Learning geometry and topology via multi-chart flows

Hanlin Yu University of Helsinki hanlin.yu@helsinki.fi

Marcelo Hartmann University of Helsinki marcelo.hartmann@helsinki.fi Søren Hauberg Technical University of Denmark sohau@dtu.dk

> Arto Klami University of Helsinki arto.klami@helsinki.fi

Georgios Arvanitidis Technical University of Denmark gear@dtu.dk

Abstract

Real world data often lie on low-dimensional Riemannian manifolds embedded in high-dimensional spaces. This motivates learning degenerate normalizing flows that map between the ambient space and a low-dimensional latent space. However, if the manifold has a non-trivial topology, it can never be correctly learned using a single flow. Instead multiple flows must be 'glued together'. In this paper, we first propose the general training scheme for learning such a collection of flows, and secondly we develop the first numerical algorithms for computing geodesics on such manifolds. Empirically, we demonstrate that this leads to highly significant improvements in topology estimation.

1 Introduction

Normalizing flows [Papamakarios et al., 2021] is a family of flexible neural network models of complex probability distribution from finite observations. They are guaranteed to be bijections allowing tractable sampling and log-density evaluations. However, due to the bijectivity requirement, the dimensionalities of the latent and ambient spaces need to be identical. As real-world datasets often lie on a low-dimensional Riemannian manifold embedded in a high-dimensional ambient space, we cannot have a bijection. As such, the modeling strategy of normalizing flows can be ill-posed.

A probability density on an unknown Riemannian manifold can be modeled using a normalizing flow [Brehmer and Cranmer, 2020, Caterini et al., 2021], by learning a mapping between a low-dimensional latent space and a high-dimensional ambient space. It was demonstrated that, for certain manifolds, the flows are able to model the dis-



Figure 1: Using multiple charts to cover a manifold gives the ability to learn the manifold's true geometry and topology.

tribution accurately. However, these methods often consider a single normalizing flow, which is justified only when the manifold is diffeomorphic to a Euclidean space. Some common manifolds,

for instance the sphere and the torus, do not satisfy this assumption. It is by definition impossible for a single-chart flow to model these manifolds correctly due to the topological mismatch.

Understanding the underlying topology and geometry of the data manifold using only finite observations is an open problem in machine learning [Hauberg, 2019, Moor et al., 2020, Acosta et al., 2023, Ross et al., 2024, Diepeveen, 2024]. Some attempts focus on capturing the Riemannian structure using the pullback metric induced by a single normalizing flow [Hoffman et al., 2019, Kruiff et al., 2024], while other works [Schonsheck et al., 2020, Kalatzis et al., 2022, Sidheekh et al., 2022] use multiple autoencoders or normalizing flows to cover the manifold. However, these approaches typically lack a principled way of choosing the correct encoder [Loaiza-Ganem et al., 2024]. Most importantly, these works did not study the geometric and topological implications of the learned generative model, but focused mainly on the density modeling task.

In this paper, we aim at learning data manifolds with the correct geometry and topology. Classic differential geometry achieves this by gluing together local "patches" of the manifold, each of which is modeled with a diffeomorphism. In similar spirit, we learn a mixture of flows, each corresponding to a "patch". In classic differential geometry, the patches are assumed to perfectly overlap in bordering regions, but such extrapolation guarantees are unrealistic with flows (see Figure 1). We propose a formalism to handle such non-overlapping patches, and further develop numerical algorithms for computing geodesics over manifolds defined through such mixture of flows. Our contributions:

1. We demonstrate that using a single normalizing flow to model a distribution on a Riemannian manifold can lead to significant misunderstanding in terms of the underlying topology and geometry.

2. We provide general training strategies for multi-chart flows based on the maximum likelihood principle and demonstrate that they can perform better than single-chart flows of similar complexities.

3. We provide specialized tools to utilize the geometry learned by multi-chart flows that crucially respect the actual topology of the underlying manifold, yielding more reliable representations. Note that beyond implementing the abstract methodology of classic differential geometry, we also account for the fact that the flows cannot perfectly overlap for data learned manifolds.

2 Preliminaries

2.1 A brief intro to Riemannian geometry

Riemannian manifolds [Do Carmo, 1992, Lee, John M., 2018] are spaces that locally resemble Euclidean space. We consider a *d*-dimensional Riemannian manifold isometrically embedded in a *D*-dimensional Euclidean space, following the typographic mnemonic that the lower-case symbol denotes the smaller dimensionality. The well-known Nash embedding theorem [Lee, John M., 2018] guarantees that such an embedding exists for all Riemannian manifolds, though it is unclear how to construct such an embedding in the general case. A *local coordinate chart* is given by the pair (U, ϕ) , where U is an open set on the manifold and ϕ is a bijection between U and an open set of *d*-dimensional Euclidean space. For simple manifolds, one might be able to construct a single chart that covers the entire manifold. However, for more complex ones, we need a collection of local charts that transition smoothly, known as an *atlas*.

A curve with zero tangential acceleration, intuitively a shortest path on the manifold, is known as a *geodesic*. For a point x on a *d*-dimensional Riemannian manifold, one can attach a *d*-dimensional Euclidean space known as the tangent space at x, containing all vectors that are tangent to the surface at x. Given an initial position x_0 and initial v_0 , one can use the exponential map $\operatorname{Exp}_{x_0}(v_0)$ to obtain the result of following along the geodesic for unit time. Given the initial position x_0 and the final position x_1 , one can use the logarithmic map $\operatorname{Log}_{x_0}(x_1)$ to obtain (one of) the initial velocity that satisfies the boundary condition. These geometric operations are also illustrated in Figure 1, where we show that these quantities are well-defined even if the manifold is covered by multiple charts.

2.2 Normalizing flows

A normalizing flow [Papamakarios et al., 2021] pushes forward a tractable distribution p(u) from a *latent space* U to a complex distribution in the input space X through the bijective transformation f, such that the resulting distribution $q_{\theta}(x)$ in X can be tractably evaluated.

For density estimation tasks on Euclidean spaces when $\dim(\mathbf{X}) = \dim(\mathbf{U}) = D$, normalizing flows are typically trained by minimizing the forward KL divergence KL $(p(\mathbf{x})|q_{\theta}(\mathbf{x}))$ [Papamakarios et al., 2021], where the density $q_{\theta}(\mathbf{x})$ can be evaluated as $q_{\theta}(\mathbf{x}) = q(\mathbf{u})|\det(\mathbf{J}_{f}(\mathbf{u}))|^{-1}|_{\mathbf{u}=f^{-1}(\mathbf{x})}$, with $\mathbf{J}_{f} \in \mathbb{R}^{D \times D}$ being the Jacobian matrix of the parameterized transformation $f: \mathbf{U} \to \mathbf{X}$; the dependency of f on θ is dropped to avoid cluttering notations. One can employ specific flow architectures, such that the log determinant of the Jacobian can be obtained efficiently. Two notable examples are RealNVP flows [Dinh et al., 2017] and Neural Spline flows [Durkan et al., 2019].

In many scenarios, one would need to model distributions that are supported on Riemannian manifolds. In the following discussions, we largely follow Brehmer and Cranmer [2020] and Caterini et al. [2021]. A normalizing flow on a Riemannian manifold is commonly parameterized as the composition of a series of d-dimensional layers g and a series of a D-dimensional layers h, with the resulting degenerate transformation given by

$$\boldsymbol{f} = \boldsymbol{h} \circ \operatorname{Pad} \circ \boldsymbol{g} = \boldsymbol{h} \circ \boldsymbol{g}, \tag{1}$$

where \circ denotes composition of functions, Pad denotes adding D - d zeros at the end, and \tilde{h} denotes $h \circ$ Pad. The resulting change of variable formula is more involved [Brehmer and Cranmer, 2020]

$$q_{\theta}(\boldsymbol{x}) = p(\boldsymbol{u}) \left| \det(\boldsymbol{J}_{\boldsymbol{f}}^{\top}(\boldsymbol{u}) \, \boldsymbol{J}_{\boldsymbol{f}}(\boldsymbol{u})) \right|^{-\frac{1}{2}} \Big|_{\boldsymbol{u} = \boldsymbol{f}^{-1}(\boldsymbol{x})},$$
(2)

where $J_f \in \mathbb{R}^{D \times d}$. Due to the degeneracy, without good reconstructions the model may learn suboptimal distributions despite having good log likelihood estimates [Brehmer and Cranmer, 2020]. For this reason, they proposed the Manifold-learning Flow, where h focuses on accurately reconstructing the data and g primarily learns the density.

While their mechanism is theoretically justified [Loaiza-Ganem et al., 2022], Caterini et al. [2021] noted that the optimization problem becomes challenging, leading to potentially inferior performances in practice: instead, it can be beneficial to explicitly optimize the KL divergence while penalizing the reconstruction loss. One can define $\tilde{h}^{\dagger} = \text{Proj} \circ h^{-1}$, where Proj is the operation that takes the *d* coordinates; this is the proper inverse function of \tilde{h} when *x* lies precisely on the hypersurface formed by \tilde{h} . The resulting loss function is given by

$$\mathbb{E}_{\boldsymbol{x} \sim p(\boldsymbol{x})} \left[-\log q_{\boldsymbol{\theta}}(\boldsymbol{x}) + \lambda \| \boldsymbol{x} - \tilde{\boldsymbol{h}} \left(\tilde{\boldsymbol{h}}^{\dagger}(\boldsymbol{x}) \right) \|^{2} \right],$$
(3)

where $\lambda > 0$ controls the trade-off between reconstruction and density estimation.

Recall that, to obtain $\log q_{\theta}(x)$, we need to compute $\log \det(J_f^{\top} J_f)$. Caterini et al. [2021] showed that this quantity can be simplified further as

$$\log \det \left(\boldsymbol{J}_{\boldsymbol{f}}^{\top} \, \boldsymbol{J}_{\boldsymbol{f}} \right) = 2 \log \det \left(\boldsymbol{J}_{\boldsymbol{g}} \right) + \log \det \left(\boldsymbol{J}_{\tilde{\boldsymbol{h}}}^{\top} \, \boldsymbol{J}_{\tilde{\boldsymbol{h}}} \right), \tag{4}$$

where $\log \det (J_g)$ is tractable as g is a square flow, and only $\log \det (J_{\tilde{h}}^{\top} J_{\tilde{h}})$ is a complex term. For optimization we only need the gradient of the log likelihood. For reasonably small d as in many practical applications, we can compute the $D \times d$ Jacobian matrix exactly using vectorization and forward mode automatic-differentiation using deep learning frameworks, while the cost of obtaining the log determinant for $d \times d$ is affordable. For larger scale problems, Caterini et al. [2021] proposed an estimator based on Hutchinson trace estimator and conjugate gradients, which can potentially speed up optimizations when the number of conjugate gradient steps is small enough and d is large enough. However, we observed that the exact estimate is largely practical and focused on it instead.

