
Latent Geodesics of Model Dynamics for Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Model-based offline reinforcement learning approaches generally rely on bounds
2 of model error. While contemporary methods achieve such bounds through an
3 ensemble of models, we propose to estimate them using a data-driven latent metric.
4 Particularly, we build upon recent advances in Riemannian geometry of generative
5 models to construct a latent metric of an encoder-decoder based forward model.
6 Our proposed metric measures both the quality of out of distribution samples as
7 well as the discrepancy of examples in the data. We show that our metric can be
8 viewed as a combination of two metrics, one relating to proximity and the other to
9 epistemic uncertainty. Finally, we leverage our metric in a pessimistic model-based
10 framework, showing a significant improvement upon contemporary model-based
11 offline reinforcement learning benchmarks.

12 1 Introduction

13 *This work focuses on leveraging Riemannian geometry of generative models in offline reinforcement*
14 *learning.*

15 **Offline reinforcement learning (offline RL)** (Levine et al., 2020), a.k.a. batch-mode reinforcement
16 learning (Ernst et al., 2005; Riedmiller, 2005; Fonteneau et al., 2013), involves learning a policy
17 from potentially suboptimal data. In contrast to imitation learning (Schaal, 1999), offline RL does
18 not rely on expert demonstrations, but rather seeks to surpass the average performance of the agents
19 that generated the data. Methodologies such as the gathering of new experience fall short in offline
20 settings, requiring reassessment of fundamental learning paradigms (Buckman et al., 2020; Wang
21 et al., 2020; Zanette, 2020).

22 **The geometry of latent generative models** has recently gained interest in unsupervised domains
23 (Chen et al., 2018; Arvanitidis et al., 2018; Chen et al., 2020a; Arvanitidis et al., 2020). There,
24 variational autoencoders (VAEs) have been shown to capture significant metrics in their latent
25 representations. The resulting manifold has been shown to capture a smooth metric of the ambient
26 output space, as well as properly capture uncertainty estimates in out of distribution (OOD) regions
27 (Arvanitidis et al., 2018).

28 In this work, we introduce the aforementioned Riemannian theory of generative models to reinforce-
29 ment learning. Specifically, we generalize previous results in VAEs to learn a Riemannian manifold
30 w.r.t. the *environment’s dynamics*. We achieve this by training a variational forward model of the
31 next state. Our latent model induces a manifold and metric which capture the natural characteristics
32 of the environment’s dynamics. Moreover, we show that this metric can be analytically decoupled
33 into metrics relating to proximity and uncertainty. Our proposed metric can be utilized in a vast range
34 of model-based approaches in reinforcement learning (e.g., offline RL, planning). Here, we show
35 how our learned metric can be leveraged in a model-based offline reinforcement learning framework.

36 Our contributions are as follows.

37 **Technical Contributions:** (1) We introduce a natural metric for forward model dynamics. The
 38 induced metric, for which we derive analytical expression for in Section 4, can be represented as
 39 a union of two metrics; namely, a metric of proximity and a metric of uncertainty. We depict the
 40 geodesics of the induced metric on a grid-like environment, suggesting our latent model captures
 41 valuable structural dependencies. (2) We integrate our metric in a model-based offline RL framework,
 42 where an agent is penalized with accordance to its distance to the data. As such, we demonstrate
 43 improved performance to contemporary offline RL approaches on several benchmarks (Section 6).

44 **Broader Impact:** Our proposed metric can be leveraged in a vast range of domains. While our work
 45 is focused on its application to offline RL, its unique characteristics can be utilized in online control,
 46 planning, and predictive analysis, as well as improve explainability of the agent and the environment.
 47 Still, using approximate models to make decisions in the real world can bring to negative social
 48 impact. Wrongful, unethical, or dangerous decisions may harm individuals affected by such actions.

49 2 Preliminaries

50 2.1 Offline Reinforcement Learning

51 We consider the standard Markov Decision Process (MDP) framework (Sutton et al., 1998) defined
 52 by the tuple $(\mathcal{S}, \mathcal{A}, r, P, \alpha)$, where \mathcal{S} is the state space, \mathcal{A} the action space, $r : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ the
 53 reward function, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ the transition kernel, and $\alpha \in (0, 1)$ is the discount factor.

54 In the online setting of reinforcement learning (RL), the environment initiates at some state $s_0 \sim \rho_0$.
 55 At any time step the environment is in a state $s \in \mathcal{S}$, an agent takes an action $a \in \mathcal{A}$ and receives a
 56 reward $r(s, a)$ from the environment as a result of this action. The environment transitions to state s'
 57 according to the transition function $P(\cdot|s, a)$. The goal of online RL is to find a policy $\pi(a|s)$ that
 58 maximizes the expected discounted return $v^\pi = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \alpha^t r(s_t, a_t) | s_0 \sim \rho_0 \right]$.

59 Unlike the online setting, the offline setup considers a dataset $\mathcal{D}_n = \{s_i, a_i, r_i, s'_i\}_{i=1}^n$ of transitions
 60 generated by some unknown agents. The objective of offline RL is to find the best policy in the test
 61 environment (i.e., real MDP) given only access to the data generated by the unknown agents.

62 2.2 Riemannian Manifolds

63 We define the Riemannian pullback metric, a fundamental component of our proposed method. We
 64 refer the reader to Carmo (1992) for further details on Riemannian geometry.

65 We are interested in studying a smooth surface M with a Riemannian metric g . A Riemannian metric
 66 is a smooth function that assigns a symmetric positive definite matrix to any point in M . At each
 67 point $z \in M$ a tangent space $T_z M$ specifies the pointing direction of vectors “along” the surface.

68 **Definition 1.** *Let M be a smooth manifold. A Riemannian metric g on M changes smoothly and
 69 defines a real scalar product on the tangent space $T_z M$ for any $z \in M$ as*

$$g_z(x, y) = \langle x, y \rangle_z = \langle x, G(z)y \rangle, \quad x, y \in T_z M,$$

70 where $G(z) \in \mathbb{R}^{d_z \times d_z}$ is the corresponding metric tensor. (M, g) is called a Riemannian manifold.

71 The Riemannian metric enables us to easily define geodesic curves. Consider some differentiable
 72 mapping $\gamma : [0, 1] \mapsto M \subseteq \mathbb{R}^{d_z}$, such that $\gamma(0) = z_0, \gamma(1) = z_1$. The length of the curve γ measured
 73 on M is given by

$$L(\gamma) = \int_0^1 \sqrt{\left\langle \frac{\partial \gamma(t)}{\partial t}, G(\gamma(t)) \frac{\partial \gamma(t)}{\partial t} \right\rangle} dt. \quad (1)$$

74 The geodesic distance $d(z_1, z_2)$ between any two points $z_1, z_2 \in M$ is then the infimum length over
 75 all curves γ for which $\gamma(0) = z_0, \gamma(1) = z_1$. That is,

$$d(z_1, z_2) = \inf_{\gamma} L(\gamma) \quad \text{s.t. } \gamma(0) = z_0, \gamma(1) = z_1.$$

76 The geodesic distance can be found by solving a system of nonlinear ordinary differential equations
 77 (ODEs) defined in the intrinsic coordinates (Carmo, 1992).

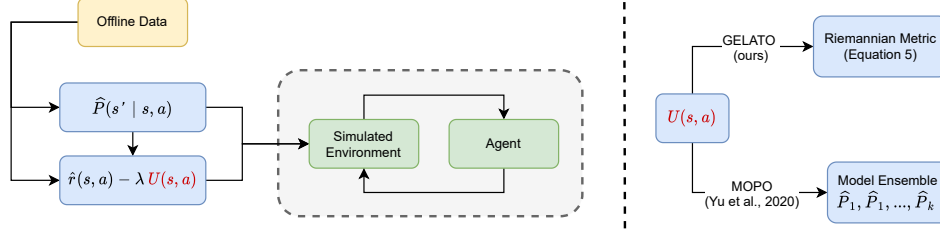


Figure 1: A reward-penalized (pessimistic) MDP is constructed from the offline data. In MOPO, the penalty is constructed using an ensemble of learned transition models. Instead, we propose to estimate the error in model dynamics through a Riemannian metric induced by a variational forward model.

78 **Pullback Metric.** Assume an ambient (observation) space \mathcal{X} and its respective Riemannian manifold
 79 $(M_{\mathcal{X}}, g_{\mathcal{X}})$. Learning $g_{\mathcal{X}}$ can be hard (e.g., learning the distance metric between images). Still, it
 80 may be captured through a low dimensional submanifold. As such, it is many times convenient to
 81 parameterize the surface $M_{\mathcal{X}}$ by a latent space $\mathcal{Z} = \mathbb{R}^{d_{\mathcal{Z}}}$ and a smooth function $f : \mathcal{Z} \mapsto \mathcal{X}$, where
 82 \mathcal{Z} is a low dimensional latent embedding space. As learning the manifold $M_{\mathcal{X}}$ can be hard, we turn
 83 to learning the immersed low dimensional submanifold $M_{\mathcal{Z}}$ (for which the chart maps are trivial,
 84 since $\mathcal{Z} = \mathbb{R}^{d_{\mathcal{Z}}}$). Given a curve $\gamma : [0, 1] \mapsto M_{\mathcal{Z}}$ we have that

$$\left\langle \frac{\partial f(\gamma(t))}{\partial t}, G_{\mathcal{X}}(f(\gamma(t))) \frac{\partial f(\gamma(t))}{\partial t} \right\rangle = \left\langle \frac{\partial \gamma(t)}{\partial t}, J_f^T(\gamma(t)) G_{\mathcal{X}}(f(\gamma(t))) J_f(\gamma(t)) \frac{\partial \gamma(t)}{\partial t} \right\rangle,$$

85 where the Jacobian matrix $J_f(z) = \frac{\partial f}{\partial z} \in \mathbb{R}^{d_{\mathcal{X}} \times d_{\mathcal{Z}}}$ maps tangent vectors in $TM_{\mathcal{Z}}$ to tangent vectors
 86 in $TM_{\mathcal{X}}$. The induced metric is thus given by

$$G_f(z) = J_f(z)^T G_{\mathcal{X}}(f(z)) J_f(z). \quad (2)$$

87 The metric G_f is known as the *pullback metric*, as it “pulls back” the metric $G_{\mathcal{X}}$ on \mathcal{X} back to G_f
 88 via $f : \mathcal{Z} \mapsto \mathcal{X}$. The pullback metric captures the intrinsic geometry of the immersed submanifold
 89 while taking into account the ambient space \mathcal{X} . The geodesic distance in ambient space is captured
 90 by geodesics in the latent space \mathcal{Z} , reducing the problem to learning the latent embedding space \mathcal{Z}
 91 and the observation function f . Indeed, learning the latent space and observation function f can be
 92 achieved through an encoder-decoder framework, such as a VAE (Arvanitidis et al., 2018).

