Regularized Latent Dynamics Prediction is a Strong Baseline For Behavioral Foundation Models

Anonymous authors

Paper under double-blind review

Abstract

1	Behavioral Foundation Models (BFMs) have seen some success recently in producing
2	agents with the capabilities to adapt to any unknown reward or task. In reality, these
3	methods are only able to produce near-optimal policies for the reward functions that are
4	in the span of some pre-existing state features. Naturally, their efficiency relies heav-
5	ily on the choice of state features used by them. As a result, these BFMs have used
6	a wide variety of potentially complex objectives to train task spanning features with
7	different inductive properties. With this work, our aim is to examine the question: are
8	these complex representation learning objectives necessary for zero-shot RL? Specifi-
9	cally, we revisit the objective of self-supervised next-state prediction for state feature
10	learning, but observe that such an objective is prone to increasing state-feature similar-
11	ity, and subsequently reducing span of reward functions that we can represent optimal
12	policies for. We show that by simply maintaining feature diversity using orthonormal
13	regularization along with next-state prediction, we can match or surpass state-of-the-art
14	complex representation learning methods for zero-shot RL.

15 1 Introduction

The reward hypothesis states that all goals and purposes can be understood as maximization of scalar reward signals. This principle has motivated development of RL algorithms that learn efficiently given a reward function. However, a large part of prior developments in RL focus on dealing with single reward function or a small subset of reward functions. But with the recent focus on generalist agents, the generalization capabilities of RL to new tasks are being tested. Still, when compared to the supervised ML counterparts, RL lags behind in showing zero-shot generalization to new tasks in an environment.

Zero shot learning has been adapted in RL (Touati et al., 2023) to learn agents that can solve *any* task in the environment without any additional training or planning, after an initial pretraining. Zero-shot RL has significant practical potential in developing generalist agents with wide applicability. For instance, robotics applications, like robotic manipulation or drone navigation, often require agents to solve a wide variety of unknown tasks. A general-purpose household robot needs to possess the capability to flexibly adapt to various household chores without explicit training for each new task.

29 Behavioral foundation models based on the idea of leveraging successor representations (Touati 30 et al., 2023; Agarwal et al., 2024) have shown promising progress towards developing algorithms that output near-optimal policies for a wide class of reward functions without additional learning 31 or training during test-time by pretraining on a dataset of reward-free interactions, i.e, zero-shot 32 RL. Such BFMs work by a) learning a state representation $\phi : s \to \mathbb{R}^d$ and b) learning a space 33 of policies parameterized by a latent vector $z \in \mathbb{R}^d$ trained to be optimal for reward defined as 34 $r(s) = \phi(s)^{\top} z$. At test time given any reward function $r^{test}(s)$, the near-optimal policy $\pi_{z_{r}test}$ 35 is obtained by projecting reward functions into the space of state-representations, and solving for 36 $z_{r^{test}}$ such that $r^{test}(s) \approx \phi(s)^{\top} \cdot z_{r^{test}}$. The near-optimal policy is then given by $\pi_{z_{stest}}$. 37

The success of modern zero-shot RL methods is often attributed to learning generalizable state-38 39 representations. State-of-the-art methods usually learn state representations that retain information suitable to represent successor measures under a wide class of policies. Successor measures are 40 41 information rich objects that capture a policy's state visitation in the environment given any starting state. Successor measures are usually learned for an explicitly defined class of policies (Agarwal 42 43 et al., 2024) or implicitly by first defining a class of reward functions (Ramesh et al., 2021; Park 44 et al., 2024) and considering optimal policies for those reward functions as the set of policies. The 45 main insight behind predicting successor measure as a target for state representation learning is 46 that representations sufficient to explain future state-visitation for a wide range of policies captures 47 features that are relevant for sequential decision making under various reward functions.

48 Unfortunately, state representation learning by estimating successor measures requires iteratively 49 applying Bellman evaluation backups or Bellman optimality backups, both of which are known to 50 result in optimization difficulties or feature collapse due to the instability inherent in bootstrapping 51 (Kumar et al., 2021). On the other hand, the dynamics learning objective is an optimization-friendly 52 objective for state representation learning that bypasses bootstrapping. However, using the learned 53 dynamics model to obtain a policy at test time would require a policy training phase with model 54 based RL algorithm going against our objective of zero-shot RL. This work aims to investigate the 55 following question:

56 Is next-state prediction enough to learn state features that enable performant zero-shot RL?

57 Our investigation is inspired by the work of Fujimoto et al. (2025), which underscored the impor-58 tance of auxiliary objective of state representation learning through dynamics prediction losses in 59 boosting performance of single-task model-free RL. Our work differs by tackling a different setting 60 - we present an empirical investigation of the simple dynamics prediction objective for learning 61 representations suitable for zero-shot RL. Unlike the single task RL setting examined by Fujimoto 62 et al. (2025), we find that in its naive form, this objective leads to a mild form of feature collapse 63 where representation of different states increase in similarity over training steps and result in poor 64 zero-shot RL performance when evaluated on a number of downstream tasks. With a simple reg-65 ularization to prevent collapse, we show that model-based representations learned via supervised 66 learning are competitive and present a scalable alternative to representations learned via complex 67 successor measure estimation methods for zero-shot RL.