2.3 Persistent homology

In order to measure whether we succeeded in learning and engaging with manifolds with non-trivial topology, we rely on the tool of topological data analysis. We briefly introduce the concept and tools for persistent homology based on Wasserman [2018]. Data can reflect the topological structures of the underlying space, and persistent homology provides a principled tool to capture these topological features, often in the form of a persistence diagram. One often care about topological features

including H_0 , H_1 and H_2 , where H_0 refers to the number of connected components, H_1 and H_2 refer to the number of one and two-dimensional holes respectively. A persistence diagram is plotted by monitoring the births and deaths of topological features as balls of varying radii are drawn around data points, with a longer life span implying a more significant feature. A persistence diagram can be generated from a matrix containing pairwise distances. While pairwise Euclidean distances can in principle be used here, the intrinsic distances may be preferred [Fernández et al., 2023].

3 Geometry of single-chart flows

Largely neglected by previous works concerning normalizing flows on Riemannian manifolds, the flow enables us to learn and utilize the underlying geometric structure of the manifold.

The learned manifold. A fundamental limitation of a single-chart flow is that it assumes that the manifold is globally diffeomorphic to Euclidean. This is an assumption that is violated by many common manifolds, including sphere and torus. Nevertheless, given a trained normalizing flow, one can reason about the learned Riemannian manifold.

Normalizing flows are by definition bijections, so a flow acts as a coordinate chart in the Riemannian sense, and naturally induces a well-defined pullback metric in the latent space U. We further observe that only \tilde{h} contributes to the manifold embedding, while g is just a reparametrization Z = g(U) of the latent space. Hence, we denote $z = \tilde{h}^{\dagger}(x)$, and the Riemannian metric in Z takes the form

$$\boldsymbol{G}_{\boldsymbol{Z}}(\boldsymbol{z}) = \left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right)^{\top} \boldsymbol{G}_{\boldsymbol{X}}(\boldsymbol{x}) \left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right) = \left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right)^{\top} \left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right), \tag{5}$$

where $G_X(x) = \mathbb{I}_D$ is the Euclidean ambient metric. Furthermore, the Jacobian $\frac{\partial x}{\partial z} \in \mathbb{R}^{D \times d}$ maps a vector from the latent space to a vector in the ambient space as $v_X(x) = \frac{\partial x}{\partial z} v_Z(z)$. In practice, the metric can be obtained through explicitly calculating the Jacobian matrices, while the transformation between the vectors can be obtained through a Jacobian vector product operation, which is reasonably cheap for modern automatic differentiation systems [Baydin et al., 2018]. Further discussions on the geometric interpretation can be found in the Appendix.

As discussed next, the flow enables us to compute geodesics, which are essential for many practical operations on a Riemannian manifold, such as the exponential map and logarithmic map, which are necessary for performing computations on the manifold. We remark that one can in principle obtain other geometric quantities, e.g. the Riemannian curvature, using the model.

Exponential maps. Using the Riemannian metric (5), one can solve the exponential map numerically by solving the geodesic equation, which is explicitly formulated below.

Theorem 1. The geodesic equation in the latent space induced by the pullback metric (5) is given by

$$\ddot{z}(t)^k = -g^{kl} \sum_{i,j,m} \frac{\partial^2 x_m}{\partial z_i \partial z_j} \frac{\partial x_m}{\partial z_l} \dot{z}^i \dot{z}^j.$$
(6)

The proof can be found in the Appendix. We remark that this is similar to the geodesic induced by Fisher information matrix with Gaussian likelihood [Song et al., 2018]. The algorithm for solving exponential maps involves explicitly forming the Jacobian matrix, calculating the Riemannian metric and its inverse, thus has complexity $O(d^3)$. However, d can be much smaller than the dimensionality of the ambient space D. Moreover, one can avoid explicitly calculating and storing the higher order gradient tensors. Further discussions on the computational aspects can be found in the Appendix.

Geodesics and Logarithmic maps. A geodesic correspond to the curve with minimum energy. As such, one can parameterize a curve using a cubic spline in the latent space of the flow and optimize its parameters to minimize the energy as observed in the ambient space [Yang et al., 2018, Detlefsen et al., 2021]. Having the optimized curve, we can estimate the distance between the two end points using discretization techniques, and the logarithmic map as the initial velocity.

Pathology of single-chart flows. A single-chart flow is, by definition, a bijection between a *d*dimensional Euclidean space and the learned manifold. As such, the learned manifold is constrained to be diffeomorphic to a Euclidean space. This is a fundamental limitation: when the underlying data manifold is not diffeomorphic to Euclidean, the flow cannot reliably represent its global structure. Perhaps the simplest example is the circle: although it can be parameterized by a single angular coordinate, the fact that *a* and $a + 2\pi$ refer to the same point for any $a \in \mathbb{R}$ means it is not Euclidean. In practice, covering the circle with a single flow inevitably introduces a small "gap" in the learned manifold. As a result, the implied geodesics are necessarily suboptimal as a geodesic needs to traverse the entire circle to connect the two sides of the gap. In addition, the implied topology is wrong.

4 Modeling using multi-chart flows

We resolve the fundamental issue mentioned above, namely, the limitations caused by using a single flow, by considering a mixture of flows that are carefully designed to properly cover the data manifold while respecting its geometric and topological structure. In particular, we aim for each flow to focus on a specific region of the manifold, and we rely on a probabilistic formulation to ensure that overlapping regions exist which is crucial when covering a manifold. We then leverage this construction to extend standard geometric computations, such as geodesics, to this data-driven manifold. We note that the main goal of the proposed generative model is to provide the suitable foundation on which geometric quantities can be computed, while preserving the intrinsic properties of the manifold.

Model. The discussions above and in Section 3 suggest that we need to use multiple normalizing flows, instead of a single one. We hence consider the following mixture model

$$\sum_{i} \log q_{\boldsymbol{\theta}}(\boldsymbol{x}_{i}) = \sum_{i} \log \left(\sum_{c=1}^{C} q_{\boldsymbol{\theta}}(\boldsymbol{x}_{i} | c) q(c) \right), \tag{7}$$

where $C \in \mathbb{Z}^+$ denotes the number of charts, with each $q_{\theta}(\boldsymbol{x} | c)$ being a separate normalizing flow with latent dimensionality d and ambient dimensionality D focusing on the chart c.

Training. With the mixture model formulation, our goal is to minimize the distance between the model distribution and the ground truth distribution. We can write the resulting KL divergence as

$$\operatorname{KL}(p(\boldsymbol{x})|q_{\boldsymbol{\theta}}(\boldsymbol{x})) = \operatorname{KL}\left(p(\boldsymbol{x})\left|\sum_{c} q_{\boldsymbol{\theta}}(\boldsymbol{x}|c)q(c)\right.\right).$$
(8)

However, as noted by previous works [Brehmer and Cranmer, 2020, Caterini et al., 2021], performing purely MLE training is not sufficient for fitting the degenerate flow to a distribution supported on the manifold. Due to degeneracy, the log-density as calculated by the normalizing flow is in fact the log-density projected onto the learned manifold. Therefore, inspired by Caterini et al. [2021], we calculate the required quantities by adding a regularization term to the log-density of the flow. The resulting log-density for a single data point is given by

$$\log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}_i | c) = \log q_{\boldsymbol{\theta}}(\boldsymbol{x}_i | c) + \lambda \| \boldsymbol{x}_i - \tilde{\boldsymbol{h}}_c \left(\tilde{\boldsymbol{h}}_c^{\dagger}(\boldsymbol{x}_i) \right) \|^2,$$
(9)

depending on the problem one may add additional regularization terms here. Upon optimality, one can expect the different flows to learn essentially the same reconstructions in overlapping regions, while modeling the target density through a mixture model. There can be potentially an infinite number of models that fit the distribution perfectly, resulting in recovering the same manifold structure.

For training a mixture model, there are generally two alternative approaches, i.e. directly performing Maximum Likelihood Estimation (MLE) and using Expectation Maximization (EM). These approaches can be employed in our case as well.

Maximum likelihood estimation (MLE). Using $\log \tilde{q}_{\theta}(x_i | c)$ as in 9, one can define

$$\log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{i} \log \sum_{c} \exp\left(\log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i} | c) + \log q(c)\right), \tag{10}$$

which can be implemented using the *logsumexp* operation, enabling direct MLE training of the flows.

Expectation-Maximization (EM). EM is arguably the most commonly used algorithm for training mixture models, and we adapt the EM algorithm to train multi-chart flows. During EM training, the E-step and the M-step are iteratively applied. In one E-step, the responsibilities of the individual flows can be calculated as

$$r_{c_i}(\boldsymbol{x}_i) = \text{stop_gradient}\left(\text{SoftMax}\left(\log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}_i \mid c_i)\right)\right), \tag{11}$$

while in one M-step, one minimizes the loss function

$$L(\boldsymbol{\theta}) = -\sum_{i} \sum_{c_i} r_{c_i}(\boldsymbol{x}_i) \log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}_i | c_i).$$
(12)

In our case, we employ a single gradient descent update for the M-step. Further discussions on the EM procedure can be found in the Appendix. Intuitively, one cannot expect a flow to correctly model data points that are far from its responsible area. As such, for an individual flow we only calculate the loss based on points where its responsibilities are above a threshold. Note that the probabilistic formulation using responsibilities allows the charts to overlap in regions of the data manifold.

General strategy. In preliminary experiments, we did not observe fundamental differences between direct MLE and EM. However, the responsibilities as provided by EM offer interpretability and ways to regularize the flows. In this work, we focus on the EM approach. For training mixture models, one common strategy is to initialize the clusters using simple clustering algorithms like K-Means [Bishop, 2006]. Similarly, we initialize the model by applying K-Means to obtain C clusters, which are then used to warm-up the charts by training each of the flows individually.

We remark that training a mixture of normalizing flows has been explored before, where both MLE and EM have been demonstrated to work [Izmailov et al., 2020, Atanov et al., 2020, Ng and Zammit-Mangion, 2023]. Nevertheless, we are unaware of works that demonstrate the successful training of mixtures of flows where both the data manifold and the distribution are unknown.

5 Geometry of multi-chart flows

The learned manifold. Similar to the case of classical differential geometry, with multi-chart flows, there may not exist a unique latent space. However, one can still get the correct geometric structure from all of them. When the flows are perfectly trained, each individual flow acts as a local chart due to its bijective nature, and we obtain a well-defined atlas that covers the entire manifold by collecting all the charts. However, in practice, we only have access to a finite and potentially noisy dataset, the modeling and optimization are most likely imperfect, and one cannot expect the flows to be perfect.

Nevertheless, we argue that still we can reason geometrically through suboptimal flows by utilizing a probabilistic treatment, which has been demonstrated useful for learning meaningful geometry in VAEs [Arvanitidis et al., 2018, Hauberg, 2019]. Denote the value of interest as e. Given C trained flows, since each flow would give us an estimate, we can collect all these estimates to form a set $\{e_c, c = 1, \ldots, C\}$, and the task is to assign the right weight to each e_c . From a Bayesian perspective, a natural approach to obtain a single estimate e is to use the predictive posterior $\mathbb{E}_{p(c|\mathbf{x})}[e_c]$.