93 3 Background: Model Error in Offline RL

94 A key element of model-based RL methods involves estimating a model $\hat{P}(s'|s, a)$. In model-based
 95 offline RL, a pessimistic MDP¹ is constructed through an upper bound on the error of the estimated
 96 model. This work builds upon MOPO, a recently proposed model-based offline RL framework (Yu
 97 et al., 2020)). Particularly, we assume access to an approximate MDP $(\mathcal{S}, \mathcal{A}, \hat{r}, \hat{P}, \alpha)$ (e.g., trained
 98 by maximizing the likelihood of the data), and define a penalized MDP $(\mathcal{S}, \mathcal{A}, \tilde{r}, \hat{P}, \alpha)$, such that

$$\tilde{r}(s, a) = \hat{r}(s, a) - \lambda d(P(\cdot|s, a), \hat{P}(\cdot|s, a)), \forall s \in \mathcal{S}, a \in \mathcal{A},$$

99 where d is a given metric (e.g., the total variation distance) and $\lambda > 0$. The offline RL problem is
 100 then solved by executing an online algorithm in the reward-penalized (simulated) MDP.

101 Unfortunately, as $P(\cdot|s, a)$ is unknown, $d(P(\cdot|s, a), \hat{P}(\cdot|s, a))$ cannot be calculated. Nevertheless,
 102 one can attempt to upper bound the distance, i.e., for some $U : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$,

$$d(P(\cdot|s, a), \hat{P}(\cdot|s, a)) \leq U(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

103 Figure 1 depicts the general framework. A question arises: how should $U(s, a)$ be chosen? In
 104 practice, MOPO learns an ensemble of models $\hat{P}_1, \dots, \hat{P}_k$ to measure the upper bound $U(s, a)$. In
 105 this work, we propose to use a naturally induced metric of a variational forward model, which we
 106 show can introduce a more effective upper bound for offline RL. In Section 4 we define this metric,
 107 and finally, we leverage it to upper bound the model error in Section 5.

¹Pessimism is a key element of offline RL algorithms (Jin et al., 2020), limiting overestimation of a trained policy due to the distribution shift between the data and the trained policy.

108 **4 Metrics of Model Dynamics**

109 We propose to measure the error in the model dynamics $d(P(\cdot|s, a), \hat{P}(\cdot|s, a))$ by embedding the
 110 offline data in a smooth Riemannian manifold, equipped with a natural metric, enabling us to measure
 111 the error of out of distribution samples. Our metric is a generalization of the metric proposed by
 112 Arvanitidis et al. (2018) for generative models.

113 **4.1 A Pullback Metric of Model Dynamics**

114 We begin by defining the immersed Riemannian submanifold and our proposed metric. The metric
 115 is defined by a latent space \mathcal{Z} and an observation function f , which will be defined later by our
 116 variational forward model.

117 **Definition 2.** We define a Riemannian submanifold $(\mathcal{M}_{\mathcal{Z}}, g_{\mathcal{Z}})$ by a differential function $f : \mathcal{Z} \mapsto \mathcal{S}$
 118 and latent space \mathcal{Z} such that

$$d_{\mathcal{Z}}(z_1, z_2) = \inf_{\gamma} \int_0^1 \left\| \frac{\partial f(\gamma(t))}{\partial t} \right\| dt \quad \text{s.t. } \gamma(0) = z_1, \gamma(1) = z_2.$$

119 A similar metric has been used in previous work on generative latent models (Chen et al., 2018;
 120 Arvanitidis et al., 2018). It states that latent codes are close w.r.t. $d_{\mathcal{Z}}$ according to the curve which
 121 induces minimal energy in ambient observation space. It is closely related to the pullback metric (see
 122 Section 2.2), as shown by the following proposition (see Appendix for proof).

123 **Proposition 1.** Let $(\mathcal{M}_{\mathcal{Z}}, g_{\mathcal{Z}})$ as defined above. Then $G_f(z) = J_f^T(z)J_f(z)$, for any $z \in \mathcal{Z}$.

124 Indeed, Proposition 1 shows us that G_f is a pullback metric. Particularly \mathcal{Z} and J_f define the
 125 structure of the ambient observation space \mathcal{S} . In what follows we characterize the submanifold $M_{\mathcal{Z}}$
 126 when $z = E(s, a)$ and $f(z) \sim P(\cdot|s, a) = P(\cdot|z)$. We show that the expected pullback metric
 127 $\mathbb{E}_{P(f)}[G_{\mathcal{Z}}(z)]$ captures notions of proximity and uncertainty and discuss how it can be utilized to
 128 measure distance to the data manifold.

129 **4.2 Metric of Proximity and Uncertainty of a Latent Forward Model**

130 We consider modeling $\hat{P}(s'|s, a)$ using a generative latent model. Specifically, we consider a latent
 131 model which consists of an encoder $E : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{B}(\mathcal{Z})$ and a decoder $D : \mathcal{Z} \mapsto \mathcal{B}(\mathcal{S})$, where
 132 $\mathcal{B}(\mathcal{X})$ is set of probability measures on the Borel sets of \mathcal{X} . While the encoder E learns a latent
 133 representation of s, a , the decoder D estimates the next state s' according to $P(\cdot|s, a)$. This model
 134 corresponds to the decomposition $P(\cdot|s, a) = D(\cdot|E(s, a))$, where here D plays the role of the
 135 observation function f , and E maps states and actions to the latent space \mathcal{Z} . Such a model can be
 136 trained by maximizing the evidence lower bound (ELBO) over the data. That is, given a prior $P(z)$,
 137 we model E_{ϕ}, D_{θ} as parametric functions and maximize the ELBO

$$\max_{\theta, \phi} \mathbb{E}_{E_{\phi}(z|s, a)} \left[\log D_{\theta}(s'|z) \right] - D_{KL}(E_{\phi}(z|s, a) || P(z))$$

138 We refer the reader to the appendix for an exhaustive overview of training VAEs by maximum
 139 likelihood and the ELBO.

140 Having trained the latent model over the data \mathcal{D}_n , we may consider the Riemannian submanifold
 141 induced by its latent space \mathcal{Z} and observation function D . Since D is stochastic, the metric $G_{\mathcal{Z}}$ also
 142 becomes stochastic, complicating analysis. Instead, Arvanitidis et al. (2018) proposed to use the
 143 expected pullback metric $\mathbb{E}[G_{\mathcal{Z}}]$, showing it is a good approximation of the underlying stochastic
 144 metric. Using Proposition 1, we have the following result (see Appendix for proof).

145 **Theorem 1.** [Arvanitidis et al. (2018)] Assume $D(\cdot|z) \sim \mathcal{N}(\mu(z), \sigma(z)I)$. Then

$$\mathbb{E}_{D(\cdot|z)} [G_D(z)] = \mathbb{E}_{D(\cdot|z)} [J_D(z)^T J_D(z)] = \underbrace{G_{\mu}(z)}_{\text{proximity}} + \underbrace{G_{\sigma}(z)}_{\text{uncertainty}}, \quad (3)$$

146 where $G_{\mu}(z) = J_{\mu}^T(z)J_{\mu}(z)$ and $G_{\sigma}(z) = J_{\sigma}^T(z)J_{\sigma}(z)$.

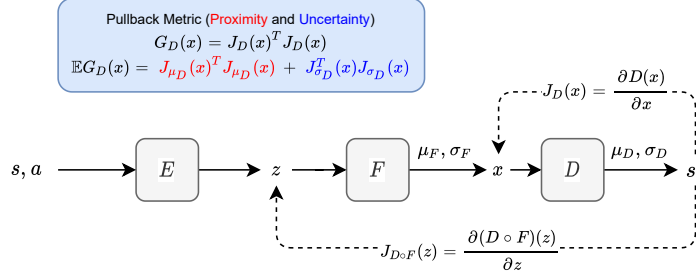


Figure 2: Plot depicts the variational latent forward model and its respective pullback metrics. Expressions for the expected pullback metrics are given in Theorems 1 and 2.

147 Given an embedded latent space \mathcal{Z} , the expected metric in Equation (3) gives us a sense of the
 148 topology of the latent space manifold induced by D . The terms $G_\mu = J_\mu^T J_\mu$ and $G_\sigma = J_\sigma^T J_\sigma$
 149 are in fact the induced pullback metrics of μ and σ , respectively. As shortest geodesics will tend to
 150 follow small values of $\|\mathbb{E}[G_D]\|$, G_μ will keep away from areas with no latent codes, whereas G_σ
 151 will remain small in regions of low uncertainty. We therefore recognize G_μ and G_σ as *metrics of*
 152 *proximity and uncertainty*, respectively.

153 **A skewed metric.** Due to the inherent decoupling between proximity and uncertainty, it may be
 154 beneficial to control the curvature of the expected metric by only focusing on one of the metrics.
 155 Denoting $\alpha_{\text{prox}} \in [0, 1]$ as the proximity coefficient, we define the skewed pullback metric of D as

$$G_D^\alpha = \alpha_{\text{prox}} G_\mu + (1 - \alpha_{\text{prox}}) G_\sigma. \quad (4)$$

156 The skewed pullback metric will become valuable in Section 6, as we carefully control the tradeoff
 157 between proximity and uncertainty in the tested domains.

158 4.3 Capturing Epistemic and Aleatoric Uncertainty

159 The previously proposed encoder-decoder model induces a metric which captures both proximity and
 160 uncertainty w.r.t. the learned dynamics. However, the decoder variance, $\sigma(z)$, does not differentiate
 161 between aleatoric uncertainty (relating to the environment dynamics) and epistemic uncertainty
 162 (relating to missing data). To bound the validity of out of distribution (OOD) samples, we wish to
 163 capture epistemic uncertainty.

164 The epistemic uncertainty can be captured by methods such as model ensembles or Monte-Carlo
 165 dropout (Gal & Ghahramani, 2016). Instead, we apply an additional forward model to our previously
 166 proposed variational model. Specifically, we assume a latent model which consists of an encoder
 167 $E : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{B}(\mathcal{Z})$, forward model $F : \mathcal{Z} \mapsto \mathcal{B}(\mathcal{X})$ and decoder $D : \mathcal{X} \mapsto \mathcal{B}(\mathcal{S})$ such that
 168 $P(\cdot|s, a) = D(\cdot|x)$, and $x \sim F(\cdot|E(s, a))$. This structure enables us to capture the aleatoric
 169 uncertainty under the forward transition model F , and the epistemic uncertainty using the decoder D .
 170 That is, $\sigma_D(z)$ is used as a measure of epistemic uncertainty, as $\sigma_F(z)$ can capture the stochasticity
 171 in model dynamics. This forward model is depicted in Figure 2.