68 2 Related Work

69 Unsupervised RL: Unsupervised RL encompasses the class of algorithms that enable learning 70 general-purpose skills and representations without relying on reward signal in the data. In this work, 71 we focus on techniques that learn representations capable of producing optimal value functions for 72 any arbitrary function reward specification.

73 Recent pre-training approaches (e.g., Ma et al. (2023); Nair et al.) borrow self-supervised tech-74 niques from computer vision—such as masked auto-encoding—to extract embeddings from large-75 scale datasets (Grauman et al. (2022)) that can be fine-tuned for downstream control. However, 76 these representations are inherently tied to the behavior policies used during data collection. These 77 policies are limited in their ability to capture the full spectrum of possible behaviors or to approxi-78 mate Q-functions for any reward functions. HILP (Park et al. (2024)) goes beyond standard masked 79 autoencoding approaches by using Hilbert-space representations to preserve temporal dynamics. 80 Auxiliary objectives, which involve complementary predictive tasks to get richer semantic or tem-81 poral structures, have also been explored in previous works (Agarwal et al. (2021), Schwarzer et al. 82 (2020)). Although representations from auxiliary objectives can accelerate policy learning, a new policy still needs to be learned from scratch for each new reward function. 83

Several works have also focused on intent or skill discovery through diversity-driven objectives. These methods consider state-visitation distribution that are defined by latents or skills. Thus, maximizing mutual information (Warde-Farley et al. (2018), Eysenbach et al. (2018), Achiam et al. (2018), Eysenbach et al. (2022)) or minimizing the Wasserstein distance (Park et al. (2023)) between latents and state-visitation distribution is used to ensure diversity.

89 Behavioral Foundation Models: Behavioral Foundation Models deals with the class of approaches

90 that can be used to train an RL agent in an unsupervised manner using task-agnostic reward-free

91 offline transitions. During inference, BFMs can approximate the optimal policy for a wide class of 92 unseen reward functions without any further training

92 unseen reward functions without any further training.

93 Forward-Backward representations (Touati & Ollivier (2021)) and PSM (Agarwal et al. (2024)) 94 provide a robust framework for BFMs based on stationary distribution, on which several succes-95 sive works are based. Fast Imitation with BFMs (Pirotta et al. (2023)) demonstrates the ability of 96 successor-measure-based BFMs to imitate new behaviors from just a few demonstrations, while Fast Adaptation with BFMs (Sikchi et al. (2025)) builds upon this by fine-tuning BFMs' latent embed-97 98 ding space, yielding 10-40% improvement over their zero-shot performance in a few of episodes. 99 Recent progress in imitation learning has led to the development of BFMs tailored for humanoid 100 control tasks (Peng et al. (2022), Won et al. (2022), Luo et al. (2023), Tirinzoni et al. (2025)) which 101 can produce diverse behaviors trained using human demonstration data.

102 **3** Preliminaries

103 We consider a reward-free Markov Decision Process (MDP) (Puterman, 2014) which is defined as 104 a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, d_0, \gamma)$, where \mathcal{S} and \mathcal{A} respectively denote the state and action spaces, P 105 denotes the transition dynamics with P(s'|s, a) indicating the probability of transitioning from s to s' by taking action a, d_0 denotes the initial state distribution and $\gamma \in (0,1)$ specifies the discount 106 107 factor. A policy π is a function $\pi : S \to \Delta(A)$ mapping a state s to probabilities of action in A. We denote by $\Pr(\cdot \mid s, a, \pi)$ and $\mathbb{E}[\cdot \mid s, a, \pi]$ the probability and expectation operators under state-108 action sequences $(s_t, a_t)_{t>0}$ starting at (s, a) and following policy π with $s_t \sim P(\cdot \mid s_{t-1}, a_{t-1})$ 109 and $a_t \sim \pi(\cdot \mid s_t)$. Given any reward function $r: S \to \mathbb{R}$, the Q-function of π for r is $Q_r^{\pi}(s, a) :=$ 110 $\sum_{t>0} \gamma^t \mathbb{E}[r(s_{t+1}) \mid s, a, \pi].$ 111

112 Successor measures based Behavioral Foundation Models: The successor measure (Dayan, 113 1993; Blier et al., 2021) of state-action (s, a) under a policy π is the (discounted) distribution of

future states obtained by taking action a in state s and following policy π thereafter:

$$M^{\pi}(s, a, X) := \sum_{t \ge 0} \gamma^{t} \Pr(s_{t+1} \in X \mid s, a, \pi) \quad \forall X \subset \mathcal{S}.$$
⁽¹⁾

Q functions can be represented using successor measures as, $Q^{\pi}(s, a) = \sum_{s^+} M^{\pi}(s, a, s^+) r(s^+)$. 115 This simple linear relationship between Q functions and Successor Measures have been exploited 116 117 by a number of works (Touati & Ollivier, 2021; Agarwal et al., 2024) to create Behavioral Foun-118 dation Models(BFMs). The BFMs parameterize their policies (and correspondingly successor mea-119 sures) using a latent $z \in \mathcal{Z}$ to pre-compute $\pi_z = \arg \max_z M_z^{\top} r$. BFMs based on successor features (Touati & Ollivier, 2021; Zheng et al., 2024) parameterize the reward functions linearly 120 121 using these latents as spans of the state features, $\phi: S \to Z$, $r = \phi z$. Hence, the inference for any 122 reward function reduces to finding this latent z from reward samples using linear regression. We 123 will be following a similar setup as these successor feature methods, where we will be representing 124 rewards as a span of the state features and learn to represent successor measures using these state features, leading to efficient computation of $M_z^{\top} r$. 125

