# Explorations in Homeomorphic Variational Auto-Encoding

**Luca Falorsi** [* 1]   **Pim de Haan** [* 1]   **Tim R. Davidson** [* 1]   **Nicola De Cao** [1]   **Maurice Weiler** [1]
**Patrick Forré** [1]   **Taco S. Cohen** [1 2]

## Abstract

The manifold hypothesis states that many kinds of high-dimensional data are concentrated near a low-dimensional manifold. If the topology of this data manifold is non-trivial, a continuous encoder network cannot embed it in a one-to-one manner without creating holes of low density in the latent space. This is at odds with the Gaussian prior assumption typically made in Variational Auto-Encoders (VAEs), because the density of a Gaussian concentrates near a blob-like manifold.

In this paper we investigate the use of manifold-valued latent variables. Specifically, we focus on the important case of continuously differentiable symmetry groups (Lie groups), such as the group of 3D rotations $SO(3)$. We show how a VAE with $SO(3)$-valued latent variables can be constructed, by extending the reparameterization trick to compact connected Lie groups. Our experiments show that choosing manifold-valued latent variables that match the topology of the latent data manifold, is crucial to preserve the topological structure and learn a well-behaved latent space.

## 1. Introduction

Many complex probability distributions can be represented more compactly by introducing latent variables. Intuitively, the idea is that there is some simple underlying latent structure, which is mapped to the observation space by a potentially complex nonlinear function. It will come as no surprise then, that most research effort has aimed at using maximally simple priors for the latent variables (e.g. Gaussians), combined with flexible likelihood functions (e.g. based on neural networks).

However, it is not hard to see (Fig. 1.1) that if the data is con-
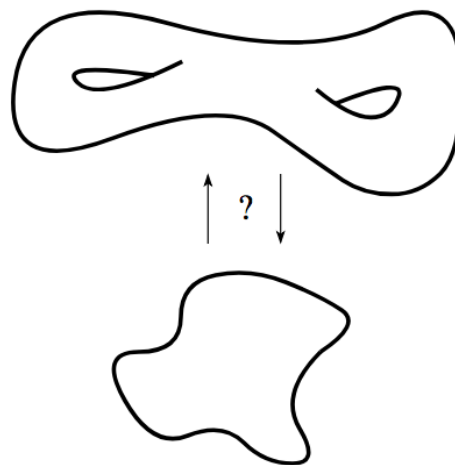


*Figure 1.1.* An example of problems that arise in mapping manifolds not diffeomorphic to each other. Notice that in the illustrated example the 'holes' in the first manifold, prevent a smooth mapping to the second.

centrated near a low-dimensional manifold with non-trivial topology, there is no continuous and invertible mapping to a blob-like manifold (the region where prior mass is concentrated). We believe that for purposes of representation learning, the embedding map (encoder) should be homeomorphic (i.e. continuous and invertible, with continuous inverse), which means that although dimensionality reduction and geometrical simplification (flattening) may be possible, the topological structure should be preserved.

Once could encode such a manifold in a higher dimensional flat space with a regular variational auto-encoder (VAE, Kingma & Welling (2013); Rezende et al. (2014)), rather than learning a homeomorphism. This has two disadvantages. The prior on the flat space will put density outside of the embedding and traversals along the extra dimensions that are normal to the manifold will either leave the decoding invariant, or move out of the data manifold. This is because at each point there will be many more degrees of freedom than the dimensionality of the manifold.

In this paper we investigate this idea for the special case of

---

[*]Equal contribution   [1]University of Amsterdam   [2]Qualcomm AI Research. Correspondence to: Luca Falorsi <luca.falorsi@gmail.com>.

Lie groups, which are symmetry groups that are simultaneously differentiable manifolds. Lie groups include rotations, translations, scaling, and other geometric transformations, which play an important role in many application domains such as robotics and computer vision. More specifically, we show how to construct[1] a VAE with latent variables that live on a Lie group, which is done by generalizing the reparameterization trick.

We will describe an approach for reparameterizing densities on $SO(3)$, the group of 3D rotations, which can be extended to general compact and connected Lie group VAEs in a straightforward manner. The primary technical difficulty in the construction of this theory is to show that the push-forward measure induced by our reparameterization has a density that is absolutely continuous w.r.t. the Haar measure. Moreover, we show how to construct the encoder such that it can learn a homeomorphic map from the data manifold to the mean parameter of the posterior. Finally, we propose a decoder that uses the group action to further encourage the latent space to respect the group structure.

We perform experiments on two types of synthetic data: $SO(3)$ embedded into a high dimensional space through its group representation, and images of 3D rotations of a single colored cube. We find that a theoretically sound architecture is capable of continuously mapping the data manifold to the latent space. On the other hand, models that do not respect topological structure, and in particular those with a standard Gaussian latent space, show discontinuities when trajectories in the latent space are visualized. To better study this phenomenon, we introduce a way to measure the continuity of the embedding based on the concept of Lipschitz continuity. We empirically demonstrate that only a manifold-valued latent variable with the required topological structure is capable of fully solving the difficult task of the more complicated experiment.

Our main contributions in this work are threefold:

1. A reparameterization trick for distributions on the $SO(3)$ group of rotations in three dimensions.

2. An encoder for the mean parameter that learns a homeomorphism between the $SO(3)$ manifold embedded in the data and $SO(3)$ itself.

3. A decoder that uses the group action to respect the group structure.

## 2. Preliminary Concepts

In this section we will first cover a number of preliminary concepts that will be used in the rest of the paper.

### 2.1. Variational Auto-Encoders

The VAE is a latent variable model, in which $\mathbf{x}$ denotes a set of observed variables, $\mathbf{z}$ stochastic latent variables, and $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ a parameterized model of the joint distribution called the *generative model*. Given a dataset $\mathbf{X} = \{\mathbf{x}_1, \cdots, \mathbf{x}_N\}$, we typically wish to maximize the average marginal log-likelihood $\frac{1}{N} \log p(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} \log \int p(\mathbf{x}_i, \mathbf{z}_i) d\mathbf{z}$, w.r.t. the parameters. However when the model is parameterized by neural networks, the marginalization of this expression is generally intractable. One solution to overcome this issue is applying variational inference in order to maximize the *Evidence Lower Bound* (ELBO) for each observation:

$$
\begin{aligned}
\log p(\mathbf{x}) &= \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\
&\geq \mathbb{E}_{q(\mathbf{z})}[\log p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z})||p(\mathbf{z})), \quad (1)
\end{aligned}
$$

where the approximate posterior $q(\mathbf{z})$ belongs to the variational family $\mathcal{Q}$. To make inference scalable an *inference network* $q(\mathbf{z}|\mathbf{x})$ is introduced that outputs a probability distribution for each data point $\mathbf{x}$, leading to the final objective

$$
\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})), \quad (2)
$$

with $\theta$ representing the parameters of $p$ and $q$. The ELBO can be efficiently approximated for continuous latent variable $\mathbf{z}$ by Monte Carlo estimates using the *reparameterization trick* of $q(\mathbf{z}|\mathbf{x})$ (Kingma & Welling, 2013; Rezende et al., 2014).

### 2.2. Lie Groups and Lie Algebras

**Lie Group** A group is a set equipped with a product that follows the four group axioms: the product is closed and associative, there exists an identity element, and every group element has an inverse. This is closely linked to symmetry transformations that leave some property invariant. For example, composing two symmetry transformations should still maintain the invariance. A *Lie group G* has additional structure, as its set is also a smooth manifold. This means that we can, at least in local regions, describe group elements continuously with parameters. The number of parameters equals the dimension of the group. We can see (connected) Lie groups as continuous symmetries where we can continuously traverse between group elements[2].

**Lie Algebra** The *Lie algebra* $\mathfrak{g}$, of a $N$ dimensional Lie group is its tangent space at the identity, which is a vector space of $N$ dimensions. We can see the algebra elements as infinitesimal generators, from which all other elements in the group can be created. For matrix Lie groups we can represent vectors $v$ in the tangent space as matrices $\mathbf{v}_\times$.

---

[1]Our implementation is available at `https://github.com/pimdh/lie-vae`.

[2]We refer the interested reader to (Hall, 2003).

**Exponential Map**    The structure of the algebra creates a map from an element of the algebra to a vector field on the group manifold. This gives rise to the *exponential map* $\exp : \mathfrak{g} \to G$ which maps an algebra element to the group element at unit length from the identity along the flow of the vector field. The zero vector is thus mapped to the identity. For compact connected Lie groups, such as $\mathrm{SO}(3)$, the exponential map is surjective.

### 2.3. The group $\mathrm{SO}(3)$

The special orthogonal Lie group of three dimensional rotations $\mathrm{SO}(3)$ is defined as:

$$\mathrm{SO}(3) := \{R \in GL(\mathbb{R}^3) : R^\top R = I \wedge det(R) = 1\} \quad (3)$$

where $GL$ is the *general linear group*, which is the set of square invertible matrices under the operation of matrix multiplication. Note that $\mathrm{SO}(3)$ is not homeomorphic to $\mathbb{R}^N$, since on $\mathbb{R}^N$ every continuous path can be continuously contracted to a point, while this is not true for $\mathrm{SO}(3)$. Consider for example a full rotation around a fixed axis.

