosen so that the solutions satisfy the position and tangency constraints,

respectively. However it has been shown that the two methods are algebraically equivalent, and the two methods produce identical positions with only the velocity approximations differing, so it has been suggested that an efficient process would use SHAKE steps wherever the velocity was not output and one RATTLE step before returning a velocity [20]. Both are second-order methods [3] and third order accurate locally. They are also both second order convergent [20].

2.3.1 Iterative methods

In order to solve the constraining equations in RATTLE and SHAKE we use iterative methods to find approximate solutions. The general form of the iterative methods used to enforce the hidden constraint given in Equation 2.22d is given by the following algorithm, which constitutes Gauss-Newton iteration. In order to solve the equation

$$\mathbf{g}(\mathbf{q}_{n+1} - \mathbf{G}(\mathbf{q}_n)^T \Lambda) = \mathbf{0}, \quad (2.24)$$

the updates can be written as [16]:

$$\Lambda^{(i+1)} = \Lambda^{(i)} + [\mathbf{G}(\mathbf{Q}^{(i)})\mathbf{G}(\mathbf{q}_n)^T] \mathbf{g}(\mathbf{Q}^{(i)}) \quad (2.25)$$

$$\mathbf{Q}^{(i)} = \mathbf{q}_{n+1} - \mathbf{G}(\mathbf{q}_n)^T \Lambda^{(i)}, \quad (2.26)$$

we can initialize this method with $\Lambda^{(0)} = \mathbf{0}$. The method can more conveniently be written as

$$\mathbf{E}^{(l)} \Delta \Lambda^{(l)} = \mathbf{g}(\mathbf{Q}^{(l)}), \quad (2.27)$$

$$\mathbf{E}^{(l)} = \mathbf{G}(\mathbf{Q}^{(l)})\mathbf{G}(\mathbf{q}_n)^T \quad (2.28)$$

where $\Delta \Lambda^{(l)} = \Lambda^{(l+1)} - \Lambda^{(l)}$. When the solution is unique, the constraints are [sufficiently] smooth and \mathbf{E} is well conditioned, the method converges rapidly; for values $\Lambda^{(i)}$ sufficiently near the solution we have

$$\|\Lambda^{(l+1)} - \Lambda_*\| \leq C \|\Lambda^{(l)} - \Lambda_*\|^2, \quad (2.29)$$

although if the initial value isn't sufficiently close to the solution there is no guarantee of convergence. The above method deals with the positional constraints in the RATTLE and SHAKE algorithms, while in RATTLE we have an additional velocity constraint Equation 2.22d which can be satisfied by Λ_v^{n+1} solving the linear system:

$$[\mathbf{G}(\mathbf{q}_n)\mathbf{G}(\mathbf{q}_n)^T] \Lambda_v^{n+1} = \mathbf{G}(\mathbf{q}_n) \left(\frac{2}{\Delta t} \mathbf{v}^{n-1/2} - \nabla_q U(\mathbf{q}_n) \right). \quad (2.30)$$

2.4 Probability measures on manifolds

Due to the assumption that \mathbf{g} is smooth and that \mathbf{G} is full row-rank everywhere, \mathcal{M} is a smooth embedded submanifold of Euclidean space \mathbb{R}^n , and it has a metric \bar{g} which is induced by the Euclidean metric in \mathbb{R}^n . Leimkuhler et al. then showed that the Riemannian measure on \mathcal{M} was equal to the surface measure, and satisfies a Poincaré inequality [19], which is needed to show that the Markov diffusion process does in fact converge to equilibrium. (X)

3 Langevin Dynamics

The Langevin equation is derived from the equations of motion of a Brownian particle suspended in a liquid, under the assumption that it experiences two forces:

1. A viscous drag, $\gamma \dot{\mathbf{q}}(t)$ representing friction.
2. A force $\eta(t)$ representing the impacts of the fluid's molecules on the particle.

where $\eta(t)$ is white noise as used in the Ornstein-Uhlenbeck process to model the collision [14]. In the case where the particle is subjected to an external potential force field, the equations of motion become [16]

$$d\mathbf{q} = \mathbf{p} dt, \quad (3.1)$$

$$d\mathbf{p} = -\nabla_{\mathbf{q}} U(\mathbf{q}) dt - \gamma \mathbf{p} dt + \sqrt{2\gamma k_B T} dW(t), \quad (3.2)$$

where $W(t)$ is a Wiener process, k_B is the Boltzmann constant and T is the temperature. We will let $\beta = (k_B T)^{-1}$ for simplicity. This is an Itô process, a stochastic differential equation (SDE) for which the solution can be represented as the sum of integrals with respect to time and the Wiener process:

$$\mathbf{q}(t) = \mathbf{q}(0) + \int_0^t \mathbf{p}(t) dt, \quad (3.3)$$

$$\mathbf{p}(t) = \mathbf{p}(0) - \int_0^t \nabla_{\mathbf{q}} U(\mathbf{q}(t) + \gamma \mathbf{p}(t) dt) + \int_0^t \sqrt{2\beta^{-1}} dW(t), \quad (3.4)$$

where the second integral in Equation 3.4 is an Itô integral. These equations can also be written as

$$\dot{\mathbf{q}} = \mathbf{p}, \quad (3.5)$$

$$\dot{\mathbf{p}} = \mathbf{F} - \gamma \mathbf{p} + \sqrt{2k_B T \gamma} \eta(t), \quad (3.6)$$

where $\eta(t)$ is a stationary zero-mean Gaussian process, which is the form that will be used in the remainder of the paper. It can be shown that the invariant distribution of \mathbf{q} is proportional to $e^{-\beta U(\mathbf{q})}$. This gives Langevin dynamics a use case outside of the physical setting of the derivation, as a means of sampling from the invariant distribution. The equations given above are called the *underdamped* Langevin diffusion equations. The *overdamped* Langevin diffusion equation refers to the SDE

$$\dot{\mathbf{x}} = -\gamma^{-1} \nabla_{\mathbf{q}} U(\mathbf{x}) + \sqrt{\gamma^{-1} \beta^{-1}} \eta(t), \quad (3.7)$$

which is the high friction limit of the underdamped equations. Clearly both have the same invariant distribution, since the invariant distribution is independent of γ . The overdamped equation, being more straightforward, is often used in data science applications [33, 31]. γ is called the *damping coefficient* for the motion; an increase in γ results in slower moving particles, and since $\sqrt{\gamma}$ also weights the random term an increase in gamma also results in larger fluctuations. The effect of an increase in β is then to amplify the effect of the random element on the dynamics, and thus alter the equilibrium state. [Should I explain how?](#)

3.1 Splitting methods for Langevin dynamics

Now we consider a family of splitting methods for the Langevin dynamics given in 3.5. The Ornstein-Uhlenbeck process alone preserves the invariant distribution, and since Hamiltonian dynamics preserve any function of the Hamiltonian, they also preserve the invariant distribution which is proportional to $\exp^{-\beta H}$. This motivates a splitting method which involves the Ornstein-Uhlenbeck process and the Hamiltonian dynamics, however since the Hamiltonian dynamics are

not integrable, they are split into the kinetic and potential term as described above. These steps are labeled \mathcal{A} , \mathcal{B} and \mathcal{O} , referring to the kinetic energy update, the potential energy update and the Ornstein-Uhlenbeck update respectively. The steps are given below [16]:

$$\hat{\mathcal{F}}_{\Delta t}^{\mathcal{A}}(\mathbf{q}, \mathbf{p}) = (\mathbf{q} + \Delta t \mathbf{p}, \mathbf{p}), \quad (3.8)$$

$$\hat{\mathcal{F}}_h^{\mathcal{B}}(\mathbf{q}, \mathbf{p}) = (\mathbf{q}, \mathbf{p} + \Delta t \mathbf{F}), \quad (3.9)$$

$$\hat{\mathcal{F}}_{\Delta t}^{\mathcal{O}}(\mathbf{q}, \mathbf{p}) = \left(\mathbf{q}, e^{-\gamma \Delta t} \mathbf{p} + \sqrt{k_b T (1 - e^{-2\gamma \Delta t})} \mathbf{R}(t) \right), \quad (3.10)$$

where $\mathbf{R}(t)$ is a vector of i.i.d normal random numbers. This splitting strategy defines a family of schemes based on compositions of the individual components. These schemes are named after the order in which the elements are composed, for example the OBABO scheme is given by

$$\hat{\mathcal{F}}_{\Delta t/2}^{\mathcal{O}} \circ \hat{\mathcal{F}}_{\Delta t/2}^{\mathcal{B}} \circ \hat{\mathcal{F}}_{\Delta t}^{\mathcal{B}} \circ \hat{\mathcal{F}}_{\Delta t/2}^{\mathcal{B}} \hat{\mathcal{F}}_{\Delta t/2}^{\mathcal{O}}, \quad (3.11)$$

noting that steps which are performed twice are given a half-step of $h/2$. Of these methods, BAOAB is considered preferable as it has shown to be stable at larger step-sizes than alternative methods [15].

3.2 Constrained Langevin dynamics

We can write the Langevin dynamics with \mathbf{q} constrained such that $\mathbf{g}(\mathbf{q}) = 0$ as:

$$\dot{\mathbf{q}} = \mathbf{p}, \quad (3.12a)$$

$$\dot{\mathbf{p}} = \mathbf{F} - \gamma \mathbf{p} + \sqrt{2\beta^{-1}\gamma} \eta(t) - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{q}), \quad (3.12b)$$

$$0 = g_j(\mathbf{q}), \quad j = 1, 2, \dots, m, \quad (3.12c)$$

$$0 = \nabla g_j(\mathbf{q})^T \mathbf{p}, \quad j = 1, 2, \dots, m, \quad (3.12d)$$

where m is the number of constraints and $\mathbf{F} = -\nabla U(\mathbf{q})$ is the force. The co-tangency constraint in this formulation, Equation 3.12d, is redundant in theory, but it will be useful in discretising the equations. It can be shown that in the constrained case above, the invariant density can be written as [17]

$$\rho^g \propto \exp(-\beta H(\mathbf{q}, \mathbf{p})) \prod_{ij} \delta_{i=1}^m [g_i(\mathbf{q})] \prod_{i=1}^m \delta[\nabla g_i(\mathbf{q})^T \mathbf{p}]. \quad (3.13)$$

The invariant distribution of \mathbf{q} is then:

$$\rho(\mathbf{q}) \propto \int_{\mathbf{p}'} \exp\left(-\beta(U(\mathbf{q}) + \frac{1}{2} \|\mathbf{p}'\|^2)\right) \prod_{i=1}^m \delta_i[g_i(\mathbf{q})] \prod_{i=1}^m \delta[\nabla g_i(\mathbf{q})^T \mathbf{p}'] d\mathbf{p}' \quad (3.14)$$

$$= \exp(-\beta U(\mathbf{q})) \prod_{ij} \delta_{i=1}^m [g_i(\mathbf{q})] C \quad (3.15)$$

for some constant C , since $\{\nabla g_i(\mathbf{q}) : i \in [1, \dots, m]\}$ are linearly independent, which gives:

$$\rho(\mathbf{q}) \propto \exp(-\beta U(\mathbf{q})) \quad (3.16)$$

for all \mathbf{q} such that $g_i(\mathbf{q}) = 0$ for all $i \in \{1, \dots, m\}$. I've used linear independence of ∇g , but I think the result should hold even if they weren't linearly independent?

3.2.1 Geodesic Integration

This section introduces a constraint-preserving integration scheme for constrained Langevin dynamics which is symplectic in the deterministic case where $\gamma = 0$. Splitting the constrained Langevin equations into potential and kinetic energies and the Ornstein-Uhlenbeck process gives [16]:

$$\begin{bmatrix} d\mathbf{q} \\ d\mathbf{p} \\ \mathbf{0} \end{bmatrix} = \underbrace{\begin{bmatrix} M^{-1}\mathbf{p} dt \\ -\mathbf{G}(\mathbf{q})^T \lambda dt \\ \mathbf{g}(\mathbf{q}) \end{bmatrix}}_A + \underbrace{\begin{bmatrix} \mathbf{0} \\ -\nabla U(\mathbf{q}) dt - \mathbf{G}(\mathbf{q})^T \lambda dt \\ \mathbf{g}(\mathbf{q}) \end{bmatrix}}_B + \underbrace{\begin{bmatrix} \mathbf{0} \\ -\gamma \mathbf{p} dt + \sqrt{2\gamma k_B T \mathbf{M}} d\mathbf{W} - \mathbf{G}(\mathbf{q})^T \lambda dt \\ \mathbf{g}(\mathbf{q}) \end{bmatrix}}_{\mathcal{O}}, \quad (3.17a)$$

where the \mathcal{A} step comes from the kinetic energy, the \mathcal{B} step from the potential and the \mathcal{O} step from the Ornstein-Uhlenbeck process. The solution of part \mathcal{A} is a geodesic curve, meaning that it has no force applied and it runs along the constraint surface, and it can be approximated using SHAKE or RATTLE [17]. Multiple steps of these algorithms can be taken within one step of the proposed algorithm, adding relatively little computational load, since they don't require an evaluation of the force function which is in practice the most expensive part of the algorithm.