126 **4 Method**

127 This method can be broadly divided into two parts - representation learning and zero-shot RL. The 128 state representation encoder is trained using dynamics prediction and orthonormality loss, enabling 129 the encoder to learn representations that will be generalizable across tasks. Leveraging these ro-130 bust state embeddings, we then pretrain a Behavioral Foundation Model (BFM) to predict successor 131 measures, enabling zero-shot inference of near-optimal policies for unseen reward functions. We refer to this method as RLDP (Regularized Latent Dynamics Prediction based Behavioral Foundation
 Policies)

134 4.1 Learning Representations with Regularized Latent Dynamics Prediction

135 Zero-shot RL based on successor measures rely on learning a state representation denoted by $\phi(s)$. 136 This state representation will define the span of reward functions that the zero-shot RL method is 137 guaranteed to output optimal policies for. Our primary representation learning objective is simple — unrolled latent dynamics prediction. We learn a state representation encoder $\phi : S \to \mathbb{R}^d$, ($\mathcal{Z} =$ 138 \mathbb{R}^d) and a latent state-action representation encoder $g: \mathbb{R}^d \times \mathcal{A} \to \mathbb{R}^d$ such that latent dynamics 139 representation remains linear in these representations $\phi(s') = q(\phi(s), a)^{\top} \mathbf{w}$ with some constant 140 141 weights w informing our loss function for representation learning. A sub-sequence of horizon H is sampled from the offline interaction dataset given by $\tau^i = \{s_0^i, a_0^i, s_1^i, a_0^i, s_1^i, a_1^i, \dots, s_{H-1}^i, a_{H-1}^i, s_H^i\}$. 142 A sequence of future latent states $h_{1:H}$ are obtained by encoding the initial state $h_0 = \phi(s_0)$ and 143 unrolling using the defined dynamics model $h_{t+1} = g(h_t, a_t)^{\top}$. w. Then the objective is to predict 144 the encoded future latent states: 145

$$\mathcal{L}_{d} = \min_{\phi,g} \mathbb{E}_{\tau \sim d^{O}} \left| \left(\sum_{t=1}^{H} g(h_{t}, a_{t})^{\top} \mathbf{w} - h_{t+1} \right)^{2} \right|$$
(2)

The idea of learning low-rank representations for dynamics prediction is inspired by prior works in linear MDP and MDP homomorphisms (Parr et al., 2008; Ravindran, 2004; Van der Pol et al., 2020) and has been shown to be successful in recent work of Fujimoto et al. (2025) where it is observed that model-free methods can be made competitive to model-based methods in sample efficiency and final performance with representations learned from dynamics prediction.

152 However, solely learning from the latent dynam-153 ics objective can lead to convergence to a col-154 lapsed solution. This is unsurprising as trivial 155 solutions of predicting a constant zero vector 156 achieves a perfect loss in Equation 2. To combat 157 this, prior works (Grill et al., 2020) have pro-158 posed the use of a semi-gradient update where a stop-gradient is used for target h_{t+1} in Equa-159 160 tion 2 along with a slowly updating target. Do 161 these techniques mitigate solution collapse? We 162 answer this question by plotting the cosine sim-163 ilarity of state representations trained via the 164 above objective on an offline dataset collected 165 by an exploration algorithm RND (Burda et al., 166 2019). Figure 1 shows that while the solutions



Figure 1: Average Cosine similarity between state-representations sampled uniformly from the training dataset: Feature similarity increases over the course of training. Shaded region shows standard deviation over 4 seeds

do not collapse, there is an increase in feature similarity over the course of learning which we refer
to as a mild form of collapse. As the space of reward functions is spanned by state features, such
an increase in feature similarity can directly reduce the class of reward functions for which we can
learn optimal policies.

171 Preventing collapse in unsupervised RL: In the setting of unsupervised RL, the dataset contains 172 purely reward-free transitions. To prevent collapse, we consider diversity regularization. Orthornor-173 mality regularizations have been widely studied in self-supervised learning (He et al., 2024b; Bansal 174 et al., 2018b). Since we are looking to span reward functions using these state features, it makes 175 sense to have these features orthogonal to each other. We project all state representations in a hypersphere: $\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$ and regularize by minimizing cosine similarity between 176 177 any two states. This technique is also referred to as orthogonal regularization and has been used in 178 self-supervised techniques for representation learning in vision and natural language (Bansal et al., 179 2018a; He et al., 2024a). Some prior unsupervised RL methods (Touati et al., 2023) use this regu-

- 180 larization as an implementation trick to stabilize training; in the case of latent dynamics prediction
- 181 this step becomes crucial to mitigate the increase in representation similarity. The orthogonal regu-

182 larization loss looks like:

$$\mathcal{L}_r = \mathbb{E}_{s \ s' \sim d^O}[\phi(s)^\top \phi(s')] \tag{3}$$

- 183 Our final loss is a weighted combination of dynamics prediction combined with orthogonal diversity
- 184 regularization : $\mathcal{L} = \mathcal{L}_d + \lambda \mathcal{L}_r \tag{4}$
- 185 where λ controls the regularization strength.

