
Value Implicit Pretraining does not learn representations suitable for Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Leveraging large-scale internet data to bootstrap representations for RL presents
2 a reasonable path forward to learning general-purpose agents. Unlike computer
3 vision and NLP, the sequential nature of the RL problem makes it unclear what the
4 learning rule should be. Goal-Conditioned RL (GCRL) presents a self-supervised
5 objective for RL that can allow agents to learn representation suitable for planning
6 to arbitrary states. Traditional GCRL objectives require action labels which are
7 usually missing from sequential data on internet (eg. videos). A recent approach,
8 Value Implicit Pretraining (VIP), presents a new objective to learn optimal value
9 functions without requiring action labels. However, our findings in this paper
10 suggest that VIP fails to learn the correct representations in simple domains. This
11 finding motivates us to conduct a detailed investigation, and through formal argu-
12 ments we establish why VIP shows this anomaly. We propose a simple yet effect
13 alternative, DERAIL, that indeed learns optimal value functions and subsequently
14 representations suitable for RL.

15 1 Introduction

16 Fields like Computer Vision and Natural Language Processing have immensely benefited from
17 utilizing large-scale vision and language data. Unfortunately, such developments in reinforcement
18 learning and robotics have been limited. Recent works like Walke et al. (2023); Padalkar et al.
19 (2023) have made significant strides in creating large-scale datasets that can be used for robotics. To
20 learn representations from such large-scale data that can be transferred easily to different domains
21 and downstream tasks requires the representations to be independent of underlying action space i.e.
22 learning representations only from observations. R3M (Nair et al., 2022), VIP (Ma et al., 2022),
23 LIV (Ma et al., 2023), ICVF (Ghosh et al., 2023) are a few successful works that use observations
24 from large-scale egocentric video data to learn value functions parameterized using encoders. These
25 encoders can then be transferred to other domains and can be used either to define reward functions
26 or as a backbone encoder for behavior cloning.

27 Learning representations for RL requires navigating a number of challenges: (1) What is the right
28 representation learning objective? (2) What are we going to use the representation for? (3) How do we
29 evaluate those representations. Unfortunately, these things are murkier in the field of reinforcement
30 learning when compared to supervised learning. A promising approach to this problem is to use
31 Goal-Conditioned RL Kaelbling (1993) as an self-supervised objective for representation learning.
32 The idea is to learn compact representations that can allow for planning from any state in the
33 environment to any other state — effectively compressing the observations into representations
34 suitable for planning. Learning representations through GCRL using offline datasets that lack actions
35 is challenging as most objectives require actions to be known (Eysenbach et al., 2021; Sikchi et al.,
36 2023a). A recently developed line of work, Value Implicit Pretraining (VIP) (Ma et al., 2022)

37 leverages Fenchel-Rockefeller duality by treating RL as a convex program and derives an action-free
 38 objective for GCRL that is later adapted to representation learning.
 39 In our work, we demonstrate that representations learnt through Value Implicit Pretraining are lacking
 40 and unsuitable for planning. VIP’s anomaly arises as a result of imposing an information bottleneck
 41 in the primal objective requiring assumptions that don’t hold in practice, to convert to a dual objective
 42 suitable for representation learning. We hypothesize that the perceived success of VIP is due to
 43 using expert-like trajectories in its offline dataset of transitions and learning representation using a
 44 time-contrastive objective between neighbouring observations. Our empirical experiments confirm
 45 our findings that the representation learned through VIP indeed fail on simple low-dimensional tasks.
 46 Next, gathering the insights from our argument, we develop a simple, principled, and action-free
 47 objective for representation learning that we then use to overcome the limitations of VIP. Our proposed
 48 objective is motivated by the recently proposed dual perspective of reinforcement learning. We show
 49 that improves representation learning from a series of qualitative and quantitative experiments.