Interestingly, these weights are precisely the responsibilities given by the EM algorithm. As such, we define the points that lie on the manifold as the set $\{\sum_{c=1,...,C} r_c(x')\tilde{h}_c(\tilde{h}_c^{\dagger}(x')), \forall x' \in \mathbb{R}^D\}$. Intuitively, this can be seen as a projection of the points x' in the ambient space onto the learned manifold which is well-defined when the points x' are sufficiently close to it. We can similarly define the tangent space at a point x as the tangent spaces of different charts weighted by r(x). When all the individual normalizing flows yield the correct estimate at the same point, these estimates correspond to the ground truth values. Otherwise, the probabilistic treatment results in estimates that are more robust. We remark that, in practice, we may set a threshold on the responsibilities and form the predictive only based on those that are above the threshold.

Exponential maps. Due to the involved definition of manifold structure as discussed above, we need specialized tools to compute the exponential map. Perhaps the most natural way is to apply an algorithm inspired by classical differential geometry, where one solves the exponential map based only on the chart with maximal responsibility, and jumps to another when the most responsible one changes. However, as the responsibilities of the charts decrease, the estimates of all the charts can

Algorithm 1: Algorithm for solving the exponential maps using a fixed T.

 $\begin{array}{l} \textbf{for } t \leftarrow 0 \textbf{ to } T - 1 \textbf{ do} \\ [c, r] \leftarrow \text{filter (resps } (\boldsymbol{x}_t), \text{resp_threshold}); \\ \textbf{for } c \in \boldsymbol{c} \textbf{ do} \\ | [\boldsymbol{x}_{t+1}^c, \boldsymbol{v}_{t+1}^c] \leftarrow \text{Exp}_{c,\Delta t} (\boldsymbol{x}_t, \boldsymbol{v}_t); \\ \boldsymbol{x}_{t+1} \leftarrow \sum_c r_c \boldsymbol{x}_{t+1}^c; \\ \boldsymbol{v}_{t+1} \leftarrow \sum_c r_c \boldsymbol{v}_{t+1}^c; \\ \boldsymbol{v}_{t+1} \leftarrow \sum_c r_c \boldsymbol{v}_{t+1}^c; \end{aligned}$

Algorithm 2: Algorithm for solving the geodesics and logarithmic maps.

 Initialize $\gamma_{\phi}(t)$;

 for $i \leftarrow 1$ to N_{iters} do

 for $t \leftarrow 0$ to T - 1 do

 $| [c, r] \leftarrow filter (resps (x_t), resp_threshold);$
 $x_t \leftarrow \sum_c r_c(\gamma_{\phi}(t))\tilde{h}_c \left(\tilde{h}_c^{\dagger}(\gamma_{\phi}(t))\right);$

 Minimize $E(\{x_t, t = 1, \dots, T - 1\})$ with respect to $\gamma_{\phi}(t)$

 Fit a new curve $\tilde{\gamma}_{\phi}(t)$ to $\sum_c r_c(\gamma_{\phi}(t))\tilde{h}_c \left(\tilde{h}_c^{\dagger}(\gamma_{\phi}(t))\right);$

 Compute v_0 from $\tilde{\gamma}_{\phi}(t)$ as the initial velocity;

degrade in different ways, such that the jump may not happen at the right moment. Practically, the numerical solvers make it even more challenging. Therefore, we argue that it is beneficial to solve the exponential maps using the responsibilities as calculated in the EM algorithm as weights.

Algorithm 1 provides an Euler-like method that addresses the issue mentioned above. We consider computing a unique curve directly on the ambient manifold, and instead of relying exclusively on the exponential map from a single chart, each update step is computed as the weighted average over multiple charts using the corresponding responsibilities. In each responsible chart, we simulate the geodesic for Δt time starting from x_t and v_t . This enhances the robustness of the resulting curve. For further ablation studies on the different solvers, we refer interested readers to the Appendix.

Geodesics and Logarithmic maps. To compute the geodesic between two points, and thus the logarithmic map, based on the solution above and Section 3, one might be tempted to parameterize a curve in the latent spaces of all the individual charts. However, a flow can be arbitrarily bad outside its responsible region, and it is unclear how to obtain the final curve based on the ones induced by the charts. Since a geodesic can naturally transition between charts, a more direct approach is to parameterize a curve across multiple latent spaces while accounting for the transitions. This leads to a particularly challenging optimization problem, and thus is prohibited for practical purposes.

Our key insight is that multi-chart flows enable us to map points onto the manifold, through the function $\mathbf{x} \to \sum_{c=1,...,C} r_c(\mathbf{x}) \tilde{\mathbf{h}}_c \left(\tilde{\mathbf{h}}_i^{\dagger}(\mathbf{x}) \right)$, which can map a curve on the manifold from a curve defined in the ambient space. Using this, we propose to parameterize in the ambient space a curve $\gamma_{\phi}(t)$, e.g., a spline, and optimize the energy of the mapped curve on the learned manifold. The logarithmic map can then be obtained as the initial velocity of the resulting curve (see Algorithm 2).

6 Experiments

We perform experiments across various datasets to demonstrate the utility of multi-chart flows to model manifold valued data while implicitly capturing the underlying geometry and topology. We compare multi-chart flows (MULT), with single-chart flows trained through both sequential manifold learning and density estimation (SINGLE- \mathcal{M}), and direct MLE training (SINGLE). For the geodesics, it is known that solutions are not unique; e.g., given two points that are approximately on the opposite of the sphere the solver may just as well find the curve completely opposite of the shortest geodesic, which results in different logarithmic map yet reasonable distance. Hence, we consider the distances instead of the logarithmic maps. We use for evaluation the reconstruction loss measured on a test set,



Figure 2: Evaluation metrics on sphere and torus, best methods are circled. MULT consistently leads to better performances, since geodesics are not unique, in some cases, it does not find the shortest.



Figure 3: Persistence diagrams for the Sphere-U data. Points further to the top left corner indicate more significant feature. MULT provides the most faithful representation of the underlying topology.

the Wasserstein distance of generated samples to a test set, the mean squared errors of the estimated exponential maps and their distances to the ground truths, and the faithfulness of the persistence diagram in capturing the topology. For fair comparisons, we keep the number of parameters for each method identical or comparable. We provide more experimental details in the Appendix.

Sphere and torus. We consider modeling distributions supported on two dimensional sphere and two dimensional torus using RealNVP flows, where the data and the ground truth geometric quantities are obtained using Geomstats [Miolane et al., 2020, 2024] and Pyro [Bingham et al., 2019]. For SINGLE- \mathcal{M} and SINGLE, we use 12 layers for g and 36 layers for h. For MULT, we use 4 charts, each chart with 3 layers for g and 9 layers for h. For both manifolds, we consider both the uniform distribution, denoted as -U, and the mixture of Vom-Mises Fisher and Bivariate Von-Mises distributions, denoted as -M, respectively. As shown in Figure 2, MULT yields the best performances in all cases considered. We additionally report the persistence diagrams computed based on the corresponding distance estimates. If these pairwise distances do not accurately reflect the topology of the underlying manifold, the resulting persistent homology may fail to recover the true topological features. Figure 3 verifies that the distances as learned using MULT allows to correctly identify the topological structure of the sphere; further diagrams can be found in the Appendix.

Triangular meshes. We use RealNVP flows with identical structures as before. The triangular meshes are provided by Trimesh [Dawson-Haggerty et al., 2025]. Inspired by Trimesh, we refine the mesh and approximate the ground truth distances between the points using a neighborhood graph; for a dense enough mesh this approximation is reasonable. The results in Figure 4 show that for these irregular shapes MULT still provides consistently better reconstructions and distance measures. Also, it enables calculating the exponential maps and logarithmic maps, even if the true ones are intractable.

Motion capture data. We use the motion capture data from MocapToolbox [Burger and Toiviainen, 2013] and generate datasets where the frame is randomly rotated. We consider both the case where the axis is fixed and the angle is uniformly at random which forms a 1-dimensional manifold (Mocap1), and the case where both the axis and the angle are uniformly at random, which forms a 3-dimensional manifold due to the connection to the well-known SO(3) manifold (Mocap3). We remark that the embeddings of the manifold may not be isometric and the ground truth exponential maps and logarithmic maps are intractable. Nevertheless, one should expect the topology to be preserved. We consider Neural Spline flows. For Mocap1, MULT uses 2 charts, each with 6 layers for g and 12 layers for h; SINGLE- \mathcal{M} and SINGLE use 12 layers for g and 24 layers for h. For Mocap3, MULT uses 5 charts, each with 4 layers for g and 8 layers for h, while SINGLE- \mathcal{M} and SINGLE use 20 layers for g and 40 layers for h. As shown in Figure 5, MULT always results in the best reconstructions and



Figure 4: Left: evaluation metrics on triangular meshes, best methods are circled. MULT consistently leads to better reconstructions and distance estimates while having competitive sample quality. Right: a trained multi-chart flow enables us to estimate the logarithmic maps.



Figure 5: Left: evaluation metrics on Mocap, best methods are circled. MULT has the best reconstructions for both Mocap1 and Mocap3, and is competitive in terms of sample quality, being the best in terms of Mocap3. Right: The Mocap frame is rotated, here along the fixed axis.

on-par or better sample quality. Additionally, we analyzed the distances as learned by the models in terms of Mocap1. Table 1 shows that only MULT can learn the correct circular structure of Mocap1.

7 Discussion

A key contribution of our paper is to handle non-overlapping charts as they are bound to happen in machine learning applications. We thus developed the first numerical algorithms, to the best of our knowledge, for computing geodesics on such multi-charted manifolds. We also demonstrated that multiple charts are essential when aiming to capture topology. While this is long-standing knowledge in mathematics, the topic seems to have been underappreciated in the machine learning community.

From the generative modeling perspective, one can make arguments that composing multiple smaller models is better than utilizing a single large model [Du and Kaelbling, 2024]. Driven by a different motivation, we make similar findings. It is an interesting question whether the success of model composition is in part due to the non-trivial manifold structure of real world datasets.

We focus on normalizing flows which enable fast and exact density evaluations, due to their ability to achieve fast inference and efficient computations of geometric quantities involving the Jacobians, along with the tractability of training degenerate normalizing flows. Flow matching [Lipman et al., 2023] is a recent family of algorithms that learns a continuous normalizing flow parameterized by a vector field, without needing numerical integrations during training. It has been extended to Riemannian manifolds [Chen and Lipman, 2024], however both the base distribution and the velocities are defined on the target manifold itself. Free-form flows [Draxler et al., 2024, Sorrenson]

Table 1: The lengths of the 5 segments between 5 points uniformly placed along the topologicallycircular learned data manifold for each method. Only MULT correctly reflects the circular nature of the data. Single-chart flows show a linear-like structure where one segment is significantly larger.

Method	Lengths (me	ean and varianc	e over three mo	odels randoml	y initialized)
Single-M Single Mult	$\begin{array}{c} 11.7 \pm 0.2 \\ 23.0 \pm 16.2 \\ 11.6 \pm 0.0 \end{array}$	$\begin{array}{c} 11.7 \pm 0.1 \\ 34.4 \pm 16.2 \\ 11.7 \pm 0.0 \end{array}$	35.1 ± 16.9 11.3 ± 0.0 11.3 ± 0.0	$\begin{array}{c} 12.0 \pm 0.3 \\ 11.8 \pm 0.0 \\ 11.8 \pm 0.0 \end{array}$	23.0 ± 16.4 11.3 ± 0.0 11.3 ± 0.0

et al., 2024] is another family of normalizing flow method that has been extended to Riemannian manifolds. However, they are not guaranteed to be bijections, making the differential geometric interpretations less grounded. In this paper we assume that the dimensionality is known beforehand. However, as shown by Zhang et al. [2023], it is possible to identify the intrinsic dimensionality while training the flow. We leave further explorations utilizing these methods as future work.