172 Next, we turn to analyze the pullback metric induced by the proposed forward transition model. As
 173 both F and D are stochastic (capturing epistemic and aleatoric uncertainty), the result of Theorem 1
 174 cannot be directly applied to their composition. The following proposition provides an analytical
 175 expression for the expected pullback metric of $D \circ F$ (see Appendix for proof).

176 **Theorem 2.** Assume $F(\cdot|z) \sim \mathcal{N}(\mu_F(z), \sigma_F(z)I)$, $D(\cdot|x) \sim \mathcal{N}(\mu_D(x), \sigma_D(x)I)$. Then, the ex-
 177 pected pullback metric of the composite function ($D \circ F$) is given by

$$\mathbb{E}_{P(D \circ F)} [G_{D \circ F}(z)] = J_{\mu_F}^T(z) \overline{G}_D(z) J_{\mu_F}(z) + J_{\sigma_F}^T(z) \text{diag}(\overline{G}_D(z)) J_{\sigma_F}(z),$$

178 where here, $\overline{G}_D(z) = \mathbb{E}_{x \sim F(\cdot|z)} [J_{\mu_D}^T(x) J_{\mu_D}(x) + J_{\sigma_D}^T(x) J_{\sigma_D}(x)]$.

179 Unlike the metric in Equation (3), the composite metric distorts the decoder metric with Jacobian
 180 matrices of the forward model statistics. The composite metric takes into account both proximity and
 181 uncertainty w.r.t. the ambient space as well as the latent forward model. As before, a skewed version
 182 of the metric can be designed, replacing \overline{G}_D with its skewed version.

Algorithm 1 GELATO: Geometrically Enriched LATent model for Offline reinforcement learning

- 1: **Input:** Offline dataset \mathcal{D}_n , RL algorithm
 - 2: Train variational latent forward model on dataset \mathcal{D}_n by maximizing ELBO.
 - 3: Construct approximate MDP $(\mathcal{S}, \mathcal{A}, \hat{r}, \hat{P}, \alpha)$
 - 4: Define $\tilde{r}_d(s, a) = \hat{r}(s, a) - \lambda \left(\frac{1}{K} \sum_{k=1}^K d_{\mathcal{Z}}(E(s, a), \text{NN}_{E(s, a)}^{(k)}) \right)$, with distance $d_{\mathcal{Z}}$ induced by pullback metric $G_{D \circ F}$ (Theorem 2).
 - 5: Train RL algorithm over penalized MDP $(\mathcal{S}, \mathcal{A}, \tilde{r}_d, \hat{P}, \alpha)$
-

183 5 GELATO: Incorporating the Metric in Offline RL

184 Having defined our metric, we are now ready to leverage it in a model based offline RL framework
 185 (see Figure 1). Specifically, provided a dataset $\mathcal{D}_n = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$ we train the variational
 186 latent forward model depicted in Figure 2. The model consists of an encoder E , which maps states
 187 and actions to a latent space \mathcal{Z} , a forward function F which maps the latent point $E(s, a)$ to a latent
 188 point $x \sim F(\cdot|E(s, a))$, and finally a decoder which maps x to the next state $s' \sim D(\cdot|x)$. The model
 189 is trained by maximizing the likelihood of state transitions in the data (a full derivation is given in the
 190 appendix). Our latent forward model induces a latent space \mathcal{Z} and a pullback metric $G_{D \circ F}(z)$ which
 191 define the distance metric $d_{\mathcal{Z}}$ (Definition 2).

192 Algorithm 1 presents GELATO, our proposed approach. In GELATO, we estimate an upper bound,
 193 $U(s, a)$, on the model error by measuring the distance of a new sample to the data manifold. That is,

$$d(P(\cdot|s, a), \hat{P}(\cdot|s, a)) \leq U(s, a) \triangleq \left(\frac{1}{K} \sum_{k=1}^K d_{\mathcal{Z}}(E(s, a), \text{NN}_{E(s, a)}^{(k)}) \right), \quad (5)$$

194 where here, $\text{NN}_{E(s, a)}^{(k)}$ is the k^{th} nearest neighbor of $E(s, a)$ in \mathcal{D}_n w.r.t. the metric $d_{\mathcal{Z}}$. Note the sum
 195 over K nearest neighbors, allowing for more robust quantification of the distance.

196 We construct the reward-penalized MDP defined in Section 3 for which the upper bound $U(s, a)$
 197 acts as a pessimistic regularizer. Finally, we train an RL algorithm over the pessimistic MDP with
 198 transition $\hat{P}(\cdot|s, a)$ and reward $r(s, a) - \lambda U(s, a)$.

199 6 Experiments

200 6.1 Metric Visualization

201 To better understand the inherent structure of our metric, we constructed a grid-world environment
 202 for visualizing our proposed latent representation and metric. The 15×15 environment, as depicted
 203 in Figure 3, consists of four rooms, with impassable obstacles in their centers. The agent, residing
 204 at some position $(x, y) \in [-1, 1]^2$ in the environment can take one of four actions: up, down, left,
 205 or right – moving the agent 1, 2 or 3 steps (uniformly distributed) in that direction. We collected a
 206 dataset of 10000 samples, taking random actions at random initializations of the environment. The
 207 ambient state space was represented by the position of the agent – a vector of dimension 2, normalized
 208 to values in $[-1, 1]$. Finally, we trained a variational latent model with latent dimension $d_{\mathcal{Z}} = 2$.
 209 We used a standard encoder $z \sim \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}(s))$ and decoder $s' \sim \mathcal{N}(\mu_{\phi}(z), \sigma_{\phi}(z))$ represented
 210 by neural networks trained end-to-end using the evidence lower bound. We refer the reader to the
 211 appendix for an exhaustive description of the training procedure.

212 The latent space output of our model is depicted by yellow markers in Figure 3a. Indeed, the latent
 213 embedding consists of four distinctive clusters, structured in a similar manner as our grid-world
 214 environment. Interestingly, the distortion of the latent space accurately depicts an intuitive notion of
 215 distance between states. As such, rooms are distinctively separated, with fair distance between each
 216 cluster. States of pathways between rooms clearly separate the room clusters, forming a topology
 217 with four discernible bottlenecks.

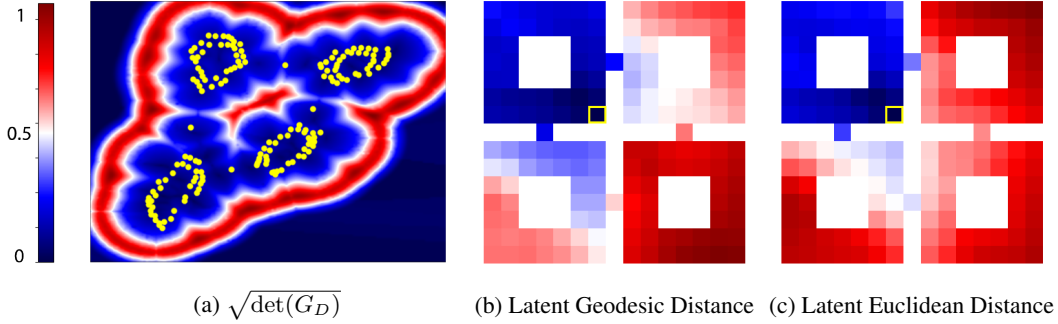


Figure 3: **(a)** The latent space (yellow markers) of the grid world environment and the geometric volume measure of the decoder-induced metric (background). **(b, c)** The geodesic distance of the decoder-induced submanifold and the Euclidean distance of latent states, as viewed in ambient space. All distances are calculated w.r.t. the yellow marked state. **Note:** colors in (a), which measure magnitude, are unrelated to colors in (b,c), which measure distance to the yellow marked state.

218 In addition to the latent embedding, Figure 3a depicts the geometric volume measure² $\sqrt{\det(G_D)}$ of
 219 the trained pullback metric induced by D . This quantity demonstrates the effective geodesic distances
 220 between states in the decoder-induced submanifold. Indeed geodesics between data points to points
 221 outside of the data manifold (i.e., outside of the red region), would attain high values as integrals over
 222 areas of high magnitude. In contrast, geodesics near the data manifold would attain low values.

223 **Comparison to Euclidean distance.** We visualize the decoder-induced geodesic distance and
 224 compare it to the latent Euclidean distance in Figures 3b and 3c, respectively. The plots depict the
 225 normalized distances of all states to the state marked by a yellow square. Evidently, the geodesic
 226 distance captures a better notion of distance in the said environment, correctly exposing the “land
 227 distance” in ambient space. As expected, states residing in the bottom-right room are farthest from
 228 the source state, as reaching them comprises of passing through at least two bottleneck states. In
 229 contrast, the latent Euclidean distance does not properly capture these geodesics, exhibiting a similar
 230 distribution of distances in other rooms. Nevertheless, both geodesic and Euclidean distances reveal
 231 the intrinsic topological structure of the environment, that of which is not captured by the extrinsic
 232 coordinates $(x, y) \in [-1, 1]^2$. Particularly, the state coordinates (x, y) would wrongly assign short
 233 distances to states across impassible walls or obstacles, i.e., measuring the “air distance”.

234 6.2 Continuous Control

235 We performed experiments to analyze GELATO on various continuous control datasets.

236 **Datasets.** We used D4RL (Fu et al., 2020) (CC BY 4.0 license) as a benchmark for all of our
 237 experiments. We tested GELATO on three Mujoco (Todorov et al., 2012) environments (Hopper,
 238 Walker2d, Halfcheetah) on datasets generated by a single policy and a mixture of two policies.
 239 Specifically, we used datasets generated by a random agent (1M samples), a partially trained agent,
 240 i.e, medium agent (1M samples), and a mixture of partially trained and expert agents (2M samples).

241 **Implementation Details.** We trained our variational model with latent dimension
 242 $\dim(\mathcal{Z}) = 32 + \dim(\mathcal{A})$. The latent model was trained for 100k steps by stochastic gradient descent
 243 with batch size of 256. We split training into two phases. First, the model was fully trained using a
 244 calibrated Gaussian decoder (Rybkin et al., 2020). Specifically, a maximum-likelihood estimate of
 245 the variance was used $\sigma^* = \text{MSE}(\mu, \hat{\mu}) \in \arg \max_{\sigma} \mathcal{N}(\hat{\mu}|\mu, \sigma^2 I)$. Then, in the second stage we fit
 246 the variance decoder network.

247 In order to practically estimate the geodesic distance in Algorithm 1, we defined a parametric curve
 248 in latent space and used gradient descent to minimize the curve’s energy. The resulting curve and
 249 pullback metric were then used to calculate the geodesic distance by a numerical estimate of the
 250 curve length (Equation (4)) (See Appendix for an exhaustive overview of the estimation method).

251 We used FAISS (Johnson et al., 2019) (MIT-license) for efficient GPU-based k -nearest neighbors
 252 calculation. We set $K = 20$ neighbors for the penalized reward (Equation (5)). Finally, we used a

²The geometric volume measure captures the volume of an infinitesimal area in the latent space.