The elements of Lie algebra $\mathfrak{so}(3)$ of group $\mathrm{SO}(3)$, are represented by the 3 dimensional vector space of the skew-symmetric $3 \times 3$ matrices. We choose a basis for the algebra:

$$L_{1,2,3} := \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (4)$$

This provides a vector space isomorphism between $\mathbb{R}^3$ and $\mathfrak{so}(3)$, written as $[\,\cdot\,]_\times : \mathbb{R}^3 \to \mathfrak{so}(3)$.

Assuming the decomposition $\mathbf{v}_\times = \theta \mathbf{u}_\times$, s.t. $\theta \in \mathbb{R}_{\geq 0}$, $\|\mathbf{u}\| = 1$, the exponential map is given by the Rodrigues rotation formula (Rodrigues, 1840):

$$\exp(\mathbf{v}_\times) = \mathbf{I} + \sin(\theta)\mathbf{u}_\times + (1 - \cos(\theta))\mathbf{u}_\times^2 \quad (5)$$

Since $\mathrm{SO}(3)$ is a compact and connected Lie group this map is surjective, however it is not injective.

## 3. Reparameterizing $\mathrm{SO}(3)$

In this section we will explain our $\mathrm{SO}(3)$ reparameterization trick by analogy to the classic version described in (Kingma & Welling, 2013; Rezende et al., 2014). An overview of the different steps and their relation to the classical case are given in Figure 2.1.

We sample from a scale-reparameterizable distribution $r(\mathbf{v}|\sigma)$ on $\mathbb{R}^3$ that is concentrated around the origin. Due to the isomorphism between $\mathbb{R}^3$ and $\mathfrak{so}(3)$ this can be identified with a sample $\mathbf{v}_\times$ from the Lie algebra $\mathfrak{so}(3)$. Next we apply the exponential map to obtain a sample $R = \exp(\mathbf{v}_\times) \sim \hat{q}(R|\sigma)$ of the group as visualized in

Figure 2.1 (a) to (b). Since the distribution $r(\mathbf{v}|\sigma)$ is concentrated around the origin, the distribution of $\hat{q}(R|\sigma)$ will be concentrated around the group identity. In order to change the location of the distribution $\hat{q}$, we left multiply $R$ by another element $R_\mu$, see Figure 2.1 (b) to (c).

To see the connection with the classical case, identify $\mathbb{R}^N$ under addition as a Lie group, with the Lie algebra isomorphic to $\mathbb{R}^N$. As the group and the algebra are in this case isomorphic, the step of taking the exponential map can be taken as the identity operation such that $r = \hat{q}$. The multiplication with a group element to change the location corresponds to a translation by $\mu$.

One critical complication is that it is not obvious that the measure we defined above through the exp map has a density function $p : \mathrm{SO}(3) \to \mathbb{R}_+$. For this to be the case we need the constructed measure to be absolutely continuous with respect to the Haar measure $\nu$, the natural measure on the Lie group. This is proven by the following theorem.

**Theorem 1.** *Let $(\mathbb{R}^3, \lambda, \mathcal{B}[\mathbb{R}^3])$ the real space, provided with the Lebesgue measure on the Borel algebra on $\mathbb{R}^3$. Let $(\mathrm{SO}(3), \nu, \mathcal{B}[\mathrm{SO}(3)])$ the group of 3 dimensional rotations, provided with the normalized Haar measure $\nu$ on the Borel $\sigma$-algebra on $\mathrm{SO}(3)$. Consider then the probability measure $\mu : \mathcal{B}[\mathbb{R}^3] \to [0, 1]$ absolutely continuous w.r.t $\lambda$, with density $r$. Consider the exponential map $\exp : \mathbb{R}^3 \to \mathrm{SO}(3)$ that is differentiable, thus continuous, thus measurable. Let then $\exp_*(\mu)$ be the pushforward of $\mu$ by the $\exp$ function. then $\exp_*(\mu)$ is absolutely continuous with respect of the Haar measure $\nu$ ($\exp_*(\mu) \ll \nu$).*

*Proof.* See Appendix B  □

As further derived in Appendix B this implies the pushforward measure on $\mathrm{SO}(3)$ to be absolutely continuous w.r.t. to the Haar measure where the density

$$\hat{q}(R|\sigma) = \sum_{k \in \mathbb{Z}} r\left(\frac{\log(R)}{\theta(R)}(\theta(R) + 2k\pi)\Big|\sigma\right) \frac{(\theta(R) + 2k\pi)^2}{3 - \mathrm{tr}(R)},$$
$$(6)$$

is defined almost everywhere. Here $R \in \mathrm{SO}(3)$ and

$$\theta(R) = \|\log(R)\| = \cos^{-1}\left(\frac{\mathrm{tr}(R) - 1}{2}\right) \quad (7)$$

Further, $\log(\cdot)$ is defined as a principal branch and maps back the group element to the unique Lie algebra element next to the origin. Notice that even if the density is singular at $R = I$, it still integrates to 1. After rotating $R$ by left multiplying with another $\mathrm{SO}(3)$ element $R_\mu$, we obtain the final sample:

$$R_z \sim q(R_z|R_\mu, \sigma) = \hat{q}(R_\mu^\top R_z|\sigma), \quad (8)$$

where the second step is valid because of the left invariance of the Haar measure.
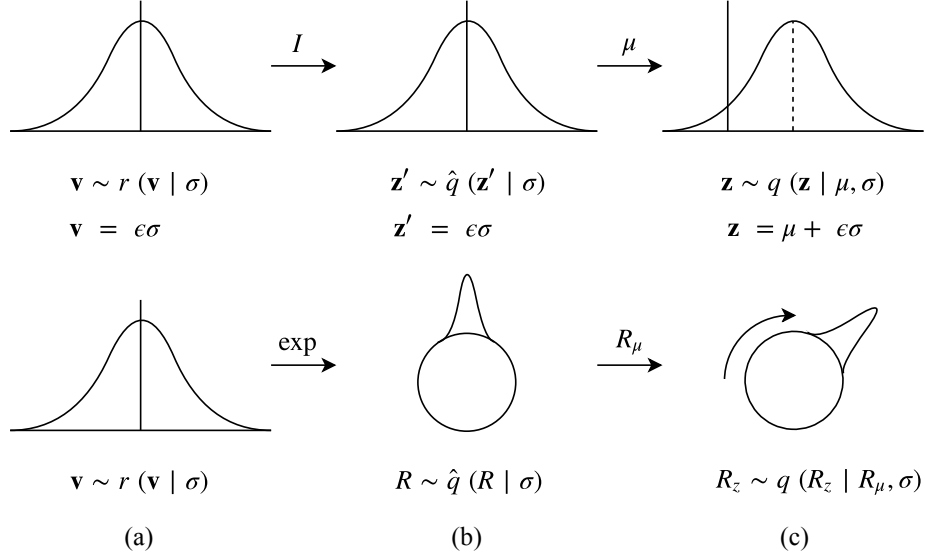
$$\mathbf{v} \sim r\left(\mathbf{v} \mid \sigma\right) \qquad \mathbf{z}' \sim \hat{q}\left(\mathbf{z}' \mid \sigma\right) \qquad \mathbf{z} \sim q\left(\mathbf{z} \mid \mu, \sigma\right)$$
$$\mathbf{v} = \epsilon\sigma \qquad \mathbf{z}' = \epsilon\sigma \qquad \mathbf{z} = \mu + \epsilon\sigma$$

$$\mathbf{v} \sim r\left(\mathbf{v} \mid \sigma\right) \qquad R \sim \hat{q}\left(R \mid \sigma\right) \qquad R_z \sim q\left(R_z \mid R_\mu, \sigma\right)$$

(a)  (b)  (c)

*Figure 2.1.* Illustration of our extended reparameterization trick in comparison to the classic reparameterization trick.

**Kullback-Leibler Divergence** The KL divergence, or *relative entropy* can be decomposed into the *entropy* and the *cross-entropy* terms, $KL(q\|p) = \mathbb{H}(q, p) - \mathbb{H}(q)$. Since the Haar measure is invariant to left multiplication, we can compute the entropy of the distribution $\hat{q}$ instead of $q$. As we have the expression of the density, the entropy can be computed using Monte Carlo samples:

$$\mathbb{H}(q) = \mathbb{H}(\hat{q}) \approx -\frac{1}{N}\sum_{i=1}^{N}\log \hat{q}(R_i|\sigma), \quad R_i \sim \hat{q}(R_i|\sigma)$$

$$= -\frac{1}{N}\sum_{i=1}^{N}\log \hat{q}(\exp(\mathbf{v}_i)|\sigma)$$

$$= -\frac{1}{N}\sum_{i=1}^{N}\log \sum_{k\in\mathbb{Z}} r\Big(\frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}(\|\mathbf{v}_i\| + 2k\pi)|\sigma\Big)\cdot$$
$$\frac{(\|\mathbf{v}_i\| + 2k\pi)^2}{2 - 2\cos(\|\mathbf{v}_i\|)}, \quad \mathbf{v}_i \sim r(\mathbf{v}_i|\sigma) \quad (9)$$

Notice that the last expression only depends on the samples taken on the Lie algebra. We found that one sample suffices when mini-batches are used. In general the cross-entropy term can be similarly approximated by MC estimates. However, in the special but important case of a uniform prior, $p$, the cross-entropy reduces to: $\mathbb{H}(p, q) = -\log\left(\frac{1}{8\pi^2}\right)$.