Limitations. Multi-chart flows are more technically involved and can be difficult to train especially for manifolds with high curvature. Similarly, computing geodesics becomes more challenging as curvature increases. However, our proposed methodology is, in principle, generic as more advanced flow models can be easily incorporated to improve the performance.

Conclusion. Single-chart flows are fundamentally incapable to respect complex geometric structures. In contrast, we showed that this is possible when using multi-chart flows which can capture the geometry and topology of the underlying space while learning the distribution. In addition, we provided specifically designed numerical methods to approximate geometric quantities, such as geodesics, on the learned manifolds. Overall, our work enables computational differential geometry on data manifolds with non-trivial topology.

Acknowledgments and Disclosure of Funding

HY, MH and AK were supported by the Research Council of Finland Flagship programme: Finnish Center for Artificial Intelligence FCAI, and by the grants 345811, 363317 and 348952. SH was supported by a research grant (42062) from VILLUM FONDEN and partly funded by the Novo Nordisk Foundation through the Center for Basic Research in Life Science (NNF20OC0062606). SH received funding from the European Research Council (ERC) under the European Union's Horizon Programme (grant agreement 101125003). GA was supported by the DFF Sapere Aude Starting Grant "GADL". The authors wish to acknowledge CSC - IT Center for Science, Finland, for computational resources.

References

- F. Acosta, S. Sanborn, K. D. Duc, M. Madhav, and N. Miolane. Quantifying Extrinsic Curvature in Neural Manifolds. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pages 610–619, 2023. doi: 10.1109/CVPRW59228.2023.00068.
- J. Ansel, E. Yang, H. He, N. Gimelshein, A. Jain, M. Voznesensky, B. Bao, P. Bell, D. Berard, E. Burovski, G. Chauhan, A. Chourdia, W. Constable, A. Desmaison, Z. DeVito, E. Ellison, W. Feng, J. Gong, M. Gschwind, B. Hirsh, S. Huang, K. Kalambarkar, L. Kirsch, M. Lazos, M. Lezcano, Y. Liang, J. Liang, Y. Lu, C. Luk, B. Maher, Y. Pan, C. Puhrsch, M. Reso, M. Saroufim, M. Y. Siraichi, H. Suk, M. Suo, P. Tillet, E. Wang, X. Wang, W. Wen, S. Zhang, X. Zhao, K. Zhou, R. Zou, A. Mathews, G. Chanan, P. Wu, and S. Chintala. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24). ACM, Apr. 2024.
- G. Arvanitidis, L. K. Hansen, and S. Hauberg. Latent Space Oddity: on the Curvature of Deep Generative Models. In *International Conference on Learning Representations*, 2018.
- A. Atanov, A. Volokhova, A. Ashukha, I. Sosnovik, and D. Vetrov. Semi-Conditional Normalizing Flows for Semi-Supervised Learning, 2020. _eprint: 1905.00505.
- A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind. Automatic Differentiation in Machine Learning: a Survey. *Journal of Machine Learning Research*, 18(153):1–43, 2018.
- K. Bennett, P. Bradley, and A. Demiriz. Constrained K-Means Clustering. Technical Report MSR-TR-2000-65, May 2000.
- E. Bingham, J. P. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. A. Szerlip, P. Horsfall, and N. D. Goodman. Pyro: Deep Universal Probabilistic Programming. J. Mach. Learn. Res., 20:28:1–28:6, 2019.

- C. M. Bishop. Pattern Recognition and Machine Learning. Information Science and Statistics. Spinger Science+Busines Media, LLC, New York, New York, USA, 2006.
- J. Brehmer and K. Cranmer. Flows for simultaneous manifold learning and density estimation. In *Advances in Neural Information Processing Systems*, volume 33, pages 442–453. Curran Associates, Inc., 2020.
- E. Brinkman. Fibonacci Lattice, 2025. URL https://github.com/erikbrinkman/fibonacci_lattice.
- B. Burger and P. Toiviainen. MoCap Toolbox A Matlab toolbox for computational analysis of movement data. In R. Bresin, editor, *Proceedings of the 10th Sound and Music Computing Conference*, pages 172–178, Stockholm, Sweden, 2013. KTH Royal Institute of Technology.
- A. L. Caterini, G. Loaiza-Ganem, G. Pleiss, and J. P. Cunningham. Rectangular Flows for Manifold Learning. In Advances in Neural Information Processing Systems, 2021.
- R. T. Q. Chen and Y. Lipman. Flow Matching on General Geometries. In *The Twelfth International Conference on Learning Representations*, 2024.

Dawson-Haggerty et al. trimesh, Mar. 2025. URL https://trimesh.org/.

- N. S. Detlefsen, A. Pouplin, C. W. Feldager, C. Geng, D. Kalatzis, H. Hauschultz, M. González-Duque, F. Warburg, M. Miani, and S. Hauberg. StochMan. *GitHub. Note:* https://github.com/MachineLearningLifeScience/stochman/, 2021.
- W. Diepeveen. Pulling back symmetric Riemannian geometry for data analysis, 2024. _eprint: 2403.06612.
- L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using Real NVP. In *International Conference on Learning Representations*, 2017.
- M. P. Do Carmo. Riemannian Geometry. Mathematics. Theory & applications. Birkhäuser, 1992.
- F. Draxler, P. Sorrenson, L. Zimmermann, A. Rousselot, and U. Köthe. Free-form Flows: Make Any Architecture a Normalizing Flow. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 2197–2205. PMLR, May 2024.
- Y. Du and L. P. Kaelbling. Position: Compositional Generative Modeling: A Single Model is Not All You Need. In *Forty-first International Conference on Machine Learning*, 2024.
- C. Durkan, A. Bekasov, I. Murray, and G. Papamakarios. Neural Spline Flows. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- X. Fernández, E. Borghini, G. Mindlin, and P. Groisman. Intrinsic Persistent Homology via Densitybased Metric Learning. *Journal of Machine Learning Research*, 24(75):1–42, 2023.
- R. Flamary, C. Vincent-Cuaz, N. Courty, A. Gramfort, O. Kachaiev, H. Quang Tran, L. David, C. Bonet, N. Cassereau, T. Gnassounou, E. Tanguy, J. Delon, A. Collas, S. Mazelet, L. Chapel, T. Kerdoncuff, X. Yu, M. Feickert, P. Krzakala, T. Liu, and E. Fernandes Montesuma. POT Python Optimal Transport. URL https://github.com/Python0T/P0T.
- A. A. Hagberg, D. A. Schult, and P. J. Swart. Exploring Network Structure, Dynamics, and Function using NetworkX. In G. Varoquaux, T. Vaught, and J. Millman, editors, *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA, 2008.
- C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. Fernández del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585:357–362, 2020.
- S. Hauberg. Only Bayes should learn a manifold (on the estimation of differential geometric structure from data), 2019. _eprint: 1806.04994.

- M. Hoffman, P. Sountsov, J. V. Dillon, I. Langmore, D. Tran, and S. Vasudevan. NeuTra-lizing Bad Geometry in Hamiltonian Monte Carlo Using Neural Transport, 2019. _eprint: 1903.03704.
- P. Houdouin, E. Ollila, and F. Pascal. Regularized EM Algorithm. In ICASSP 2023 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5, 2023.
- P. Izmailov, P. Kirichenko, M. Finzi, and A. G. Wilson. Semi-Supervised Learning with Normalizing Flows. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4615–4630. PMLR, July 2020.
- D. Kalatzis, J. Z. Ye, A. Pouplin, J. Wohlert, and S. Hauberg. Density estimation on smooth manifolds with normalizing flows, 2022. _eprint: 2106.03500.
- F. d. Kruiff, E. Bekkers, O. Öktem, C.-B. Schönlieb, and W. Diepeveen. Pullback Flow Matching on Data Manifolds, 2024. _eprint: 2410.04543.
- Lee, John M. Introduction to Riemannian Manifolds. Graduate Texts in Mathematics. Springer International Publishing AG, Cham, Switzerland, 2nd edition, 2018.
- J. Levy-Kramer. k-means-constrained, Apr. 2018. URL https://github.com/joshlk/ k-means-constrained.
- H. Li, K. Zhang, and T. Jiang. The regularized EM algorithm. In *Proceedings of the 20th National Conference on Artificial Intelligence Volume 2*, AAAI'05, pages 807–812. AAAI Press, 2005. Place: Pittsburgh, Pennsylvania.
- Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow Matching for Generative Modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- G. Loaiza-Ganem, B. L. Ross, J. C. Cresswell, and A. L. Caterini. Diagnosing and Fixing Manifold Overfitting in Deep Generative Models. *Transactions on Machine Learning Research*, 2022.
- G. Loaiza-Ganem, B. L. Ross, R. Hosseinzadeh, A. L. Caterini, and J. C. Cresswell. Deep Generative Models through the Lens of the Manifold Hypothesis: A Survey and New Connections. *Transactions on Machine Learning Research*, 2024.
- N. Miolane, N. Guigui, A. L. Brigant, J. Mathe, B. Hou, Y. Thanwerdas, S. Heyder, O. Peltre, N. Koep, H. Zaatiti, H. Hajri, Y. Cabanes, T. Gerald, P. Chauchat, C. Shewmake, D. Brooks, B. Kainz, C. Donnat, S. Holmes, and X. Pennec. Geomstats: A Python Package for Riemannian Geometry in Machine Learning. *Journal of Machine Learning Research*, 21(223):1–9, 2020.
- N. Miolane, L. F. Pereira, S. Utpala, N. Guigui, A. L. Brigant, Hzaatiti, Y. Cabanes, J. Mathe, N. Koep, elodiemaignant, ythanwerdas, xpennec, tgeral68, Christian, T. M. Nguyen, O. Peltre, J. Harvey, pchauchat, julesdeschamps, Q. Barthélemy, mortenapedersen, Maya95assal, Abdellaoui-Souhail, A. Myers, F. Ambellan, Florent-Michel, S. Heyder, S. Talbar, Y. d. Mont-Marin, and Marius. geomstats/geomstats: Geomstats v2.8.0, Sept. 2024.
- M. Moor, M. Horn, B. Rieck, and K. Borgwardt. Topological Autoencoders. In Proceedings of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 7045–7054. PMLR, July 2020.
- K. P. Murphy. Probabilistic Machine Learning: Advanced Topics. MIT Press, 2023.
- C. T. Nathaniel Saul. Scikit-TDA: Topological Data Analysis for Python, 2019. URL https://doi.org/10.5281/zenodo.2533369.
- T. L. J. Ng and A. Zammit-Mangion. Mixture Modeling with Normalizing Flows for Spherical Density Estimation, 2023. _eprint: 2301.06404.
- G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan. Normalizing Flows for Probabilistic Modeling and Inference. *Journal of Machine Learning Research*, 22(57): 1–64, 2021.

- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and É. Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- B. L. Ross, G. Loaiza-Ganem, A. L. Caterini, and J. C. Cresswell. Neural Implicit Manifold Learning for Topology-Aware Density Estimation. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- S. Schonsheck, J. Chen, and R. Lai. Chart Auto-Encoders for Manifold Structured Data, 2020. _eprint: 1912.10094.
- S. Sidheekh, C. B. Dock, T. Jain, R. Balan, and M. K. Singh. VQ-Flows: Vector quantized local normalizing flows. In *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 1835–1845. PMLR, Aug. 2022.
- Y. Song, J. Song, and S. Ermon. Accelerating Natural Gradient with Higher-Order Invariance. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4713–4722. PMLR, July 2018.
- P. Sorrenson, F. Draxler, A. Rousselot, S. Hummerich, L. Zimmermann, and U. Koethe. Lifting Architectural Constraints of Injective Flows. In *The Twelfth International Conference on Learning Representations*, 2024.
- V. Stimper, D. Liu, A. Campbell, V. Berenz, L. Ryll, B. Schölkopf, and J. M. Hernández-Lobato. normflows: A PyTorch Package for Normalizing Flows. *Journal of Open Source Software*, 8(86): 5361, June 2023.
- P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17: 261–272, 2020.
- L. Wasserman. Topological Data Analysis. *Annual Review of Statistics and Its Application*, 5(Volume 5, 2018):501–532, 2018. Publisher: Annual Reviews Type: Journal Article.
- T. Yang, G. Arvanitidis, D. Fu, X. Li, and S. Hauberg. Geodesic clustering in deep generative models. In *arXiv preprint*, 2018.
- M. Zhang, Y. Sun, C. Zhang, and S. Mcdonagh. Spread Flows for Manifold Modelling. In *Proceedings* of The 26th International Conference on Artificial Intelligence and Statistics, volume 206 of Proceedings of Machine Learning Research, pages 11435–11456. PMLR, Apr. 2023.

A Proof of Theorem 1

We restate Theorem 1, and provide the proof.

Theorem 1. The geodesic equation in the latent space induced by the pullback metric is given by

$$\ddot{z}(t)^{k} = -g^{kl} \sum_{i,j,m} \frac{\partial^{2} x_{m}}{\partial z_{i} \partial z_{j}} \frac{\partial x_{m}}{\partial z_{l}} \dot{z}^{i} \dot{z}^{j}.$$
(13)

Proof. Recall that the pullback metric is given by

$$\boldsymbol{g} = \left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right)^{\top} \frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}},\tag{14}$$

$$g_{ij} = \sum_{m} \frac{\partial x_m}{\partial z_i} \frac{\partial x_m}{\partial z_j}.$$
(15)

The Christoffel symbols and the geodesic equation are given by [Lee, John M., 2018] (Equation 5.8 and Equation 4.16)

$$\Gamma_{ij}^{k} = \frac{1}{2} g^{kl} \left(\partial_{i} g_{jl} + \partial_{j} g_{il} - \partial_{l} g_{ij} \right), \tag{16}$$

$$\ddot{z}(t)^{k} + \dot{z}^{i}(t)\dot{z}^{j}(t)\Gamma^{k}_{ij}(z(t)) = 0.$$
(17)

We have

$$\partial_l g_{ij} = \sum_m \frac{\partial^2 x_m}{\partial z_i \partial z_l} \frac{\partial x_m}{\partial z_j} + \frac{\partial x_m}{\partial z_i} \frac{\partial^2 x_m}{\partial z_j \partial z_l},\tag{18}$$

and

$$\frac{1}{2} \left(\partial_i g_{jl} + \partial_j g_{il} - \partial_l g_{ij} \right) \tag{19}$$

$$=\frac{1}{2}\left(\sum_{m}\frac{\partial^{2}x_{m}}{\partial z_{j}\partial z_{i}}\frac{\partial x_{m}}{\partial z_{l}}+\frac{\partial x_{m}}{\partial z_{j}}\frac{\partial^{2}x_{m}}{\partial z_{l}\partial z_{i}}+\frac{\partial^{2}x_{m}}{\partial z_{i}\partial z_{j}}\frac{\partial x_{m}}{\partial z_{l}}+\frac{\partial x_{m}}{\partial z_{i}}\frac{\partial^{2}x_{m}}{\partial z_{l}\partial z_{j}}\right)$$
(20)

$$-\frac{\partial^2 x_m}{\partial z_i \partial z_l} \frac{\partial x_m}{\partial z_j} - \frac{\partial x_m}{\partial z_i} \frac{\partial^2 x_m}{\partial z_j \partial z_l} \right)$$
(21)

$$= \frac{1}{2} \left(\sum_{m} 2 \frac{\partial^2 x_m}{\partial z_i \partial z_j} \frac{\partial x_m}{\partial z_l} \right) = \sum_{m} \frac{\partial^2 x_m}{\partial z_i \partial z_j} \frac{\partial x_m}{\partial z_l}.$$
 (22)

As such,

$$\dot{z}^{i}(t)\dot{z}^{j}(t)\Gamma_{ij}^{k}(z(t)) = \dot{z}^{i}(t)\dot{z}^{j}(t)\frac{1}{2}g^{kl}\left(\partial_{i}g_{jl} + \partial_{j}g_{il} - \partial_{l}g_{ij}\right)$$
(23)

$$= \dot{z}^{i}(t)\dot{z}^{j}(t)g^{kl}\sum_{m}\frac{\partial^{2}x_{m}}{\partial z_{i}\partial z_{j}}\frac{\partial x_{m}}{\partial z_{l}} = g^{kl}\sum_{m}\frac{\partial^{2}x_{m}}{\partial z_{i}\partial z_{j}}\frac{\partial x_{m}}{\partial z_{l}}\dot{z}^{i}\dot{z}^{j},$$
(24)

and

$$\ddot{z}(t)^{k} = -g^{kl} \sum_{i,j,m} \frac{\partial^{2} x_{m}}{\partial z_{i} \partial z_{j}} \frac{\partial x_{m}}{\partial z_{l}} \dot{z}^{i} \dot{z}^{j}.$$
(25)

The above expression explicitly gives the form of the geodesic accelerations and affords tractable implementations using Automatic Differentiation systems, as discussed in detail in Section C.2.

B EM training

We briefly explain the entire EM procedure, using as references Murphy [2023] and Bishop [2006]. We provide an outline of the general procedure of using the EM algorithm to train the model, and provide the expressions for the precise EM training procedure in our case.

Denote $f(c_i)$ as an arbitrary set of distributions over c for each x_i . One can write

$$\sum_{i} \log q(\boldsymbol{x}_{i}) = \sum_{i} \log \left(\sum_{c_{i}} q(\boldsymbol{x}_{i}, c_{i}) \right) = \sum_{i} \log \left(\sum_{c_{i}} f(c_{i}) \frac{q(\boldsymbol{x}_{i}, c_{i})}{f(c_{i})} \right).$$
(26)

Using Jensen's inequality, we have

$$\sum_{i} \log(\sum_{c_i} f(c_i) \frac{q(\boldsymbol{x}_i, c_i)}{f(c_i)}) \ge \sum_{i} \sum_{c_i} f(c_i) \log\left(\frac{q(\boldsymbol{x}_i, c_i)}{f(c_i)}\right),$$
(27)

which is the Evidence Lower BOund (ELBO). In order to achieve maximum likelihood training, we would like to maximize the above quantity.

For the E step, observe that

$$\sum_{c_i} f(c_i) \log\left(\frac{q(\boldsymbol{x}_i, c_i)}{f(c_i)}\right) = \sum_{c_i} f(c_i) \left(\log\frac{q(c_i | \boldsymbol{x}_i)q(\boldsymbol{x}_i)}{f(c_i)}\right)$$
(28)

$$= \sum_{c_i} f(c_i) \left(\log \frac{q(c_i | \boldsymbol{x}_i)}{f(c_i)} \right) + \sum_{c_i} f(c_i) \log q(\boldsymbol{x}_i) = -\mathrm{KL} \left(f(c_i) | q(c_i | \boldsymbol{x}_n) \right) + \log q(\boldsymbol{x}_i).$$
(29)

As such, we can set $f(c_i)$ to $q(c_i | \mathbf{x}_i)$ to maximize it. Given $q(c_i)$ and $q(\mathbf{x}_i | c_i)$, $q(c_i | \mathbf{x}_i)$ can be obtained in closed-form using the Bayes rule as

$$q(c_i | \boldsymbol{x}_i) = \frac{q(\boldsymbol{x}_i | c_i)q(c_i)}{q(\boldsymbol{x}_i)} = \frac{q(\boldsymbol{x}_i | c_i)q(c_i)}{\sum_c q(\boldsymbol{x}_i | c)q(c)}.$$
(30)

Note that $f(c_i)$ thus has an intuitive interpretation of *responsibilities* of the charts on the point x_i . We thus refer to it accordingly, and denote it as $r_{c_i}(x_i)$.

For the M step, we use $r_{c_i}(\boldsymbol{x}_i)$ as derived in the E step and optimize for the best $q(\boldsymbol{x}_i | c_i)$ (and possibly $q(c_i)$). In other words, we now need to maximize the following

$$\sum_{i} \sum_{c} r_{c_i}(\boldsymbol{x}_i)(\boldsymbol{x}_i) \log \left(q(\boldsymbol{x}_i, c_i)\right) = \sum_{i} \sum_{c} r_{c_i}(\boldsymbol{x}_i) \log q(\boldsymbol{x}_i | c_i) + \sum_{i} \sum_{c} r_{c_i}(\boldsymbol{x}_i) \log q(c_i).$$
(31)

Observe that the optimization for $q(c_i)$ can be carried out in closed-form, while the optimization for $q(\mathbf{x}_i | c_i)$ affords standard gradient based training.

In the specific case of our multi-chart flows, we directly fix $q(c_i)$ to be uniform over the charts, and train the individual normalizing flows to model $q(x_i | c_i)$. However, since we would like to enforce the flows to perform reconstructions as well, following Caterini et al. [2021] we additionally add a term penalizing the reconstruction error, resulting in

$$\log \tilde{q}_{\boldsymbol{\theta}}(\boldsymbol{x}_i | c_i) = \log q_{\boldsymbol{\theta}}(\boldsymbol{x}_i | c_i) + \lambda \| \boldsymbol{x}_i - \tilde{\boldsymbol{h}_{c_i}}\left(\tilde{\boldsymbol{h}_{c_i}}^{\dagger}(\boldsymbol{x}_i) \right) \|^2.$$
(32)

Following Brehmer and Cranmer [2020], depending on the problem, we may add additional terms, e.g. a term that penalizes the hidden variables $\tilde{h_{c_i}}^{\dagger}(x_i)$.

To summarize, when training the mixture model using EM, we alternate between the E step, which sets $r_{c_i}(\boldsymbol{x}_i)$ based on the flows' log-density evaluations, and the M step, which trains each individual flow based on its responsibilities of the data points.

Additionally, sometimes it is beneficial to include a term penalizing the deviation of the mean responsibilities of the charts averaged over the data points. Since we set q(c) to be uniform over the charts, we expect the mean responsibilities to be close to uniform as well especially when the batch size is larger. We therefore add a loss term in the form of

$$\lambda \sum_{c} \left(\frac{\sum_{i} r_{i}}{N} - \frac{1}{N} \right)^{2}.$$
(33)

We remark that it has been demonstrated useful to add regularization terms in EM in general [Li et al., 2005, Houdouin et al., 2023].