For the solution of the \mathcal{B} -step, we get [17]

$$0 = \mathbf{G}(\mathbf{q})\mathbf{p}, \quad (3.18)$$

$$\implies 0 = \frac{d}{dt} \mathbf{G}(\mathbf{q})\mathbf{p} + \mathbf{G}(\mathbf{q}) \frac{d}{dt} \mathbf{p} \quad (3.19)$$

$$\implies \mathbf{G}(\mathbf{q}) \frac{d}{dt} \mathbf{p} = 0 \quad (3.20)$$

$$\implies \lambda = -[\mathbf{G}\mathbf{G}^T]^{-1} \mathbf{G} \nabla U(\mathbf{q}) \quad (3.21)$$

Where we used \mathbf{G} as shorthand for $\mathbf{G}(\mathbf{q})$. This implies that the \mathcal{B} -step can be written as

$$\frac{d}{dt} \mathbf{p} = -\mathbf{\Pi} \nabla U(\mathbf{q}), \quad (3.22)$$

where $\mathbf{\Pi}$ is a projector onto the cotangent space, $\mathbf{\Pi} : \mathcal{M} \rightarrow T\mathcal{M}$, defined by

$$\mathbf{\Pi} = \mathbf{I} - \mathbf{G}^T [\mathbf{G}\mathbf{G}^T]^{-1} \mathbf{G}. \quad (3.23)$$

Therefore the \mathcal{B} -step can be solved exactly as

$$\mathbf{q}(t) = \mathbf{q}(0) \quad (3.24)$$

$$\mathbf{p}(t) = \mathbf{p}(0) - t \mathbf{\Pi}(\mathbf{q}(0)) \nabla U(\mathbf{q}) \quad (3.25)$$

which is called a projected “kick” of the force onto the tangent space of \mathcal{M} .

Similarly for the \mathcal{O} -step we get an SDE of the form

$$\frac{d}{dt} \mathbf{p} = \gamma \mathbf{\Pi} \mathbf{p} + \sqrt{2\gamma k_B T \mathbf{\Pi}} \eta(t), \quad (3.26)$$

which can be weakly solved by

$$\mathbf{p}(t) = e^{-\gamma t \mathbf{\Pi}} \mathbf{p}(0) + \sqrt{k_B T} [I - e^{-2\gamma t \mathbf{\Pi}}]^{1/2} \mathbf{\Pi} \mathbf{R}(t), \quad (3.27)$$

where $\mathbf{R}(t)$ is a vector of i.i.d normal random variables. Since $\boldsymbol{\Pi}^2 = \boldsymbol{\Pi}$, we have

$$e^{-\gamma t \boldsymbol{\Pi}} = \mathbf{I} + (e^{\gamma t} - 1)\boldsymbol{\Pi}, \quad (3.28)$$

$$\implies [I - e^{-2\gamma t \boldsymbol{\Pi}}]^{1/2} = \sqrt{1 - e^{-2\gamma t}}\boldsymbol{\Pi} \quad (3.29)$$

and if we choose the initial velocity to satisfy the tangency condition so that $\boldsymbol{\Pi}\mathbf{p}(0) = \mathbf{p}(0)$, then we have

$$\mathbf{p}(t) = e^{-\gamma t}\mathbf{p}(0) + \sqrt{\beta^{-1}(1 - e^{-2\gamma t})}\boldsymbol{\Pi}\mathbf{R}(t). \quad (3.30)$$

We can now use these solutions of the individual steps to construct integrators based on their compositions.

3.2.2 g-BAOAB

The geodesic BAOAB method uses a similar splitting to that used in unconstrained BAOAB, but ensures that both constraints are satisfied for each step. The name geodesic BAOAB comes from the fact that the solution of the \mathcal{A} step is a geodesic along the surface of the manifold. The algorithm is defined as follows:

Algorithm 1 g-BAOAB Algorithm

```

 $\mathbf{p} \leftarrow \mathbf{p} + \frac{dt}{2}F(\mathbf{q}) + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^*\mathcal{M}$   $\triangleright \frac{1}{2}\mathcal{B}-\text{step}$ 
for  $k$  from 1 to  $K_r$  do
     $(\mathbf{q}, \mathbf{p}) \leftarrow A(\mathbf{q}, \mathbf{p}, \frac{dt}{2K_r})$ 
end for  $\triangleright \frac{1}{2}\mathcal{A}-\text{step}$ 
 $\mathbf{p} \leftarrow a_2\mathbf{p} + b_2\mathbf{M}^{1/2}R + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^*\mathcal{M}$   $\triangleright \mathcal{O}-\text{step}$ 
for  $k$  from 1 to  $K_r$  do
     $(\mathbf{q}, \mathbf{p}) \leftarrow A(\mathbf{q}, \mathbf{p}, \frac{dt}{2K_r})$ 
end for  $\triangleright \frac{1}{2}\mathcal{A}-\text{step}$ 
 $\mathbf{p} \leftarrow \mathbf{p} + \frac{dt}{2}F(\mathbf{q}) + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^*\mathcal{M}$   $\triangleright \frac{1}{2}\mathcal{B}-\text{step}.$ 

```

The solution of the \mathcal{A} -step is calculated using the RATTLE algorithm with a force of zero using a total of K_r steps. As shown by Leimkuhler and Matthews [17], a low value of K_r is sufficient for g-BAOAB and larger values don't provide a significant difference in stability or error. This will be studied in [the next section].

3.3 Example

The g-BAOAB algorithm was used to generate the samples used in Figure 6. The potential function used was $U(x, y, z) = (x^2 - 1)^2 + (y^2 - 1)^2 + (z^2 - 1)^2$. The samples used in Figure 6 were generated using 50,000 steps for 6 different initial points on the surface of the sphere, for a total of 300,000 samples.

In Figure 7 we can see that the algorithm converges to the correct theoretical values of the expectations of the observables. The speed of convergence is slower and the bias is larger, somewhat counter intuitively, for lower stepsizes. This is particularly surprising when we remember that due to the fact that the x -axis represents time rather than iterations, meaning that with lower step-sizes we can converge more slowly to the true value than a larger step-size even while taking 1,000 times more steps. This is a phenomenon that has been seen before in Langevin dynamics - [Is there a reference for this?](#).

The number of RATTLE iterations performed per \mathcal{A} -step, K_r , was seen to have effects on how accurately the solutions satisfied the constraints, as shown in Figure 8. The errors decrease rapidly with higher values of K_r , with $K_r = 2$ already bringing the error in the constraint below machineprecision. The iterations in Figure 8 were run with a step size of 0.2, $\gamma = 1$ and

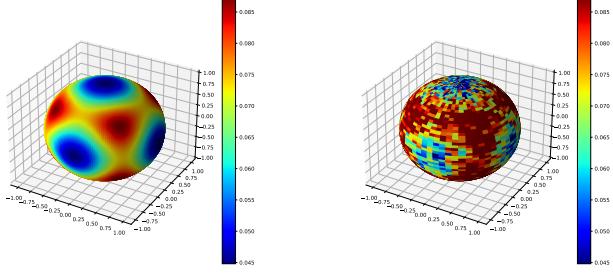


Figure 6: Density compared to g-BAOAB for the function $U(x, y, z) = (x^2 - 1)^2 + (y^2 - 1)^2 + (z^2 - 1)^2$

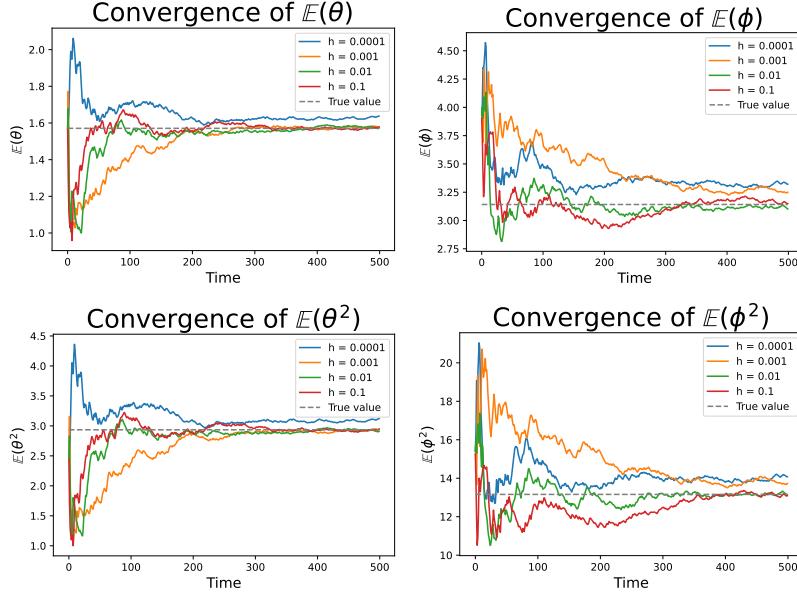


Figure 7: Convergence of ϕ , θ , ϕ^2 and θ^2 for different stepsizes h . (note plot with math h)

$\beta = 1$. While in general it's preferable to run the Newton iterations with an error tolerance as a stopping criterion, this may not always be practical.

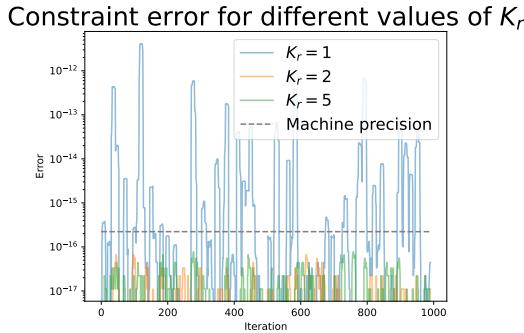


Figure 8: 10 point rolling average of error per iteration for different values of K_r

Again slightly counter intuitively, we see from Figure 9 that convergence is also much slower for lower values of γ . There are two possible causes for this stemming from the fact that γ is a factor in two terms in the equations. A low γ in the damping term could mean that the particles end up moving too quickly and essentially overshoot high density areas. On the other hand, the lower weighting on the random term could mean that the particles get trapped in local maxima in low-density areas. In this case, the latter are not present in the density function, so we can

conclude that it is due to the former explanation.

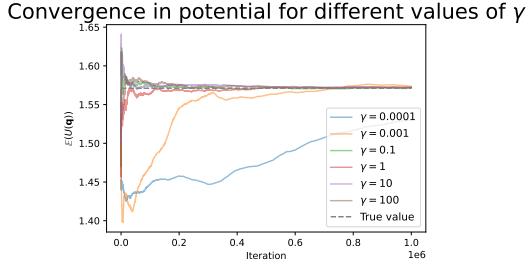


Figure 9: Convergence in expected potential $U(\mathbf{q})$ for different values of γ

4 Score based diffusion

Now that we have seen how the g-BAOAB integrator can be used for solving SDEs subject to holonomic constraints, we move into the area of score-based diffusion, where g-BAOAB can be used to noisify data on the constraint manifold, and reverse the noising process while remaining on the constraint manifold. Key to this process is approximating the score function of the marginal distribution.

4.1 Score approximation

For an SDE of the form given in Equations 3.12, our aim is to construct an estimator of the score function of the marginal distribution of \mathbf{q} , $\hat{s}_\theta(\mathbf{q}, t) \approx \nabla_{\mathbf{q}} \log \rho_t(\mathbf{q})$, where ρ_t is the distribution at time t in the noising process. The score function is intractable, and needs to be approximated rather than calculated exactly. The fact that the score-based model $s_\theta(\mathbf{x})$ is independent of the normalizing constant allows us to use a variety of models, since the normalizing constant does not need to be tractable. Usually, training the model score function involves minimizing the Fisher divergence [32]:

$$\mathcal{L}(\hat{s}_\theta) = \int_0^T \mathbb{E}_\rho [\nabla_{\mathbf{q}} \log \rho_t(\mathbf{x}) - \hat{s}_\theta(\mathbf{x}, t)]_2^2 dt, \quad (4.1)$$

where \mathbb{E}_ρ signifies that the expectation is taken over $\rho(\mathbf{x})$. Although the Fisher divergence includes the *true* score function which is intractable, Hyvärinen showed that it can remarkably be written purely in terms of the *model* score function, plus a constant. In the following theorem we deal a score model which is not time-dependent, and write $\hat{s}_\theta(\mathbf{x})$ rather than $\hat{s}_\theta(\mathbf{x}, t)$, but the time-dependent version will follow.

Theorem 4.1. [12] *Assume that the model score function $\hat{s}(\mathbf{x})$ is differentiable, the data distribution is differentiable, the expectation of the squared model density and the expectation of the squared norm of the score are finite. Also assume that $\rho(\mathbf{q})\hat{\rho}(\mathbf{q}) \rightarrow 0$ as $\|\mathbf{q}\| \rightarrow \infty$. Then the objective in Equation 4.1 can be expressed as*

$$\mathcal{L}(\hat{s}_\theta) = \mathbb{E}_\rho \left[\text{Tr}(\nabla_x \hat{s}_\theta(\mathbf{x})) + \frac{1}{2} \|\hat{s}_\theta(\mathbf{x})\|_2^2 \right] + C, \quad (4.2)$$

where $\nabla_x \hat{s}_\theta(\mathbf{x})$ is the Jacobian matrix of the score model evaluated at \mathbf{x} .