## 4. Encoder and Decoder networks

Having defined the reparameterizable density $q(R_z|R_\mu, \sigma)$, we need to design encoder networks which map elements from the input space $\mathcal{X}$ to the reparameterization parameters $R_\mu$, $\sigma$ and decoder networks which map group elements to the output prediction.

### 4.1. Homeomorphic Encoder

We split the encoder network in two parts $\text{enc}^\mu$ and $\text{enc}^\sigma$, which predict reparameterization parameters $R_\mu$ and $\sigma$ respectively. Since $\sigma$ are parameters of a distribution in $\mathbb{R}^3$, the corresponding network $\text{enc}^\sigma$ does not pose any problems and can be chosen similarly as in classical VAEs. However, special attention needs to be paid to designing $\text{enc}^\mu$ which predicts a group element $R_\mu \in \text{SO}(3)$.

We consider the data as lying in a lower dimensional manifold $\mathcal{M}$, embedded in the input space $\mathcal{X}$. In our particular problem the manifold $\mathcal{M}$ is assumed to be generated by $\text{SO}(3)$, acting on a canonical object and a subsequent projection into ambient space (e.g. pixel space) $\mathcal{X}$ which, for simplicity we assume to be injective. This means that we can make the simplifying assumption that $\text{SO}(3)$ can be recovered from its image in $\mathcal{X}$, i.e. that the map $\text{SO}(3) \to \mathcal{M}$ is a homeomorphism. The encoder is now meant to learn the inverse map, i.e. to learn a map from $\mathcal{X}$ to $\text{SO}(3)$, which when restricted to $\mathcal{M}$ is a homeomorphism and thus preserves the topological structure of $\text{SO}(3)$.

In general there is no standard way to define and parameterize the class of functions which are guaranteed to have these properties by design via a neural network. Instead we will give a general way to build $\text{enc}^\mu$ *capable* of learning such a mapping. We divide the encoder network in two functions: $\text{enc}^\mu = \pi \circ f$, where $f : \mathcal{X} \to \mathcal{Y}$, for some space $\mathcal{Y}$, is parametrized by a neural network, and $\pi : \mathcal{Y} \to \text{SO}(3)$ is a fixed surjective function. Not any space $\mathcal{Y}$ or function $\pi$ is suited: since neural networks can only model continuous functions, a necessary condition on $\mathcal{Y}$ for $\text{enc}^\mu$ to be able to learn to be a homeomorphism (when its domain

$\mathcal{X}$ is restricted to $\mathcal{M}$), is that there exists an embedding $i : \mathrm{SO}(3) \to \mathcal{Y}$. Then by definition a function $\pi$ exist, such that $\pi|_{i(\mathrm{SO}(3))} = i^{-1}|_{i(\mathrm{SO}(3))}$ is a homeomorphism. Any extension of $\pi|_{i(\mathrm{SO}(3))}$ to $\mathcal{Y}$ is a suitable candidate. Moreover, if we choose $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathbb{R}^m$ for some $n, m$, then some continuous $f : \mathcal{X} \to \mathcal{Y}$ exists (which we can approximate with neural networks) such that an appropriate $\mathrm{enc}^\mu = \pi \circ f$ exist. Several choices for $\mathcal{Y}$ and $\pi$ are proposed in Appendix D and investigated in the experimental section 6.

### 4.2. Group Action Decoder

Our decoder $p(\mathbf{x}|R_z)$ must be capable to map a group element and optionally additional latent structure back to the original high dimensional input space. When the factor of variation in the input is the pose of an object, and we learn a latent variable $R_z \in \mathrm{SO}(3)$, we desire that a transformation $R \in \mathrm{SO}(3)$ applied to a latent object representation $\mathbf{z}$ results in a corresponding transformation of the pose of the decoded object. The task of the decoder is thus to learn a three dimensional representation of the object, to rotate it according to the latent variable and finally to project it back to the two dimensional frame of the input image. A naive approach could be to simply provide the 9 elements of the rotation matrix to a neural network. However, although it may learn to reconstruct the input well, it provides no guarantee that the latent space accurately reflects the pose variations of the object. Therefore, we like to make the method more explicit.

Hypothetically, one could learn a vector-valued signal on the sphere, $f : S^2 \to \mathbb{R}^N$, to represent the three dimensional object in the input, as 3D shapes can be well represented by its information projected to a sphere (Cohen et al., 2018a). The decoder can rotate this signal with the latent variable, before projecting it back to pixel space. A major downside of this approach is that parameterizing and projecting a function on the sphere is highly non-trivial.

Alternatively, we propose a method based on the representation theory of groups (Hall, 2003). Rather than learning the function $f$, we directly learn its (band-limited) Fourier modes, which form a simple vector space. It can be shown (Chirikjian & Kyatkin, 2000) that rotations of a signal on the sphere correspond to a linear transformation of the Fourier modes. The transformed Fourier modes are subsequently fed through an image generative network, and the linear transformation is the Wigner-D-matrix, which is a function of the $\mathrm{SO}(3)$ element. Technically, the Wigner-D-matrices form representations of the group. This means that as mapping to the linear transformations is a homomorphism, it preserves the group structure: $D(g)D(g') = D(gg')$, for $g, g' \in \mathrm{SO}(3)$ and $D(g)$ a Wigner-D-matrix. This method encourages the latent space to represent the actual pose of
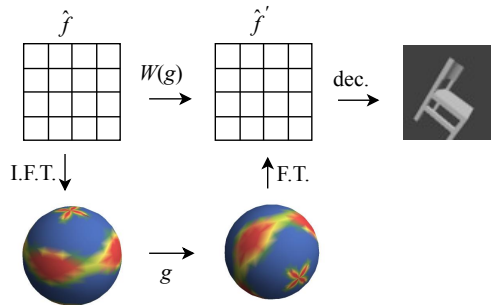


*Figure 4.1.* The encoder infers the $R \in \mathrm{SO}(3)$ and Fourier modes $\hat{f}$ if working with multiple objects, otherwise $\hat{f}$ is a parameter. Shown is the commutative diagram between taking the Inverse Fourier Transform, rotating the result and taking the Fourier Transform, and acting with the group representation $W$ on $\hat{f}$. The decoder maps the transformed $\hat{f}'$ to pixels.

the input, while only requiring the construction of the $W$ matrices and performing a linear transformation. We refer to Figure 4.1 and Appendix F for details.

An additional advantage is that this decoder allows for disentangling of content and pose, as it is forced to encode the pose in a meaningful way. The Fourier modes are in that case also generated by the encoder. We leave this for future work.

## 5. Related Work

As VAEs utilize VI to recover some distribution on a latent manifold responsible for generating the observed data, the majority of extensions is focused on increasing the flexibility of the prior and approximate posterior. Although the majority of approaches make use of a normal Gaussian prior, recently there has been a surge to provide additional options to offset some of this distribution's perceived limitations. Tomczak & Welling (2017) propose to directly tie the prior to the approximate posterior and learn it as a mixture over approximate posteriors. Nalisnick & Smyth (2017) introduce a non-parametric prior applying a truncated stick-breaking method. Research to support discrete latent variables was done in Jang et al. (2017); Maddison et al. (2017), while in Naesseth et al. (2017); Figurnov et al. (2018) recently novel techniques were introduced to reparameterize a suite of continuous distributions. In (Davidson et al., 2018), the reparameterization technique of Naesseth et al. (2017) is extended to explore the properties of the hyperspherical von Mises-Fisher distribution to better capture intrinsically hyperspherical data. This is done in the context of avoiding *manifold mismatches*, and as such is closely related to the motivation of this work.
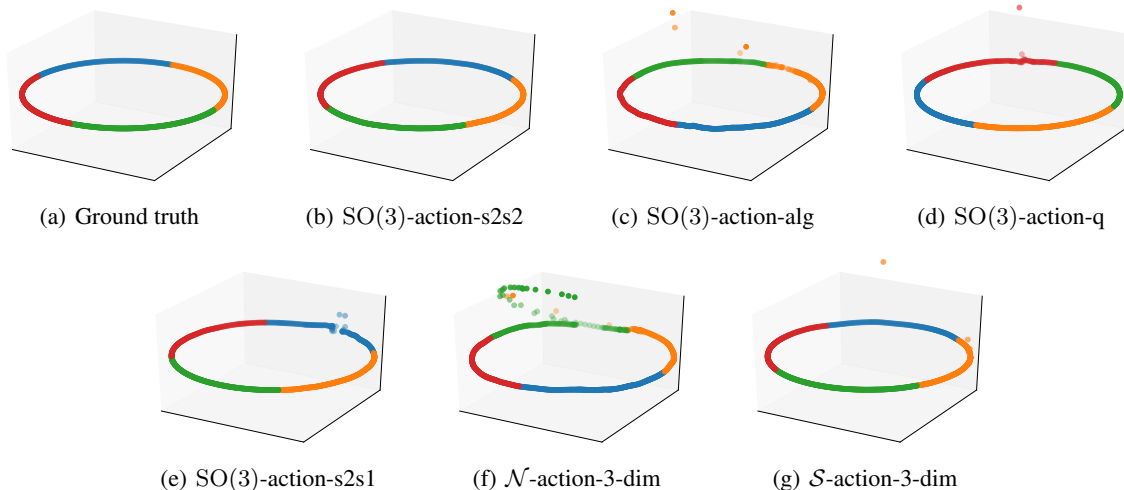
(a) Ground truth    (b) SO(3)-action-s2s2    (c) SO(3)-action-alg    (d) SO(3)-action-q

(e) SO(3)-action-s2s1    (f) $\mathcal{N}$-action-3-dim    (g) $\mathcal{S}$-action-3-dim

*Figure 4.2.* The latent encoding of a $S^1$ trajectory in the Toy data set. The $SO(3)$ elements are mapped to $\mathbb{R}^9$ by taking the rotation matrix elements and are subsequently mapped to 3D by Principal Component Analysis.