186 4.2 Zero-shot RL with Regularized Latent Dynamics Predictive Representations

187 We construct a Behavioral Foundation Model using the state-encoder ϕ trained by minimizing \mathcal{L} (Equation 4). We define the class of reward functions as: $\mathcal{T} = \{\mathbf{r}_z | \mathbf{r}_z = (\phi^{\top} \phi)^{-1} \phi^{\top} z \text{ for some } z \in \mathbb{C} \}$ 188 \mathbb{S}^{d-1} and define π_z to be the optimal policy for r_z . Successor Measures, M^{π_z} , are parameterized 189 as $M^{\pi_z}(s, a, s^+) = \psi^{\pi_z}(s, a)\phi(s^+)$ with $\psi^{\pi_z}(s, a)$ (or $\psi(s, a, z)$) being the successor features 190 for the state feature ϕ i.e. $\psi^{\pi_z}(s,a) = \mathbb{E}_{\pi_z} [\sum_{t=0}^{\infty} \gamma^t \phi(s)]$. Since ϕ is already obtained using 191 Equation 4, learning M^{π_z} would mean optimizing for ψ . Following Touati & Ollivier (2021); Touati 192 et al. (2023); Agarwal et al. (2024), we use a contrastive objective to train M^{π_z} , parameterized as 193 194 $M^{\pi_z}(s, a, s^+) = \psi(s, a, z)\phi(s^+),$

$$\mathcal{L}_{BFM} = -\mathbb{E}_{s,a,s'\sim d^O}[M_z^{\pi}(s,a,s')] + \frac{1}{2}\mathbb{E}_{s,a,s'\sim d^O,s^+\sim d^O}[(M_z^{\pi}(s,a,s^+) - \gamma \bar{M}_z^{\pi}(s',\pi_z(s'),s^+))^2].$$
(5)

195 Equation 5 requires samples from π_z . Hence the policy is optimized to maximize Q^{π_z} ,

$$\pi_z(s) = \max_a Q^{\pi_z}(s, a) = \max_a \psi(s, a, z).z$$
(6)

During inference, given a reward function r, we obtain the near-optimal policy by finding the corresponding z via linear regression,

$$z_R = \min \mathbb{E}_{d^{\mathcal{O}}}[(r(s) - (\phi^\top \phi)^{-1} \phi(s) \cdot z)]^2 \implies z_R = \mathbb{E}_{d^{\mathcal{O}}}[\phi(s) \cdot r(s)]$$
(7)

198 5 Experiments

The goal of our experiments is to perform an empirical study of suitability of state representations learned by latent next-state prediction objective when compared to other methods that employ more complex strategies. We perform several empirical ablations to understand our design choices. In particular, we aim to answer the following questions: a) Under an apples-to-apples setting of keeping all other learning factors similar, how does our method compare to baselines in enabling generalization to unseen reward functions? b) How does our method fare in the setting of large observation space of pixels? c) What design decisions are crucial to the success of our method?

206 Setup: We consider continuous control tasks from DeepMind control suite - Pointmass, Cheetah, 207 Walker, Quadruped under a similar setup considered by prior works in zero-shot RL. A dataset is 208 collected in these experiments using an exploration algorithm RND (Burda et al., 2019) without 209 specifying any reward functions. With such data, we pretrain a BFM using the method outlined in 210 our method section 4. Each algorithm is given the same budget of gradient steps during pretraining, 211 controlling the state representation dimension and the final performance is obtained by taking the 212 model obtained at the end of learning and querying it for different task-rewards for 50 episodes. All 213 of our results are aggregated across 4 seeds.

Baselines: We compare this method against commonly used state-of-the-art baselines in rewardfunction spanning/representation such as: FB, HILP, PSM, and Laplacian. The Laplacian approach (Wu et al., 2018) learns state representation using eigenvectors of graph-Laplacian induced by a random-walk operator. FB (Touati et al., 2023) intertwines learning of state-representation with policy learning step where state-representations are learned such that they can represent successor