Thus we can train the score-based model by calculating its norm and divergence on samples of the data, without any explicit knowledge of the true score function. This fact is key to score matching and thus score-based generative modelling. The proof of the theorem is based on

integration by parts; we have that

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) = \int \rho(\mathbf{x}) \left[\frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{x})\|^2 + \frac{1}{2} \|\mathbf{s}(\mathbf{x})\|^2 - \hat{\mathbf{s}}_\theta(\mathbf{x})^T \mathbf{s}(\mathbf{x}) \right] d\mathbf{x} \quad (4.3)$$

where the integration domain has been omitted for simplicity. Then it must be shown that

$$\int \rho(\mathbf{q}) \hat{\mathbf{s}}_\theta^T \mathbf{s}(\mathbf{q}) d\mathbf{q} = \int \rho(\mathbf{q}) \nabla \hat{\mathbf{s}}_\theta d\mathbf{q}. \quad (4.4)$$

The partial integration is based on a generalization of the following statement to multiple dimensions: Given a one dimensional density function ρ and a function f , we have

$$\int \rho(\mathbf{q}) (\log \rho)'(\mathbf{q}) f(\mathbf{q}) d\mathbf{q} = - \int \rho(\mathbf{q}) f'(\mathbf{q}) d\mathbf{q}, \quad (4.5)$$

under the assumption that f is differentiable, and that the limit goes to zero at infinity as stated in the theorem. Using this theorem, we can use the following estimator of the Fisher divergence to train the score model:

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\text{Tr}(\nabla_x \hat{\mathbf{s}}_\theta(\mathbf{x})) + \frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{x})\|_2^2 \right]. \quad (4.6)$$

Note that this theorem implies that the minimum value for the loss is $-C = -\mathbb{E}_\rho(\frac{1}{2} \|\mathbf{s}(\mathbf{q})\|^2)$, so losses can and generally should be negative but we won't know what the minimum value, or even what is an acceptable value, in advance. Finally, we can generalize this result to a continuous time-dependent vector field $\hat{\mathbf{s}}_\theta$ which aims to minimize the loss function in Equation 4.1, where we get the loss function

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) \approx \frac{1}{T} \sum_{j=1}^M \frac{1}{N} \sum_{i=1}^N \left[\text{Tr}(\nabla_{\mathbf{x}_i} \hat{\mathbf{s}}_\theta(\mathbf{x}, t_j)) + \frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{x}_i, t_j)\|_2^2 \right] \quad (4.7)$$

where the times t_j , $j \in [1, \dots, M]$ are sampled uniformly.

4.1.1 Fisher divergence on Riemannian manifolds

It can be shown that on a Riemannian manifold, we get an expression which is directly analogous to the Fisher divergence used in Euclidean space. If we define the loss function as

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) = \int_{\mathcal{M}} \frac{1}{2} \|\nabla_{\mathbf{q}} \log \rho(\mathbf{q}) - \hat{\mathbf{s}}_\theta(\mathbf{q})\| \quad (4.8)$$

Theorem 4.2. [8] *Under sufficient regularity of $\rho_{t|s}(\mathbf{q}_t | \mathbf{q}_s) \mathbf{s}$ for all $t, s \in (0, T]$ with $t > s$, we have*

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) = \int_{\mathcal{M}} \left(\frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{q})\|^2 + \Delta \hat{\mathbf{s}}_\theta(\mathbf{q}) \right) \rho(\mathbf{q}) d\mathbf{q} \quad (4.9)$$

The proof involves expanding $\rho_t(\mathbf{q}_t)$ as $\rho_t(\mathbf{q}_t) = \int_{\mathcal{M}} \nabla_{\mathbf{q}} \log \rho_{t|s}(\mathbf{q}_t | \mathbf{q}_s) \rho_s(\mathbf{q}_s) d\mathbf{q}_s$, and therefore requires regularity conditions on $\rho_{t|s}$ which are discussed in Appendix C. This theorem in combination with a relationship between the Riemannian and the Euclidean divergence will, under certain conditions, let us express the loss function in terms of the Euclidean divergence. First, we denote by $\mathcal{X}(\mathcal{M})$ the space of smooth tangent vector fields to \mathcal{M} . Then:

Theorem 4.3. [27] *If $u \in \mathcal{X}(\mathbb{R}^n)$, and the restriction of u on \mathcal{M} , $u|_{\mathcal{M}}$ is infinitesimally constant in the normal directions of \mathcal{M} , then for $\mathbf{q} \in \mathcal{M}$, $\Delta u(\mathbf{q}) = \Delta_E u(\mathbf{q})$, where Δ_E is the standard Euclidean divergence.*

In other words, given a vector field which is tangent to the manifold at each point, $\mathbf{u}(\mathbf{q}) \in T_{\mathbf{q}}\mathcal{M} \forall \mathbf{q} \in \mathcal{M}$, its divergence is equal to its Euclidean divergence, and since we have an explicit projection matrix $\mathbf{\Pi}$ onto $T_{\mathbf{q}}\mathcal{M}$, we can construct such a vector field using projections.

4.2 Perturbing data

For regions of low-density data, it can be difficult to accurately estimate the score. Since score matching aims to minimize the Fisher divergence (4.1), the low density regions of the [space] naturally contribute relatively little to the overall loss, since the norms are weighted by $\rho(\mathbf{x})$. To remedy this we can perturb the data points with noise, which will populate low density areas with points. Generally multiple scales of Gaussian noise have been used to ensure that we cover low-density areas [33]. A noise-conditional score-based model is then trained on the noisy data. The training objective for the noise-conditional score-based model is given by a weighted sum of Fisher divergences for different noise scales:

$$\sum_{i=1}^L \zeta(i) \mathbb{E}_{\rho_i} \mathbf{x} [\nabla_{\mathbf{x}} \log \rho_i(\mathbf{x}) - \hat{s}_{\theta}(\mathbf{x}, i)]_2^2 \quad (4.10)$$

where ρ_i is the distribution of x perturbed with noise with standard deviation σ_i and $\zeta(i)$ is some positive weighting function, usually taken to be $\zeta(i) = \sigma_i$.

However recently, instead of using this discrete noising process, an SDE has been used to perturb the data, effectively generalizing the number of noise scales to infinity [33]. In perturbing the data using an SDE, we must bear in mind that the invariant distribution of the SDE used should be known, since we will need to sample from the invariant distribution as the *initial* distribution of the reverse process.

Choosing a suitable distribution on manifolds can present challenges. The uniform distribution seems a sensible choice, however a uniform distribution cannot be defined if the manifold is not closed. In that case, we would need to choose a potential function $U(\mathbf{q})$ so that the Langevin dynamics give a stationary distribution we can sample from. Two distributions proposed by De Bortoli et al. are the “Riemannian normal” and the “exponential wrapped” Gaussian [8]. For a closed manifold, the uniform distribution on the manifold is the invariant distribution given by Langevin dynamics with a force term of zero. In the overdamped limit this equation is a Brownian motion on the manifold and is given by

$$\dot{\mathbf{q}} = \sqrt{2k_B T \gamma} \mathbf{M}^{1/2} \eta(t) - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{q}), \quad (4.11)$$

$$g_i(\mathbf{q}) = 0, \quad i = 1, \dots, m. \quad (4.12)$$

However we can also use a system of underdamped Langevin dynamics, and since the invariant distribution is proportional to $e^{-U(\mathbf{q})}$, letting $U(\mathbf{q}) = 0$ gives $p(\mathbf{q}) \propto e^{-0} = 1$ which is a uniform distribution on the manifold. This gives us the familiar noising process

$$\dot{\mathbf{q}} = \mathbf{M}^{-1} \mathbf{p} \quad (4.13)$$

$$\dot{\mathbf{p}} = -\gamma \mathbf{p} + \sqrt{2k_B T \gamma} \mathbf{M}^{1/2} \eta(t) - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{q}) \quad (4.14)$$

$$0 = g_j(\mathbf{q}), \quad j = 1, 2, \dots, m \quad (4.15)$$

$$0 = \nabla g_j(\mathbf{q})^T \mathbf{M}^{-1} \mathbf{p}, \quad j = 1, \dots, m \quad (4.16)$$

from which we can generate samples using the g-BAOAB algorithm.

4.3 Time Reversal

Central to the idea of score-based diffusion is the reverse SDE [33], which allows us to move backwards in time from a known distribution to the unknown data distribution. Now that we have a noising process using constrained underdamped Langevin dynamics, we aim to construct the reverse SDE corresponding to our formulation of constrained Langevin dynamics. **This will allow us to move backwards from noised data to data which follows the data distribution.**

4.3.1 Time reversal of constrained Langevin dynamics

As shown by [4], any SDE in Euclidean space

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) + \mathbf{g}(t)dW \quad (4.17)$$

yields a time-reversed SDE, given by

$$d\mathbf{x} = [\mathbf{f}(x, t) - \nabla_x \log p_t(x)] dt + dW, \quad (4.18)$$

where can avoid explicitly reversing the time by imagining $d\mathbf{x}$ as another process moving forward through time, where $p_t(\mathbf{x})$ is the distribution of \mathbf{x} at time t in the *denoising* process [4].

In order to perform score-based diffusion on the manifold however, we must find the reverse-time formulation of the constrained Langevin dynamics SDE given in Equations 3.12. As shown in by Lelievre et al., if \mathcal{G}_t^f is the generator of the forward dynamics, then the backward dynamics can be defined through its generator by first reversing the momentum in the initial condition, then reversing it *back* after the time evolution [22]:

$$\mathcal{G}_{t'}^b = \mathcal{R}\mathcal{G}_{T-t'}^f\mathcal{R} \quad (4.19)$$

where $\mathcal{R} : \phi \rightarrow \phi \circ S$ where S is an operator which flips the momentum, $S(\mathbf{q}, \mathbf{p}) = S(\mathbf{q}, -\mathbf{p})$

This means that constrained Langevin dynamics of the form given in Equations 3.12 yield the reverse equations:

$$\dot{\mathbf{q}} = -\mathbf{M}^{-1}\mathbf{p} \quad (4.20)$$

$$\dot{\mathbf{p}} = -\mathbf{F}(\mathbf{q}) - \gamma\mathbf{p} + \sqrt{2k_B T \gamma} \eta(t) - \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{q}) \quad (4.21)$$

$$0 = g_j(\mathbf{q}), \quad j = 1, 2, \dots, m \quad (4.22)$$

$$0 = \nabla g_j(\mathbf{q})^T M^{-1} \mathbf{p}, \quad j = 1, 2, \dots, m \quad (4.23)$$

Modifying the g-BAOAB algorithm for the reversed dynamics, we have that the \mathcal{B} -step can be solved exactly as

$$\mathbf{q}(t) = \mathbf{q}(0), \quad \mathbf{p}(t) = \mathbf{p} - t\Pi(\mathbf{q}(0))\mathbf{F}. \quad (4.24)$$

The \mathcal{O} -step doesn't change, and the \mathcal{A} -step is simply given by reversing the sign of the velocity. For completeness, the algorithm is given below:

5 Connection with ODEs

One issue with score-based diffusion is that it doesn't give a likelihood for the density of the model, and it also doesn't necessarily maximize the likelihood of the observations. In order to link the score-based diffusion SDE to the likelihood, we use continuous normalizing flows (CNFs), which are ODEs describing the rate of change of a probability density function.

Algorithm 2 Reverse g-BAOAB algorithm

$\mathbf{p} \leftarrow \mathbf{p} - \frac{dt}{2} F(\mathbf{q}) + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^* \mathcal{M}$	$\triangleright \frac{1}{2} \mathcal{B}-\text{step}$
for k from 1 to K_r do	
$(\mathbf{q}, -\mathbf{p}) \leftarrow A(\mathbf{q}, -\mathbf{p}, \frac{dt}{2K_r})$	$\triangleright \frac{1}{2} \mathcal{A}-\text{step}$
end for	
$\mathbf{p} \leftarrow a_2 \mathbf{p} + b_2 \mathbf{M}^{1/2} R + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^* \mathcal{M}$	$\triangleright \mathcal{O}-\text{step}$
for k from 1 to K_r do	
$(\mathbf{q}, -\mathbf{p}) \leftarrow A(\mathbf{q}, -\mathbf{p}, \frac{dt}{2K_r})$	$\triangleright \frac{1}{2} \mathcal{A}-\text{step}$
end for	
$\mathbf{p} \leftarrow \mathbf{p} - \frac{dt}{2} F(\mathbf{q}) + \sum_j \mu_j G_j(\mathbf{q}) \in T_q^* \mathcal{M}$	$\triangleright \frac{1}{2} \mathcal{B}-\text{step}$

5.1 Normalizing flows

Normalizing flows provide a mechanism for defining probability distributions based on a base distribution. In unconstrained problems, we express points \mathbf{x} as a transformation of a real vector \mathbf{u} sampled from the base distribution $\rho_u(\mathbf{u})$. For a flow based model, T must be invertible and T^{-1} and T must be differentiable. Then we can use the change of variables formula

$$\rho_x(\mathbf{x}) = \rho_u(T^{-1}(\mathbf{x})) |\det J_T(\mathbf{u})|^{-1}, \quad (5.1)$$

where J_T is the Jacobian matrix of T , using a change of variables [25]. The name *normalizing* flow refers to the fact that the change of variables gives a normalized density.