The predominant procedure to generate a more complex approximate posterior is through *normalizing flows* (Rezende & Mohamed, 2015), in which a class of invertible transformations is applied sequentially to a reparameterizable density. This general idea has later been extended Kingma et al. (2016); Berg et al. (2018), to improve flexibility even further. As this framework does not hold any specific distributional requirements on the prior besides being reparameterizable, it would be interesting to investigate possible applications to SO(3) in future work.

The problem of defining distributions on homogeneous spaces, including Lie groups, was investigated in (Chirikjian & Kyatkin, 2000; Chirikjian, 2010; Chirikjian & Kyatkin, 2016; Chirikjian, 2012). Cohen & Welling (2015) devised harmonic exponential families which are a powerful family of distributions defined on homogeneous spaces. These works did not concentrate on making the distributions reparameterizable.

Rendering complex scenes from multiple poses has been explored in (Eslami et al., 2018). However, this work assumes access to ground truth poses and does not do unsupervised pose learning as in the presented framework.

The idea of incorporating prior knowledge on mathematical structures into machine learning models has proven fruitful in many works. Cohen et al. (2018a) adapt convolutional networks to operate on spherical and SO(3) valued data. Equivariant networks, investigated in Cohen & Welling (2016; 2017); Worrall et al. (2017); Weiler et al. (2018); Cohen et al. (2018b) reduce the complexity of a learning task by taking a quotient over group orbits which explain a subset of dimensions of the data manifold.
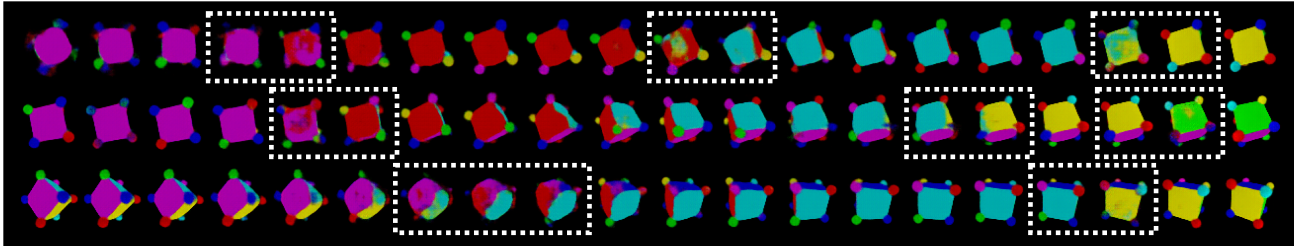
## 6. Experiments

We perform two experiments to investigate the importance of using a homeomorphic parameterization of the VAE in recovering the original underlying SO(3) manifold. In both experiments we explore three main axes of comparison: (1) manifold topology, (2) decoder architecture, and (3) specifically for the SO(3) models we compare different mean parameterizations as discussed in section 4.1. For each model we compute a tight bound on the negative log likelihood (NLL) through importance sampling following Burda et al. (2016).

For manifold topology we examine VAEs with the Gaussian parameterization ($\mathcal{N}$-VAE), the hyperspherical parameterization of Davidson et al. (2018) ($\mathcal{S}$-VAE), and the SO(3) latent variable discussed above. The two decoder variants are a simple MLP versus the group action decoder described in section 4.2. Lastly we explore mean parameterizations through unit Quaternions (q), the Lie algebra (alg), $\mathcal{S}^2 \times \mathcal{S}^1$ (s2s1), and $\mathcal{S}^2 \times \mathcal{S}^2$ (s2s2). These parameterizations are chosen to be either valid ($\mathcal{S}^2 \times \mathcal{S}^2$) or invalid (q, alg, $\mathcal{S}^2 \times \mathcal{S}^1$) for the purpose of investigating the soundness of our theoretical considerations and to compare their behaviour. Details and derivations on the properties of these different parameterizations can be found in Appendix D.

### 6.1. Toy experiment

The simplest way of creating a non-linear embedding of SO(3) in a high dimensional space is through the representation as discussed in Section 4.2. The data is created in the following way: a fixed representation $W : SO(3) \to \mathbb{R}^{N \times N}$ is chosen, in this experiment this is three copies

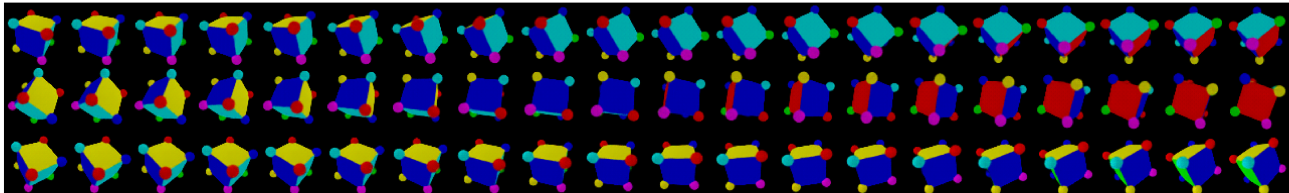(a) 10 dimensional Normal latent space with MLP decoder



(b) SO(3) latent space with S2S2 mean and action decoder

*Figure 6.1.* Three interpolations of two models. Discontinuities in the reconstructions of the Normal model are outlines by a dashed line.

of the direct sum of the Wigner-D-matrices up to order 3, making the embedding space $\mathbb{R}^{64}$. Subsequently a single element $v$ of $\mathbb{R}^{64}$ is generated. The data set now consists of the vectors $W(R)v$, where $R \in \mathrm{SO}(3)$ sampled uniformly. Since the representation is faithful ($W$ is injective) the data points lie on a 3 dimensional submanifold of $\mathbb{R}^{64}$ homeomorphic to $SO(3)$.

In order to verify the ability of the models to correctly learn to encode from the embedded manifold to the manifold itself, we learn various variational and non-variational auto-encoders on this data set. The encoder is a 3 layer MLP, and for the decoder we use the group action decoder of Section 4.2. The same representation $W$ is used as in the data generation, but we learn $v$. In addition to the SO(3) models, we use a 3 dimensional normal, which we map to SO(3) using the ZYZ-Euler angles, and a $S^3$ von Mises-Fisher, which we map to SO(3) by identifying $S^3$ as the unit quaternions.

**Results** The quantitative results are shown in Table 1. We observe that the choice for the mean parametrization significantly impacts the ability of the model to correctly learn the manifold. The $\mathcal{S}^2 \times \mathcal{S}^2$ method strongly outperforms the competing methods in the non-variational Auto-Encoder achieving near-perfect reconstructions. Additionally, the metric indicating the continuity of the encoder, which we define in Appendix E, shows it is the only model that does not have discontinuities in the latent space. These results are in line with our theoretical prediction outlined in Section 4.1 and Appendix D.

The qualitative results in Figure 4.2 and Figures A.1, A.2 in Appendix A tell a similar story. These plots are created by

| Algorithm | VAE | | | AE | |
| --- | --- | --- | --- | --- | --- |
| | NLL | recon | KL | recon | disc. |
| SO(3)-q | 10.9 | 2.32 | **9.16** | 0.29 | .992 |
| SO(3)-alg | 13.4 | 6.24 | 9.36 | 4.02 | 1. |
| SO(3)-s2s1 | 11.0 | 2.12 | 9.41 | 0.29 | 1. |
| SO(3)-s2s2 | **10.7** | 1.81 | 9.21 | **0.01** | 0. |
| $\mathcal{N}$-3-dim | 18.9 | 9.91 | 10.3 | 14.7 | 1. |
| $\mathcal{S}$-3-dim | 13.6 | **1.79** | 11.8 | 0.27 | 1. |

*Table 1.* Summary of results for the toy experiment for Variational Auto-Encoder (VAE) and Auto-Encoder (AE), including the discontinuity metric (disc.).

taking a $S^1$ subgroup of $SO(3)$ and making a $S^1$ submanifold in the data space using the same process with which the data was generated. This embedded trajectory is then encoded and reconstructed. The trajectory is divided in four equally sized partitions, each shown in a different color. We clearly see that only the $\mathcal{S}^2 \times \mathcal{S}^2$ method is able to learn a continuous latent space.