	Task	Laplace	FB	HILP	PSM	RLDP
Walker	Stand Run Walk Flip	$\begin{array}{c} 243.70 \pm 151.40 \\ 63.65 \pm 31.02 \\ 190.53 \pm 168.45 \\ 48.73 \pm 17.66 \end{array}$	$\begin{array}{c} 902.63 \pm 38.94 \\ 392.76 \pm 31.29 \\ 877.10 \pm 81.05 \\ 206.22 \pm 162.27 \end{array}$	$\begin{array}{c} 607.07 \pm 165.28 \\ 107.84 \pm 34.24 \\ 399.67 \pm 39.31 \\ 277.95 \pm 59.63 \end{array}$	$\begin{array}{c} 872.61 \pm 38.81 \\ 351.50 \pm 19.46 \\ 891.44 \pm 46.81 \\ 640.75 \pm 31.88 \end{array}$	$\begin{array}{c} 890.40 \pm 27.33 \\ 334.26 \pm 49.69 \\ 779.768 \pm 137.156 \\ 492.94 \pm 22.79 \end{array}$
	Average(*)	136.65	594.67	348.13	689.07	624.34
Cheetah	Run Run Backward Walk Walk Backward	$\begin{array}{c} 96.32 \pm 35.69 \\ 106.38 \pm 29.40 \\ 409.15 \pm 56.08 \\ 654.29 \pm 219.81 \end{array}$	$\begin{array}{c} 257.59\pm 58.51\\ 307.07\pm 14.91\\ 799.83\pm 67.51\\ 980.76\pm 2.32 \end{array}$	$\begin{array}{c} 68.22 \pm 47.08 \\ 37.99 \pm 25.16 \\ 318.30 \pm 168.42 \\ 349.61 \pm 236.29 \end{array}$	$\begin{array}{c} 244.38 \pm 80.00 \\ 296.44 \pm 20.14 \\ 984.21 \pm 0.49 \\ 979.01 \pm 7.73 \end{array}$	$\begin{array}{c} 157.12 \pm 29.92 \\ 170.52 \pm 15.30 \\ 592.92 \pm 104.66 \\ 821.51 \pm 50.62 \end{array}$
	Average(*)	316.53	586.31	193.53	626.01	435.52
Quadruped	Stand Run Walk Jump	$\begin{array}{c} 854.50 \pm 41.47 \\ 412.98 \pm 54.03 \\ 494.56 \pm 62.49 \\ 642.84 \pm 114.15 \end{array}$	$\begin{array}{c} 740.05 \pm 107.15 \\ 386.67 \pm 32.53 \\ 566.57 \pm 53.22 \\ 581.28 \pm 107.38 \end{array}$	$\begin{array}{c} 409.54 \pm 97.59 \\ 205.44 \pm 47.89 \\ 218.54 \pm 86.67 \\ 325.51 \pm 93.06 \end{array}$	$\begin{array}{c} 842.86 \pm 82.18 \\ 431.77 \pm 44.69 \\ 603.97 \pm 73.67 \\ 596.37 \pm 94.23 \end{array}$	$\begin{array}{c} 794.94 \pm 43.25 \\ 457.41 \pm 74.70 \\ 465.40 \pm 185.29 \\ 733.322 \pm 55.304 \end{array}$
	Average(*)	601.22	568.64	289.75	618.74	612.77
Pointmass	Top Left Top Right Bottom Left Bottom Right	$\begin{array}{c} 713.46 \pm 58.90 \\ 581.14 \pm 214.79 \\ 689.05 \pm 37.08 \\ 21.29 \pm 42.54 \end{array}$	$\begin{array}{c} 897.83 \pm 35.79 \\ 274.95 \pm 197.90 \\ 517.23 \pm 302.63 \\ 19.37 \pm 33.54 \end{array}$	$\begin{array}{c} 944.46 \pm 12.94 \\ 96.04 \pm 166.34 \\ 192.34 \pm 177.48 \\ 0.17 \pm 0.29 \end{array}$	$\begin{array}{c} 831.43 \pm 69.51 \\ 730.27 \pm 58.10 \\ 451.38 \pm 73.46 \\ 43.29 \pm 38.40 \end{array}$	$\begin{array}{c} 890.406 \pm 60.791 \\ 795.469 \pm 21.103 \\ 805.172 \pm 20.443 \\ 193.381 \pm 167.633 \end{array}$
	Average(*)	501.23	427.34	308.25	514.09	671.10

Table 1: Comparison (over 4 seeds) of zero-shot RL performance between different methods. RLDP demonstrates a marked improvement over prior methods.

219 measures for a class of reward-optimal policies. HILP (Park et al., 2024) learns state representation

220 using a goal reaching objective which is subsequently used for zero-shot RL. PSM (Agarwal et al.,

221 2024) learns state representation to represent the successor measures for a class of policies defined

with a discrete codebook. These baselines represent a set of diverse and strong approaches in the area of zero-shot RL.

Implementation and Evaluation: To evaluate the different zero-shot RL methods we take the pretrained policies and query them on a variety of tasks. For each environment, we consider 4 tasks similar to prior works (Touati et al., 2023; Park et al., 2024; Agarwal et al., 2024).

227 5.1 Benchmarking Zero-Shot RL for Continuous Control

We conduct our experiments across two axis - a) Table 1 compares against representation dimensions suggested by authors for different methods with the same number of gradient updates for pretraining each BFM. b) Table 2 pretrains all the BFMs on same number of representation dimensions (512) and gradient steps (3 million). We also conducted experiments for pixel-based inputs, the results for which are in Table 3. We consider four environments – Walker, Cheetah, Quadruped, and Pointmass and use the ExoRL suite (Yarats et al. (2022)) for obtaining exploratory datasets using RND (Burda et al. (2019)).

Overall, RLDP fares competitively to baselines across the environments despite its simplicity. Abla tions studying the effects of the orthogonality loss and encoding horizon in representation learning
 are presented in the next section.

238 5.2 What matters for supervising representation suitable for control?

In this section, we aim to ablate components of our method and understand which factors have a strong effect on final performance. To that end, we consider ablating diversity regularization as well as encoding horizon - the two design choices we make in our method. We begin by keeping the encoding horizon constant (H = 5) while we change the orthogonality regularization coefficient. For the second study, we keep the orthogonality regularizer constant ($\lambda = 1.0$) while changing the encoding horizon. The results for these ablations are in Figure 2.

We observe that for zero regularization, the performance takes a steep dip compared to $\lambda > 0$. This shows that orthogonality regularizer is critical to the representation learning part of our algorithm.