	Hopper			Walker2d			Halfcheetah		
Method	Rand	Med	Med-Expert	Rand	Med	Med-Expert	Rand	Med	Med-Expert
Data Score	299	1021	1849	1	498	1062	-303	3945	8059
GELATO	685	1676	574	412	1269	1515	116	5168	6449
GELATO-unc	481	1158	879	290	487	1473	23	3034	7130
GELATO-prox	240	480	920	158	571	1596	-28	3300	7412
MOPO	677	1202	1063	396	518	1296	4114	4974	7594
MBPO	444	457	2105	251	370	222	3527	3228	907
SAC	664	325	1850	120	27	-2	3502	-839	-78
Imitation	615	1234	3907	47	193	329	-41	4201	4164

Table 1: Performance of GELATO and its variants in comparison to contemporary model-based methods on D4RL datasets. Scores correspond to the return, averaged over 5 seeds (standard deviation removed due to space constraints and is given in the appendix). Results for MOPO, MBPO, SAC, and imitation are taken from Yu et al. (2020). Mean score of dataset added for reference. Bold scores show an improved score w.r.t other methods.

253 variant of Soft Learning, as proposed by Yu et al. (2020) as our RL algorithm, trained for 1M steps.
254 Each experiment was run on a single GPU, RTX 2080 (see Appendix for more details).

255 **Proximity vs. Uncertainty.** To test GELATO we constructed two variants, trading off proximity
256 and uncertainty through our latent-induced metric. Particularly, we denote by GELATO-UNC and
257 GELATO-PROX variants which implement the skewed metric (see Equation (4)), with $\alpha_{\text{prox}} = 0.1$
258 and $\alpha_{\text{prox}} = 0.9$, respectively. We compared GELATO and its variants to contemporary model-based
259 offline RL approaches; namely, MOPO (Yu et al., 2020) and MBPO (Janner et al., 2019), as well as
260 the standard baselines of SAC (Haarnoja et al., 2018) and imitation (behavioral cloning).

261 Results for all of the tested domains are shown in Table 1. For the non-skewed version of GELATO
262 (i.e., $\alpha_{\text{prox}} = 0.5$) we found performance increase on most domains, and most significantly in the
263 medium domain, i.e., partially trained agent. We believe this to be due to the inherent nature of our
264 metric to take into account both proximity and uncertainty, allowing the agent to leverage proximity
265 to the data in areas of high uncertainty. Since the medium dataset contained average behavior, mixing
266 proximity benefited the agent’s overall performance.

267 In most of the tested datasets we found an increase in performance for at least one of the GELATO
268 variants. The med-expert datasets showed better performance for the proximity-oriented metric.
269 These suggest flexibility of our metric to increase performance when the quality of the dataset is
270 known, a reasonable assumption in many domains. Moreover, the non-skewed version of GELATO,
271 showed consistency over all datasets, favorably leveraging the strengths of proximity and uncertainty.

272 6.3 RBF Networks.

273 A question arises as to how to represent $\sigma_D(z)$. In general, neural networks may result in a poor
274 measure of uncertainty, due to uncontrolled extrapolations of the neural network to arbitrary regions
275 of the latent space, i.e., areas of little data. However, Arvanitidis et al. (2017) showed that the inverse
276 variance $\beta(z) = (\sigma^2(z))^{-1}$ with a positive Radial Basis Function (RBF) network achieves a reliable
277 uncertainty estimate, with well-behaved extrapolations in latent space. Formally the RBF network is
278 defined by $\sigma(z) = \sqrt{(\beta(z))^{-1}}$ where $\beta(z) = W^T \phi(z)$, $\phi_i(z) = \exp\left(-\frac{1}{2}\lambda_i \|z - c_i\|_2^2\right)$, W are the
279 positive learned weights of the network, λ_i the bandwidth, and c_i the centers trained using k -means
280 over the learned latent representations of the offline data.

281 We tested GELATO on D4RL Mujoco benchmarks with an RBF decoder network. Specifically, we
282 followed a similar training procedure with two training phases. In the first training phase we trained
283 our variational model with a calibrated decoder as before. In the second training phase we used
284 k -means to cluster our dataset to 128 clusters, after which an RBF network was trained for a second
285 phase of 50000 iterations.

	Hopper		
Method	Random	Medium	Med-Expert
GELATO	685 ± 15	1676 ± 223	574 ± 16
GELATO-RBF	613 ± 24	1700 ± 319	498 ± 55

Table 2: Performance of GELATO using an RBF decoder compared to standard decoder.

286 Results for GELATO with RBF decoder networks for the Hopper environment are presented in
287 Table 2. We did not find significant improvement in using RBF networks over decoder variance. We
288 believe this is due to the smoothness in ambient space, allowing for well behaved extrapolations of
289 uncertainty. We conjecture RBF networks may show improved performance on higher dimensional
290 problems (e.g., images), yet we leave this for future work, as these may involve utilizing more
291 involved variational models (Vahdat & Kautz, 2020).

292 7 Related Work

293 **Offline Reinforcement Learning.** The field of offline RL has recently received much attention as
294 several algorithmic approaches were able to surpass standard off-policy algorithms. Value-based
295 online algorithms do not perform well under highly off-policy batch data (Fujimoto et al., 2019;
296 Kumar et al., 2019; Fu et al., 2019; Fedus et al., 2020; Agarwal et al., 2020), much due to issues with
297 bootstrapping from out-of-distribution (OOD) samples. These issues become more prominent in the
298 offline setting, as new samples cannot be acquired.

299 Several works on offline RL have shown improved performance on standard continuous control
300 benchmarks (Laroche et al., 2019; Kumar et al., 2019; Fujimoto et al., 2019; Chen et al., 2020b;
301 Swazinna et al., 2020; Kidambi et al., 2020; Yu et al., 2020; Kumar et al., 2020). In this work we were
302 specifically interested in model-based approaches (Yu et al., 2020; Kidambi et al., 2020), in which
303 the agent is incentivized to remain close to areas of low uncertainty. Our work focused on controlling
304 uncertainty estimation in high dimensional environments. Our methodology utilized recent success
305 on the geometry of deep generative models (Arvanitidis et al., 2018, 2020), proposing an alternative
306 approach to uncertainty estimation.

307 **Representation Learning.** Representation learning seeks to find an appropriate representation of
308 data for performing a machine-learning task (Goodfellow et al., 2016). Variational Auto Encoders
309 (Kingma & Welling, 2013; Rezende et al., 2014) have been a popular representation learning
310 technique, particularly in unsupervised learning regimes (Kingma et al., 2014; Sønderby et al.,
311 2016; Chen et al., 2016; Van Den Oord et al., 2017; Hsu et al., 2017; Serban et al., 2017; Engel
312 et al., 2017; Bojanowski et al., 2018; Ding et al., 2020), though also in supervised learning and
313 reinforcement learning (Hausman et al., 2018; Li et al., 2019; Petangoda et al., 2019; Hafner et al.,
314 2019, 2020). Particularly, variational models have been shown able to derive successful behaviors in
315 high dimensional benchmarks (Hafner et al., 2020).

316 Various representation techniques in reinforcement learning have also proposed to disentangle
317 representation of both states (Engel & Mannor, 2001; Littman & Sutton, 2002; Stooke et al., 2020;
318 Zhu et al., 2020), and actions (Tennenholtz & Mannor, 2019; Chandak et al., 2019). These allow
319 for the abstraction of states and actions to significantly decrease computation requirements, by e.g.,
320 decreasing the effective dimensionality of the action space (Tennenholtz & Mannor, 2019). Unlike
321 previous work, GELATO is focused on measuring proximity and uncertainty for the purpose of
322 mitigating the OOD problem and enhancing offline reinforcement learning performance.

323 8 Discussion and Future Work

324 This work presented a metric for model dynamics and its application to offline reinforcement learning.
325 While our metric showed supportive evidence of improvement in model-based offline RL we note
326 that it was significantly slower – comparably, 5 times slower than using the decoder’s variance for
327 uncertainty estimation. The apparent slowdown in performance was mostly due to computation of the
328 geodesic distance. Improvement in this area may utilize techniques for efficient geodesic estimation
329 (Chen et al., 2018, 2019).

330 We conclude by noting possible future applications of our work. In Section 6.1 we demonstrated the
331 inherent geometry our model had captured, its corresponding metric, and geodesics. Still, in this
332 work we focused specifically on metrics related to the decoded state. In fact, a similar derivation
333 to Theorem 2 could be applied to other modeled statistics, e.g., Q-values, rewards, future actions,
334 and more. Each distinct statistic would induce its own unique metric w.r.t. its respective probability
335 measure. Particularly, this concept may benefit a vast array of applications in continuous or large
336 state and action spaces.