Moreover, the worst performing models are the 3 dimensional $\mathcal{N}$ and SO(3) algebra mean models. Interestingly, these share that at one intermediate point SO(3) is represented by $\mathbb{R}^3$. This indicates that using flat space to represent a non-trivial manifold results in a poorly structured latent space and worse reconstruction performance and Log Likelihoods.

### 6.2. Sphere-Cube

For this experiment we learn auto-encoders on renderings of a cube. The cube is made highly asymmetrical through the

colors of the faces and the colored spheres at the vertices. This should make it easier for the encoder to detect the orientation. This *sphere-cube* is then rotated by applying uniformly sampled group elements from $SO(3)$, to create a training set of 1M images. Ideally the model learns to correctly represent these encodings in the latent space.

The encoder consists of 5 convolutional layers, followed by one of the mean encoders and reparameterization methods. The decoder uses either the group action or a 3 layer MLP, both followed by a 5 deconvolutional layers. In order to balance reconstruction and the KL divergence in a controlled manner, we follow Burgess et al. (2018) and replace the negative KL term in the original VAE loss with a squared difference of the computed KL value and a target value. We found that a target value of 7 early in training to 15 at the end of the training gave good results. This allows the model to first organize the space and later become more certain of its predictions. We found that two additional regularizing loss terms were needed to correctly learn the latent space. Details can be found in Appendix G.

**Results**   Quantitative results comparing the best performing $SO(3)$ parameterization to $\mathcal{N}$-VAEs of diff dimensionality are shown in Table 2. Although the higher dimensional $\mathcal{N}$-VAEs are able to achieve competitive metrics compared to the best $SO(3)$ model, they only learn to embed $SO(3)$ in a high dimensional space in an unstructured fashion. As can be seen in in 6.1, the $SO(3)$ latent space with $\mathcal{S}^2 \times \mathcal{S}^2$ mean parameterization learns a nearly perfect encoding, while the 10 dimensional Normal learns disconnected patches of the data manifold.[3]

It can be seen in Table 3 that the results from the Toy experiment extend to this more complicated task. We observe that only the continuous encoding, $\mathcal{S}^2 \times \mathcal{S}^2$, achieves low log likelihood and reconstruction losses compared to the other mean parameterizations.

Lastly, we observe that the group action decoder yields significantly higher performance than the MLP decoder. This is in line with the hypotheses that using the group action encourages structure in the latent space.

## 7. Discussion & Conclusion

In this paper we explored the use of manifold-valued latent variables, by proposing an extension of the reparameterization trick to compact connected Lie groups. We worked out the implementation details for the specific case of $SO(3)$, and highlighted the various subtleties that must be taken into account to ensure a successful parameterization of the VAE. Through a series of experiments, we showed the importance

---

[3]Animated interpolations can be found at `https://sites.google.com/view/lie-vae`.

| Algorithm | NLL | ELBO | recon. |
|---|---|---|---|
| SO(3)-MLP-s2s2 | 123.6 | 144.6 | 129.6 |
| SO(3)-action-s2s2 | **46.90** | **63.35** | **48.35** |
| $\mathcal{N}$-MLP 3-dim | 140.7 | 157.7 | 142.7 |
| $\mathcal{N}$-MLP 10-dim | 64.02 | 80.80 | 65.80 |
| $\mathcal{N}$-MLP 30-dim | 55.7 | 74.37 | 59.37 |

*Table 2.* Results on sphere-cube of $SO(3)$ encodings and $\mathbb{R}^N$ embedding encodings. The $SO(3)$ models employ both regularizers, the $\mathbb{R}^N$ models neither. This achieved the best respective performance.

| Algorithm | NLL | ELBO | recon. |
|---|---|---|---|
| SO(3)-MLP-q | 111.8 | 140.1 | 135.1 |
| SO(3)-MLP-alg | 218.9 | 316.7 | 301.7 |
| SO(3)-MLP-s2s1 | 106.0 | 144.5 | 129.5 |
| SO(3)-MLP-s2s2 | 123.6 | 144.6 | 129.6 |
| SO(3)-action-q | 378.6 | 471.2 | 456.2 |
| SO(3)-action-alg | 241.2 | 333.2 | 318.2 |
| SO(3)-action-s2s1 | 128.5 | 173.0 | 158.0 |
| SO(3)-action-s2s2 | **46.90** | **63.35** | **48.35** |

*Table 3.* Summary of results comparing $SO(3)$ mean parameterization and model decoder on sphere-cubes. The models employ both regularizers.

of matching the topology of the latent data manifold with that of the latent variables to induce a continuous, well-behaved latent space. Additionally we demonstrated the improvement in learned latent space structure by using a group action decoder, and the need for care in choosing an embedding space for the posterior distribution's mean parameter.

We believe that the use of $SO(3)$ and other well-known manifold-valued latent variables could present an interesting addition to tackling problems in such fields as model based RL and computer vision. Moving forward we thus aim to extend this theory to other Lie groups such as $SE(3)$. A limitation of the current work, and reparameterizing distributions on specific manifolds in general, is that it relies on the assumption of *a priori* knowledge about the observed data's latent structure. Hence in future work our ambition is to find a general theory to learn arbitrary manifolds not known in advance.

# References

Berg, Rianne van den, Hasenclever, Leonard, Tomczak, Jakub M, and Welling, Max. Sylvester normalizing flows for variational inference. *UAI*, 2018.

Burda, Yuri, Grosse, Roger, and Salakhutdinov, Ruslan. Importance weighted autoencoders. *ICLR*, 2016.

Burgess, Christopher P, Higgins, Irina, Pal, Arka, Matthey, Loic, Watters, Nick, Desjardins, Guillaume, and Lerchner, Alexander. Understanding disentangling in *beta*-vae. *arXiv preprint arXiv:1804.03599*, 2018.

Chirikjian, Gregory S. Information-theoretic inequalities on unimodular lie groups. *Journal of geometric mechanics*, 2(2):119, 2010.

Chirikjian, Gregory S. *Stochastic Models, Information Theory, and Lie Groups*. 2012.

Chirikjian, Gregory S and Kyatkin, Alexander B. *Engineering applications of noncommutative harmonic analysis: with emphasis on rotation and motion groups*. CRC press, 2000.

Chirikjian, Gregory S and Kyatkin, Alexander B. *Harmonic Analysis for Engineers and Applied Scientists: Updated and Expanded Edition*. Courier Dover Publications, July 2016.

Cohen, Taco and Welling, Max. Group equivariant convolutional networks. In *ICML*, pp. 2990–2999, 2016.

Cohen, Taco S and Welling, Max. Harmonic exponential families on manifolds. *ICML*, 2015.

Cohen, Taco S and Welling, Max. Steerable cnns. *ICLR*, 2017.

Cohen, Taco S., Geiger, Mario, Köhler, Jonas, and Welling, Max. Spherical CNNs. *ICLR*, 2018a.

Cohen, Taco S, Geiger, Mario, and Weiler, Maurice. Intertwiners between induced representations (with applications to the theory of equivariant neural networks). March 2018b.

Davidson, Tim R., Falorsi, Luca, Cao, Nicola De, Kipf, Thomas, and Tomczak, Jakub M. Hyperspherical Variational Auto-Encoders. *UAI*, 2018.

Eslami, S. M. Ali, Jimenez Rezende, Danilo, Besse, Frederic, Viola, Fabio, Morcos, Ari S., Garnelo, Marta, Ruderman, Avraham, Rusu, Andrei A., Danihelka, Ivo, Gregor, Karol, Reichert, David P., Buesing, Lars, Weber, Theophane, Vinyals, Oriol, Rosenbaum, Dan, Rabinowitz, Neil, King, Helen, Hillier, Chloe, Botvinick, Matt, Wierstra, Daan, Kavukcuoglu, Koray, and Hassabis,

Demis. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. ISSN 0036-8075. doi: 10.1126/science.aar6170. URL http://science.sciencemag.org/content/360/6394/1204.

Figurnov, Michael, Mohamed, Shakir, and Mnih, Andriy. Implicit reparameterization gradients. *arXiv preprint arXiv:1805.08498*, 2018.

Hall, B. *Lie Groups, Lie Algebras, and Representations: An Elementary Introduction*. Graduate Texts in Mathematics. Springer, 2003. ISBN 9780387401225.

Jang, Eric, Gu, Shixiang, and Poole, Ben. Categorical reparameterization with gumbel-softmax. *ICLR*, abs/1611.01144, 2017.

Kingma, Diederik P. and Welling, Max. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.

Kingma, Diederik P, Salimans, Tim, Jozefowicz, Rafal, Chen, Xi, Sutskever, Ilya, and Welling, Max. Improved variational inference with inverse autoregressive flow. In *NIPS*, pp. 4743–4751, 2016.

Maddison, Chris J, Mnih, Andriy, and Teh, Yee Whye. The concrete distribution: A continuous relaxation of discrete random variables. *ICLR*, 2017.

Naesseth, Christian, Ruiz, Francisco, Linderman, Scott, and Blei, David. Reparameterization gradients through acceptance-rejection sampling algorithms. In *AISTATS*, pp. 489–498, 2017.