C Geometry

C.1 Details on the geometric operations

Consider \tilde{h} , \tilde{h}^{\dagger} , x and z as defined in the main paper. For points that lie precisely on the learned manifold, \tilde{h} and \tilde{h}^{\dagger} are a pair of bijective functions that are the inverse of each other.

Given a value of z, one can obtain the unique corresponding $x = \tilde{h}(z)$. The pullback metric can also be tractably evaluated as $J_{\tilde{h}}^{\top} J_{\tilde{h}}$. Given v_z , one can also obtain the corresponding v_x using the formula $v_x = \frac{\partial \tilde{h}}{\partial z} v_z$. One way to see that is to note that $v_x = \frac{\partial x}{\partial t} = \frac{\partial x}{\partial z} \frac{\partial z}{\partial t} = \frac{\partial x}{\partial z} v_z$. These operations are generally well-defined.

However, due to the degeneracy, multiple x can correspond to the same z. \tilde{h}^{\dagger} acts as a projection operation that projects data points onto the manifold, thus naturally induces quotient structures. When x lies exactly on the learned manifold and v_x lies exactly on the tangent space, one can obtain the corresponding v_z using $v_z = \frac{\partial \tilde{h}^{\dagger}}{\partial x} v_x$. Due to the way that \tilde{h}^{\dagger} is composed, we have

$$\frac{\partial \tilde{\boldsymbol{h}}^{\dagger}}{\partial \boldsymbol{x}} = \frac{\partial \left(\operatorname{Proj} \cdot \boldsymbol{h}^{-1} \right)}{\partial \boldsymbol{x}} = \frac{\partial \operatorname{Proj}}{\partial \boldsymbol{h}^{-1}(\boldsymbol{x})} \frac{\partial \boldsymbol{h}^{-1}}{\partial \boldsymbol{x}}.$$
(34)

Projection is a linear operation by definition, and its Jacobian matrix is simply given by a rectangular matrix with ones at the entries where the coordinates are retained. As such, when two distinct values x_1 and x_2 correspond to the same z, the Jacobian matrices are the same when the corresponding rows of $\frac{\partial h^{-1}}{\partial x_1}$ and $\frac{\partial h^{-1}}{\partial x_2}$ are the same.

C.2 Implementing the geodesic equation

Here we discuss how to implement an efficient and numerically stable version of the geodesic equation using modern Automatic Differentiation (AD) systems.

Recall the form of the geodesic equation; at each step we need $-g^{kl}\sum_{i,j,m}\frac{\partial^2 x_m}{\partial z_i \partial z_j}\frac{\partial x_m}{\partial z_i}\dot{z}^i\dot{z}^j$. One can divide it into three parts: the inverse of the Riemannian metric g^{kl} , the Jacobian $\frac{\partial x_m}{\partial z_i}$ and the quadratic form $\frac{\partial^2 x_m}{\partial z_i \partial z_j}\dot{z}^i\dot{z}^j$. Noting that the Riemannian metric is precisely given by the outer product of the Jacobian, we only need to obtain the Jacobian and the quadratic form using AD.

Jacobian. In general, for a function $f : \mathbb{R}^n \to \mathbb{R}^m$, the cost to obtain its Jacobian is kn for forward mode AD and km for reverse mode AD, where k is some constant dependent on the function [Baydin et al., 2018]. In our case, we need to obtain the Jacobian of h, which maps from \mathbb{R}^d to \mathbb{R}^D , with d < D. As such, one may prefer forward mode AD. In practice, the batched Jacobian matrices can be implemented by a vectorization of the function that obtains the individual Jacobian over the batch, where for each sample we obtain the Jacobian matrix using forward mode AD.

Quadratic form. Forward mode AD step efficiently calculates Jacobian-vector product [Baydin et al., 2018]. As such, one can first calculate $\frac{\partial x_m}{\partial z_i} \dot{z}^i$, and then calculate $\frac{\partial^2 x_m}{\partial z_i \partial z_j} \dot{z}^i \dot{z}^j$, employing Jacobian-vector product in each. This is to some extent known in the AD community; see e.g. https://github.com/jax-ml/jax/discussions/8456#discussioncomment-1586157.

We remark that we do not claim to have the most efficient implementation, as the focus of the paper is on obtaining the correct geometric quantities instead of having the fastest speed in doing so.

D Experimental details

For open source libraries we provide the license information inside brackets after the citations.

We mainly use PyTorch [Ansel et al., 2024] (custom license), NumPy [Harris et al., 2020] (custom license) and SciPy [Virtanen et al., 2020] (BSD-3 license) for the experiments, with the normalizing flows implementations largely based on Stimper et al. [2023] (MIT license). We use Scikit-TDA [Nathaniel Saul, 2019] (MIT license) to plot the persistence diagrams.

For all models, we tune the learning rates among [1e - 4, 3e - 4, 1e - 3], and tune the reconstruction factors among [100, 1000, 10000]. We tune the hyperparameters using a single run for each configuration, then for the best configuration run two additional runs and report the results obtained using all three runs. For SINGLE- \mathcal{M} , in the first half of all epochs the model is trained solely based on reconstruction errors as induced by h, while in the second half of all epochs the model is trained solely based on log probabilities as induced by g.

For SINGLE and MULT, we use early stopping solely in terms of the reconstruction loss; for SINGLE- \mathcal{M} , in the manifold learning phase we use early stopping in terms of the reconstruction loss, while in the density estimation phase we use early stopping in terms of \mathcal{W} . This may lead to some advantages for SINGLE- \mathcal{M} in terms of \mathcal{W} .

For sphere, torus and triangular meshes experiments, we use RealNVP flows [Dinh et al., 2017], where the MLPs use Tanh activation, and *double* data type. We remark that it is important to use a smooth activation function when intending to solve the IVP, due to the need to differentiate through the flow. We train all models for a maximum of 1000 epochs, using a total of 12000 data points and batch size 256. For MULT, 200 epochs are used for pretraining. The pretraining datasets for the charts are based on K-Means clusters obtained through Scikit-Learn [Pedregosa et al., 2011] (BSD-3-Clause license). We use early stopping with patience 50, starting validation after the model has been trained for 100 epochs in the current stage.

For SINGLE and MULT, we use as validation loss the reconstruction loss on the val set. For SINGLE- \mathcal{M} , we use sequential training, where the first phase uses validation loss the reconstruction loss, and the second phase uses validation loss the Wasserstein distance.

For solving the exponential maps, most numerical integrations are carried out using *scipy.integrate.solve_ivp* from SciPy. For solving the logarithmic maps, our implementation is largely based on Stochman [Detlefsen et al., 2021] (Apache-2.0 license), where we optimize a parameterized spline using gradient descent.

We use Python Optimal Transport [Flamary et al.] (MIT license) to calculate the Wasserstein distances between batches of data and batches of samples drawn using the models, where the subsamplings are performed for 5 times and each batch has 1024 samples.

D.1 Sphere and torus

We consider four distributions on the two dimensional sphere and the two dimensional torus. Specifically, they are: uniform distribution on the sphere (Sphere-U), mixture of Von-Mises Fisher (VMF) distribution on the sphere (Sphere-M), uniform distribution on the torus (Torus-U) and mixture of Bivariate Von-Mishes (BVM) distribution on the torus (Torus-M).

The datasets are generated mainly using as tools Geomstats [Miolane et al., 2020, 2024] (MIT license) and Pyro [Bingham et al., 2019] (Apache-2.0 license).

For the mixture of VMF distribution, we use four components, with the means given by $\left[\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}\right], \left[-\frac{1}{\sqrt{3}}, -\frac{1}{\sqrt{3}}, \frac{1.0}{\sqrt{3}}\right], \left[-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, -\frac{1}{\sqrt{3}}\right], \left[\frac{1}{\sqrt{3}}, -\frac{1}{\sqrt{3}}, -\frac{1}{\sqrt{3}}\right]$ and $\kappa = 5$.

For the mixture of BVM distribution, we use four components specified in angular coordinates. The means are given by $[0,0], [\pi,0], [0,\pi], [\pi,\pi]$. Each component has concentration 1, with the correlation between the two dimensions 0.

When performing evaluations of the exponential maps, the logarithmic maps and distances, based on preliminary experiments, we observe that it is beneficial to choose the evaluation points as distant to each other as possible to make the resulting persistence diagrams clear. As such, on the sphere, we use Fibonacci lattice as implemented by Brinkman [2025] (MIT license). On the torus, since it is the product manifold formed by two circles, we analytically draw samples on the circles using angular coordinates and form the final samples using Cartesian products. For exponential maps, we use 100 data points, resulting in 4950 evaluations. For logarithmic maps and distances we use 225 data points, resulting in 25200 evaluations. Note that the distance of x_1 to x_2 is naturally the same as the distance of x_2 to x_1 , and the distance of x_1 to itself is naturally 0. The exponential maps and logarithmic maps are evaluated based on the pairs determined by *tril_indices* function in NumPy and PyTorch with offset -1.

D.2 Triangular meshes

The triangular meshes experiments are largely based on Trimesh [Dawson-Haggerty et al., 2025] (MIT license). We use 7_8ths_cube.stl, busted.STL and rock.obj.bz2 as provided by Trimesh, while refining the meshes partly to make them smooth and the resulting distance approximations reasonable. Additionally, the meshes are normalized.

For the evaluation of distances, we use 100 data points. Following the reasoning in Section D.1, there are a total of 4950 distance evaluations to be carried out. Due to the lack of a principled approach to draw maximally distant samples on the meshes, the data points are simply drawn uniformly. Following Dawson-Haggerty et al. [2025], each *ground truth* distance between two points is approximated as

the shortest distance between the closest vertices to the points on the graph induced by the mesh and solved using NetworkX [Hagberg et al., 2008] (3-clause BSD license).

D.3 Motion capture data

We use the initial frame from *dance1* as provided by MocapToolbox [Burger and Toiviainen, 2013] (GNU general public license). We construct two variants of Mocap datasets, *Mocap1* and *Mocap3*.

In Mocap1, the initial frame is randomly rotated around the axis given by $[1/\sqrt{14}, 2/\sqrt{14}, 3/\sqrt{14}]$, with the rotation angle uniformly sampled between 0 and 2π . We use 30000 points as the train set, 10000 points as the val set and 10000 points as the test set. In Mocap3, the initial frame is rotated using a random element from SO(3). We use 50000 points as the train set, 10000 points as the test set. For both datasets, we standardize the data based on the train set. In Mocap1 the batch size is 256, while for Mocap3 in preliminary experiments we observe that using a larger batch size stabilizes training especially for MULT, so we employ a batch size of 1024.