It's been shown by Papamakarios et al. that for any pair of sufficiently well behaved distributions τ and ρ , there is exists a diffeomorphism ϕ which maps τ to ρ , where the requirements on the distributions are only that they are differentiable and strictly positive [25].

5.1.1 Kullback-Leibler divergence

The Kullback-Leibler (KL) divergence is one of the most popular objective functions used in flow-based models (cite uses); it is given by the expectation of the logarithmic difference between probabilities. The forward KL divergence between the target distribution ρ and the modelled distribution $\hat{\rho}$ is given by:

$$D_{KL}[\rho|\hat{\rho}] = \mathbb{E} \left[\log \left(\frac{\rho}{\hat{\rho}} \right) \right] \quad (5.2)$$

$$= -\mathbb{E}_\rho [\log \rho \mathbf{x}] + C \quad (5.3)$$

$$= \mathbb{E}_\rho [\log \rho_u(T^{-1}(\mathbf{x}))] + \log |\det J_{T^{-1}}(\mathbf{x})| \quad (5.4)$$

The minimizer of the KL divergence can be approximated by minimizing the loss function

$$\mathcal{L}(\rho, \tau) = -\frac{1}{N} \sum_{n=1}^N (\log \tau(T^{-1}(\mathbf{x}_n)) + \log |\det J_{T^{-1}}(\mathbf{x})|) \quad (5.5)$$

which can be minimized in practice using gradient descent, using its gradients with respect to the parameters of ρ , θ_ρ and the parameters of τ, θ_τ :

$$\nabla_{\theta_\rho} \mathcal{L}(\rho, \tau) = -\frac{1}{N} \sum_{n=1}^N \nabla_{\theta_\rho} \log \tau(T^{-1}(\mathbf{x}_n)) + \nabla_{\theta_\rho} \log |\det J_{T^{-1}}(\mathbf{x})| \quad (5.6)$$

$$\nabla_{\theta_\tau} \mathcal{L}(\rho, \tau) = -\frac{1}{N} \sum_{n=1}^N \log \tau(T^{-1}(\mathbf{x}_n)). \quad (5.7)$$

The KL divergence is useful since we can estimate it using samples of the target distribution, and it can be shown that maximising the likelihood is asymptotically equivalent to minimizing the KL divergence, ie. they are equivalent as the number of samples goes to infinity.

5.1.2 Continuous normalizing flows

Now considering a continuous time version of the transformation gives an ODE

$$\dot{\mathbf{z}} = \mathbf{g}(t, \mathbf{z}) \quad (5.8)$$

where the function \mathbf{g} is Lipschitz continuous, and continuous in t . The transformed variable \mathbf{x} is then given by $\mathbf{x} = \mathbf{u} + \int_{t=0}^{t=1} g_\phi(t, \mathbf{z}_t) dt$. We define the *pushforward* map so that the pushforward of the density by the transformation is the density of the transformed variable, ie. $\mathbf{z} = \phi(\mathbf{z}_0)$ with $\mathbf{z}_0 \sim \tau \implies \mathbf{z} \sim \phi_*\tau$ [24]. The ODE can then be solved using, for example, a simple ODE discretisation such as Euler's method:

$$\mathbf{z}_{t+h} \simeq \mathbf{z}_t + hg(t, \mathbf{z}_t) \quad (5.9)$$

It can be shown that the rate of change of the log-density can be characterized in terms of the trace of the Jacobian matrix [25]:

$$\frac{d}{dt} \log \rho(\mathbf{z}_t) = -\text{Tr}(J_{g_\phi(t)}(\mathbf{z}_t)), \quad (5.10)$$

which is obtained from a Fokker-Planck equation with a diffusion term of zero [25]. Thus we can evaluate both the transformation and the log-density by integrating:

$$\begin{bmatrix} \mathbf{x} \\ \log \rho(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \log \rho(\mathbf{u}) \end{bmatrix} + \int_{t=0}^{t=1} \begin{bmatrix} g(t, \mathbf{z}_t) \\ -\text{Tr}(J_{g_\phi(t)}(\mathbf{z}_t)) \end{bmatrix} dt \quad (5.11)$$

using ODE solvers and then update the parameters based on a backward pass, in order to minimize the KL divergence.

5.2 Riemanninan normalizing flows

Continuous normalizing flows can also be constructed on Riemannian manifolds. Consider a particle $\mathbf{q}(t)$ lying on the manifold \mathcal{M} . Then let f_θ be a vector field, $f_\theta : \mathcal{M} \times \mathbb{R} \rightarrow T\mathcal{M}$, where $T_{\mathbf{q}}\mathcal{M}$ is the tangent space at \mathbf{q} and $T\mathcal{M} = \bigcap_{\mathbf{q} \in \mathcal{M}} T_{\mathbf{q}}\mathcal{M}$ is the tangent bundle. Then the particle's time-evolution according to f_θ is given by [24]:

$$\frac{d}{dt} \mathbf{q}(t) = f_\theta(\mathbf{q}(t), t). \quad (5.12)$$

Then when starting at an initial position $\mathbf{q}(0) = \mathbf{q}_0$, the flow operator $\phi : \mathcal{M} \times \mathbb{R} \rightarrow \mathcal{M}$ is defined so that

$$\phi(\mathbf{q}_0, t) = \mathbf{q}(t) \quad (5.13)$$

gives the position of the particle at time t . Mathieu and Nickel give a generalization of the log-probability flow given in Equation 5.33 for the flow on a Riemannian manifold in the following theorem:

Theorem 5.1. [24] *Let $\mathbf{q}(t)$ be a continuous manifold-valued random variable which is described by the ODE from 5.12 with probability density $\rho_\theta(\mathbf{q}(t))$. Then the change in log-probability follows a differential equation given by*

$$\frac{d}{dt} \log \rho_\theta(\mathbf{q}(t)) = -\Delta_{\mathcal{M}}(f_\theta(\mathbf{q}(t), t)). \quad (5.14)$$

Therefore we can also solve for the change in log-density between two manifold-valued random variables, and notably we have that

$$\log\left(\frac{\rho}{\tau}\right) = - \int_0^1 \Delta f_\theta(\mathbf{q}(t), t) dt. \quad (5.15)$$

This will allow us to evaluate the transformation by integrating using an ODE solver, and then to minimize the log-difference, or KL divergence, by backpropogating and updating the parameters.

5.3 Score-based flow

In the unconstrained case, score-based diffusions can be much more efficient to train than CNFs, because minimizing the maximum likelihood while training CNFs requires running an expensive ODE solver for each optimization step. However an issue with score-based diffusion is that training the model doesn't necessarily maximize the likelihood, which may be desirable in certain cases. However, we can model the flow of the marginal distribution of the SDE, as shown by Song [31]. For a stochastic differential equation given by

$$\dot{\mathbf{x}} = f(\mathbf{x}, t) + \mathbf{g}(t)\eta(t). \quad (5.16)$$

There is a deterministic process sharing the same marginal densities as the SDE, which is an example of a CNF, given by:

$$\dot{\mathbf{x}} = f(\mathbf{x}, t) - \frac{1}{2}g(t)^2\hat{\mathbf{s}}_\theta(\mathbf{x}, t) \quad (5.17)$$

called the *probability flow ODE*. The ODE is a continuous normalizing flow (CNF), and it can be used to generate samples using numerical ODE solvers starting with a sample from the base distribution. This CNF also uses the score function, which can then either be taken from the diffusion model in order to calculate its likelihood, or trained using the maximum likelihood objective of the CNF. The main result from Song's paper was the bounding of the likelihood of *score-based* models through their relationship with the probability flow ODE. The results are as follows:

Theorem 5.2. *Let $\rho(\mathbf{x})$ be the unknown data distribution and $\tau(\mathbf{x})$ be a known base distribution. If $d\mathbf{x} = f(\mathbf{x}, t)dt + g(t)dW$, then (under certain regularity conditions)*

$$D_{KL}(\rho||\rho_{SDE}) \leq \mathcal{J}_{SM}(\theta; g(\cdot)^2) + D_{KL}(\rho_T|\pi) \quad (5.18)$$

where \mathcal{J}_{SM} is the score-matching loss:

$$\mathcal{J}_{SM}(\theta; \zeta) = \frac{1}{2} \int_0^T \mathbb{E}_{\rho_t} [\zeta(t) \|\nabla_x \log \rho_t - \mathbf{s}_\theta(\mathbf{x}, t)\|_2^2] dt, \quad (5.19)$$

for a some positive weighting function ζ .

The regularity conditions required are given in Appendix C. This theorem shows that training with a weighted combination of score-matching losses similarly to Equation 4.10 rather than training based on the maximum likelihood leads to a bounding of the KL divergence, of the score-based density. Note that the function g used in the score matching loss is the time-dependent weighting of the stochastic term in the SDE; if we have $\frac{d}{dt}g = 0$, then the score matching loss reduces to the Fisher divergence. The next theorem shows that under certain conditions, this bound becomes an equality. Let ρ_θ^{SDE} be the marginal distribution of \mathbf{q}_0 obtained from the score-based diffusion, ie. our approximation of the data distribution, and ρ_θ^{ODE} be that obtained from the score-based flow model. Then:

Theorem 5.3. Suppose $\rho_1(\mathbf{x})$ and $\rho_2(\mathbf{x})$ have continuous second-order derivatives and finite second moments. Then if $d\mathbf{x} = f(\mathbf{x}, t)dt + g(t)dW$, then under the conditions $q_T = \pi$ and $s_\theta(\mathbf{x}, t) = \nabla \log q_t(\mathbf{x})$ for all $t \in [0, T]$ we have the following equivalence in distributions

$$\hat{\rho}^T SDE = \hat{\rho}_{ODE} = q \quad (5.20)$$

and

$$D_{KL}(\rho || \rho_{SDE}) = \mathcal{J}_{SM}(\theta; g(\cdot)^2) + D_{KL}(\rho_T || \pi) \quad (5.21)$$

Therefore, in the Euclidean overdamped setting, if we set our prior distributions to τ , then the marginal distribution obtained from both methods is equal to the true distribution, if the conditions are held. However, as noted in the paper, the score function will rarely be estimated exactly, so in practice there is no guarantee that minimizing one likelihood will decrease the other.

5.3.1 Score-based flow on Riemannian manifolds

Song's derivation of the probability flow ODE comes from the Fokker-Planck equation for the SDE. In order to examine the probability flow for underdamped Langevin dynamics subject to constraints, we will derive the **Fokker-Planck equation**. In order to sample from the uniform distribution \mathcal{M} , De Bortoli et al. used the overdamped Langevin dynamics with a force $\mathbf{F} = \mathbf{0}$ to noise the data, which constitutes a Brownian motion on the manifold [8]. Using a Brownian motion yields a probability-flow ODE of

$$\frac{d}{dt}\rho_t(\mathbf{q}) = \frac{1}{2}\Delta_{\mathcal{M}}\rho_t(\mathbf{q}), \quad (5.22)$$

which is a generalization of the Euclidean Fokker-Planck equation for Brownian motion. In the case of underdamped Langevin dynamics, in order to converge to the uniform distribution we can still use a force of $\mathbf{F} = \mathbf{0}$, but the \mathbf{q} dynamics are no longer a Brownian motion. Therefore the probability flow will take a different form, which we will discuss below.