Nalisnick, Eric and Smyth, Padhraic. Stick-breaking variational autoencoders. *ICLR*, 2017.

Rezende, Danilo and Mohamed, Shakir. Variational inference with normalizing flows. *ICML*, 37:1530–1538, 2015.

Rezende, Danilo Jimenez, Mohamed, Shakir, and Wierstra, Daan. Stochastic backpropagation and approximate inference in deep generative models. *ICML*, pp. 1278–1286, 2014.

Rodrigues, Olinde. *Des lois géométriques qui régissent les déplacements d'un système solide dans l'espace: et de la variation des cordonnées provenant de ces déplacements considérés indépendamment des causes qui peuvent les produire*. 1840.

Tomczak, Jakub M and Welling, Max. VAE with a VampPrior. *AISTATS*, 2017.

Weiler, Maurice, Hamprecht, Fred A, and Storath, Martin. Learning steerable filters for rotation equivariant CNNs. In *CVPR*, 2018.

Worrall, Daniel E, Garbin, Stephan J, Turmukhambetov, Daniyar, and Brostow, Gabriel J. Harmonic networks: Deep translation and rotation equivariance. In *CVPR*, 2017.
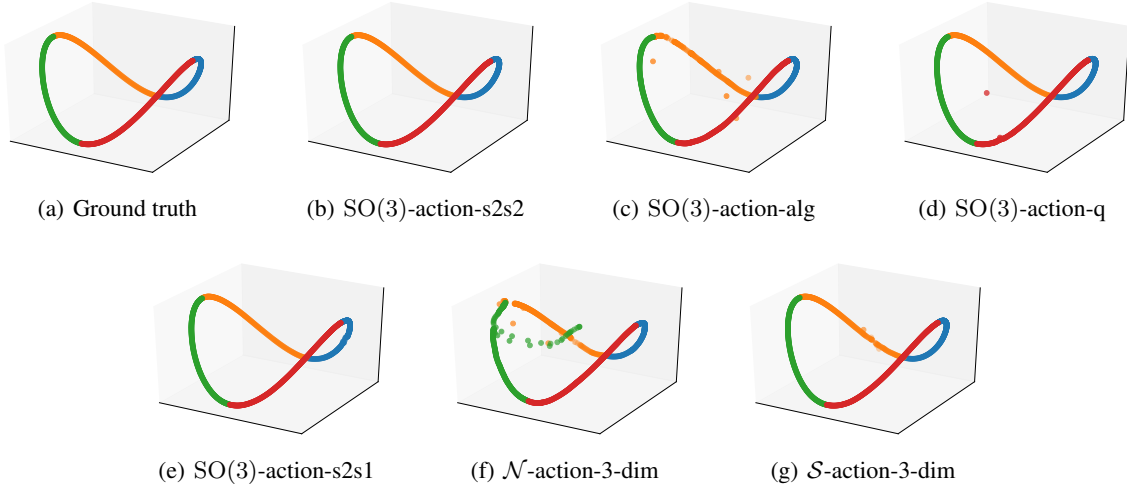
# A. Additional Figures



(a) Ground truth    (b) SO(3)-action-s2s2    (c) SO(3)-action-alg    (d) SO(3)-action-q

(e) SO(3)-action-s2s1    (f) $\mathcal{N}$-action-3-dim    (g) $\mathcal{S}$-action-3-dim

*Figure A.1.* Reconstructions of a $S^1$ trajectory in the Toy data set. The $\mathbb{R}^{64}$ elements are mapped to 3D by Principal Component Analysis. See Section 6.1 for details.
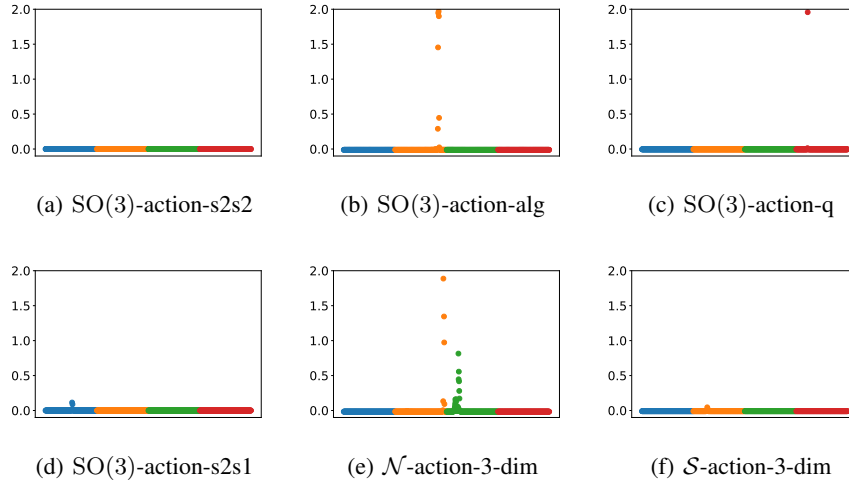


(a) SO(3)-action-s2s2    (b) SO(3)-action-alg    (c) SO(3)-action-q

(d) SO(3)-action-s2s1    (e) $\mathcal{N}$-action-3-dim    (f) $\mathcal{S}$-action-3-dim

*Figure A.2.* Discontinuities in the latent space along a $S^1$ trajectory in Toy data set. Shown is $\|f(x_{i+1}) - f(x_i)\|^2$ for encoder $f$ along the trajectory. See Section 6.1 for details.

# B. Pushforward Measure SO(3)

**Theorem 1.** *Let* $(\mathbb{R}^3, \lambda, \mathcal{B}[\mathbb{R}^3])$ *the real space provided with the Lebesgue measure on the Borel algebra on* $\mathbb{R}^3$. *Let* $(\mathrm{SO}(3), \nu, \mathcal{B}[\mathrm{SO}(3)])$ *The group of 3 dimensional rotations provided with the normalized Haar measure* $\nu$ *on the Borel algebra on* $\mathrm{SO}(3)$. *Consider then the probability measure* $\mu : \mathcal{B}[\mathbb{R}^3] \to [0,1]$ *absolutely continuous w.r.t* $\lambda$ *with density* $r$. *Consider the exponential map* $\exp : \mathbb{R}^3 \to \mathrm{SO}(3)$ *that is differentiable thus continuous thus measurable. Let then* $\exp_*(\mu)$ *the pushforward of* $\mu$ *by the* $\exp$ *function. then* $\exp_*(\mu)$ *is absolutely continuous with respect of the Haar measure* $\nu$. $(\exp_*(\mu) \ll \nu)$

*Proof.* Define the sets:

$$A_k = \{x \in \mathbb{R}^3 : \|x\| \in (k\pi, (k+1)\pi)\} \quad B_k = \{x \in \mathbb{R}^3 : \|x\| = k\pi\} \quad k \in \mathbb{N} \tag{10}$$

Then note that:

$$\mathbb{R}^3 = \left( \dot{\bigcup_{k \in \mathbb{N}}} A_k \right) \dot{\cup} \left( \dot{\bigcup_{k \in \mathbb{N}}} B_k \right) \tag{11}$$

And since $m(B_k) = \mu(B_k) = 0$, then $\mu_{B_k}(E) = 0 \quad \forall E \in \mathcal{B}[\mathbb{R}^3] \, k \in \mathbb{N}$. Therefore:

$$\mu(E) = \sum_{k \in \mathbb{N}} \mu_{A_k}(E) \quad \forall E \in \mathcal{B}[\mathbb{R}^3] \tag{12}$$

Then we have $\mu = \sum_{k \in \mathbb{N}} \mu_{A_k}$.
Now consider the pushforward measure $\exp_*(\mu)$, we then have that:

$$(\exp_*(\mu))(E) = \mu(\exp^{-1}(E)) = \sum_{k \in \mathbb{N}} \mu_{A_k}(\exp^{-1}(E)) = \sum_{k \in \mathbb{N}} (\exp_*(\mu_{A_k}))(E) = \tag{13}$$

$$\sum_{k \in \mathbb{N}} ((\exp_{|A_k})_*(\mu))(E) \quad \forall E \in \mathcal{B}[\mathrm{SO}(3)] \tag{14}$$

Then we have $\exp_*(\mu) = \sum_{k \in \mathbb{N}} \exp_*(\mu_{A_k}) = \sum_{k \in \mathbb{N}} ((\exp_{|A_k})_*(\mu))$. Where $\exp_{|A_k}$ is the exp function restricted to $A_k$. Moreover notice that $\exp_{|A_k}$ is a injective, therefore we can apply the change of variable formula:

$$((\exp_{|A_k})_*(\mu))(E) = \int_{\exp_{|A_k}^{-1}(E)} r \, d\lambda = \int_E (r \circ \exp_{|A_k}^{-1}) \cdot |J_{\exp_{|A_k}^{-1}}| \, d\nu \tag{15}$$

Then $(\exp_{|A_k})_*(\mu) \ll \nu$ and since $\exp_*(\mu) = \sum_{k \in \mathbb{N}} ((\exp_{|A_k})_*(\mu))$ then $\exp_*(\mu) \ll \nu$. $\qquad \square$