50 2 Related Works

51 Representation Learning in RL can be broadly classified into two categories: (a) offline pretraining
 52 (Ma et al., 2022; Nair et al., 2022) and (b) using an auxilliary loss (Schwarzer et al., 2021; Agarwal
 53 et al., 2021; Agarwal et al.) over the RL loss. The goal of both these paradigms is to induce an
 54 inductive bias on the representation space by using reward-free interaction data. There are a variety
 55 of auxilliary objectives that can be added to produce desired properties in the representations majority
 56 of them being contrastive objectives (Schwarzer et al., 2021; Agarwal et al.; Srinivas et al., 2020)
 57 aiming for sample efficiency (Schwarzer et al., 2021; Srinivas et al., 2020), generalization (Agarwal
 58 et al., 2021; Agarwal et al.) and temporal consistency (Zhao et al., 2023). While these methods
 59 do introduce some interesting properties in the representation space and observes gains in sample
 60 efficieny and generalization, these do not look into pre-training task-agnostic generalizable encoders
 61 from offline data.
 62 With the availability of large scale datasets like Ego4D (Grauman et al., 2022) and Epic Kitchens
 63 (Damen et al., 2018), several methods have been developed that look learn representations from
 64 large-scale pretraining trying to bridge the gap between RL and fields like computer vision and NLP.
 65 RRL (Shah & Kumar, 2021) and VC1 (Majumdar et al., 2024) are some methods that have attempted
 66 using classical computer vision techniques for representation learning. However, these works do not
 67 take into account the sequential nature and temporal data. R3M (Nair et al., 2022) used Ego4D to
 68 learn representations using Time Constrastive Loss on the trajectories of Ego4D while VIP (Ma et al.,
 69 2022) and the follow up work LIV (Ma et al., 2023) introduced this temporaral consistency implicitly
 70 by learning a goal-conditioned value function for the trajectories in the large-scale dataset.

71 3 Preliminaries

72 We consider a learning agent in a Markov Decision Process (MDP) (Puterman, 2014; Sutton & Barto,
 73 2018) which is defined as a tuple: $\mathcal{M} = (\mathcal{O}, \mathcal{A}, P, R, \gamma, d_0)$ where \mathcal{O} and \mathcal{A} denote the state and action
 74 spaces respectively, P denotes the transition function with $P(s'|s, a)$ indicating the probability of
 75 transitioning from s to s' taking action a ; R denotes the reward function and $\gamma \in (0, 1)$ specifies the
 76 discount factor. We use $o \in \mathcal{O}$ to denote the space of observations, where an observation is generated
 77 as a stochastic function of underlying state. The reinforcement learning objective is to obtain a policy
 78 $\pi : \mathcal{O} \rightarrow \Delta()$ that maximizes expected return: $\pi \sum_{t=0}^{\infty} \gamma^t r(o_t, a_t)$, where we use \mathbb{E}_{π} to denote the
 79 expectation under the distribution induced by $a_t \sim \pi(\cdot|o_t)$, $o_{t+1} \sim p(\cdot|o_t, a_t)$ and $\Delta()$ denotes a
 80 probability simplex supported over \mathcal{W} .

81 **Dual formulation of RL** Dual RL (Sikchi et al., 2023b), also called Distribution Correction
 82 Estimation (DICE) (Nachum & Dai, 2020) present a family of principled off-policy algorithms that
 83 can leverage data from arbitrary sources to learn optimal policy. Dual RL works by first considering
 84 the following two convex program formulations of regularized reinforcement learning in the form of
 85 **primal-Q** and **primal-V**:

$$\begin{aligned}
 \max_{\pi} J(\pi) &= \max_{\pi} \left[\max_d \mathbb{E}_{d(o,a)}[r(o,a)] - \alpha d(o,a)d^O(o,a) \right] \\
 \text{s.t } d(o,a) &= (1-\gamma)d_0(o).\pi(a|o) + \gamma \sum_{s',a'} d(o',a')p(o'|o',a')\pi(a'|o), \forall o \in \mathcal{O}, a \in \mathcal{A}.
 \end{aligned} \tag{1}$$