337 **References**

- 338 Agarwal, R., Schuurmans, D., and Norouzi, M. An optimistic perspective on offline reinforcement
339 learning. In *International Conference on Machine Learning*, pp. 104–114. PMLR, 2020.
- 340 Arvanitidis, G., Hansen, L. K., and Hauberg, S. Maximum likelihood estimation of riemannian
341 metrics from euclidean data. In *International Conference on Geometric Science of Information*, pp.
342 38–46. Springer, 2017.
- 343 Arvanitidis, G., Hansen, L. K., and Hauberg, S. Latent space oddity: On the curvature of deep
344 generative models. In *6th International Conference on Learning Representations, ICLR 2018*,
345 2018.
- 346 Arvanitidis, G., Hauberg, S., and Schölkopf, B. Geometrically enriched latent spaces. *arXiv preprint*
347 *arXiv:2008.00565*, 2020.
- 348 Bojanowski, P., Joulin, A., Lopez-Pas, D., and Szlam, A. Optimizing the latent space of generative
349 networks. In *International Conference on Machine Learning*, pp. 600–609, 2018.
- 350 Buckman, J., Gelada, C., and Bellemare, M. G. The importance of pessimism in fixed-dataset policy
351 optimization. *arXiv preprint arXiv:2009.06799*, 2020.
- 352 Carmo, M. P. d. *Riemannian geometry*. Birkhäuser, 1992.
- 353 Chandak, Y., Theodorou, G., Kostas, J., Jordan, S., and Thomas, P. Learning action representations
354 for reinforcement learning. In *International Conference on Machine Learning*, pp. 941–950, 2019.
- 355 Chen, N., Klushyn, A., Kurle, R., Jiang, X., Bayer, J., and Smagt, P. Metrics for deep generative
356 models. In *International Conference on Artificial Intelligence and Statistics*, pp. 1540–1550.
357 PMLR, 2018.
- 358 Chen, N., Ferroni, F., Klushyn, A., Paraschos, A., Bayer, J., and van der Smagt, P. Fast approximate
359 geodesics for deep generative models. In *International Conference on Artificial Neural Networks*,
360 pp. 554–566. Springer, 2019.
- 361 Chen, N., Klushyn, A., Ferroni, F., Bayer, J., and van der Smagt, P. Learning flat latent manifolds
362 with vaes. 2020a.
- 363 Chen, X., Kingma, D. P., Salimans, T., Duan, Y., Dhariwal, P., Schulman, J., Sutskever, I., and Abbeel,
364 P. Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*, 2016.
- 365 Chen, X., Zhou, Z., Wang, Z., Wang, C., Wu, Y., and Ross, K. Bail: Best-action imitation learning
366 for batch deep reinforcement learning. *Advances in Neural Information Processing Systems*, 33,
367 2020b.
- 368 Ding, Z., Xu, Y., Xu, W., Parmar, G., Yang, Y., Welling, M., and Tu, Z. Guided variational autoencoder
369 for disentanglement learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
370 *and Pattern Recognition*, pp. 7920–7929, 2020.
- 371 Dinh, L., Krueger, D., and Bengio, Y. Nice: Non-linear independent components estimation. *arXiv*
372 *preprint arXiv:1410.8516*, 2014.
- 373 Engel, J., Hoffman, M., and Roberts, A. Latent constraints: Learning to generate conditionally from
374 unconditional generative models. *arXiv preprint arXiv:1711.05772*, 2017.
- 375 Engel, Y. and Mannor, S. Learning embedded maps of markov processes. In *Proceedings of the*
376 *Eighteenth International Conference on Machine Learning*, pp. 138–145, 2001.
- 377 Ernst, D., Geurts, P., and Wehenkel, L. Tree-based batch mode reinforcement learning. *Journal of*
378 *Machine Learning Research*, 6(Apr):503–556, 2005.
- 379 Fedus, W., Ramachandran, P., Agarwal, R., Bengio, Y., Larochelle, H., Rowland, M., and Dabney, W.
380 Revisiting fundamentals of experience replay. In *International Conference on Machine Learning*,
381 pp. 3061–3071. PMLR, 2020.

- 382 Fonteneau, R., Murphy, S. A., Wehenkel, L., and Ernst, D. Batch mode reinforcement learning based
383 on the synthesis of artificial trajectories. *Annals of operations research*, 208(1):383–416, 2013.
- 384 Fu, J., Kumar, A., Soh, M., and Levine, S. Diagnosing bottlenecks in deep q-learning algorithms. In
385 *International Conference on Machine Learning*, pp. 2021–2030, 2019.
- 386 Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4rl: Datasets for deep data-driven
387 reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 388 Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration.
389 In *International Conference on Machine Learning*, pp. 2052–2062. PMLR, 2019.
- 390 Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty
391 in deep learning. In *international conference on machine learning*, pp. 1050–1059, 2016.
- 392 Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT press
393 Cambridge, 2016.
- 394 Guss, W. H., Houghton, B., Topin, N., Wang, P., Codel, C., Veloso, M., and Salakhutdinov, R.
395 MineRL: A large-scale dataset of Minecraft demonstrations. *Twenty-Eighth International Joint*
396 *Conference on Artificial Intelligence*, 2019. URL <http://miner1.io>.
- 397 Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy
398 deep reinforcement learning with a stochastic actor. In *International Conference on Machine*
399 *Learning*, pp. 1861–1870. PMLR, 2018.
- 400 Hafner, D., Lillicrap, T., Ba, J., and Norouzi, M. Dream to control: Learning behaviors by latent
401 imagination. In *International Conference on Learning Representations*, 2019.
- 402 Hafner, D., Lillicrap, T., Norouzi, M., and Ba, J. Mastering atari with discrete world models. *arXiv*
403 *preprint arXiv:2010.02193*, 2020.
- 404 Hausman, K., Springenberg, J. T., Wang, Z., Heess, N., and Riedmiller, M. Learning an embedding
405 space for transferable robot skills. In *International Conference on Learning Representations*, 2018.
- 406 Hsu, W.-N., Zhang, Y., and Glass, J. Unsupervised learning of disentangled and interpretable
407 representations from sequential data. In *Advances in neural information processing systems*, pp.
408 1878–1889, 2017.
- 409 Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust your model: Model-based policy
410 optimization. *arXiv preprint arXiv:1906.08253*, 2019.
- 411 Jin, Y., Yang, Z., and Wang, Z. Is pessimism provably efficient for offline rl? *arXiv preprint*
412 *arXiv:2012.15085*, 2020.
- 413 Johnson, J., Douze, M., and Jégou, H. Billion-scale similarity search with gpus. *IEEE Transactions*
414 *on Big Data*, 2019.
- 415 Kalatzis, D., Eklund, D., Arvanitidis, G., and Hauberg, S. Variational autoencoders with riemannian
416 brownian motion priors. In *Proceedings of the 37th International Conference on Machine Learning*
417 *(ICML)*. PMLR, 2020.
- 418 Kidambi, R., Rajeswaran, A., Netrapalli, P., and Joachims, T. Morel: Model-based offline reinforce-
419 ment learning. *arXiv preprint arXiv:2005.05951*, 2020.
- 420 Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*
421 *arXiv:1412.6980*, 2014.
- 422 Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,
423 2013.
- 424 Kingma, D. P., Mohamed, S., Jimenez Rezende, D., and Welling, M. Semi-supervised learning with
425 deep generative models. *Advances in neural information processing systems*, 27:3581–3589, 2014.

- 426 Klushyn, A., Chen, N., Kurle, R., Cseke, B., and van der Smagt, P. Learning hierarchical priors in
427 vaes. In *Advances in Neural Information Processing Systems*, pp. 2870–2879, 2019.
- 428 Kumar, A., Fu, J., Soh, M., Tucker, G., and Levine, S. Stabilizing off-policy q-learning via bootstrap-
429 ping error reduction. In *Advances in Neural Information Processing Systems*, pp. 11784–11794,
430 2019.
- 431 Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative q-learning for offline reinforcement
432 learning. *arXiv preprint arXiv:2006.04779*, 2020.
- 433 Larocche, R., Trichelair, P., and Des Combes, R. T. Safe policy improvement with baseline bootstrap-
434 ping. In *International Conference on Machine Learning*, pp. 3652–3661. PMLR, 2019.
- 435 Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and
436 perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 437 Li, M., Wu, L., Jun, W., and Ammar, H. B. Multi-view reinforcement learning. In *Advances in neural
438 information processing systems*, pp. 1420–1431, 2019.
- 439 Littman, M. L. and Sutton, R. S. Predictive representations of state. In *Advances in neural information
440 processing systems*, pp. 1555–1561, 2002.
- 441 Petangoda, J. C., Pascual-Diaz, S., Adam, V., Vrancx, P., and Grau-Moya, J. Disentangled skill
442 embeddings for reinforcement learning. *arXiv preprint arXiv:1906.09223*, 2019.
- 443 Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference
444 in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- 445 Riedmiller, M. Neural fitted q iteration—first experiences with a data efficient neural reinforcement
446 learning method. In *European Conference on Machine Learning*, pp. 317–328. Springer, 2005.
- 447 Rybkin, O., Daniilidis, K., and Levine, S. Simple and effective vae training with calibrated decoders.
448 *arXiv preprint arXiv:2006.13202*, 2020.
- 449 Schaal, S. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3(6):
450 233–242, 1999.
- 451 Senge, R., Bösner, S., Dembczyński, K., Haasenritter, J., Hirsch, O., Donner-Banzhoff, N., and
452 Hüllermeier, E. Reliable classification: Learning classifiers that distinguish aleatoric and epistemic
453 uncertainty. *Information Sciences*, 255:16–29, 2014.
- 454 Serban, I., Sordani, A., Lowe, R., Charlin, L., Pineau, J., Courville, A., and Bengio, Y. A hierarchical
455 latent variable encoder-decoder model for generating dialogues. In *Proceedings of the AAAI
456 Conference on Artificial Intelligence*, volume 31, 2017.
- 457 Sønderby, C. K., Raiko, T., Maaløe, L., Sønderby, S. K., and Winther, O. Ladder variational
458 autoencoders. In *Advances in neural information processing systems*, pp. 3738–3746, 2016.
- 459 Stooke, A., Lee, K., Abbeel, P., and Laskin, M. Decoupling representation learning from reinforce-
460 ment learning. *arXiv preprint arXiv:2009.08319*, 2020.
- 461 Sutton, R. S., Barto, A. G., et al. *Introduction to reinforcement learning*, volume 135. MIT press
462 Cambridge, 1998.
- 463 Swazinna, P., Udluft, S., and Runkler, T. Overcoming model bias for robust offline deep reinforcement
464 learning. *arXiv preprint arXiv:2008.05533*, 2020.
- 465 Tennenholtz, G. and Mannor, S. The natural language of actions. In *International Conference on
466 Machine Learning*, pp. 6196–6205, 2019.
- 467 Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012
468 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE,
469 2012.

- 470 Vahdat, A. and Kautz, J. Nvae: A deep hierarchical variational autoencoder. *arXiv preprint*
471 *arXiv:2007.03898*, 2020.
- 472 Van Den Oord, A., Vinyals, O., et al. Neural discrete representation learning. In *Advances in Neural*
473 *Information Processing Systems*, pp. 6306–6315, 2017.
- 474 Wang, R., Foster, D. P., and Kakade, S. M. What are the statistical limits of offline rl with linear
475 function approximation? *arXiv preprint arXiv:2010.11895*, 2020.
- 476 Yu, T., Thomas, G., Yu, L., Ermon, S., Zou, J., Levine, S., Finn, C., and Ma, T. Mopo: Model-based
477 offline policy optimization. *arXiv preprint arXiv:2005.13239*, 2020.
- 478 Zanette, A. Exponential lower bounds for batch reinforcement learning: Batch rl can be exponentially
479 harder than online rl. *arXiv preprint arXiv:2012.08005*, 2020.
- 480 Zhu, J., Xia, Y., Wu, L., Deng, J., Zhou, W., Qin, T., and Li, H. Masked contrastive representation
481 learning for reinforcement learning. *arXiv preprint arXiv:2010.07470*, 2020.