5.3.2 Score-based flow for constrained (underdamped) Langevin dynamics with zero force.

In order to discuss the probability flow of constrained Langevin dynamics, we first introduce some new concepts and notation for constrained systems. The Poisson bracket $\{\cdot, \cdot\}$ of two smooth functions f_1, f_2 acting on (\mathbf{q}, \mathbf{p}) is defined as

$$\{f_1, f_2\} = (\nabla_{\mathbf{q}} f_1)^T \nabla_{\mathbf{p}} f_2 - (\nabla_{\mathbf{p}} f_1)^T \nabla_{\mathbf{q}} f_2 \quad (5.23)$$

Let $\{\cdot, \cdot\}_\Xi$ be the Poisson bracket for the constrained system, defined by

$$\{\rho_1, \rho_2\} = \nabla \rho_1^T J_G \nabla \rho_2 \quad (5.24)$$

where J_G is the constrained symplectic matrix given by

$$J_G = J - J \nabla \Xi \left[\nabla \Xi^T J \nabla \Xi \right]^{-1} \nabla \Xi^T J, \quad (5.25)$$

J is the symplectic matrix $J \in \mathbb{R}^{m \times m}$

$$J = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}, \quad (5.26)$$

and $\Xi(\mathbf{q}, \mathbf{p}) \in \mathbb{R}^m$ is a vector of constraints. In our case

$$\Xi(\mathbf{q}, \mathbf{p}) = \begin{bmatrix} \mathbf{g}(\mathbf{q}) \\ \mathbf{G}(\mathbf{q})\mathbf{p} \end{bmatrix} \quad (5.27)$$

Now we will make use of the following theorem which gives the Markov generator of the constrained underdamped Langevin dynamics:

Theorem 5.4. [21] *For constraints given by [reference them], the Markov generator of the constrained process 3.12 where the matrix $\mathbf{G}(\mathbf{q})$ is invertible for all \mathbf{q} can be written*

$$\mathcal{G}_\Xi = \mathcal{G}_\Xi^{ham} + \mathcal{G}_\Xi^{thm} \quad (5.28)$$

where the Hamiltonian generator is

$$\mathcal{G}_\Xi^{ham} = \{-, H\}_\Xi \quad (5.29)$$

and the fluctuation-dissipation part can be written

$$\mathcal{G}_\Xi^{thm} = \frac{1}{\beta} \Delta_{\mathbf{p}} (e^{-\beta H} \Pi \nabla_{\mathbf{p}}). \quad (5.30)$$

Now to examine the flow of the system with a force of $\mathbf{F} = \mathbf{0}$, consider the case where $\nabla_{\mathbf{q}} U = \mathbf{0}$. We suppose that ρ is a probability distribution defined only over \mathbf{q} , meaning that $\nabla_{\mathbf{p}} \rho = \mathbf{0}$. Then we have that

$$\frac{d}{dt} \rho(\mathbf{q}) = \mathcal{G}_\Xi^{ham} = \{\rho; H\}_\Xi \quad (5.31)$$

$$= \{\rho, H\} - \{\rho, \Xi\} \Gamma^{-1} \{\Xi, H\}, \quad (5.32)$$

Then we have that $\{\rho, \Xi\} = \mathbf{0}$, however the second term in Equation 5.32 does not simplify, so we are unable to use the probability flow in the same way as Song et al. and De Bortoli et al. where the flow can be written in terms of the score function [31, 8].

5.4 Moser Flows

Moser flows on Riemannian manifolds have been proposed to improve upon some issues with Riemannian CNFs, being that they can be slow to train, and they also can be slow to converge [27]. The idea behind Moser flows is to construct an interpolant between the source and target distributions, ie. a map $\alpha : [0, 1] \times \mathcal{M} \rightarrow \mathbb{R}_{>0}$ such that $\alpha_0 = \tau$, $\alpha_1 = \rho$ and $\int_{\mathcal{M}} \alpha_t dV = 1$ for all $t \in [0, 1]$. In this way the Moser flow will generate a CNF, but by limiting the *flow space*, we can express the learned probability in a desirable way as the divergence of a vector field. Then the flow of the CNF is defined by

$$\frac{d}{dt} \phi(\mathbf{q}) = v_t(\phi_t) \quad (5.33)$$

where v is a time dependent vector field, $v_t \in T_{\mathbf{q}} \mathcal{M}$, defined so that the flow satisfies the continuous normalization equation:

$$\Phi_t^* \alpha_0 = \alpha_t, \quad (5.34)$$

meaning that the way the pullback transforms densities agrees with the interpolant we defined on the density space. For $t = 1$, we get the map Φ_1 such that $\hat{\rho} = \Phi_1^*(\tau)$. Now the following theorem shows that if we construct a vector field $u_t \in T\mathcal{M}$ such that

$$\Delta(u_t) = -\frac{d}{dt} \alpha_t \quad (5.35)$$

and define the vector field v_t as

$$v_t = \frac{u_t}{\alpha_t} \quad (5.36)$$

then the diffeomorphism Φ defined by u_t and v_t satisfies the normalization equation as required.

Theorem 5.5. *The diffeomorphism $\Phi = \Phi_1$, defined by the ODE in 5.33 and the vector field v_t in equation 5.36 solves the normalization equation 5.34.*

This theorem is derived from works by Moser, after whom the Moser flow is named. Essentially, it tells us the vector field which defines the flow induced by an interpolation in the probability distribution space. As an example, the simplest choice of interpolant α_t is

$$\alpha_t = (1-t)\tau + t\rho, \quad (5.37)$$

which means that the divergence of the vector field u_t should be constant:

$$\text{div}(u) = -\frac{d}{dt}\alpha_t = \tau - \rho. \quad (5.38)$$

In this case the vector field v_t is given by

$$v_t = \frac{u}{(1-t)\rho + t\pi}, \quad (5.39)$$

and the model density $\hat{\rho}$ is defined by

$$\hat{\rho} = \tau - \Delta u. \quad (5.40)$$

Rozen et al. give a description of how Moser flows can be used for generative modelling over Euclidean submanifolds [27]. Specifically, the following theorem speaks to the universality of Moser flows for Euclidean submanifolds. By parameterizing the vector field u as the projection onto the $T_x\mathcal{M}$ of an ambient vector field, we construct our network in the same way as in Section 6.2, and define

$$\mathbf{u}(\mathbf{q}) = \mathbf{\Pi}v_\theta \quad (5.41)$$

Then from Theorem 4.3 since $u_t(\mathbf{x}) \in T_{\mathbf{x}}\mathcal{M}$ for all $\mathbf{x} \in \mathcal{M}$, the Riemannian divergence is given by the Euclidean divergence for all $\mathbf{x} \in \mathcal{M}$. The following theorem validates this approach:

Theorem 5.6. [27] *Given an orientable, compact, boundaryless, connected, differentiable n -dimensional submanifold $\mathcal{M} \in \mathbb{R}^d$, with $n < d$, and a continuous target density $\rho : \mathcal{M} \rightarrow \mathbb{R}_{>0}$, there exists for each $\epsilon > 0$ a vector field \mathbf{u} so that $\hat{\rho}$ defined by 5.40 satisfies:*

$$\max_{\mathbf{q} \in \mathcal{M}} |\rho(\mathbf{q}) - \hat{\rho}(\mathbf{q})| < \epsilon. \quad (5.42)$$

It's worth noting here that not all manifolds given by the set of constraints ?? are connected, so the usefulness of this theorem is slightly limited in our case. Therefore we will aim to compare Moser flows to score-based diffusion on compact, connected manifolds, but score-based diffusion can be employed on a wider range of manifolds, as is done in Section 6.2.

5.4.1 Training the model

The paper notes that \mathbf{u} integrates to 1 over \mathcal{M} by construction, since we know that ρ integrates to 1 over the manifold, and Stokes' theorem proves that the integral of $\Delta \mathbf{u}$ over \mathcal{M} is zero. In order to ensure *positivity* over the manifold we choose a small value ϵ and then define

$$\hat{\mu}_+(\mathbf{x}) = \max\{\epsilon, \hat{\mu}(\mathbf{x})\} \quad (5.43)$$

$$\hat{\mu}_-(\mathbf{x}) = \min\{\epsilon, \hat{\mu}(\mathbf{x})\} \quad (5.44)$$

so that $\hat{\mu}_+ + \hat{\mu}_- = \hat{\mu}$. Here $\hat{\mu}_-$ is used to [encourage] $\hat{\rho} > \epsilon$, as the term is zero for all $\hat{\rho} \leq \epsilon$. Then the loss can be given by

$$\mathcal{L}(\theta) = -\mathbf{E}_{\mu} \log \hat{\mu}(\mathbf{x}) + \lambda \int_{\mathcal{M}} \hat{\mu}(\mathbf{x}) dV, \quad (5.45)$$

where the first term is the negative log-likelihood of the observations and can be approximated by

$$\mathbb{E}_{\rho} \log \hat{\rho}_+(\mathbf{x}) \quad (5.46)$$

and the integral in the second term is approximated by a Monte Carlo estimate over some distribution τ over \mathcal{M} :

$$\int_{\mathcal{M}} \bar{\mu}_- dV \sim \frac{1}{l} \sum_{i=1}^l \mu_i. \quad (5.47)$$

In this way training seeks to maximize the expectation of the observations using the first term, while ensuring that the density is strictly positive using the second term.

Note: this needs compactness. It is noted in [27] that the additional term does not change the fact that the unique minimum of the loss is the true density ρ , for $\lambda > 1$ and sufficiently small ϵ , and since the loss used is the KL divergence, this implies that $\hat{\rho}$ is the maximum likelihood estimator within the Moser flow framework.

As noted by [8], both score-based generative models and Moser flow models interpolate between the base density π and the [data] density ρ . The key difference is that while Moser flows interpolate on the density space by defining a map $MF(\rho, \pi)$, the score-based generative model instead interpolates by defining a path on $\mathcal{M} \times \mathcal{M} \rightarrow \mathcal{M}$. This path then generates an interpolation of the densities given by $\rho_t = \int_{\mathcal{M}} \rho_0(\mathbf{y}) \rho$. In this way, the stochastic interpolant approach discussed in Section ??

6 Experiments

Human pose estimation is the active research area in computer vision concerned with predicting the joint locations of a human body, usually based on 2D images. There are a wide range of applications for 3D pose estimation including autonomous driving, human-computer interaction and biomechanics [36].

6.1 3D Human Poses Estimation

This experiment aims to estimate the 3D coordinates of human skeleton joints based on a 2D image captured by a camera with known coordinates and focal lengths. The constraint manifold is defined to ensure that the lengths of the limbs remain constant:

$$\|q^i - q^j\|_2^2 = l_{i,j}^2 \quad \forall (i, j) \in \mathcal{J}, \quad (6.1)$$

where \mathbf{q}^i is the 3D position of joint i , $l_{i,j}$ is the length of the limb, and \mathcal{J} is the set of limbs. Assuming that we are given noisy 2D locations x^i of the joints and the matrix A of camera parameters, Brubaker et al. use a linear pose model based on eigenposes obtained from the training data [30, 6]. Essentially the model fit to the distribution is a normal distribution on the principal component analysis (PCA) of the data. The eigenposes are obtained from the PCA as the columns of the matrix U in the singular value decomposition of the data matrix X , where each column of X is a pose q in the training data, which has been de-meaned, meaning that the mean of each joint coordinate is subtracted from the data before the SVD is applied.

The eigenposes capture the most dominant directions of deviation from the mean of the

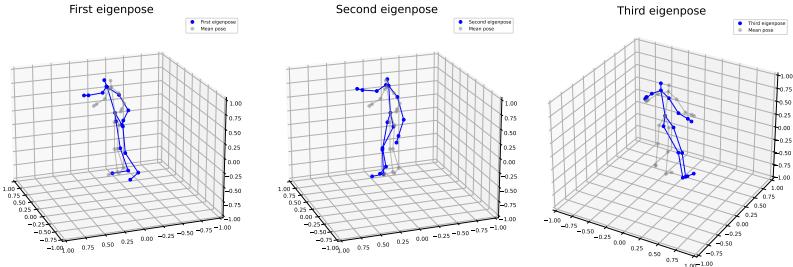


Figure 10: First, second and third eigenpose.

poses. The density is then defined as

$$\begin{aligned} \pi(\mathbf{q}|\mathbf{x}) &\propto p(\mathbf{x}|\mathbf{q}) \cdot p(\mathbf{q}) \\ &\propto \prod_{i=1}^N \exp(-\|x(q^i) - x^i\|^2/\sigma_m^2) \prod_{i=1}^{3N} \exp(-(P_j^T(q - q_0))^2/\sigma_j^2), \end{aligned} \quad (6.2)$$

where $x(q^i)$ is the 2D projection of joint i , P_j is the j th eigenpose and σ_j^2 is the associated eigenvalue. The first term in the density function describes the probability of a 2D observation \mathbf{x} given a 3D pose \mathbf{q} . Since we have added Gaussian noise to our observations this is a simple Gaussian distribution. The second term describes the prior distribution that was fit by Brubaker et al. to the HumanEva data [6]. Consider the matrix X with columns given by the data poses, $X = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m]$. Let \hat{X} be the de-meanned observations, so that $\hat{X} + \mu \otimes [1 \ 1 \ \dots \ 1] = X$ where μ is the vector of the row-wise means. Then the PCA is applied using an SVD of \hat{X} , $\hat{X} = U\Sigma V^T$. The transformed data is then modelled using a multivariate Gaussian distribution. **Since the principal components are independent, there is no covariance and we get the density given above. Since the PCA applies a linear transformation to the data, this is equivalent to a multivariate Gaussian model on the original data.**

The data used to train and test the model was generated according to the variational auto-encoder (VAE) used by Graham et al [9]. It was trained using the *PosePrior* motion-capture dataset, and generates angles according to a log-normal model that was fit to the data. For the purposes of this paper, the generated lengths were kept constant (at the mean of the model fit to the data), although there is some variation in the *perceived* bone lengths since the skeleton we use does not include all of the joints used to generate the data. The data was centred with the pelvis at $(0, 0, 2)$ so that the manifold was closed.