86 and,

$$\begin{aligned} & \max_{d \geq 0} \mathbb{E}_{d(o,a)}[r(o,a)] - \alpha d(o,a) d^O(o,a) \\ & \text{s.t. } \sum_{a \in \mathcal{A}} d(o,a) = (1-\gamma)d_0(o) + \gamma \sum_{(o',a') \in \mathcal{O} \times \mathcal{A}} d(o',a') p(o|o',a'), \forall o \in \mathcal{O}. \end{aligned} \quad (2)$$

87 The constraints above represent the Bellman flow conditions that any valid visitation distribution
88 should satisfy. That is, the visitation distribution should be induced by some policy under the
89 dynamics of the environment. Applying Lagrangian duality and using convex conjugates result in
90 respective unconstrained optimization problems for solving regularized RL:

$$\max_{\pi} \min_Q (1-\gamma)o \sim d_0, a \sim \pi(o)Q(o,a) + \alpha(o,a) \sim d^O f^*([Q(o,a) - Q(o,a)]/\alpha), \quad (3)$$

91 and,

$$\min_V (1-\gamma)o \sim d_0 V(o) + \alpha(o,a) \sim d^O f_p^*([\mathcal{T}V(o,a) - V(o)])/\alpha, \quad (4)$$

92 where denotes Bellman operator with policy π and reward function r such that $Q(o,a) = r(o,a) +$
93 $\gamma o' \sim p(\cdot|o,a), a' \sim \pi(\cdot|o')Q(o',a')$ and $\mathcal{T}V(o,a) = r(o,a) + \gamma s' \sim p(\cdot|s,a)V(o')$.

94 4 Relating Value Implicit Pretraining to Optimal Value Function Learning

95 In this section, elucidate why a classical implementation of VIP objective does not quite learn optimal
96 value functions. Our analysis below identifies why VIP value functions perform worse and later
97 propose guidance on how to train representations with a dual objective.

98 4.1 What objective is VIP trying to learn?

99 VIP learns goal-conditioned value functions by leveraging the GCRL as a convex optimization
100 problem. Following the DICE framework, it constructs a convex objective by regularizing the distri-
101 bution matching problem with linear constraints on visitation distribution to generate a corresponding
102 dual objective which is action-free. Additionally, VIP introduces an information bottleneck on
103 observations, thus maximizing the expected return of reaching various goals in the environment.

$$\begin{aligned} & \max_{d,\phi} \mathbb{E}_d[r(o;g)] - D_{KL}(d(\phi(o),a;\phi(g))||d^D(\phi(o),a;\phi(g))) \\ & \text{s.t. } \sum_a d(\phi(o),a;\phi(g)) = (1-\gamma)\mu_0(o,g) + \gamma \sum_{o',a'} T(o|o',a')d(\phi(o'),a';\phi(g)) \end{aligned} \quad (5)$$

104 The idea is that ϕ encodes sufficient statistics about observations and goals in such a way that still
105 allows it to solve the GCRL problem as well as possible. We ask the question *if VIP indeed succeeds*
106 *in maximizing this objective below.*

107 4.2 Issues with VIP under the representation learning objective

108 To identify the issues with the representation learning framework presented by VIP, we shall go
109 through the broad derivation of the representation learning objective from the dual objective. The
110 Lagrangian dual of the above primal representation learning objective with some algebraic rearrange-
111 ment is shown below. Although VIP uses KL divergence, we will be presenting the derivation for any
112 general f -Divergence.

$$\begin{aligned} & \min_V \max_{\phi,d \geq 0} (1-\gamma)\mathbb{E}_{\mu_0(o;g)}[V(\phi(o),\phi(g))] + \mathbb{E}_{d(\phi(o),a;\phi(g))}[r(o;g) + \gamma V(\phi(o'),\phi(g')) \\ & \quad - V(\phi(o),\phi(g))] - D_f(d(\phi(o),a;\phi(g))||d^D(\phi(o),a;\phi(g))) \end{aligned} \quad (6)$$