Figure 2: Evaluating the impact of Encoding Horizon and Orthogonality Regularization

247 For fixed encoding horizon, we see that orthogonality regularizer $\lambda = 1$ performs best. To further 248 understand the role of the orthogonality regularizer in representation learning and how it helps pre-249 vent feature collapse, we refer to Section 4.1 and Section 6.2, where we look at the cosine similarity 250 between representations. For the second part of our ablations where we see how different encoding 251 horizons can affect the average reward obtained by the agent under constant orthonormality regular-252 ization ($\lambda = 1$), we observe that encoding horizons 5 and 20 perform better than the others, with 253 5 performing best on average across all environments. For the purpose of our experiments, we use 254 encoding horizon 5.

255 6 Conclusion

256 This paper introduces 257 RLDP, a method to de-258 couple representation 259 learning and reinforcement 260 learning, which allows 261 for learning generalizable 262 representations. Our ob-263 jective takes the simple 264 form of regularized latent-265 dynamics prediction, an 266 objective that does not 267 require any reconstruc-268 tion, making it able to 269 handle high-dimensional 270 observation space and 271 does not require Bellman backups, making it more 272 273 amenable to optimization. 274 We show that simply using 275 latent-dynamics prediction 276 leads to a mild form of 277 feature collapse where 278 the state-representation

-				
	Task	FB	PSM	RLDP
Walker	Stand Run Walk Flip	$\begin{array}{c} 918.29 \pm 28.83 \\ 381.31 \pm 17.32 \\ 779.29 \pm 63.60 \\ 977.08 \pm 2.76 \end{array}$	$\begin{array}{c} 899.54 \pm 30.73 \\ 450.57 \pm 28.95 \\ 875.61 \pm 33.44 \\ 621.36 \pm 75.62 \end{array}$	$\begin{array}{c} 890.40 \pm 27.33 \\ 334.26 \pm 49.69 \\ 779.768 \pm 137.16 \\ 492.94 \pm 22.79 \end{array}$
	Average(*)	763.99	711.77	624.34
Cheetah	Run Run Backward Walk Walk Backward	$\begin{array}{c} 129.39 \pm 37.63 \\ 142.41 \pm 36.77 \\ 604.54 \pm 80.51 \\ 630.40 \pm 144.23 \end{array}$	$\begin{array}{c} 181.85\pm54.17\\ 158.64\pm18.56\\ 576.98\pm209.45\\ 817.92\pm98.86\end{array}$	$\begin{array}{c} 157.12 \pm 29.92 \\ 170.52 \pm 15.30 \\ 592.92 \pm 104.66 \\ 821.51 \pm 50.62 \end{array}$
	Average(*)	376.69	433.85	435.52
Quadruped	Stand Run Walk Jump	$\begin{array}{c} 732.59 \pm 101.33 \\ 425.15 \pm 52.02 \\ 492.91 \pm 17.55 \\ 567.27 \pm 48.90 \end{array}$	$\begin{array}{c} 708.03 \pm 34.99 \\ 404.32 \pm 23.26 \\ 523.94 \pm 52.13 \\ 549.57 \pm 15.86 \end{array}$	$\begin{array}{c} 794.94 \pm 43.25 \\ 457.41 \pm 74.70 \\ 465.40 \pm 185.29 \\ 733.322 \pm 55.30 \end{array}$
	Average(*)	554.48	546.46	612.77
Pointmass	Top Left Top Right Bottom Left Bottom Right	$\begin{array}{c} 943.85 \pm 17.31 \\ 550.84 \pm 282.41 \\ 672.28 \pm 153.06 \\ 272.97 \pm 274.99 \end{array}$	$\begin{array}{c} 924.20 \pm 10.64 \\ 666.00 \pm 133.15 \\ 800.93 \pm 15.62 \\ 123.44 \pm 138.82 \end{array}$	$\begin{array}{c} 890.41 \pm 60.79 \\ 795.47 \pm 21.10 \\ 805.17 \pm 20.44 \\ 193.38 \pm 167.63 \end{array}$
	Average(*)	461.77	488.05	671.11

Table 2: Comparison (over 4)	seeds) of zero-shot RL performance
between FB, PSM, and RLDP	with representation size of $d = 512$.
RLDP results are unchanged.	

similarity increases over time. To combat this issue, we propose using orthogonal regularization, a well-known technique to prevent feature collapse. Using our method enables learning generalizable, stable, and robust representations that can achieve competitive performance compared to established successor measure-based techniques without relying on reinforcement-driven signals. In this work, we present initial investigation results. Future research directions include qualitatively examining the learned representations, alternative regularization strategies, further scaling these methods to complex pixel-based observations, and extending the applicability to real-world robotics and control tasks. This work, thus, paves the way for simpler yet effective approaches to developing behavioralfoundation models.