482 **Checklist**

- 483 1. For all authors...
- 484 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
485 contributions and scope? [Yes]
- 486 (b) Did you describe the limitations of your work? [Yes] computational limitation (discus-
487 sion section).
- 488 (c) Did you discuss any potential negative societal impacts of your work? [Yes] broader
489 impact in introduction section
- 490 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
491 them? [Yes]
- 492 2. If you are including theoretical results...
- 493 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 494 (b) Did you include complete proofs of all theoretical results? [Yes]
- 495 3. If you ran experiments...
- 496 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
497 mental results (either in the supplemental material or as a URL)? [Yes]
- 498 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
499 were chosen)? [Yes]
- 500 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
501 ments multiple times)? [Yes] standard deviation in Table 1
- 502 (d) Did you include the total amount of compute and the type of resources used (e.g., type
503 of GPUs, internal cluster, or cloud provider)? [Yes]
- 504 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 505 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 506 (b) Did you mention the license of the assets? [Yes]
- 507 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 508 (d) Did you discuss whether and how consent was obtained from people whose data you're
509 using/curating? [No] FAISS is under MIT license. Our usage of D4RL dataset does
510 not require consent according to CC BY 4.0 license.
- 511 (e) Did you discuss whether the data you are using/curating contains personally identifiable
512 information or offensive content? [No] Data is of standard Mujoco benchmarks
- 513 5. If you used crowdsourcing or conducted research with human subjects...
- 514 (a) Did you include the full text of instructions given to participants and screenshots, if
515 applicable? [N/A]
- 516 (b) Did you describe any potential participant risks, with links to Institutional Review
517 Board (IRB) approvals, if applicable? [N/A]
- 518 (c) Did you include the estimated hourly wage paid to participants and the total amount
519 spent on participant compensation? [N/A]

Method	Hopper			Walker2d			Halfcheetah		
	Random	Medium	Med-Expert	Random	Medium	Med-Expert	Random	Medium	Med-Expert
Data Score	299 ± 200	1021 ± 314	1849 ± 1560	1 ± 6	498 ± 807	1062 ± 1576	-303 ± 79	3945 ± 494	8059 ± 4204
GELATO	685 ± 15	1676 ± 223	574 ± 16	412 ± 85	1269 ± 549	1515 ± 379	116 ± 43	5168 ± 849	6449 ± 2790
GELATO-unc	481 ± 29	1158 ± 423	879 ± 153	290 ± 79	487 ± 289	1473 ± 389	23 ± 35	3034 ± 585	7130 ± 3780
GELATO-prox	240 ± 22	480 ± 15	920 ± 249	158 ± 35	571 ± 326	1596 ± 416	-28 ± 31	3300 ± 613	7412 ± 3166
MOPO	677 ± 13	1202 ± 400	1063 ± 193	396 ± 76	518 ± 560	1296 ± 374	4114 ± 312	4974 ± 200	7594 ± 4741
MBPO	444 ± 193	457 ± 106	2105 ± 1113	251 ± 235	370 ± 221	222 ± 99	3527 ± 487	3228 ± 2832	907 ± 1185
SAC	664	325	1850	120	27	-2	3502	-839	-78
Imitation	615	1234	3907	47	193	329	-41	4201	4164

Table 3: Results of GELATO as presented in Table 1 with added std for each run, averaged over 5 seeds.

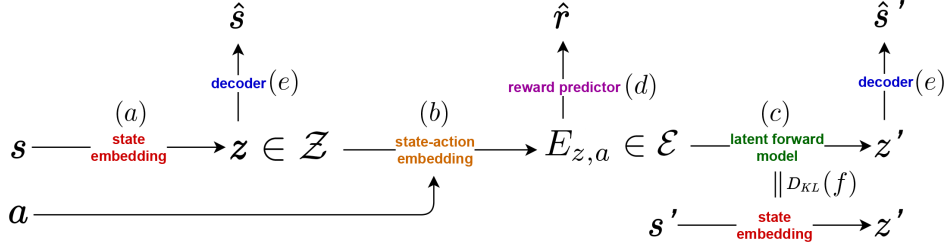


Figure 4: A graphical representation of our latent variable model. **(a)** The state s is embedded via the state embedding function (i.e., approximate posterior) $z \sim q(\cdot|s)$. **(b)** The action and embedded state pass through an invertible embedding function E to produce the state-action embedding $E_{z,a}$. **(c,d)** The state-action embedding is passed through a reward predictor and latent forward model, $\hat{r} \sim P(\cdot|E_{z,a})$ and $z' \sim P(\cdot|z, a)$, respectively. **(e)** The next latent state z' is decoded back to observation space to generate $\hat{s}' \sim P(\cdot|z')$. **(f)** Finally, during training, the target state s' is embedded and compared to z' (by the KL-divergence term in Equation (7)), preserving the consistency of the latent space \mathcal{Z} .

520 Appendix

521 9 Variational Latent Model

522 We begin by describing our variational forward model. The model, based on an encoder, latent forward
523 function, and decoder framework assumes the underlying dynamics and reward are governed by a
524 state-embedded latent space $\mathcal{Z} \subseteq \mathbb{R}^{d_z}$. The probability of a trajectory $\tau = (s_0, a_0, r_0, \dots, s_h, a_h, r_h)$
525 is given by

$$P(\tau) = \int_{z_0, \dots, z_h} P(z_0) \prod_{i=0}^h P(s_i|z_i) \pi(a_i|s_i) P(r_i|E_{z_i, a_i}) \prod_{j=1}^h P(z_j|E_{z_{j-1}, a_{j-1}}) dz_0 \dots dz_h, \quad (6)$$

526 where $E : \mathcal{Z} \times \mathcal{A} \mapsto \mathcal{E} \subseteq \mathbb{R}^{d_{\mathcal{E}}}$ is a deterministic, invertible embedding function which maps pairs
527 (z, a) to a state-action-embedded latent space \mathcal{E} . $E_{z,a}$ is thus a sufficient statistic of (z, a) . The
528 proposed graphical model is depicted in Figure 4. We note that an extension to the partially observable
529 setting replaces s_t with $h_t = (s_0, a_0, \dots, s_t)$, a sufficient statistic of the unknown state.

530 Maximizing the log-likelihood $\log P(\tau)$ is hard due to intractability of the integral in Equation (6).
531 We therefore introduce the approximate posterior $q(z|s)$ and maximize the evidence lower bound.

532 To clear notations we define $E_{z_{-1}, a_{-1}} = 0$, so that we can rewrite the above expression as

$$P(s_0, a_0, r_0, \dots, s_h, a_h, r_h) = \int_{z_0, \dots, z_h} \prod_{i=0}^h P(s_i|z_i) \pi(a_i|s_i) P(r_i|E_{z_i, a_i}) P(z_j|E_{z_{-1}, a_{-1}})$$

533 Introducing $q(z_i | s_i)$ we can write

$$\begin{aligned}
& \log \int_{z_0, \dots, z_h} \prod_{i=0}^h \frac{q(z_i | s_i)}{q(z_i | s_i)} \prod_{i=0}^h P(s_i | z_i) \pi(a_i | s_i) P(r_i | E_{z_i, a_i}) P(z_j | E_{z_{-1}, a_{-1}}) \\
& \geq \int_{z_0, \dots, z_h} \prod_{i=0}^h q(z_i | s_i) \left(\sum_{i=0}^h \log \left(\frac{P(s_i | z_i) \pi(a_i | s_i) P(r_i | E_{z_i, a_i}) P(z_j | E_{z_{-1}, a_{-1}})}{q(z_i | s_i)} \right) \right) \\
& = \sum_{i=0}^H \mathbb{E}_{q(z_i | s_i)} \left[\log(P(s_i | z_i) \pi(a_i | s_i) P(r_i | E_{z_i, a_i})) \right] \\
& \quad - \sum_{i=0}^{H-1} \mathbb{E}_{q(z_i | s_i)} \left[D_{KL}(q(z_{i+1} | s_{i+1}) || P(z_{i+1} | E_{z_i, a_i})) \right] \\
& \quad - D_{KL}(q(z_0 | s_0) || P(z_0)).
\end{aligned}$$

534 Hence,

$$\begin{aligned}
& \sum_{i=0}^h \mathbb{E}_{z_i \sim q(z_i | s_i)} \left[\log(P(s_i | z_i) \pi(a_i | s_i) P(r_i | E_{z_i, a_i})) \right] \\
& \quad - \sum_{i=0}^{h-1} \mathbb{E}_{z_i \sim q(z_i | s_i)} \left[D_{KL}(q(z_{i+1} | s_{i+1}) || P(z_{i+1} | E_{z_i, a_i})) \right] \\
& \quad - D_{KL}(q(z_0 | s_0) || P(z_0)). \tag{7}
\end{aligned}$$

535 The distribution parameters of the approximate posterior $q(z|s)$, the likelihoods
536 $P(s|z)$, $\pi(a|s)$, $P(r|E_{z,a})$, and the latent forward model $P(z'|E_{z,a})$ are represented by neu-
537 ral networks. The invertible embedding function E is represented by an invertible neural network,
538 e.g., affine coupling, commonly used for normalizing flows (Dinh et al., 2014). Though various latent
539 distributions have been proposed (Klushyn et al., 2019; Kalatzis et al., 2020), we found Gaussian
540 parametric distributions to suffice for all of our model’s functions. Particularly, we used two outputs
541 for every distribution, representing the expectation μ and variance σ . All networks were trained
542 end-to-end to maximize the evidence lower bound in Equation (7).

543 Our latent variable model is designed to capture both the epistemic and aleatoric uncertainty (Senge
544 et al., 2014). The variance output of the decoder captures epistemic uncertainty, while stochasticity
545 of the latent forward model $P(z'|E_{z,a})$ captures aleatoric uncertainty. For the purpose of offline RL,
546 we will focus on the epistemic uncertainty of our model.

547 We tested the quality of our variational model on datasets of two tasks in Minecraft (Guss et al.,
548 2019); namely, a navigation task (150k examples) and a tree chopping task (250k examples), both
549 generated by human players. The variational model was trained only on the navigation task. We

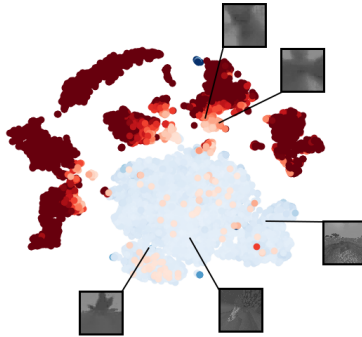


Figure 5: TSNE projection of latent space \mathcal{Z} for navigation dataset (blue) and tree chopping dataset (red) in Minecraft (Guss et al., 2019). Darker colors correspond to higher decoder variance.

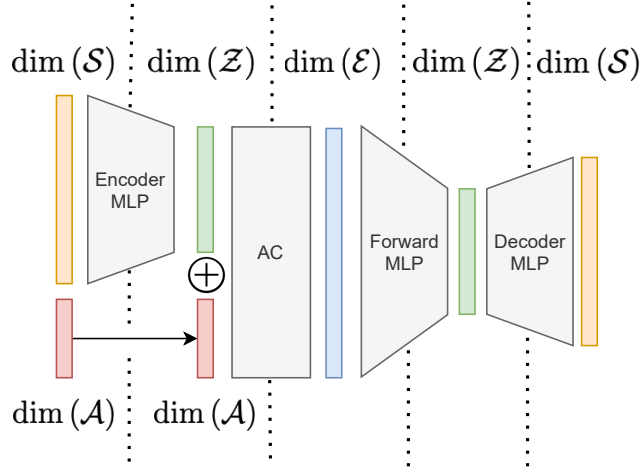


Figure 6: Latent model architecture (does not depict reward MLP).