A known camera matrix was used to project the 3D data to a 2D image. The 2D image was then perturbed randomly by adding a random variable with a standard deviation of 0.1 to each joint's coordinates. The integrator was initialized at a random point in the training dataset for each frame of the test dataset, and run for 2,000 steps with a stepsize of 0.01.

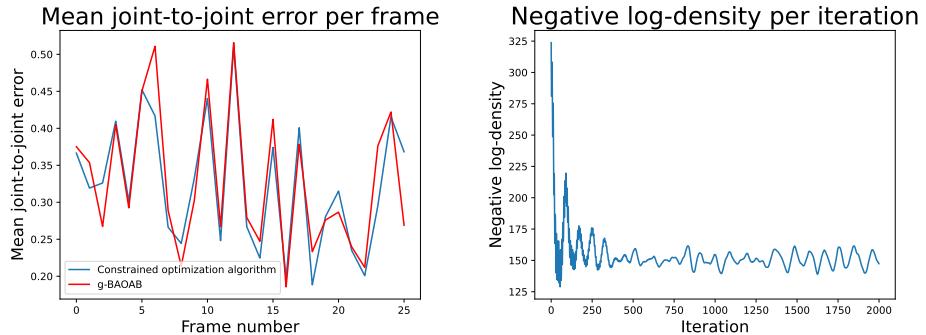


Figure 11: Mean joint-to-joint error per frame and negative log density per iteration

The g-BAOAB estimates of the 3d poses were compared to estimates obtained from a

gradient-based constrained optimization algorithm which minimized the negative log-density, as was done by De Bortoli et al. [8]. The error was measured as the average Euclidean distance between the joints. The algorithm used was sequential least-squares, a method similar to Newton’s method but which can minimize nonlinear functions subject to constraints [should I cite the package used \(Scipy\)?](#). The g-BAOAB estimate was taken to be the *intrinsic* mean of the positions generated by the algorithm, since the Euclidean mean will not necessarily lie on the manifold - this was obtained by converting to intrinsic coordinates, taking the average, and converting the average back to Euclidean coordinates. As shown in Figure 11, the g-BAOAB estimate was roughly as effective as the sequential least-squares estimate, underperforming the optimization algorithm for some frames but outperforming it for others.

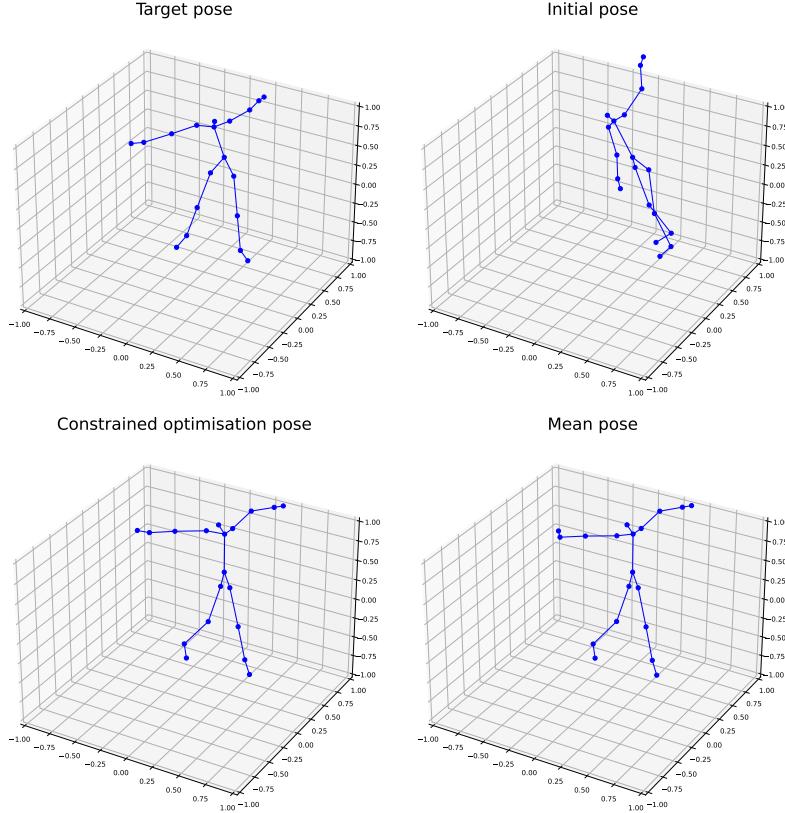


Figure 12: Target pose, constrained optimization algorithm pose and g-BAOAB mean pose after 2,000 g-BAOAB iterations.

6.2 3D Human Poses Generation

Next, we aim to use the score-based diffusion method described in Section ?? to noisify 3D poses from our dataset in order to train a score-based model to de-noisify poses, in order to sample from the pose distribution.

6.2.1 Adding noise

We noisify the data by constrained underdamped Langevin dynamics using the g-BAOAB algorithm, as described in . We fix one joint of the poses so that our constraint manifold is closed. This lets us perturb the data using the g-BAOAB algorithm with a force $\mathbf{F} = \mathbf{0}$, which will give a uniform stationary distribution. We can sample from the uniform distribution on the manifold by sampling the angle of each bone from a uniform spherical distribution, which will make every

point on the manifold equally likely:

$$\tau(\mathbf{q}) = \mathbb{P}(\mathbf{q}_1 = \mathbf{q}_1^*)\mathbb{P}(\mathbf{q}_2 = \mathbf{q}_2^*|\mathbf{q}_1^*) \cdots \mathbb{P}(\mathbf{q}_n = \mathbf{q}_n^*|\mathbf{q}_1^*, \mathbf{q}_2^* \dots \mathbf{q}_{n-1}^*) \quad (6.3)$$

$$\implies \tau(\mathbf{q}) = \prod_{i=2}^n \frac{1}{4\pi l_i^2} \quad (6.4)$$

$$(6.5)$$

where i indexes the bones and l_i is the length of bone i . Note that since we have fixed the initial point, $\mathbb{P}(\mathbf{q}_1 = \mathbf{q}_1^*) = 0$. As a note on the uniform spherical distribution, it is incorrect to uniformly sample $\phi \in [0, 2\pi]$ and $\theta \in [0, \pi]$. Due to the fact that the area element is given by $d\Omega = \sin \theta d\theta d\phi$, this method would not give a uniform distribution and points would instead more likely be sampled nearer to the poles. Points are instead generated by sampling $\phi \in [0, 2\pi)$ and $\theta = \cos^{-1}(2v - 1)$, where v is uniformly distributed on $[0, 1)$.

The data were perturbed using the g-BAOAB algorithm, with $h = 0.01$, $\gamma = \beta = 1$, and $T = 10$ giving 1000 steps. Figure 13 shows how perturbing using the g-BAOAB algorithm changes the pose. As the time approaches infinity, this will approach the distribution defined above.

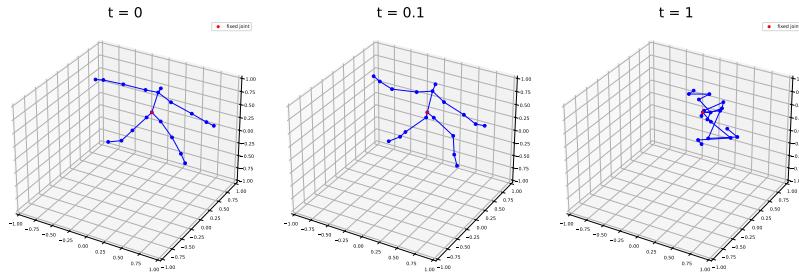


Figure 13: Initial pose being perturbed using the g-BAOAB algorithm, with a step size of 0.01.

6.2.2 Score Matching

Similarly to De Bortoli et al.'s experiments in Riemannian score-based modelling [8], a 5-layer multilayer perceptron was used to parameterize the score based network. The time was encoded in a manner similar to Song's method using Gaussian projections [33], which have been used to enable MLPs to learn high-frequency functions [34].

As shown in Section 4.1.1, we can train the network by minimizing the loss function:

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) = \int_{\mathcal{M}} \frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{q}, t)\|^2 + \Delta \hat{\mathbf{s}}_\theta(\mathbf{q}, t) d\mathbf{q} + C. \quad (6.6)$$

We approximate the minimizer of this loss by

$$\mathcal{L}(\hat{\mathbf{s}}_\theta) \simeq \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{q}, t)\|_2^2 + \Delta \hat{\mathbf{s}}_\theta(\mathbf{q}, t) \right). \quad (6.7)$$

where x_t are the data from the distribution. In this case, x_t are the noisified data points.

The neural network is constructed as shown by De Bortoli et al., using the projection onto the tangent space [8]:

$$\hat{\mathbf{s}}_\theta(\mathbf{q}, t) = \mathbf{\Pi}(\mathbf{q}) \tilde{\mathbf{s}}_\theta(\mathbf{q}, t) \quad (6.8)$$

where $\tilde{\mathbf{s}} : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$ is an ambient vector field. Then using Theorem ??, the Riemannian divergence $\Delta \hat{\mathbf{s}}_\theta(\mathbf{q}, t)$ is equal to the Euclidean divergence $\Delta_E \hat{\mathbf{s}}_\theta$, the Euclidean divergence. This allows us to calculate the loss without knowledge of the Riemannian metric, which is useful for

constraints of the form given in Equation 3.12c.

By an application of Stokes' lemma, if our neural network is remaining tangent to the manifold at each point and the divergence is calculated accurately, we should have that

$$\int_{\mathcal{M}} \Delta \hat{\mathbf{s}}(\mathbf{q}) d\mathbf{q} = 0. \quad (6.9)$$

In order to check that this is satisfied, we use Monte Carlo integration using a uniform distribution over the manifold. Then we should have that

$$\sum_{\mathbf{q} \in \mathcal{M}} \Delta \hat{\mathbf{s}}(\mathbf{q}) \approx 0 \quad (6.10)$$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N \Delta \hat{\mathbf{s}}(\mathbf{q}_i) \approx 0 \quad (6.11)$$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N (1 - \Delta \hat{\mathbf{s}}(\mathbf{q}_i)) \approx 1 \quad (6.12)$$

so we can validate that the neural net is behaving as expected using this final value. Indeed, as shown in 14, the approximation approaches 1 as the number of samples is increased. The

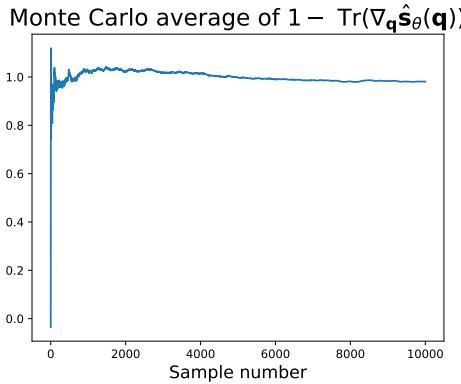


Figure 14: Monte Carlo integration of $(1 - \nabla \hat{\mathbf{s}})$

specific construction of the neural network is shown in Figure 15. The main structure of an MLP was similar to that used by De Bortoli et al., while the structure of the inclusion of the time embeddings was taken from Song et al. [8, 33]. The time was mapped to 58 different nodes, each taking a value of either $\sin(2\pi\omega t)$ or $\cos(2\pi\omega t)$ for a random variable ω for each node which was fixed at initialisation. The arrows from the time embedding indicate element-wise addition (after the application of the activation function). 4 main hidden layers were used rather than 5 in order to both speed up training and reduce the risk of overfitting (due to having less computational resources and so a lower number of samples).

In each epoch, for each sample from the data distribution, a time t was randomly chosen, the data was noised from 0 to t , and the loss was updated based on the score function given at time t . The full experiment algorithms are given below:

The model was trained on 2,000 initial data points, although each datapoint was noised 1,000 times, so in effect the model was trained on 2,000,000 different (\mathbf{q}, \mathbf{t}) values. The learning rate used was 10^{-4} , and training was stopped after 3 successive epochs without an improvement in the loss on the test data. Even with a relatively low number of samples, the reverse diffusion model generated realistic looking samples.