The proof then tells us how to compute the Radon-Nikodym derivative of the pushforward with respect to the Haar measure. In fact:

$$\frac{d(\exp_{|A_k})_*(\mu)}{d\nu} = (r \circ \exp_{|A_k}^{-1}) \cdot |J_{\exp_{|A_k}^{-1}}| \quad , \quad \frac{d \exp_*(\mu)}{d\nu} = \sum_{k \in \mathbb{N}} \frac{d(\exp_{|A_k})_*(\mu)}{d\nu} \tag{16}$$

Defining $\hat{q} := \frac{d \exp_*(\mu)}{d\nu}$ we then have:

$$\hat{q}(R) = \sum_{k \in \mathbb{N}} (r \circ \exp_{|A_k}^{-1}(R)) \cdot |J_{\exp_{|A_k}^{-1}}(R)| = \sum_{v \in \exp^{-1}(R)} r(v) \cdot |J_{\exp}(v)|^{-1} \tag{17}$$

From (Chirikjian, 2010) we have that

$$|J_{\exp}(v)| = \frac{2 - 2\cos\|v\|}{\|v\|^2} \tag{18}$$

We then have:

$$\hat{q}(R) = \sum_{v \in \exp^{-1}(R)} r(v) \frac{\|v\|^2}{2 - 2\cos(\|v\|)} \tag{19}$$

To then have an expression explicitly dependent on $R$ consider that

$$\exp_{|A_k}^{-1}(R) = \frac{\exp_{|A_0}^{-1}(R)}{\|\exp_{|A_0}^{-1}(R)\|}(\|\exp_{|A_0}^{-1}(R)\| + 2k\pi) = \frac{\log(R)}{\|\log(R)\|}(\|\log(R)\| + k\pi) \quad \text{if } k \text{ is even} \tag{20}$$

$$\exp_{|A_k}^{-1}(R) = \frac{\exp_{|A_0}^{-1}(R)}{\|\exp_{|A_0}^{-1}(R)\|}(\|\exp_{|A_0}^{-1}(R)\| + 2k\pi) = \frac{\log(R)}{\|\log(R)\|}(\|\log(R)\| + -(k+1)\pi) \quad \text{if } k \text{ is odd} \tag{21}$$

Where we have defined $\log := \exp_{|A_0}^{-1}(R)$. Moreover we then have:

$$|J_{\exp_{|A_k}^{-1}}(R)| = \frac{\|\exp_{|A_k}^{-1}(R)\|^2}{2 - 2\cos(\|\exp_{|A_k}^{-1}(R)\|)} = \frac{(\|\log(R)\| + k\pi)^2}{2 - 2\cos(\|\log(R)\|)} \quad \text{if } k \text{ is even} \tag{22}$$

$$|J_{\exp_{|A_k}^{-1}}(R)| = \frac{\|\exp_{|A_k}^{-1}(R)\|^2}{2 - 2\cos(\|\exp_{|A_k}^{-1}(R)\|)} = \frac{(\|\log(R)\| - (k+1)\pi)^2}{2 - 2\cos(\|\log(R)\|)} \quad \text{if } k \text{ is odd} \tag{23}$$

Putting everything together:

$$\hat{q}(R) = \sum_{k \in \mathbb{Z}} r\left(\frac{\log(R)}{\|\log(R)\|}(\|\log(R)\| + 2k\pi)\right) \frac{(\|\log(R)\| + 2k\pi)^2}{2 - 2\cos(\|\log(R)\|)} \tag{24}$$

Where from (Chirikjian, 2010) we have:

$$\log(R) = \frac{\theta(R)}{2\sin(\theta(R))}(R - R^\top) \quad \theta(R) = \cos^{-1}\left(\frac{\text{tr}(R) - 1}{2}\right) \tag{25}$$

This gives us the final expression:

$$\hat{q}(R|\sigma) = \sum_{k \in \mathbb{Z}} r\left(\frac{\log(R)}{\theta(R)}(\theta(R) + 2k\pi)\right) \frac{(\theta(R) + 2k\pi)^2}{3 - \text{tr}(R)} \tag{26}$$

## C. Entropy computation

We oprimize MC estimates of the Entropy:

$$\mathbb{H}(q) = \mathbb{H}(\hat{q}) \approx -\frac{1}{N}\sum_{i=1}^{N} \log\hat{q}(R_i), \quad R_i \sim \hat{q}(R_i)$$

(Where we dropped dependency on the parameters for simplicity) Then using, Equation (24):

$$\mathbb{H}(q) = \mathbb{H}(\hat{q}) \approx -\frac{1}{N}\sum_{i=1}^{N} \log\sum_{k \in \mathbb{Z}} r\left(\frac{\log(R_i)}{\|\log(R_i)\|}(\|\log(R_i)\| + 2k\pi)\right) \frac{(\|\log(R_i)\| + 2k\pi)^2}{2 - 2\cos(\|\log(R_i)\|)}, \quad R_i \sim \hat{q}(R_i)$$

In the way we defined $\hat{q}$ we obtain samples from it in the following way:

$$R_i = \exp(\mathbf{v}_i) \quad \mathbf{v}_i \sim r(\mathbf{v}_i) \tag{27}$$

Substituting it in in the previous expression we get:

$$\mathbb{H}(q) \approx -\frac{1}{N}\sum_{i=1}^{N} \log\hat{q}(\exp(\mathbf{v}_i))$$

$$= -\frac{1}{N}\sum_{i=1}^{N} \log\sum_{k \in \mathbb{Z}} r\left(\frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}(\|\mathbf{v}_i\| + 2k\pi)\right) \cdot \frac{(\|\mathbf{v}_i\| + 2k\pi)^2}{2 - 2\cos(\|\mathbf{v}_i\|)}, \quad \mathbf{v}_i \sim r(\mathbf{v}_i)$$

Notice that this expression depends only on the samples from $r$ in the lie algebra

Assuming the density $r$ decays quickly enough to zero, the above infinite summation can be truncated. This is always the case for exponentially decaying distributions, like the Normal. The truncated summation can then can be computed using the *logsumexp* trick:

$$\mathbb{H}(q) \approx -\frac{1}{N}\sum_{i=1}^{N} \text{logsumexp}_k\left(\log r\left(\frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}(\|\mathbf{v}_i\| + 2k\pi)\right) + \log\frac{(\|\mathbf{v}_i\| + 2k\pi)^2}{2 - 2\cos(\|\mathbf{v}_i\|)}\right), \quad \mathbf{v}_i \sim r(\mathbf{v}_i) \tag{28}$$

# D. Mean parameterization

As discussed above, some requirements exist on $\pi : \mathcal{Y} \to \mathrm{SO}(3)$ for the encoder $\mathrm{enc}^m u$ to correctly represent the data manifold.

We split $\pi$ in the composition of $\phi : \mathcal{Y} \to \mathcal{Y}'$ and $\psi : \mathcal{Y}' \to \mathrm{SO}(3)$, both generally discontinuous. We assume $\mathcal{Y} = \mathbb{R}^m$ to be a neural network output. The functions $\psi$ below are known ways to surjectively map to $\mathrm{SO}(3)$. $\phi$ are constructed to map from $\mathbb{R}^m$ to the domain of $\psi$.

We discuss the existence of a map $i : \mathrm{SO}(3) \to \mathcal{Y}'$ such that it is a right inverse of $\psi$ ($\psi \circ i = \mathrm{id}_{\mathrm{SO}(3)}$), which is necessary for the correct encoder to exist.

1. **Algebra** with $\mathcal{Y} = \mathcal{Y}' = \mathbb{R}^3$, $\phi = \mathrm{id}$. This method simply uses the exponential map:

$$\pi : \mathbb{R}^3 \to \mathrm{SO}(3) \tag{29}$$
$$\mathbf{v} \mapsto \exp(\mathbf{v}_\times) \tag{30}$$

   It's inverses are the branches of the log map. However, a path in $\mathrm{SO}(3)$ that is a full rotation around a fixed axis is continuous in $\mathrm{SO}(3)$ but discontinuous in the algebra, when mapped with the log map. Thus the log map is not continuous.

2. **Quaternions** with $\mathcal{Y} = \mathbb{R}^4$, $\mathcal{Y}' = S^3$, $\phi(x) = x/\|x\|$. The unit Quaternions, which are homeomorphic to $S^3$ are a 'double cover' of $\mathrm{SO}(3)$, which means that a continuous surjective projection $\pi : S^3 \to \mathrm{SO}(3)$ exists that is two-to-one. The projection map can be found in Chirikjian & Kyatkin (2000, Eqn. (5.60)). Using the theory of Fiber Bundles (recognizing $S^3$ as a non-trivial principle bundle with base space $\mathrm{SO}(3)$), one can show that no embedding $i$ exists.