550 embedded the data from both datasets using our trained model, and measured the decoder variance for
 551 all samples. Figure 5 depicts a TSNE projection of the latent space \mathcal{Z} , coloring in blue the navigation
 552 task and in red the tree chopping task. Light colors correspond to low variance (i.e., sharp images),
 553 whereas dark colors correspond to large variance (i.e. OOD samples). We found that our variational
 554 model was able to properly distinguish between the two tasks, with some overlap due to similarity
 555 in state space features. Additionally, we noticed a clear transition in decoding variance as samples
 556 farther away from the trained latent data attained larger variance, suggesting our variational model
 557 was properly able to distinguish OOD samples.

558 We refer the reader to the appendix for further analysis and approaches of uncertainty quantification
 559 in variational models. In our experiments, we found that the standard decoder variance sufficed for
 560 all of the tested domains.

561 10 Specific Implementation Details

562 As a preprocessing step rewards were normalized to values between $[-1, 1]$. We trained our variational
 563 model with latent dimensions $\dim(\mathcal{Z}) = 32$ and $\dim(\mathcal{E}) = \dim(\mathcal{Z}) + \dim(\mathcal{A})$. All domains were
 564 trained with the same hyperparameters. Specifically, we used a 2-layer Multi Layer Perceptron (MLP)
 565 to encode \mathcal{Z} , after which a 2-layer Affine Coupling (AC) (Dinh et al., 2014) was used to encode \mathcal{E} .
 566 We also used a 2-layer MLP for the forward, reward, and decoder models. All layers contained 256
 567 hidden layers.

568 The latent model was trained in two separate phases for 100k and 50k steps each by stochastic gradient
 569 descent and the ADAM optimizer (Kingma & Ba, 2014). First, the model was fully trained using a
 570 calibrated Gaussian decoder (Rybkin et al., 2020). Specifically, a maximum-likelihood estimate of
 571 the variance was used $\sigma^* = \text{MSE}(\mu, \hat{\mu}) \in \arg \max_{\sigma} \mathcal{N}(\hat{\mu} | \mu, \sigma^2 I)$. Finally, in the second stage we
 572 fit the variance decoder network. We found this process of to greatly improve convergence speed
 573 and accuracy, and mitigate posterior collapse. We used a minimum variance of 0.01 for all of our
 574 stochastic models.

575 To further stabilize training we used a momentum encoder. Specifically we updated a target encoder
 576 as a slowly moving average of the weights from the learned encoder as

$$\bar{\theta} \leftarrow (1 - \tau)\bar{\theta} + \tau\theta$$

577 Hyperparameters for our variational model are summarized in Table 4. The latent architecture is
 578 visualized in Figure 6.

Parameter	Value	Parameter	Value
$\dim(\mathcal{Z})$	32	LEARNING RATE	10^{-3}
$\dim(\mathcal{E})$	$32 + \dim(\mathcal{A})$	BATCH SIZE	128
ENCODER MLP HIDDEN	256, 256	TARGET UPDATE τ	0.01
FORWARD MLP HIDDEN	256, 256	TARGET UPDATE INTERVAL	1
DECODER MLP HIDDEN	256, 256	PHASE 1 UPDATES	100000
REWARD MLP HIDDEN	256, 256	PHASE 2 UPDATES	50000

Table 4: Hyper parameters for variational model

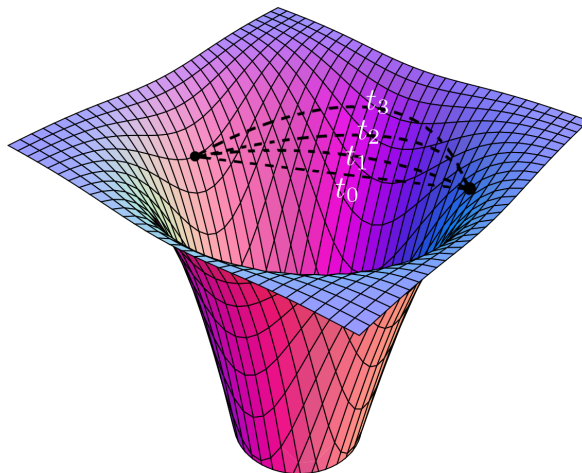


Figure 7: Illustration of geodesic curve optimization in Algorithm 2.

579 10.1 Geodesic Distance Estimation

580 In order to practically estimate the geodesic distance between two points $e_1, e_2 \in \mathcal{E}$ we defined a
581 parametric curve in latent space and used gradient descent to minimize the curve’s energy³. The
582 resulting curve and pullback metric were then used to calculate the geodesic distance by a numerical
583 estimate of the curve length (Equation (4)).

584 Pseudo code for Geodesic Distance Estimation is shown in Algorithm 2. Our curve was modeled
585 as a cubic spline with 8 coefficients. We used SGD (momentum 0.99) to optimize the curve energy
586 over 20 gradient iterations with a grid of 10 points and a learning rate of 10^{-3} . An illustration of the
587 convergence of such a curve is illustrated in Figure 7

588 10.2 RL algorithm

589 Our learning algorithm is based on the Soft Learning framework proposed in Algorithm 2 of Yu et al.
590 (2020). Pseudo code is shown in Algorithm 3. Specifically we used two replay buffers $\mathcal{R}_{\text{model}}, \mathcal{R}_{\text{data}}$,
591 where $|\mathcal{R}_{\text{model}}| = 50000$ and $\mathcal{R}_{\text{data}}$ contained the full offline dataset. In every epoch an initial state
592 s_0 was sampled from the offline dataset and embedded using our latent model to generate $z_0 \in \mathcal{Z}$.
593 During rollouts of π , embeddings $E_{z,a} \in \mathcal{E}$ were then generated from z and used to (1) sample
594 next latent state z' , (2) sample estimated rewards r , and (3) compute distances to $K = 20$ nearest
595 neighbors in embedded the dataset.

596 We used Algorithm 2 to compute the geodesic distances, and FAISS (Johnson et al., 2019) for
597 efficient nearest neighbor computation on GPUs. To stabilize learning, we normalized the penalty
598 $\frac{1}{K} \sum_{k=1}^K d_k$ according to the maximum penalty, ensuring penalty lies in $[0, 1]$ (recall that the latent
599 reward predictor was trained over normalized rewards in $[-1, 1]$). For the non-skewed version of

³Other methods for computing the geodesic distance include solving a system of ODEs (Arvanitidis et al., 2018), using graph based geodesics (Chen et al., 2019), or using neural networks (Chen et al., 2018).

Algorithm 2 Geodesic Distance Estimation

Input: forward latent F , decoder D , learning rate λ , number of iterations T , grid size n , eval points e_0, e_1
Initialize: parametric curve $\gamma_\theta : \gamma_\theta(0) = e_0, \gamma_\theta(1) = e_1$
for $t = 1$ **to** T **do**
 $L_\mu(\theta) \leftarrow \sum_{i=1}^n \mu_D(\mu_F(\gamma_\theta(\frac{i}{n}))) - \mu_D(\mu_F(\gamma_\theta(\frac{i-1}{n})))$
 $L_\sigma(\theta) \leftarrow \sum_{i=1}^n \sigma_D(\mu_F(\gamma_\theta(\frac{i}{n}))) - \sigma_D(\mu_F(\gamma_\theta(\frac{i-1}{n})))$
 $L(\theta) \leftarrow L_\mu(\theta) + L_\sigma(\theta)$
 $\theta \leftarrow \theta - \lambda \nabla_\theta L(\theta)$
end for
 $G_{D \circ F} = J_{\mu_F}^T \bar{G}_D J_{\mu_F} + J_{\sigma_F}^T \text{diag}(\bar{G}_D) J_{\sigma_F}$
 $\forall i, \Delta_i \leftarrow \gamma_\theta(\frac{i}{n}) - \gamma_\theta(\frac{i-1}{n})$
 $d(e_0, e_1) \leftarrow \sum_{i=1}^n \left(\frac{\partial \gamma_\theta}{\partial t} \Big|_{\frac{i}{n}} \right)^T G_{D \circ F}(\gamma_\theta(\frac{i}{n})) \left(\frac{\partial \gamma_\theta}{\partial t} \Big|_{\frac{i}{n}} \right) \Delta_i$
Return: $d(e_0, e_1)$

Algorithm 3 GELATO with Soft Learning

Input: Reward penalty coefficient λ , rollout horizon h , rollout batch size b , training epochs T , number of neighbors K .
Train variational latent forward model on dataset \mathcal{D} by maximizing ELBO (Equation (7))
Construct embedded dataset $\mathcal{D}_{\text{embd}} = \{E_i\}_{i=1}^n$ using latent model to initialize KNN.
Initialize policy π and empty replay buffer $\mathcal{R}_{\text{model}} \leftarrow \emptyset$.
for epoch = 1 **to** T **do**
 for $i = 1$ **to** b (in parallel) **do**
 Sample state s_1 from \mathcal{D} for the initialization of the rollout and embed using latent model to produce z_1 .
 for $j = 1$ **to** h **do**
 Sample an action $a_j \sim \pi(\cdot | z_j)$.
 Embed $(z_j, a_j) \rightarrow E_{z_j, a_j}$ using latent model
 Sample z_{j+1} from latent forward model $F(E_{z_j, a_j})$.
 Sample r_j from latent reward model $R(E_{z_j, a_j})$.
 Use Algorithm 2 to compute K nearest neighbors $\left\{ \text{NN}_{z_j, a_j}^{(k)} \right\}_{k=1}^K$ and their distances $\{d_k\}_{k=1}^K$ to E_{z_j, a_j} .
 Compute $\tilde{r}_j = r_j - \lambda \left(\frac{1}{K} \sum_{k=1}^K d_k \right)$
 Add sample $(z_j, a_j, \tilde{r}_j, z_{j+1})$ to $\mathcal{R}_{\text{model}}$.
 end for
 end for
 Drawing samples from $\mathcal{R}_{\text{data}} \cup \mathcal{R}_{\text{model}}$, use SAC to update π .
end for

600 GELATO, we used $\lambda = 1$ as our reward penalty coefficient and $\lambda = 2$ for the skewed versions. We
601 used rollout horizon of $h = 5$, and did not notice significant performance improvement for different
602 values of h .