We use Neural Spline flows [Durkan et al., 2019] with ReLU activation and *float* data type. For pretraining as in MULT, we use Constrained K-Means clustering [Bennett et al., 2000] as implemented in Levy-Kramer [2018] (BSD-3-Clause license) to obtain the clusters, where the clusters are constrained to be between 0.9t and 1.1t, where t is the target size given by dividing the total number of samples by the number of charts. We use early stopping with infinite patience and perform validations throughout the process. For Mocap1, SINGLE is trained for the first 20 epochs solely based on reconstructions, while during pretraining of MULT which has a total of 30 epochs the first 20 epochs are solely based on reconstructions and the following 10 epochs also account for log probabilities. For Mocap3, the first 50 epochs of SINGLE are solely based on reconstructions and the later 20 epochs also account for log probabilities. Following Brehmer and Cranmer [2020], a loss term with weight 1e - 3 is added to constrain the latent representations. For MULT, an additional loss term with weight 10 is added to encourage the mean responsibilities of the charts to be more uniform.

In Mocap1, the ground truth latent space is by definition 1. As such, we perform an experiment where we calculate the pairwise distance between 5 rotations of the original frame whose angles are uniformly distributed. While it is not clear whether the ground truth distances should be the same for these pairs, we should expect them to reflect the circular nature of the latent space.

D.4 Compute resources

We mainly use a compute cluster to perform the experiments. The CPU is Intel Xeon Gold 6230 with 192GB memory, the GPU is Nvidia Volta V100 with 32GB memory.

Generally, for each CPU job we use 4 CPU cores, while each GPU job uses one GPU and 10 CPU cores. For jobs that involved training the flows, for Mocap datasets the jobs are run on GPU, otherwise they were run on CPU. The running times of the jobs may vary, but they are generally within several hours. SINGLE- \mathcal{M} generally results in faster training than SINGLE and MULT. Experiments involving evaluating the geometric and topological quantities were run on CPU, whose running times may vary and may need more than a day for those with ill-defined geometries. Some preliminary experiments were run, whose results did not make it to the paper, which naturally took up additional compute resources. The total amount of compute used is thus most likely thousands of CPU hours and hundreds of GPU hours.

E Additional experimental results

In the result tables, we use Recons to denote the reconstruction errors, W to denote the Wasserstein distances, Exps to denote the qualities of exponential maps, Logs to denote the qualities of logarithmic maps and Dists to denote the qualities of distances.

Algorithm 3: Hard_switch algorithm for solving the exponential maps. Exp_c denotes the exponential map induced by the *c*th flow.

while $\arg \max (resps(x_t)) = c$ and t < 1 do $\mid x_t \leftarrow \operatorname{Exp}_c(x_t, v_t)$ end

Algorithm 4: Ambient algorithm for solving the exponential maps. Geo_c denotes the function that takes in position and velocity and returns velocity and acceleration of the geodesic induced by the *c*th flow. Integral_{Δt} denotes integrating the ODE for time Δt .

 $\begin{array}{l} \text{for } t \leftarrow 0 \text{ to } T_{auto} - 1 \text{ do} \\ & \left[c, r \right] \leftarrow \text{filter} \left(\text{resps} \left(x_t \right), \text{resp_threshold} \right); \\ & \text{for } c \text{ in } c \text{ do} \\ & \left| v_{t+1}^c, a_{t+1}^c \leftarrow \text{Geo}_c \left(x_t, v_t \right); \\ v_{t+1} \leftarrow \sum_c r_c v_{t+1}^c, a_{t+1} \leftarrow \sum_c r_c a_{t+1}^c; \\ x_{t+1}, v_{t+1} \leftarrow \text{Integral}_{\Delta t} (x_t, v_t, v_{t+1}, a_{t+1}); \end{array} \right.$

// Get filtered responsibilities

E.1 Different exponential map solvers for multi-chart flows

Algorithm 1 in the main paper introduces one algorithm for solving for exponential maps based on multi-chart flows. We refer to that algorithm as Euler, and discuss two alternative algorithms, Hard_switch and Ambient.

Hard_switch. The naive algorithm for solving the exponential maps involves following along the geodesic in one chart at a time, and is presented in Algorithm 3. We refer to this algorithm as *Hard_switch*. In terms of numerical implementations, this can be achieved by specifying an event that terminates the current integration when the most responsible chart changes and jumps to another chart upon the event happens. With *scipy.integrate.solve_ivp*, we can specify the event in two ways: it can be a *discrete* event checking whether the most responsible chart is the current one or a *continuous* event checking the value of the difference between responsibilities of the current chart and the second most responsible one.

Ambient. In the Euler algorithm, there is a hyperparameter T that needs to be tuned. One interesting question to be asked is, what if we take the infinite limit? This question leads to an alternative algorithm where one averages over the velocities and accelerations, so as to directly perform the integrations of the geodesics in the ambient space. We derive the acceleration of geodesics when viewed in the ambient space.

Theorem 2. The acceleration of a geodesic can be expressed in the ambient space as

$$a_X = \frac{\partial J}{\partial z} v_Z v_Z + J a_Z, \qquad (35)$$

where $a_{\mathbf{Z}}$ is the acceleration as calculated in Theorem 1.

Proof. Recall that the velocity in the latent space and the velocity in the ambient space are related by $v_X = J v_Z$. Taking gradient with respect to t on both sides, we have

$$\boldsymbol{a}_{\boldsymbol{X}} = \frac{\partial \left(\boldsymbol{J} \, \boldsymbol{v}_{\boldsymbol{Z}}\right)}{\partial t} = \frac{\partial \, \boldsymbol{J}}{\partial t} \, \boldsymbol{v}_{\boldsymbol{Z}} + \boldsymbol{J} \, \boldsymbol{a}_{\boldsymbol{Z}} = \frac{\partial \, \boldsymbol{J}}{\partial \, \boldsymbol{z}} \, \boldsymbol{v}_{\boldsymbol{Z}} \, \boldsymbol{v}_{\boldsymbol{Z}} + \boldsymbol{J} \, \boldsymbol{a}_{\boldsymbol{Z}} \,. \tag{36}$$

Interestingly, observe that the first term also appears in Theorem 1. Moreover, after a_Z is calculated, we naturally have access to J. As such, the acceleration in the ambient space can be calculated with little extra computational overhead. The resulting algorithm is outlined in Algorithm 4.

Table 2: The errors of the exponential map integrators in the form of mean \pm std, lower is better. The best are highlighted in bold, and the ones with failed runs are italic. *hard switch* is notably worse than the other two.

Data	Euler	Ambient	Hard_switch
Sphere-U	$9.55\!\cdot\!10^{-3}\pm 6.31\!\cdot\!10^{-3}$	$6.5 \!\cdot\! \mathbf{10^{-3}} \pm 2.73 \!\cdot\! 10^{-3}$	$1.39\!\cdot\!10^{-1}\pm2.67\!\cdot\!10^{-2}$
Torus-U	$7.31 \cdot 10^{-2} \pm 1.26 \cdot 10^{-2}$	$3.5 \cdot 10^{-2} \pm 3.37 \cdot 10^{-3}$	$5.69 \cdot 10^{-1} \pm 1.05 \cdot 10^{-1}$
Sphere-M	$1.12 \cdot 10^{-2} \pm 6.33 \cdot 10^{-3}$	$8.79 \cdot \mathbf{10^{-3}} \pm 2.44 \cdot 10^{-4}$	$2.24 \cdot 10^{-1} \pm 4.45 \cdot 10^{-2}$
Torus-M	$1.07 \cdot 10^{0} \pm 2.39 \cdot 10^{-1}$	$2.14 \cdot 10^{16} \pm 2.14 \cdot 10^{16}$	$1.87{\cdot}10^0\pm1.63{\cdot}10^{-1}$

Table 3: Evaluation metrics on circle in the form of mean \pm std, lower is better. The best are highlighted in bold, and the ones with failed runs are italic.

Method	Recons	W
Single- <i>M</i> Single	$\begin{array}{c} 3.65 \cdot 10^{-3} \pm 3.66 \cdot 10^{-3} \\ 4.68 \cdot 10^{-3} \pm 5.63 \cdot 10^{-3} \end{array}$	$\begin{array}{c} 6.87 \cdot 10^{-1} \pm 4.3 \cdot 10^{-1} \\ 1.1 \cdot 10^{-1} \pm 2.49 \cdot 10^{-2} \end{array}$
Mult-Direct Mult	$7.73 \cdot 10^{-7} \pm 4.32 \cdot 10^{-7} 5.48 \cdot 10^{-7} \pm 1.2 \cdot 10^{-7}$	$\frac{8.24 \cdot 10^{-2} \pm 2.88 \cdot 10^{-2}}{8.56 \cdot 10^{-2} \pm 2.79 \cdot 10^{-2}}$

Comparisons. We compare the *discrete* variant of *Hard_switch*, *Euler* and *Ambient*. The results are shown in Table E.1. In practice, we observe that Euler as outlined in Algorithm 1 offers reasonable performances and stabilities. Ambient yields good performances in some cases, but is unstable in some others. As expected, Hard_switch consistently performs worse than Euler.

E.2 Results on circle

The evaluation metrics of the different algorithms trained to model the uniform distribution on a circle are shown in Figure 8, and the numerical results can be found in Table E.2. We do not observe a qualitative difference between direct MLE and EM. We additionally report example samples generated by the models in Figure 6 and example persistence diagrams computed based on the models in Figure 7. Both MULT-DIRECT and MULT are able to capture the topology well.

E.3 Results on sphere and torus

We report the numerical results of the different models on sphere and torus in Table E.3, Table E.3



Figure 6: From left to right: SINGLE- \mathcal{M} , SINGLE, MULT-DIRECT and MULT. While single-chart flows fail to cover the circle, multi-chart flows are able to provide samples across the circle.



Figure 7: From left to right: SINGLE- \mathcal{M} , SINGLE, MULT-DIRECT and MULT. Using a single chart results in pathological failures, where the model fails to cover the entire circle. Only when using multiple charts the model learns the correct topology.



Figure 8: Evaluation metrics of different methods. Generally, MULT-DIRECT and MULT yield roughly the same level of performances while clearly outperform SINGLE-M and SINGLE.



Figure 9: An example of geodesic solver failure on torus. The solver finds the path on the opposite side, resulting in a drastically different logarithmic map and a longer distance.