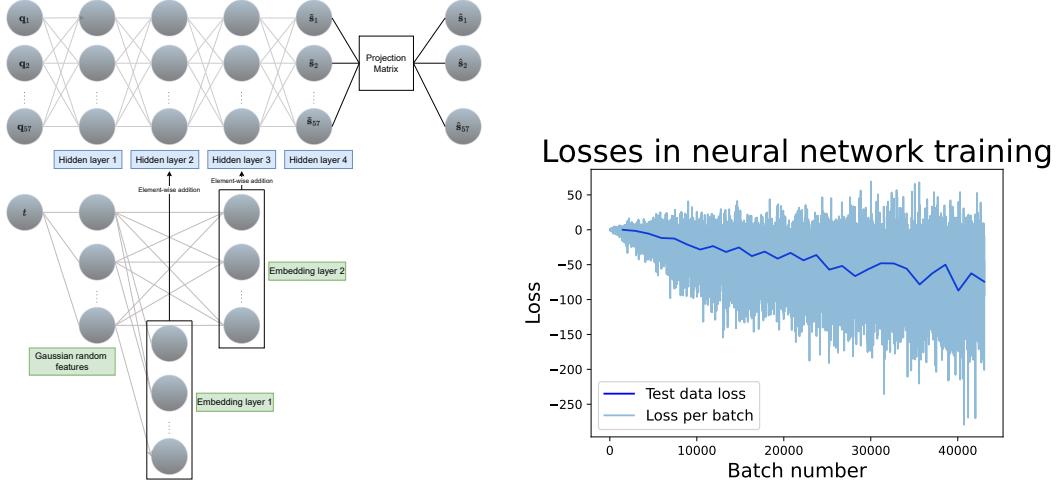


Figure 15: Score-based diffusion neural network architecture and training loss [progression]

Algorithm 3 Score-based training using g-BAOAB

```

for  $n \in \{0, \dots, N_{\text{epochs}} - 1\}$  do
    for  $\mathbf{x} \in \text{dataset}$  do                                 $\triangleright$  A batch is selected from the shuffled data
         $t \sim \mathcal{U}(h, 1)$   $\triangleright$  A batch of  $ts$  are sampled uniformly and independently between  $h$  and 1
         $v = \mathbf{0}$                                           $\triangleright$  Initializing the velocities at zero
        for  $i \in \{0, \dots, \text{Int}(\frac{t}{h}) - 1\}$  do
             $\mathbf{x}, \mathbf{v} = \text{g-BAOAB}(\mathbf{x}, \mathbf{v}, t, \text{Force} = \mathbf{0})$            $\triangleright$  noisify the data using g-BAOAB
             $\mathcal{L}(\hat{\mathbf{s}}_\theta) = \frac{1}{2} \|\hat{\mathbf{s}}_\theta(\mathbf{x}, t)\|^2 + \Delta \hat{\mathbf{s}}_\theta(\mathbf{q}, t)$ 
             $\hat{\mathbf{s}} = \text{optimizer.update}(\hat{\mathbf{s}})$ 
        end for
    end for
end for

```

Algorithm 4 Score-based diffusion time-reversal using g-BAOAB

```

 $\mathbf{y} \sim \pi$                                       $\triangleright$  Sampling from the base (uniform) distribution
 $\mathbf{v} = \mathbf{0}$                                   $\triangleright$  Initializing the velocity
 $\mathbf{y}, \mathbf{v} = \text{reverse g-BAOAB}(\mathbf{y}, \mathbf{v}, \frac{T}{h}, \text{Force} = -\hat{\mathbf{s}}(\mathbf{y}, T - t))$   $\triangleright$  de-noisify the data using g-BAOAB

```

6.3 Results

The results of the g-BAOAB score-based diffusion sampling process are presented here, compared with samples from the true data distribution.

Retraining with 2x the samples to see if I can improve these a bit.

7 Comparison with Moser flow

As a comparison to the score-based diffusions using g-BAOAB, a Moser flow model and a score-based generative model were trained and used to sample from the data distribution for the toy problem considered in Section 3.3. As shown above, Moser flows are a class of CNFs which are well suited to problems on certain Riemannian manifolds. Both models were trained using 20,000 data points from a g-BAOAB simulation, taken from every 50th step of a 1,000,000 step run. A histogram of the data are given in ??.

7.1 Moser flow model

A neural network, similar to that used in the score-based model was constructed based on a 5-layer MLP; the only difference in construction between this network and that used for the

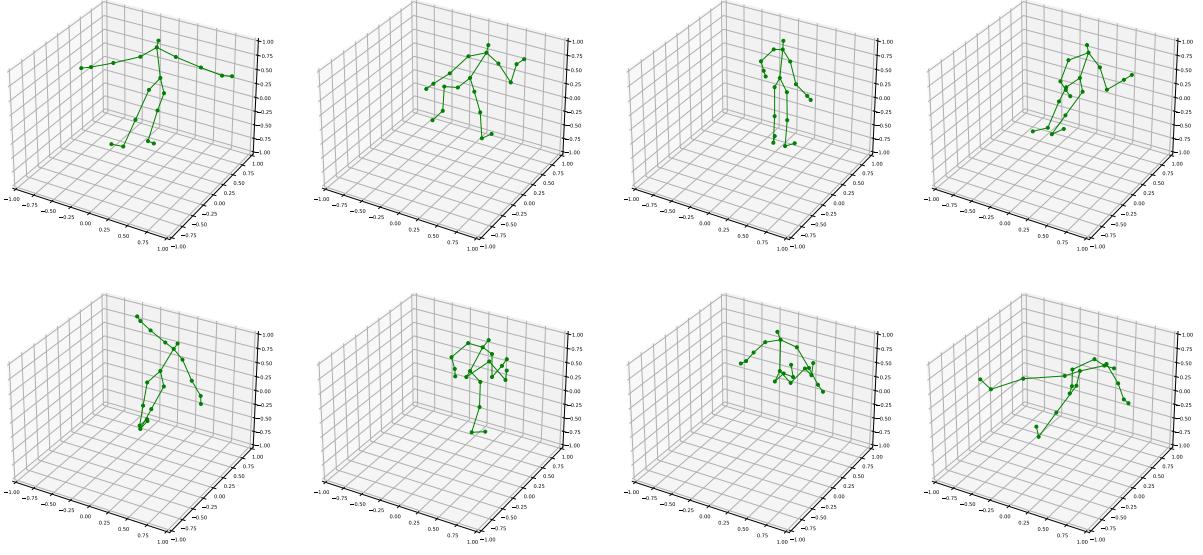


Figure 16: Samples from the data distribution

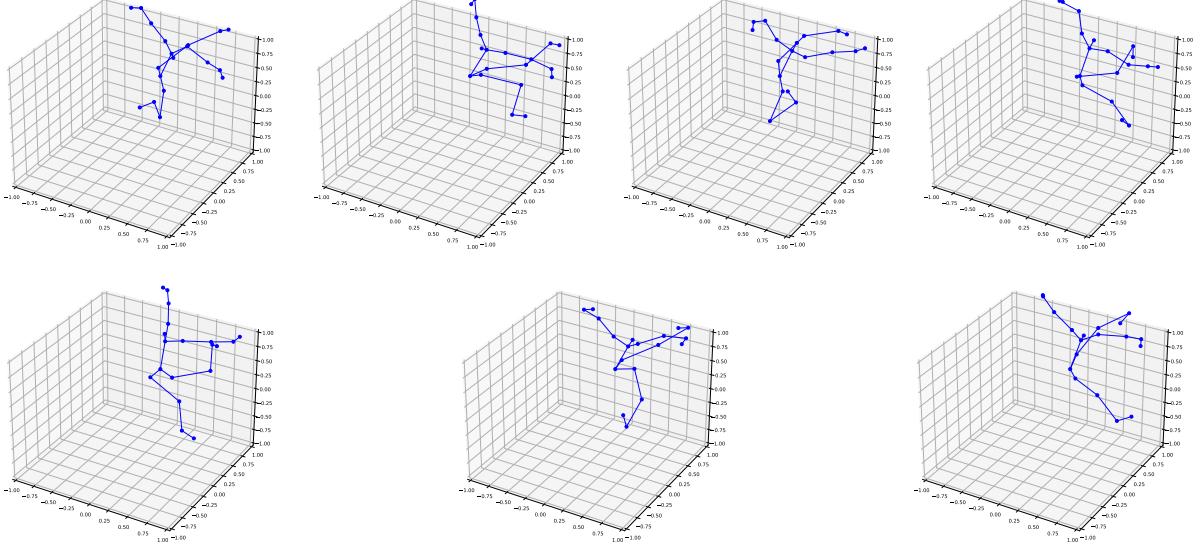


Figure 17: Samples from a score-based diffusion model using g-BAOAB

experiments in Section 6.2 was that this neural network did not take time as an input. The network architecture is shown in Figure ?? for completeness. The loss function used was

$$\mathcal{L}(\hat{\mu}) = +\frac{1}{l} \sum_1^l \frac{\hat{u}_-(\mathbf{y}_j)}{v(\mathbf{y}_j)} \quad (7.1)$$

as described in Section 5.4.1, where v was chosen to be the uniform density over the sphere, with each point taking a probability of $v(\mathbf{q}) = \frac{1}{4\pi}$. The number of samples used for the Monte Carlo integration was taken to be $l = 1,000$, which was high enough to diminish the effect of sampling error on the loss. The neural network was trained using a batch size of 64.

7.2 Score-based model

For the score-based model, the neural network used was identical to that used in Section 6.2 apart from the input and output layers which were naturally three-dimensional. The data were noised using a step-size of 0.01 for a time of $T = 1$ using the g-BAOAB algorithm with a force

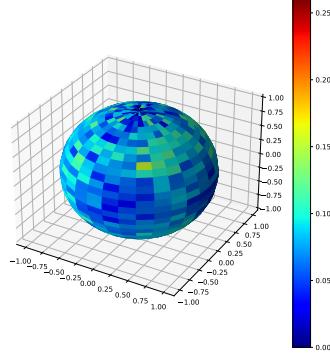


Figure 18: Training data for the Moser flow and SBG models (Note that the colours have been normalized differently than in Figure 6 in order to match the scale in ??)

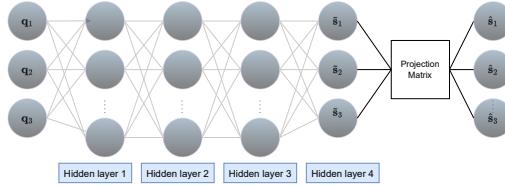


Figure 19: Moser flow neural network architecture

of zero to move towards the uniform distribution on the sphere.

7.3 Comparison / Results

Both models were trained on the same dataset so that a direct comparison could be made, however it is worth noting that training the Moser flow model is a much faster process, since it eliminates the need for the noising of the dataset. Since both models are trained based on the divergence of a vector field, the noising process was the only significant difference in training times.

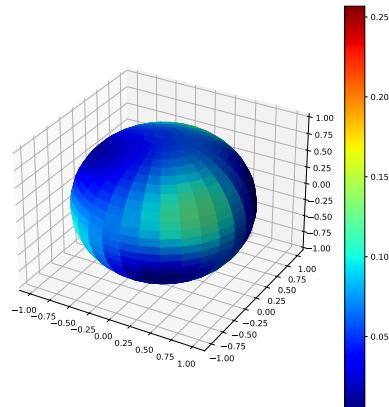


Figure 20: Result of a Moser flow model fit and a score-based generative model fit to data generated by g-BAOAB waiting for the score-based model to train

8 Related work

8.1 Stochastic interpolants

Stochastic interpolants have been proposed as a framework that can unify flows and diffusions. Similarly to score-based diffusion, they involve interpolation on the event space. Given two densities ρ_0 and $\rho_1 : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$, a stochastic interpolant is a stochastic process

$$\mathbf{x}_t = I(t, x_0, x_1) \quad (8.1)$$

for $t \in [0, 1]$, where $I(0, x_0, x_1) = x_0$, $I(1, x_0, x_1) = x_1$, $x_0 \sim \rho_0$, $x_1 \sim \rho_1$ and [1] where I is twice continuously differentiable in t and \mathbf{x} , and satisfies the conditions $I(0, \mathbf{x}_0, \mathbf{x}_1) = \mathbf{x}_0$ and $I(1, \mathbf{x}_0, \mathbf{x}_1) = \mathbf{x}_1$, as well as

$$\exists C_1 < \infty : |\delta_t I(t, x_0, x_1)| \leq C_1 |x_0 - x_1| \quad \forall (t, x_0, x_1) \in [0, 1] \times \mathbb{R}^n \times \mathbb{R}^n. \quad (8.2)$$

The first condition on I is simply that it does what we would expect of an interpolant, while the second requires that its output doesn't move too quickly relative to its input.

The velocity is then defined as the vector field v_t such that:

$$\dot{\mathbf{x}}_t = v_t(\mathbf{x}_t). \quad (8.3)$$

The stochastic interpolant \mathbf{x}_t defined above with $I_t(\mathbf{x}_0, \mathbf{x}_1)$ satisfying the above conditions has a probability density that satisfies the continuity equations

$$\partial_t \rho_t + \nabla \cdot (v_t \rho_t) = 0, \quad (8.4)$$

which maps from ρ_0 at time $t = 0$ to ρ_1 at time $t = 1$, with its velocity given by the unique minimizer of

$$\mathcal{L}(\hat{v}) = \mathbb{E} [| \hat{v}_t I_t(x_0, x_1) |^2 - 2 \partial_t I_t(\mathbf{x}_0, \mathbf{x}_1) \cdot \hat{v}_t(I_t(\mathbf{x}_0, \mathbf{x}_1))]. \quad (8.5)$$

Thus the velocity can be parameterised as a neural network and learned, as was done with Moser flows and SBGM.