3. **s2s1**($S^2 \times S^1$) with $\mathcal{Y} = \mathbb{R}^3 \times \mathbb{R}^2$, $\mathcal{Y}' = S^2 \times S^1$, $\phi(x, y) = (x/\|x\|, y/\|y\|)$. This is the map from an axis in $S^2$ and angle in $S^1$.

$$\pi : \mathcal{S}^2 \times \mathcal{S} \to \mathrm{SO}(3) \tag{31}$$
$$(\mathbf{u}, \mathbf{v}) \mapsto \mathbf{I} + v_2 \mathbf{u}_\times + (1 - v_1)\mathbf{u}_\times^2 \tag{32}$$

   For $i$, to be continuous, its image $i(\mathrm{SO}(3))$ must be closed (as it is a compact subset of a Hausdorff space). Thus so must the set $A = i(\mathrm{SO}(3)) \cap S^2 \times \{\pi\}$. However, as $\psi(\mu, \pi) = \psi(-\mu, \pi)$ for $\mu \in S^2$, $A$ is a hemisphere (times a point) that does not contain its entire boundary, thus it is not closed and $i$ is not continuous.

4. **s2s2**($S^2 \times S^2$) with $\mathcal{Y} = \mathbb{R}^3 \times \mathbb{R}^3$, $\mathcal{Y}' = S^2 \times S^2$, $\phi(x, y) = (x/\|x\|, y/\|y\|)$. This method creates two orthonormal vectors.

$$\pi : \mathcal{S}^2 \times \mathcal{S}^2 \to \mathrm{SO}(3) \tag{33}$$
$$(\mathbf{u}, \mathbf{v}) \mapsto \mathrm{concat}(\mathbf{w_1}, \mathbf{w_2}, \mathbf{w_3}) \tag{34}$$
$$\text{Where:} \tag{35}$$
$$\mathbf{w_1} = \mathbf{u} \tag{36}$$
$$\mathbf{w_2}' = \mathbf{v} - \langle \mathbf{u}, \mathbf{v} \rangle \mathbf{u} \tag{37}$$
$$\mathbf{w_2} = \frac{\mathbf{w_2}'}{\|\mathbf{w_2}'\|} \tag{38}$$
$$\mathbf{w_3} = \mathbf{w_1} \times \mathbf{w_2} \tag{39}$$

   Notice that there exists a continuous and injective map $i : SO(3) \to S^2 \times S^2$. It simply consists of taking the first two rows of the matrix representation of the $SO(3)$ element (The third row is the vector product between the first two, so it can always be recovered). Moreover we have that $\pi \circ i = \mathrm{Id}_{SO(3)}$

# E. Continuity Metric

Consider a map $f : X \to Y$ where $X, Y$ are metric spaces with metrics $d_X$ and $d_Y$ respectively. In order to compute the proposed continuity metric we take a *continuous* path $(x_i)_{i \in [N]}$, defined as $N$ pairwise close points, and compute the

relative distances

$$L_i = \frac{d_Y(f(x_{i+1}), f(x_i))}{d_X(x_{i+1}, x_i)} \tag{40}$$

From this we further compute the quantities

$$M := \max_i L_i \quad \text{and} \quad P_\alpha := \alpha\text{-th percentile of } \{L_i : i \in [N-1]\}. \tag{41}$$

By comparing these two values, we want to discover whether there is at least one outlier in the set of $L_i$. Such outliers corresponds to a transition with a *big* jump, signalling a discontinuity point. We define a path to be discontinuous if $M > \gamma P_\alpha$.

In order to capture stochastic effects we repeat the above procedure with several paths. The final score is the fraction of discontinuous paths. In the practical implementation we chose 1000 paths, using $\gamma = 10$ and $\alpha = 90$ (90th percentile).

## F. Group Action

For each degree $l \in \mathbb{Z}_{\geq 0}$, the Wigner-D-matrix can be expressed in a real basis as $D^l : \mathrm{SO}(3) \to \mathbb{R}^{(2l+1)\times(2l+1)}$. We choose the $n_l$ copies of each degree $l$ and stack the matrices in block-diagonal form to create our representation, which amounts to taking the direct sum of the representations.

The Wigner-D-matrices represent rotations of the Fourier modes of a signal on the sphere, which provides an interpretation for using the group action. We consider a real signal one the sphere $f : S^2 \to \mathbb{R}$. It has a generalized Fourier transformation (Chirikjian & Kyatkin, 2000):

$$f(s) = \sum_{l=0}^{\infty} (2l+1) \sum_{m,n=-l}^{l} \hat{f}_m^l D_{m0}^l(\alpha_s, \beta_s, 0)$$

where $\hat{f}$ are the Fourier components and $D$ is the Wigner-D-matrix. We use identity $D_{m0}^l(\alpha, \beta, 0) = Y_m^l(\alpha, \beta)$, where $\alpha, \beta$ are the first two Euler angles, to write the spherical harmonics that are the basis functions of the Fourier modes as Wigner-D-matrices. Then for a rotation $g \in \mathrm{SO}(3)$, using the homomorphism property:

$$f(g(s)) = \sum_{l=0}^{\infty} (2l+1) \sum_{m=-l}^{l} \hat{f}_m^l \sum_{r=-l}^{l} D_{mr}^l(g) D_{r0}^l(\alpha_s, \beta_s, 0)$$

$$= \sum_{l=0}^{\infty} (2l+1) \sum_{m=-l}^{l} \left( \sum_{r=-l}^{l} D_{mr}^l(g) \hat{f}_m^l \right) D_{r0}^l(\alpha_s, \beta_s, 0)$$

where $g(s)$ corresponds to rotating a point on the sphere.

We see that our method of using representations in the decoder corresponds to having the content latent code represent the Fourier coefficients of a virtual signal on the sphere.

## G. Regularizers

Even when an appropriate mean parametrization is selected and proper behaviour of the decoder is encouraged by the group action decoder, the network can still learn a discontinuous latent space. To encourage it to learn the data manifold correctly, we employ two additional loss terms that act as regularizers. An ablative analysis of the effectiveness of these regularizers is shown in Table 4.

### G.1. Equivariance regularizer

If the 3D object on which $\mathrm{SO}(3)$ acts is centered in the frame, then a $S^1$ subgroup $H$ of $\mathrm{SO}(3)$ exist such that its action corresponds to the rotations whose axis is orthogonal to the camera frame. For any $R \in \mathrm{SO}(3)$, angle $\theta$, decoder

| Regularizer | NLL | ELBO | recon. . |
|---|---|---|---|
| Neither | 235.24 | 246.2 | 231.2 |
| Equivariance | 75.62 | 93.76 | 78.76 |
| Continuity | 87.36 | 125.6 | 110.6 |
| Both | **45.18** | **60.62** | **45.62** |

*Table 4.* Ablative analysis of the regularizers. The has model uses $\mathrm{SO}(3)$ latent space, the group action decoder and the S2S2 mean.
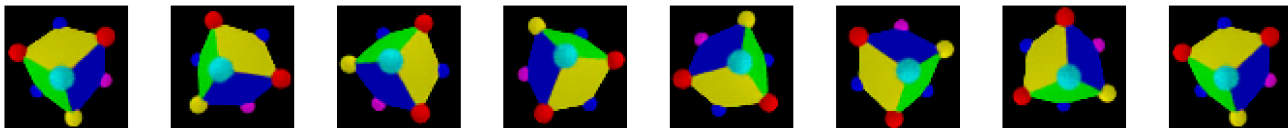


*Figure G.1.* Decodings of an orbit of the equivariance subgroup $S^1$ whose axis is orthogonal to the camera frame. This orbit corresponds to planar rotations of the pixels.

$g : \mathrm{SO}(3) \to \mathbb{R}^N$, action of subgroup $H$ on the latent space, $\psi_\theta : \mathrm{SO}(3) \to \mathrm{SO}(3)$ and action on the pixels through planar rotations $\phi_\theta : \mathbb{R}^N \to \mathbb{R}^N$, we have equivariance relationship:

$$g(\psi_\theta(R)) = \phi_\theta(g(R)) \tag{42}$$

This equivariance is shown in Figure G.1. The relationship is exact if the object is centered, $\mathrm{SO}(3)$ acts on all pixels and if the camera is orthographic (located infinitely far away from the subject). If the object is off center, the pixel rotation can be performed around a learned center point. If the images have a rotation-invariant background, a learned mask can be applied. If the camera is not orthographic, the equivariance relationship is not exact, but approximate. The decoder is regularized by enforcing Equation (42) through a mean squared error loss on the pixels for uniformly sampled $R \in \mathrm{SO}(3)$ and $\theta \in S^1$. We choose $\psi_\theta$ to correspond to rotation around the $z$-axis.

This regularizer helps align all rotations in each $S^1$ orbit, but does not help in correctly aligning the orbits among each other. Thus we reduce the problem from aligning $\mathrm{SO}(3)$ to aligning $\mathrm{SO}(3)/S^1 \cong S^2$, since the cosets of $\mathrm{SO}(3)$ after the orbit are identified, are homeomorphic to the sphere.

## G.2. Continuity regularizer

If the learner is provided with pairs images that are nearby with respect to the manifold metric, the encoder can be regularized by penalizing differences in the encodings of the two inputs. This is done by penalizing the mean squared error of the Frobenius norms of the two encoded rotation matrices, which is a proper metric on the $\mathrm{SO}(3)$ manifold.

This simplifies the problem from unsupervised learning on i.i.d. samples to learning a VAE on two frame samples from random trajectories of data lying on the $\mathrm{SO}(3)$ manifold.