113 The inner maximization w.r.t d assumes a closed-form solution, denoted by f_{VIP}^* which reduces the
114 representation learning objective to a two-player game between V and ϕ .

$$\begin{aligned} & \max_{\phi} \min_V (1-\gamma)\mathbb{E}_{\mu_0(o;g)}[V(\phi(o),\phi(g))] + \mathbb{E}_{d(\phi(o),a;\phi(g))}[f_{VIP}^*(r(o;g) + \gamma V(\phi(o'),\phi(g'))) \\ & \quad - V(\phi(o),\phi(g)))] \end{aligned} \quad (7)$$

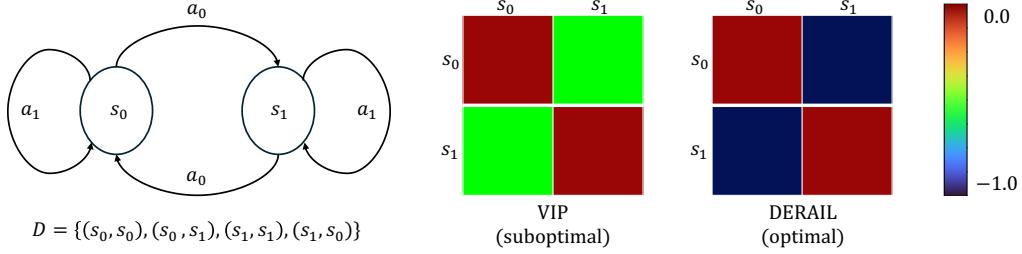


Figure 1: (left) A simple 2-state 2-action MDP with a reward, and the dataset containing all 2 state trajectories obtained from the MDP. (right) A comparison of the goal conditioned value functions $V(s_1, s_2)$ learnt by VIP and DERAIL for the goal conditioned reward, $r(s, g) = 0$ if $s = g$ and -1 otherwise. *VIP produces suboptimal value functions while DERAIL learns the optimal value function.*

115 To ease learning, it is assumed that the optimal value function for any ϕ has the following structure
 116 $V^*(o, g) = -\|\phi(o) - \phi(g)\|$.

$$\max_{\phi} (1 - \gamma) \mathbb{E}_{\mu_0(o; g)} [V^*(\phi(o), \phi(g))] + \mathbb{E}_{d(\phi(o), a; \phi(g))} [f_{VIP}^*(r(o; g) + \gamma V^*(\phi(o'), \phi(g')) - V^*(\phi(o), \phi(g)))] \quad (8)$$

117 We use a simple symmetric MDP example below with two states where one state transitions to another
 118 with a particular action or remains at the same location with the other action. The use of symmetric
 119 MDP is to ensure that the optimal value functions are representable by the structural assumption
 120 of VIP. We will use this MDP as a proof by counterexample to demonstrate VIP’s failure to learn
 121 optimal value functions.

122 *Issue 1: VIP ignores the positivity constraint $d \geq 0$ changes the fixed point of optimization*

123 The inner maximization with respect to the visitation distribution d admits a closed-form solution.
 124 This closed-form solution can differ significantly if the positivity constraint of $d \geq 0$ is ignored. We
 125 compare the functional forms of the conjugate function used in VIP vs the true conjugate function
 126 under positivity constraints in Table 1.

| Divergence Name | Generator $f(x)$ | Conjugate in VIP $f_{VIP}^*(y)$ | True Conjugate $f_p^*(y)$ |
|-------------------|--------------------|---------------------------------|--|
| Reverse KL | $x \log x$ | $\log x e^{y-1}$ | $e^{(y-1)}$ |
| Squared Hellinger | $(\sqrt{x} - 1)^2$ | – | $\frac{y}{1-y}$ |
| Pearson χ^2 | $(x - 1)^2$ | $\frac{(y+1)^2}{2}$ | $\max(\frac{y}{2} + 1, 0)y - (\max(\frac{y}{2} + 1, 0) - 1)^2$ |

Table 1: VIP’s use of conjugate functions vs the true conjugate under positivity constraints. ‘–’ denotes divergences not discussed in VIP.