603 **11 Missing Proofs**

604 **11.1 Proof of Proposition 1**

605 For any curve γ , we have that

$$\begin{aligned}
 \int_0^1 \left\| \frac{\partial f(\gamma(t))}{\partial t} \right\| dt &= \int_0^1 \left\| \frac{\partial f(\gamma(t))^T}{\partial \gamma(t)} \frac{\partial \gamma(t)}{\partial t} \right\| dt \\
 &= \int_0^1 \left\| J_f(\gamma(t))^T \frac{\partial \gamma(t)}{\partial t} \right\| dt \\
 &= \int_0^1 \sqrt{\left\langle \frac{\partial \gamma(t)}{\partial t}, J_f^T(\gamma(t)) J_f(\gamma(t)) \frac{\partial \gamma(t)}{\partial t} \right\rangle} dt \\
 &= \int_0^1 \sqrt{\left\langle \frac{\partial \gamma(t)}{\partial t}, G_f(\gamma(t)) \frac{\partial \gamma(t)}{\partial t} \right\rangle} dt.
 \end{aligned}$$

606 This completes the proof.

607 **11.2 Proof of Theorems 1 and 2**

608 We begin by restating the theorems.

609 **Theorem 1.** [Arvanitidis et al. (2018)] Assume $D(\cdot|z) \sim \mathcal{N}(\mu(z), \sigma(z)I)$. Then

$$\mathbb{E}_{D(\cdot|z)} [G_D(z)] = \mathbb{E}_{D(\cdot|z)} [J_D(z)^T J_D(z)] = \underbrace{G_\mu(z)}_{\text{proximity}} + \underbrace{G_\sigma(z)}_{\text{uncertainty}}, \quad (3)$$

610 where $G_\mu(z) = J_\mu^T(z) J_\mu(z)$ and $G_\sigma(z) = J_\sigma^T(z) J_\sigma(z)$.

611 **Theorem 2.** Assume $F(\cdot|z) \sim \mathcal{N}(\mu_F(z), \sigma_F(z)I)$, $D(\cdot|x) \sim \mathcal{N}(\mu_D(x), \sigma_D(x)I)$. Then, the expected pullback metric of the composite function $(D \circ F)$ is given by

$$\mathbb{E}_{P(D \circ F)} [G_{D \circ F}(z)] = J_{\mu_F}^T(z) \overline{G}_D(z) J_{\mu_F}(z) + J_{\sigma_F}^T(z) \text{diag}(\overline{G}_D(z)) J_{\sigma_F}(z),$$

613 where here, $\overline{G}_D(z) = \mathbb{E}_{x \sim F(\cdot|z)} [J_{\mu_D}^T(x) J_{\mu_D}(x) + J_{\sigma_D}^T(x) J_{\sigma_D}(x)]$.

614 Notice that Theorem 1 is a special case of Theorem 2 with F being the trivial identity function.
 615 Additionally, a complete proof of Theorem 1 can be found in Arvanitidis et al. (2018). We turn to
 616 prove Theorem 2.

617 We begin by proving the following auxiliary lemma.

618 **Lemma 1.** Let $\epsilon \sim \mathcal{N}(0, I_K)$, $f : \mathbb{R}^d \mapsto \mathbb{R}^K$, $A \in \mathbb{R}^{K \times K}$. Denote $S_i = \text{diag}\left(\frac{\partial f^1}{\partial z_i}, \frac{\partial f^2}{\partial z_i}, \dots, \frac{\partial f^K}{\partial z_i}\right)$
 619 for $1 \leq i \leq d$ and

$$B = [S_1 \epsilon, S_2 \epsilon, \dots, S_d \epsilon]_{K \times d}.$$

620 Then $\mathbb{E} [B^T A B] = J_f^T \text{diag}(A) J_f$.

621 *Proof.* We have that

$$\begin{aligned}
 \mathbb{E} [B^T A B] &= \mathbb{E} \left[\begin{bmatrix} \epsilon^T S_1^T \\ \epsilon^T S_2^T \\ \vdots \\ \epsilon^T S_d^T \end{bmatrix} A [S_1 \epsilon, S_2 \epsilon, \dots, S_d \epsilon] \right] \\
 &= \begin{bmatrix} \mathbb{E} [\epsilon^T S_1^T A S_1 \epsilon], \mathbb{E} [\epsilon^T S_1^T A S_2 \epsilon], \dots, \mathbb{E} [\epsilon^T S_1^T A S_d \epsilon], \\ \mathbb{E} [\epsilon^T S_2^T A S_1 \epsilon], \mathbb{E} [\epsilon^T S_2^T A S_2 \epsilon], \dots, \mathbb{E} [\epsilon^T S_2^T A S_d \epsilon], \\ \vdots \\ \mathbb{E} [\epsilon^T S_d^T A S_1 \epsilon], \mathbb{E} [\epsilon^T S_d^T A S_2 \epsilon], \dots, \mathbb{E} [\epsilon^T S_d^T A S_d \epsilon] \end{bmatrix}.
 \end{aligned}$$

622 Finally notice that for any matrix M

$$\mathbb{E}[\epsilon^T M \epsilon] = \sum_{i=1}^d \sum_{j=1}^d M_{ij} \mathbb{E}[\epsilon_i \epsilon_j] = \sum_{i=1}^d M_{ii} = \text{trace}(M).$$

623 Then,

$$\mathbb{E}[B^T AB] = \begin{bmatrix} \text{trace}(S_1^T AS_1), \text{trace}(S_1^T AS_2), \dots, \text{trace}(S_1^T AS_d) \\ \text{trace}(S_2^T AS_1), \text{trace}(S_2^T AS_2), \dots, \text{trace}(S_2^T AS_d) \\ \vdots \\ \text{trace}(S_d^T AS_1), \text{trace}(S_d^T AS_2), \dots, \text{trace}(S_d^T AS_d) \end{bmatrix}.$$

624

□

625 Next, note that

$$\text{trace}(S_i AS_j) = \sum_{k=1}^K \frac{\partial f^k}{\partial z_i} \frac{\partial f^k}{\partial z_j} A_{kk}.$$

626 Therefore,

$$\mathbb{E}[B^T AB] = J_f^T \text{diag}(A) J_f.$$

627 We are now ready to prove Theorem 2.

628 *Proof of Theorem 2.* We can write $z' = \mu_F(z) + \sigma_F(z) \odot \epsilon_F$ and $s' = \mu_D(z') + \sigma_D(z') \odot \epsilon_D$ where
 629 $\epsilon_F \sim \mathcal{N}(0, I_d)$, $\epsilon_D \sim \mathcal{N}(0, I_K)$, $\mu_F: \mathbb{R}^d \mapsto \mathbb{R}^\ell$, $\mu_D: \mathbb{R}^\ell \mapsto \mathbb{R}^K$ and $\sigma_F: \mathbb{R}^d \mapsto \mathbb{R}^\ell$, $\sigma_D: \mathbb{R}^\ell \mapsto$
 630 \mathbb{R}^K .

631 Applying the chain rule we get

$$J_{D \circ F} = \frac{\partial(D \circ F)}{\partial z} = J_{\mu_D} J_{\mu_F} + J_{\mu_D} B_{\epsilon_F} + B_{\epsilon_D} J_{\mu_F} + B_{\epsilon_D} B_{\epsilon_F}$$

632 where

$$\begin{aligned} B_{\epsilon_F} &= (S_{F,1}\epsilon_F, S_{F,2}\epsilon_F, \dots, S_{F,d}\epsilon_F)_{d \times d}, \\ S_{F,i} &= \text{diag}\left(\frac{\partial \sigma_F^1}{\partial z_i}, \frac{\partial \sigma_F^2}{\partial z_i}, \dots, \frac{\partial \sigma_F^d}{\partial z_i}\right)_{d \times d}, \quad \text{and} \\ B_{\epsilon_D} &= (S_{D,1}\epsilon_D, S_{D,2}\epsilon_D, \dots, S_{D,d}\epsilon_D)_{K \times d}, \\ S_{D,i} &= \text{diag}\left(\frac{\partial \sigma_D^1}{\partial z'_i}, \frac{\partial \sigma_D^2}{\partial z'_i}, \dots, \frac{\partial \sigma_D^d}{\partial z'_i}\right)_{K \times K}. \end{aligned}$$

633 This yields

$$\begin{aligned} \mathbb{E}[J_{F \circ D}^T J_{F \circ D}] &= \mathbb{E}\left[(J_{\mu_D} J_{\mu_F} + J_{\mu_D} B_{\epsilon_F} + B_{\epsilon_D} J_{\mu_F} + B_{\epsilon_D} B_{\epsilon_F})^T (J_{\mu_D} J_{\mu_F} + J_{\mu_D} B_{\epsilon_F} + B_{\epsilon_D} J_{\mu_F} + B_{\epsilon_D} B_{\epsilon_F})\right] \\ &= J_{\mu_F}^T J_{\mu_D}^T J_{\mu_D} J_{\mu_F} + \mathbb{E}[B_{\epsilon_F}^T J_{\mu_D}^T J_{\mu_D} B_{\epsilon_F}] + \mathbb{E}[J_{\mu_F}^T B_{\epsilon_D}^T B_{\epsilon_D} J_{\mu_F}] + \mathbb{E}[B_{\epsilon_F}^T B_{\epsilon_D}^T B_{\epsilon_D} B_{\epsilon_F}] \\ &= J_{\mu_F}^T (J_{\mu_D}^T J_{\mu_D} + \mathbb{E}[B_{\epsilon_D}^T B_{\epsilon_D}]) J_{\mu_F} + \mathbb{E}[B_{\epsilon_F}^T (J_{\mu_D}^T J_{\mu_D} + B_{\epsilon_D}^T B_{\epsilon_D}) B_{\epsilon_F}] \end{aligned}$$

634 where in the second equality we used the fact that ϵ_D, ϵ_F are independent and $\mathbb{E}[B_{\epsilon}] = 0$. By
 635 Lemma 1 we have

$$\mathbb{E}[B_{\epsilon_D}^T B_{\epsilon_D}] = J_{\sigma_D}^T J_{\sigma_D}.$$

636 Similarly,

$$\mathbb{E}[B_{\epsilon_F}^T B_{\epsilon_D}^T B_{\epsilon_D} B_{\epsilon_F}] = \mathbb{E}[\mathbb{E}[B_{\epsilon_F}^T B_{\epsilon_D}^T B_{\epsilon_D} B_{\epsilon_F} | \epsilon_F]] = \mathbb{E}[B_{\epsilon_F}^T J_{\sigma_D}^T J_{\sigma_D} B_{\epsilon_F}] = J_{\sigma_F}^T \text{diag}(J_{\sigma_F}^T J_{\sigma_F}) J_{\sigma_F}$$

637 Finally,

$$\mathbb{E}[B_{\epsilon_F}^T J_{\mu_D}^T J_{\mu_D} B_{\epsilon_F}] = J_{\sigma_F}^T \text{diag}(J_{\mu_D}^T J_{\mu_D}) J_{\sigma_F}.$$

638 This completes the proof. □