Table 4: Evaluation metrics concerning modeling of SINGLE- \mathcal{M} on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	W
Sphere-U	$8.61\!\cdot\!10^{-5}\pm4.9\!\cdot\!10^{-5}$	$1.29{\cdot}10^{-1}\pm1.24{\cdot}10^{-2}$
Torus-U	$3.52 \cdot 10^{-2} \pm 4.63 \cdot 10^{-2}$	$3.11 \cdot 10^{-1} \pm 5.2 \cdot 10^{-2}$
Sphere-M	$1.73 \cdot 10^{-4} \pm 1.93 \cdot 10^{-4}$	$1.52 \cdot 10^{-1} \pm 1.56 \cdot 10^{-2}$
Torus-M	$2.58 \cdot 10^{-2} \pm 1.49 \cdot 10^{-2}$	$1.42 \cdot 10^0 \pm 2.32 \cdot 10^0$

Table 5: Evaluation metrics concerning modeling of SINGLE on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	W
Sphere-U	$6.81 \cdot 10^{-5} \pm 4.66 \cdot 10^{-5}$	$1.35 \cdot 10^{-1} \pm 9.43 \cdot 10^{-3}$
Torus-U	$4.19 \cdot 10^{-3} \pm 2.25 \cdot 10^{-3}$	$4.75 \cdot 10^{-1} \pm 2.28 \cdot 10^{-1}$
Sphere-M	$9.43 \cdot 10^{-5} \pm 6.66 \cdot 10^{-5}$	$1.7 \cdot 10^1 \pm 3.1 \cdot 10^1$
Torus-M	$3.67 \cdot 10^{-2} \pm 2.72 \cdot 10^{-2}$	$6.34{\cdot}10^{-1}\pm2.98{\cdot}10^{-1}$

Table 6: Evaluation metrics concerning modeling of MULT on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	\mathcal{W}
Sphere-U	$2.22{\cdot}10^{-6}\pm7.44{\cdot}10^{-7}$	$1.29{\cdot}10^{-1}\pm1.29{\cdot}10^{-2}$
Torus-U	$2.25 \cdot 10^{-5} \pm 1.99 \cdot 10^{-6}$	$2.46 \cdot 10^{-1} \pm 2.55 \cdot 10^{-2}$
Sphere-M	$1.58 \cdot 10^{-6} \pm 2.42 \cdot 10^{-7}$	$1.34 \cdot 10^{-1} \pm 1.66 \cdot 10^{-2}$
Torus-M	$6.9{\cdot}10^{-5} \pm 2.96{\cdot}10^{-5}$	$3.15\!\cdot\!10^{-1}\pm5.93\!\cdot\!10^{-2}$

Table 7: Evaluation metrics concerning geometry of SINGLE- \mathcal{M} on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Exps	Logs	Dists
Sphere-U	$1.86{\cdot}10^3\pm2.62{\cdot}10^3$	$8.58\!\cdot\!10^{-1}\pm3.45\!\cdot\!10^{-1}$	$2.35{\cdot}10^{-2}\pm7.66{\cdot}10^{-3}$
Torus-U	$2.01 \cdot 10^{6} \pm 2.01 \cdot 10^{6}$	$2.79 \cdot 10^1 \pm 6.68 \cdot 10^0$	$7.98 \cdot 10^0 \pm 3.09 \cdot 10^0$
Sphere-M	$2.1 \cdot 10^{-1} \pm 7.73 \cdot 10^{-2}$	$1.19 \cdot 10^0 \pm 2.27 \cdot 10^{-1}$	$4.91 \cdot 10^{-2} \pm 1.61 \cdot 10^{-2}$
Torus-M	$1.23 \cdot 10^5 \pm 1.68 \cdot 10^5$	$5.3{\cdot}10^2\pm6.75{\cdot}10^2$	$5.0{\cdot}10^2\pm 6.68{\cdot}10^2$

Table 8: Evaluation metrics concerning geometry of SINGLE on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Exps	Logs	Dists
Sphere-U	$1.91 \!\cdot\! 10^{24} \pm 1.91 \!\cdot\! 10^{24}$	$1.09{\cdot}10^0\pm2.0{\cdot}10^{-1}$	$5.22{\cdot}10^{-2}\pm6.04{\cdot}10^{-3}$
Torus-U	$5.63 \cdot 10^0 \pm 6.62 \cdot 10^0$	$1.9 \cdot 10^1 \pm 3.85 \cdot 10^0$	$4.06 \cdot 10^0 \pm 2.31 \cdot 10^0$
Sphere-M	$1.15 \cdot 10^2 \pm 1.0 \cdot 10^2$	$1.59 \cdot 10^0 \pm 1.36 \cdot 10^0$	$1.51 \cdot 10^{-1} \pm 1.93 \cdot 10^{-1}$
Torus-M	$1.98\!\cdot\!10^{11}\pm2.8\!\cdot\!10^{11}$	$1.49 \cdot 10^3 \pm 2.03 \cdot 10^3$	$1.45 \cdot 10^3 \pm 2.01 \cdot 10^3$

Table 9: Evaluation metrics concerning geometry of MULT on sphere and torus in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Exps	Logs	Dists
Sphere-U	$2.28{\cdot}10^{-2}\pm1.52{\cdot}10^{-2}$	$1.32 \cdot 10^{-1} \pm 2.86 \cdot 10^{-2}$	$3.77{\cdot}10^{-4}\pm1.27{\cdot}10^{-4}$
Torus-U	$1.17 \cdot 10^{-1} \pm 3.99 \cdot 10^{-2}$	$6.41 \cdot 10^0 \pm 1.02 \cdot 10^0$	$4.34 \!\cdot\! 10^{-1} \pm 1.38 \!\cdot\! 10^{-1}$
Sphere-M	$9.19 \cdot 10^{-3} \pm 5.65 \cdot 10^{-3}$	$1.67 \cdot 10^{-1} \pm 5.72 \cdot 10^{-2}$	$8.2 \cdot 10^{-4} \pm 4.1 \cdot 10^{-4}$
Torus-M	$1.5 \cdot 10^0 \pm 3.82 \cdot 10^{-1}$	$2.32 \cdot 10^0 \pm 6.07 \cdot 10^{-1}$	$5.13 \cdot 10^{-2} \pm 1.71 \cdot 10^{-2}$



Figure 10: From top to bottom: persistence diagrams induced by different models on Sphere-M, Torus-U and Torus-M.

Table 10: Evaluation metrics of SINGLE- \mathcal{M} on triangular meshes in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	\mathcal{W}	Dists
7-8 cube Busted Rock	$\begin{array}{c} 1.67 \cdot 10^{-4} \pm 5.68 \cdot 10^{-5} \\ 5.14 \cdot 10^{-4} \pm 2.18 \cdot 10^{-4} \\ 3.43 \cdot 10^{-4} \pm 5.39 \cdot 10^{-5} \end{array}$	$\begin{array}{c} 1.98 \cdot 10^{-1} \pm 2.28 \cdot 10^{-2} \\ 2.0 \cdot 10^{-1} \pm 1.76 \cdot 10^{-2} \\ 1.74 \cdot 10^{-1} \pm 1.49 \cdot 10^{-2} \end{array}$	$\begin{array}{c} 2.45 \cdot 10^{-1} \pm 2.06 \cdot 10^{-1} \\ 2.17 \cdot 10^{-1} \pm 9.59 \cdot 10^{-2} \\ 3.93 \cdot 10^{-1} \pm 6.47 \cdot 10^{-2} \end{array}$

E.4 Results on triangular meshes

We report the numerical results of SINGLE- \mathcal{M} , SINGLE and MULT on triangular meshes in Table E.4, Table E.4 and Table E.4. We additionally report example persistence diagrams obtained under these settings in Figure 11. Due to the lack of a principled way to obtain maximally uniformly distributed samples from the meshes and the relatively small number of data points, the diagrams may not that clearly reflect the underlying topology.

E.5 Results on Mocap data

The numerical results of SINGLE- \mathcal{M} , SINGLE and MULT on Mocap data are reported in Table E.5, Table E.5 and Table E.5, respectively.

Table	11: Evaluation metrics of SING	LE on triangular mes	shes in the form	of mean \pm std,	lower 1s
better.	The ones with failed runs are ita	lic.			

Manifold	Recons	\mathcal{W}	Dists
7-8 cube Busted	$8.37 \cdot 10^{-4} \pm 5.76 \cdot 10^{-4} \\ 6.06 \cdot 10^{-4} \pm 3.92 \cdot 10^{-4} \\ 1.7 \cdot 10^{-4} \pm 4.41 \cdot 10^{-5}$	$2.23 \cdot 10^{-1} \pm 4.53 \cdot 10^{-2}$ $2.15 \cdot 10^{-1} \pm 3.26 \cdot 10^{-2}$ $1.01 \cdot 10^{-1} \pm 1.07 \cdot 10^{-2}$	$\frac{1.71 \cdot 10^{-1} \pm 6.49 \cdot 10^{-2}}{2.64 \cdot 10^{-1} \pm 9.69 \cdot 10^{-2}}$
Rock	$1.7 \cdot 10^{-4} \pm 4.41 \cdot 10^{-6}$	$1.91 \cdot 10^{-1} \pm 1.97 \cdot 10^{-2}$	$2.4 \cdot 10^{-1} \pm 2.28 \cdot 10^{-1}$

Table 12: Evaluation metrics of MULT on triangular meshes in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	\mathcal{W}	Dists
7-8 cube Busted Rock	$\begin{array}{c} 1.44 \cdot 10^{-5} \pm 1.75 \cdot 10^{-6} \\ 5.65 \cdot 10^{-5} \pm 7.79 \cdot 10^{-6} \\ 1.24 \cdot 10^{-5} \pm 2.77 \cdot 10^{-6} \end{array}$	$\begin{array}{c} 1.97 \cdot 10^{-1} \pm 2.36 \cdot 10^{-2} \\ 2.18 \cdot 10^{-1} \pm 3.39 \cdot 10^{-2} \\ 1.79 \cdot 10^{-1} \pm 2.89 \cdot 10^{-2} \end{array}$	$\begin{array}{c} 9.22 \cdot 10^{-2} \pm 1.46 \cdot 10^{-3} \\ 4.9 \cdot 10^{-2} \pm 5.97 \cdot 10^{-3} \\ 6.21 \cdot 10^{-2} \pm 8.56 \cdot 10^{-3} \end{array}$



Figure 11: From top to bottom: persistence diagrams induced by different models on 7-8th cube, Busted and Rock.

Table 13: Evaluation metrics of SINGLE- \mathcal{M} on Mocap in the form of mean \pm std, lower is better. The ones with failed runs are italic.

Manifold	Recons	\mathcal{W}
Mocap1 Mocap3	$\begin{array}{c} 4.06 \cdot 10^{-2} \pm 4.95 \cdot 10^{-2} \\ 2.24 \cdot 10^{0} \pm 8.98 \cdot 10^{-2} \end{array}$	$\begin{array}{c} 7.63 \cdot 10^{-1} \pm 2.25 \cdot 10^{-1} \\ 2.7 \cdot 10^{0} \pm 5.01 \cdot 10^{-2} \end{array}$

Table 14: Evaluation metrics of SINGLE on Mocap in the form of mean \pm std, lower is better. The ones with failed <u>runs are italic</u>.

Manifold	Recons	\mathcal{W}
Mocap1 Mocap3	$\begin{array}{c} 9.48 \cdot 10^{-3} \pm 5.63 \cdot 10^{-3} \\ 2.26 \cdot 10^{0} \pm 2.8 \cdot 10^{-1} \end{array}$	$\begin{array}{c} 7.13 \cdot 10^{-1} \pm 2.08 \cdot 10^{-1} \\ 2.95 \cdot 10^{0} \pm 1.51 \cdot 10^{-1} \end{array}$

Table 15: Evaluation metrics of MULT on Mocap in the form of mean \pm std, lower is better. The ones with failed <u>runs are italic</u>.

Manifold	Recons	\mathcal{W}
Mocap1 Mocap3	$\begin{array}{c} 6.99 \cdot 10^{-4} \pm 1.3 \cdot 10^{-5} \\ 3.98 \cdot 10^{-3} \pm 1.96 \cdot 10^{-4} \end{array}$	$\begin{array}{c} 7.16 \cdot 10^{-1} \pm 1.8 \cdot 10^{-1} \\ 2.48 \cdot 10^{0} \pm 8.05 \cdot 10^{-2} \end{array}$