127 *Issue 2: Assumes a structure on the optimal value function that will be true regardless of the
 128 representations ϕ to simplify optimization*

129 In two-player game or bilevel optimization where one variable depends on the value of another
 130 variable (in our case $V(\phi(o), \phi(g))$ is a function of phi), assuming an analytical relation on the fixed
 131 point (eq. $V^*(\phi(o), \phi(g)) = -\|\phi(o) - \phi(g)\|$) can lead to substantially different fixed point solution.
 132 Figure 1 shows the value function learned by VIP and value function learned by the method we
 133 propose later in this work. VIP converges to an incorrect value function.

134 In general, without imposing structural assumption on the value function, VIP objective remains a
 135 two-player game (Eq 7) bringing optimization challenges. Furthermore, the sampling distribution
 136 required in VIP is in the space of encoded representations and contributes to gradient updates.
 137 The traditional implementation ignores this by assuming no conflicting encoded observations (no
 138 embedding collisions i.e $\phi(o) \neq \phi(o') \forall o, o' \in \mathcal{O}$). This is additionally also a requirement when the
 139 dynamics are stochastic as the constraints in Eq 5 use transitions defined in the space of unencoded

140 observations. This effectively means that if there were a simple two state MDP as in Figure 1 with
 141 D observations of each state, VIP would necessitate compression to $\log(2D)$ bits instead of the
 142 sufficient 1-bit representation of the MDP, thus losing the compression capability afforded by the
 143 structure of MDP. Our work aims to get away with these assumptions and propose a representation
 144 learning objective that retains the benefits of a single-player learning objective.

145 5 DERRAIL: Learning optimal value functions with Dual-V Learning

146 Understanding the limitations of VIP, motivated by the dual framework, we turn to presenting a
 147 simple action-free objective for representation learning. Our key insight that the issues of learning
 148 optimal value function can be mitigated by first deriving a mathematically sound dual objective and
 149 then imposing a information bottleneck as opposed to imposing an information bottleneck on the
 150 primal and using approximation to derive the dual objective.

151 We consider the dual-RL objective Sikchi et al. (2023b) reformulated for goal-conditioned RL:

$$\max_d \mathbb{E}_d[r(o; g)] - D_f(d(o, a; g) \| d^D(o, a; g)) \quad (9)$$

152 where d represents the visitation distribution and is subject to the traditional bellman flow constraints
 153 that ensure d is induced by some policy respecting dynamics of the environment. Here the constraints
 154 become:

$$\sum_a d(o, a; g) = (1 - \gamma)\mu_0(o, g) + \gamma \sum_{o', a'} T(o|o', a')d(o', a'|g) \quad (10)$$

155 Computing the dual of the problem is easy, and can be done following the same steps as in Dual-
 156 RL (Sikchi et al., 2023b). Using straightforward algebraic manipulations we can rewrite the above
 157 equation as:

$$\begin{aligned} \min_V \max_{d \geq 0} & (1 - \gamma)\mathbb{E}_{\mu_0(o; g)}[V(o, g)] + \mathbb{E}_{d(o, a; g)}[r(o; g) + \gamma \sum_{o'} T(o'|o, a, g)V(o', g) - V(o, g)] \\ & - D_f(d(o, a; g) \| d^D(o, a; g)) \end{aligned} \quad (11)$$

158 The distributions in the above equation are all conditioned on a particular goal. The inner maximiza-
 159 tion problem w.r.t. d has a analytical solution and leads us to the final optimization objective:

$$\min_V (1 - \gamma)\mathbb{E}_{\mu_0(o; g)}[V(o, g)] + \mathbb{E}_{d(o, a; g)}[f_p^*(r(o; g) + \gamma \sum_{o'} T(o'|o, a, g)V(o', g) - V(o, g))] \quad (12)$$

160 With the action-free GCRL objective in hand, we
 161 now directly impose an information bottleneck on the
 162 observations by enforcing value function predictions
 163 to only use the encoded representations ($\phi(o)$) of
 164 observations.