288 **References**

- Joshua Achiam, Harrison Edwards, Dario Amodei, and Pieter Abbeel. Variational option discovery
 algorithms. *arXiv preprint arXiv:1807.10299*, 2018.
- Siddhant Agarwal, Aaron Courville, and Rishabh Agarwal. Behavior predictive representations
 for generalization in reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021. URL
 https://openreview.net/forum?id=b5PJaxS6Jxq.
- Siddhant Agarwal, Harshit Sikchi, Peter Stone, and Amy Zhang. Proto successor measure:
 Representing the space of all possible solutions of reinforcement learning. *arXiv preprint* arXiv:2411.19418, 2024.
- Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regularizations in training deep networks? *Advances in Neural Information Processing Systems*, 31, 2018a.
- Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regular izations in training deep cnns? *CoRR*, abs/1810.09102, 2018b. URL http://arxiv.org/
 abs/1810.09102.
- Léonard Blier, Corentin Tallec, and Yann Ollivier. Learning successor states and goal-dependent
 values: A mathematical viewpoint. *arXiv preprint arXiv:2101.07123*, 2021.
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
 distillation. In *International Conference on Learning Representations*, 2019. URL https:
 //openreview.net/forum?id=H11JJnR5Ym.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representa tion. *Neural computation*, 5(4):613–624, 1993.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need:
 Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. The information geometry of
 unsupervised reinforcement learning. In *International Conference on Learning Representations*,
 2022. URL https://openreview.net/forum?id=3wU2UX0voE.
- Scott Fujimoto, Pierluca D'Oro, Amy Zhang, Yuandong Tian, and Michael Rabbat. Towards
 general-purpose model-free reinforcement learning. *arXiv preprint arXiv:2501.16142*, 2025.
- Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in
 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18995–19012, 2022.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena
 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Junlin He, Jinxiao Du, and Wei Ma. Preventing dimensional collapse in self-supervised learning via orthogonality regularization. *arXiv preprint arXiv:2411.00392*, 2024a.
- Junlin He, Jinxiao Du, and Wei Ma. Preventing dimensional collapse in self-supervised learning via orthogonality regularization, 2024b. URL https://arxiv.org/abs/2411.00392.

- Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron Courville, George Tucker, and Sergey Levine.
 Dr3: Value-based deep reinforcement learning requires explicit regularization. *arXiv preprint arXiv:2112.04716*, 2021.
- Zhengyi Luo, Jinkun Cao, Josh Merel, Alexander Winkler, Jing Huang, Kris Kitani, and Weipeng
 Xu. Universal humanoid motion representations for physics-based control. *arXiv preprint arXiv:2310.04582*, 2023.
- Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy
 Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training.
 In *The Eleventh International Conference on Learning Representations*, 2023. URL https:
 //openreview.net/forum?id=YJ7o2wetJ2.
- Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A univer sal visual representation for robot manipulation. In *6th Annual Conference on Robot Learning*.
- Seohong Park, Oleh Rybkin, and Sergey Levine. Metra: Scalable unsupervised rl with metric-aware
 abstraction. *arXiv preprint arXiv:2310.08887*, 2023.
- Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert represen tations. In *Forty-first International Conference on Machine Learning*, 2024. URL https:
 //openreview.net/forum?id=LhNsSaAKub.
- Ronald Parr, Lihong Li, Gavin Taylor, Christopher Painter-Wakefield, and Michael L Littman. An
 analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*,
 pp. 752–759, 2008.
- Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale
 reusable adversarial skill embeddings for physically simulated characters. *ACM Transactions On Graphics (TOG)*, 41(4):1–17, 2022.
- Matteo Pirotta, Andrea Tirinzoni, Ahmed Touati, Alessandro Lazaric, and Yann Ollivier. Fast imita tion via behavior foundation models. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
 Wiley & Sons, 2014.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
 and Ilya Sutskever. Zero-shot text-to-image generation. *CoRR*, abs/2102.12092, 2021. URL
 https://arxiv.org/abs/2102.12092.
- Balaraman Ravindran. An algebraic approach to abstraction in reinforcement learning. University
 of Massachusetts Amherst, 2004.
- Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bach man. Data-efficient reinforcement learning with self-predictive representations. *arXiv preprint arXiv:2007.05929*, 2020.
- Harshit Sikchi, Andrea Tirinzoni, Ahmed Touati, Yingchen Xu, Anssi Kanervisto, Scott Niekum,
 Amy Zhang, Alessandro Lazaric, and Matteo Pirotta. Fast adaptation with behavioral foundation
 models. *arXiv preprint arXiv:2504.07896*, 2025.
- Andrea Tirinzoni, Ahmed Touati, Jesse Farebrother, Mateusz Guzek, Anssi Kanervisto, Yingchen
 Xu, Alessandro Lazaric, and Matteo Pirotta. Zero-shot whole-body humanoid control via behav ioral foundation models. *arXiv preprint arXiv:2504.11054*, 2025.
- Ahmed Touati and Yann Ollivier. Learning one representation to optimize all rewards. Advances in
 Neural Information Processing Systems, 34:13–23, 2021.

- 374 Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist?
- In The Eleventh International Conference on Learning Representations, 2023. URL https:
 //openreview.net/forum?id=MYEap_OcQI.
- Elise Van der Pol, Thomas Kipf, Frans A Oliehoek, and Max Welling. Plannable approximations to
 mdp homomorphisms: Equivariance under actions. *arXiv preprint arXiv:2002.11963*, 2020.
- David Warde-Farley, Tom Van de Wiele, Tejas Kulkarni, Catalin Ionescu, Steven Hansen, and
 Volodymyr Mnih. Unsupervised control through non-parametric discriminative rewards. *arXiv preprint arXiv:1811.11359*, 2018.
- Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Physics-based character controllers using conditional vaes. *ACM Transactions on Graphics (TOG)*, 41(4):1–12, 2022.
- Yifan Wu, George Tucker, and Ofir Nachum. The laplacian in rl: Learning representations with
 efficient approximations. *arXiv preprint arXiv:1810.04586*, 2018.
- Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro Lazaric,
 and Lerrel Pinto. Don't change the algorithm, change the data: Exploratory data for offline
- 388 reinforcement learning. *arXiv preprint arXiv:2201.13425*, 2022.
- Chongyi Zheng, Jens Tuyls, Joanne Peng, and Benjamin Eysenbach. Can a misl fly? analysis and
 ingredients for mutual information skill learning. *arXiv preprint arXiv:2412.08021*, 2024.