In the more general case proposed later [1], the interpolant is generalized to include a *latent variable* and we have

$$\mathbf{x}_t = I(t, \mathbf{x}_0, \mathbf{x}_1) + \gamma(t) \mathbf{z}, \quad (8.6)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\gamma : [0, 1] \rightarrow \mathbb{R}$ is continuously differentiable. The addition of the latent variable smooths the density and velocity, and also allows stochastic interpolants to be more closely related to score-based diffusion. The latent variable allows the score of the marginal distribution to be written as

$$\nabla_{\mathbf{x}} \log \rho_t(\mathbf{x}) = -\gamma(t)^{-1} \mathbb{E}(\mathbf{z} | \mathbf{x}_t = \mathbf{x}) \quad (8.7)$$

and the score associated with the model is the unique continuously differentiable minimizer of

$$\mathcal{L}(\hat{s}) = \int_0^1 \mathbb{E} \left(\frac{1}{2} |\hat{s}(t, \mathbf{x}_t)|^2 + \gamma^{-1}(t) \hat{s}(t, \mathbf{x}_t) \right) \quad (8.8)$$

In particular, it has been shown that minimizing the loss functions given in Equation 8.8 and Equation 8.5 maximizes the likelihood of a stochastic generative model with force given by a weighted sum of the score and velocity estimates

$$\mathbf{F}(\mathbf{x}, t) = \hat{\mathbf{b}}(\mathbf{x}, t) + \kappa \hat{\mathbf{s}}(\mathbf{x}, t) \quad (8.9)$$

for a constant κ .

8.2 Pushforward of Euclidean normalizing flows

In this paper I focused on methods for constructing normalizing flows directly on manifolds in order to compare with constrained score-based diffusion, however another way to [construct] a normalizing flow on a d -dimensional manifold which is embedded in \mathbb{R}^n is to use the embedding map $T : \mathbb{R}^d \rightarrow \mathbb{R}^n$. Then the embedding map induces a metric $G()$ on the tangent space of \mathcal{M} , which gives a volume $d\nu = \sqrt{\det G(\mathbf{u})} d\mathbf{u}$ on \mathcal{M} . Then in a manner not dissimilar to the case outlined in Section ??, the inverse mapping of $T|_{\mathcal{M}}$ as $T^{-1} : \mathcal{M} \rightarrow \mathbb{R}^n$, we arrive at the density on the manifold [25]:

$$\rho(\mathbf{q}) = \tau(T^{-1}(\mathbf{q})) [\det G(T^{-1}(\mathbf{q}))]^{-1/2}, \quad (8.10)$$

which is a generalization of Equation ?? in the case where $\mathcal{M} = \mathbb{R}^n$. However, the diffeomorphism T above can only exist if the manifold \mathcal{M} is topologically equivalent to Euclidean space, and without this map any embedding map will create singularities, which lead to numerical issues.

9 Limitations

As discussed in Section 5.3, the main limitation of using constrained overdamped Langevin dynamics in SBD model is the lack of an accompanying probability-flow model when compared to overdamped Langevin dynamics

10 Conclusion

We have explored aspects of a score-based diffusion model based which uses underdamped rather than overdamped Langevin dynamics.

Appendices

A Computational cost

A.1 Score-based generative modelling

The methods used in approximating the vector fields in the SBGM experiments defined the vector field using the projection onto the tangent space. Automatic differentiation in order to minimize the objective function then requires the *divergence* which is the trace of the Jacobian matrix of the model. Since the projection matrix is defined using the Jacobian matrix of the constraints, this adds complexity to the computational graph. As noted by Song et al., a naive algorithm for computing the divergence requires n times more backward passes than computing the gradient [32], which may be too slow for large dimensions n . Since the constraints are known in advance, the computational graph could be simplified slightly by explicitly evaluating the Jacobian and Hessian matrices of the constraints rather than computing them via backpropogation, in order to maximize efficiency, however the other elements of the neural network would still need to be differentiated using backpropogation. To aid the calculation of the divergence, Song et al. suggested a “sliced score-matching” approach, which would make the process more feasable in large dimensions [32].

A.2 Moser flows

The computational cost of Moser flows is much lower than for general CNFs, however it still requires the divergence of the vector field to be computed, which similarly would need to be approximated in higher dimensions [27]. Compared to CNFs which use ODE solvers...

B Activation functions

Throughout the paper various vector fields are constructed using neural networks. In Section 5.4, the neural networks are required to be $l - \text{finite}$, since for a smooth target function it has been proven that there exists an MLP with l -finite smooth activation function that uniformly approximates the first m derivatives of f with an error of at most ϵ for a given ϵ . Therefore the activations chosen are **Tanh** and **Softplus**, which are l -finite for $l \geq 1$ and $l \geq 2$, respectively. A function is said to be l -finite if it is l times continuously differentiable, and satisfies $\int_{-\infty}^{\infty} |f^{(l)}(t)| \leq \infty$. **Swish**, **ReLU** and **LeakyReLU** are not $l - \text{finite}$ for any l . **ReLU** has been widely shown to be a [stable] activation function which helps to mitigate issues with exploding or vanishing gradients, however for [this reason] these functions were not used in this context. For the score-based diffusion model,

C Assumptions

C.1 Theorem ??

The regularity conditions imposed on $\rho_{t|s}(\mathbf{x}_t|\mathbf{x}_s)$ are simply that

1. $\rho_{t|s}(\mathbf{q}_t|\mathbf{q}_s)\mathbf{s}_t(\mathbf{q})$ is a vector field in C_1 for all \mathbf{x}_s
2. $|\rho_{t|s}(\mathbf{q}_t|\mathbf{q}_s)| \in L_1 \forall \mathbf{q}_s \in \mathcal{M}$
3. $\Delta [\rho_{t|s}(\mathbf{q}_t|\mathbf{q}_s)\mathbf{s}_t(\mathbf{q})] \in L_1 \forall \mathbf{q}_s \in \mathcal{M}$

[Need to define \$L_1\$ in the same way as in the appendix](#) where C_1 is the set of all continuously differentiable vector fields, and $f \in L_1 \iff \int_{\mathcal{M}} |f| < \infty$. The continuous differentiability requirement will be dealt with in practice when we consider activation functions in B. The proof is straightforward and similar to Hyvärinen's for compact manifolds, where Stokes' theorem can be used, and for non-compact manifold it uses a divergence theorem.

C.2 Theorems 5.2 and 5.3

The assumptions made by Song et al. and required for either Theorem 5.2 (C.2, C.2 and C.2) or Theorem 5.3 are given below.

- $\rho(\mathbf{x}), \tau(\mathbf{x}) \in \mathcal{C}^2$ and $\mathbb{E}_{\rho} [||\mathbf{x}||_2^2], \mathbb{E}_{\tau} [||\mathbf{x}||_2^2]$
- $\mathbf{f}(\cdot, t) \in \mathcal{C}^1$ for all $t \in [0, T]$, and there exists a constant $C > 0$ such that for all $\mathbf{x} \in \mathbb{R}^n$ and each $t \in [0, T]$, $||\mathbf{f}(\mathbf{x}, t)||_2 \leq C(1 + ||\mathbf{x}||_2)$.
- There exists a constant $C > 0$ such that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $||\mathbf{f}(\mathbf{x}, t) - \mathbf{f}(\mathbf{y}, t)||_2 \leq C||\mathbf{x} - \mathbf{y}||_2$.

References

- [1] M. S. Albergo, N. M. Boffi, and E. Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions, 2023.
- [2] M. S. Albergo and E. Vanden-Eijnden. Building normalizing flows with stochastic interpolants, 2023.
- [3] H. C. Andersen. Rattle: A “velocity” version of the shake algorithm for molecular dynamics calculations. *Journal of Computational Physics*, 52(1):24–34, 1983.
- [4] B. D. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.

- [5] D. Berthelot, T. Schumm, and L. Metz. Began: Boundary equilibrium generative adversarial networks, 2017.
- [6] M. A. Brubaker, M. Salzmann, and R. Urtasun. In N. D. Lawrence and M. A. Girolami, editors, *AISTATS*, volume 22 of *JMLR Proceedings*, pages 161–172. JMLR.org, 2012.
- [7] X. Cheng, N. S. Chatterji, P. L. Bartlett, and M. I. Jordan. Underdamped langevin mcmc: A non-asymptotic analysis. In S. Bubeck, V. Perchet, and P. Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 300–323. PMLR, 06–09 Jul 2018.
- [8] V. De Bortoli, E. Mathieu, M. Hutchinson, J. Thornton, Y. W. Teh, and A. Doucet. Riemannian score-based generative modelling. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 2406–2422. Curran Associates, Inc., 2022.
- [9] M. M. Graham and A. J. Storkey. Asymptotically exact inference in differentiable generative models. *Electronic Journal of Statistics*, 11(2):5105 – 5164, 2017.
- [10] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models, 2020.
- [11] E. P. Hsu. *Stochastic Analysis on Manifolds*. Graduate Studies in Mathematics. American Mathematical Society, 2001.
- [12] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24):695–709, 2005.
- [13] A. Jolicoeur-Martineau, K. Li, R. Piché-Taillefer, T. Kachman, and I. Mitliagkas. Gotta go fast when generating data with score-based models, 2021.
- [14] B. Leimkuhler and C. Matthews. *The Langevin Equation: With Applications to Stochastic Problems in Physics, Chemistry and Electrical Engineering*. World Scientific Series in Contemporary Chemical Physics. World Scientific, September 2012.
- [15] B. Leimkuhler and C. Matthews. Robust and efficient configurational molecular sampling via Langevin dynamics. *The Journal of Chemical Physics*, 138(17):174102, 05 2013.
- [16] B. Leimkuhler and C. Matthews. *Molecular Dynamics: With Deterministic and Stochastic Numerical Methods*. Interdisciplinary Applied Mathematics. Springer, May 2015.
- [17] B. Leimkuhler and C. Matthews. Efficient molecular dynamics using geodesic integration and solvent–solute splitting. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 472(2189):20160138, 2016.
- [18] B. Leimkuhler and S. Reich. *Simulating Hamiltonian Dynamics*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, September 2004.
- [19] B. Leimkuhler, T. Vlaar, T. Pouchon, and A. Storkey. Better training using weight-constrained stochastic dynamics, 2021.
- [20] B. J. Leimkuhler and R. D. Skeel. Symplectic numerical integrators in constrained hamiltonian systems. *Journal of Computational Physics*, 112(1):117–125, 1994.
- [21] T. Lelièvre, M. Rousset, and G. Stoltz. *Free Energy Computations: A Mathematical Perspective*, volume 31 of *Springer Series in Computational Mathematics*. Imperial College Press, second edition, 2010.
- [22] T. Lelièvre, M. Rousset, and G. Stoltz. Langevin dynamics with constraints and computation of free energy differences, 2011.

- [23] L. Liao, Y. Zhang, S. J. Maybank, Z. Liu, and X. Liu. Image recognition via two-dimensional random projection and nearest constrained subspace. *Journal of Visual Communication and Image Representation*, 25(5):1187–1198, 2014.
- [24] E. Mathieu and M. Nickel. Riemannian continuous normalizing flows, 2020.
- [25] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan. Normalizing flows for probabilistic modeling and inference, 2021.
- [26] S. Pascual, G. Bhattacharya, C. Yeh, J. Pons, and J. Serrà. Full-band general audio synthesis with score-based diffusion, 2022.
- [27] N. Rozen, A. Grover, M. Nickel, and Y. Lipman. Moser flow: Divergence-based generative modeling on manifolds, 2021.
- [28] J.-P. Ryckaert, G. Ciccotti, and H. J. Berendsen. Numerical integration of the cartesian equations of motion of a system with constraints: molecular dynamics of n-alkanes. *Journal of Computational Physics*, 23(3):327–341, 1977.
- [29] M. Shapovalov and R. Dunbrack. A smoothed backbone-dependent rotamer library for proteins derived from adaptive kernel density estimates and regressions. *Structure*, 19(6):844–858, 2011.
- [30] J. D. Skufca, E. M. Boltt, R. Pilkar, and C. J. Robinson. Eigenposes: Using principal components to describe body configuration for analysis of postural control dynamics, 2010.
- [31] Y. Song, C. Durkan, I. Murray, and S. Ermon. Maximum likelihood training of score-based diffusion models, 2021.
- [32] Y. Song, S. Garg, J. Shi, and S. Ermon. Sliced score matching: A scalable approach to density and score estimation. In *Uncertainty in Artificial Intelligence*, pages 574–584. PMLR, 2020.
- [33] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *CoRR*, abs/2011.13456, 2020.
- [34] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *CoRR*, abs/2006.10739, 2020.
- [35] A. Vahdat and J. Kautz. NVAE: A deep hierarchical variational autoencoder. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- [36] J. Wang, S. Tan, X. Zhen, S. Xu, F. Zheng, Z. He, and L. Shao. Deep 3d human pose estimation: A review. *Computer Vision and Image Understanding*, 210:103225, 2021.
- [37] L. Yang, Z. Zhang, Y. Song, S. Hong, R. Xu, Y. Zhao, W. Zhang, B. Cui, and M.-H. Yang. Diffusion Models: A Comprehensive Survey of Methods and Applications. *arXiv e-prints*, page arXiv:2209.00796, Sept. 2022.