165 **Lemma 5.1.** *Value function learning with the DE-*
 166 *RAIL objective converges to optimal value function*
 167 *under sufficient representation capacity of the infor-*
 168 *mation bottleneck.*

Algorithm 1 DERRAIL

Init $V_\phi(s, g)$, conservatism λ
 Let $\mathcal{D} = \hat{\rho} = \{(s, a, s')\}$ be an offline
 dataset.
for $t = 1..T$ iterations **do**
 Train V_ϕ via Orthogonal gradient update
 on Eq. 13
end for
 return ϕ

169 *Proof.* The derivation follows from leveraging the strong duality argument from Section B.1.4
 170 in Sikchi et al. (2023b) along with the assumption of lemma that sufficient representation capability
 171 allows representing all value functions.

172 The objective for representation learning using obser-
 173 vational data, DERAIL, can be written as:

$$\min_{\phi} (1-\gamma) \mathbb{E}_{\mu_0(o;g)} [V_{\phi}(o,g)] + \mathbb{E}_{d(o,a;g)} [f_p^*(r(o;g) + \gamma \sum_{o'} T(o'|o,a,g) V_{\phi}(o',g) - V_{\phi}(o,g))] \quad (13)$$

174 Interestingly, the objective we obtain indicates an al-
 175 most opposite behavior than the VIP objective which
 176 maximizes a monotonic function of bellman error, whereas DERAIL objective minimizes it.

177 5.1 Practical Algorithm

178 Learning representations from offline datasets require tuning conservatism akin to offline RL algo-
 179 rithms. Following Sikchi et al. (2023b), we incorporate conservatism by a linear weighting between
 180 the two terms in the objective:

$$\min_{\phi} (1-\lambda)(1-\gamma) \mathbb{E}_{\mu_0(o;g)} [V_{\phi}(o,g)] + \lambda \mathbb{E}_{d(o,a;g)} [f_p^*(r(s;g) + \gamma \sum_{s'} T(o'|s,a,g) V_{\phi}(o',g) - V_{\phi}(o,g))] \quad (14)$$

181 We instantiate our algorithm (Algorithm 1) using the Pearson χ^2 divergence for which the f_p^* takes
 182 the following closed form:

$$f_p^*(y) = \max\left(\frac{x}{2} + 1, 0\right)x - \left(\max\left(\frac{x}{2} + 1, 0\right) - 1\right)^2 \quad (15)$$

183 Substituting the above form of f_p^* in Eq. 13 gets us the practical objective we use in this work. To
 184 optimize Eq 13 we use orthogonal gradient updates that have been shown to be more effective in
 185 practice in finding the fixed point of the objective Mao et al. (2024) compared to semi-gradient
 186 updates Sikchi et al. (2023b). Prior works have found feature co-adaptation between features of the
 187 current state and the next state, leading to gradients of $V(o',g)$ and $V(o,g)$ canceling out. Orthogonal
 188 gradient updates fix this by considering the projection of the gradient of the next observation in the
 189 orthogonal direction to the gradient of the current observation.

190 We parameterize value functions by considering two representation bottlenecks in this work: (a)
 191 L2/Eucleadian ($V_{\phi}(o,g) = -\|\phi(o) - \phi(g)\|$) used in Ma et al. (2022), and (b) Multilinear used
 192 in Ghosh et al. (2023) ($V_{\phi}(s,g) = \phi_1(s)\phi_2(g)\phi_3(g)$) where $\phi_1(s) \in \mathbb{R}^d, \phi_2(g) \in \mathbb{R}^{d \times d}$ and
 193 $\phi_3(g) \in \mathbb{R}^d$. For multilinear representations we use ϕ_1 as the resulting representation encoder.
 194 The Eucleadian bottleneck linearizes the value function in representation of observation but has the
 195 downside of enforcing symmetric value functions $V(o,g) = V(g,o)$, a condition often violated in
 196 practice. For this reason, we consider Multilinear representation that allow for assymetric value
 197 function learning while still imposing an information bottleneck.