391 Supplementary Materials 392 The following content was not necessarily subject to peer review.

394 Appendix

395 6.1 Pseudocode for the update

Algorithm 1 UPDATE RLDP (PRETRAINING)

```
Require: Replay buffer \mathcal{D}; encoder \phi with target \overline{\phi}; forward model \psi with target \overline{\psi}; actor \pi
Require: Pretraining steps T_{repr}, Polyak factor \tau
   for num train steps do
      if step < representation steps then
           Representation Learning (Section 4.1):
           Sample segment batch \{s_{0:H}, a_{0:H}\} \sim \mathcal{D}
           L_{\text{repr}} \leftarrow \lambda L_{\text{orth}} + L_{\text{dyn}} (Equation 2, 3, 4)
          if step\%250 == 0 then
              \bar{\phi} \leftarrow \phi
           end if
       else
           Sample transitions \{(s, a, s', done)\} \sim D
           \gamma \leftarrow \gamma_0 (1 - \text{done})
           Sample z \sim \text{MixUniform} \cup \phi(s') {mix random prior + goal-encoded}
           Value update (Section 4.2):
           a' \sim \pi(s', z)
           F^{\star} \leftarrow \overleftarrow{\psi}(s', z, a'), \quad B^{\star} \leftarrow \overline{\phi}(s')
           M_{ij}^{\star} = F_i^{\star} \cdot B_j^{\star}
           F \leftarrow \psi(s, z, a), \quad B \leftarrow \phi(s')
           M_{ij} = F_i \cdot B_j
           L_{FB} = \frac{1}{2} \mathbb{E}[\|(M - \gamma M^{\star}) \circ \mathbf{1}_{i \neq j}\|^2]
           Update \psi \leftarrow \psi - \eta \nabla_{\psi} L_{FB}
          Policy update (Section 4.2):
          a = \pi(s, z)
           Q = \psi(s, z, a) \cdot z
           L_{\pi} = -\mathbb{E}[Q]
           Update \pi \leftarrow \pi - \eta \nabla_{\pi} L_{\pi}
           Target-network sync:
           \bar{\psi} \leftarrow \tau \psi + (1 - \tau) \bar{\psi}
       end if
   end for
```

396 6.2 Additional Results

Learning Zero-shot Policies for Continuous Control Table 3 shows the performance of RLDP on
 pixel inputs.

Understanding the role of orthogonality regularization Figure 3 shows the impact of changing orthogonality regularization while keeping a constant encoding horizon (H = 5). Specifically, we look at how the cosine similarity between representations changes during training for different regularization coefficients.

We observe that for a regularization coefficient of 0, the cosine similarity increases, indicating that all states are getting mapped to similar representations, i.e., representation collapse. For any regularization coefficient greater than 0 ($\lambda > 0$), we observe that the cosine similarity drops below 0.1,

Dataset	Environment	Task	FDM	FB	HILP	RLDP
	XX7 11	Flip	282 ± 52	62 ± 57	232 ± 41	242.6 ± 33.5
		Run	146 ± 60	42 ± 25	126 ± 8	106.6 ± 12.3
	walker	Stand	557 ± 97	172 ± 111	496 ± 73	464.1 ± 67.7
		Walk	452 ± 52	104 ± 82	376 ± 52	229.8 ± 55.4
	Average(*)		359.25	95.0	307.5	149.0
RND	Cheetah	Run	178 ± 41	221 ± 15	276 ± 46	101.1 ± 30.3
		Run Backward	126 ± 9	171 ± 123	297 ± 46	31.4 ± 5.5
		Walk	470 ± 182	535 ± 251	895 ± 33	439.7 ± 133.7
		Walk Backward	441 ± 107	535 ± 440	927 ± 35	137.1 ± 27.2
	Average(*)		303.8	365.5	598.8	177.3
		Jump	273 ± 66	224 ± 149	244 ± 122	430.5 ± 94.9
	Quadruped	Run	192 ± 29	158 ± 82	148 ± 19	323.4 ± 37.1
		Stand	374 ± 85	347 ± 191	327 ± 126	507.6 ± 76.8
		Walk	199 ± 63	162 ± 92	163 ± 45	285.7 ± 29.9
	Average(*)		259.5	222.8	220.5	386.8

Table 3: Performance comparison on the pixel-based ExORL benchmark across different environments and tasks.

406 close to 0.0, indicating that the states are being mapped to different representations. This highlights the importance of the regularization coefficient in preventing representation collapse.



Figure 3: Evaluating the impact of Orthogonality Regularization on representations learned across four environments: Cheetah (top left), Pointmass (top right), Quadruped (bottom left), and Walker (bottom right).

407