198 6 Experiments

199 Our experiments aim to validate the arguments made in the paper about the failure of VIP in learning
 200 correct representations for reinforcement learning and demonstrate the effectiveness of the proposed
 201 method DERAIL. To this end, we consider a number of MuJoCo tasks with freely accessible datasets Fu
 202 et al. (2020). We do not use the Ego4D dataset used in as the version of the dataset used in the
 203 paper is not available to public. Our results are not limited by the consideration of simulated tasks,
 204 as any representation learning method should learn meaningful representation invariant of domain
 205 being used. Indeed, recent work () shows that representations learned on states (already a compact
 206 representation) can speed up RL. Simulated tasks gives us the ability to do a more detailed analysis
 207 of OOD generalization capabilities of the learned representation. Our experiments below evaluate the
 208 learned representations from VIP and DERAIL both qualitatively and quantitatively.

209 **Datasets:** We use the D4RL datasets and consider the problem of learning representations that allow
 210 the agent to plan from any state in the dataset to any other state. Our preliminary results investigate
 211 representations learned on the following datasets — halfcheetah-medium-expert, hopper-medium-
 212 expert, walker-medium-expert, ant-medium-expert.

213 **6.1 Nature of Learned Representations**

214 In this section, we investigate whether the representations learned by are meaningful. To do so, we
 215 train each method to convergence for 100k gradient updates on all the datasets and use the encoder
 216 (the information bottleneck of the value function) to generate 16 and 32 dimensional representation
 217 of states in the environment. For the multilinear representation we discard ψ which encodes the task
 218 information of which goal to reach and only use ϕ . Figure 2 and Figure 3 plot an MDS projection of
 219 the representations in 2D for L2 and Multilinear representations respectively. Our choice of MDS
 220 projection is motivated by its distance preserving nature even after projection.

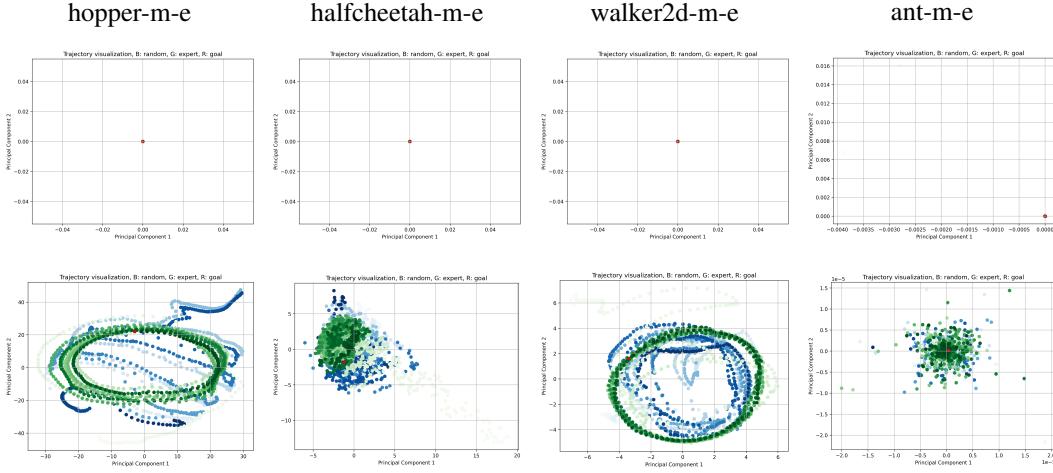


Figure 2: MDS plot of in-distribution representations with L2 bottleneck: We sample an expert trajectory that is not a part of the medium-expert data

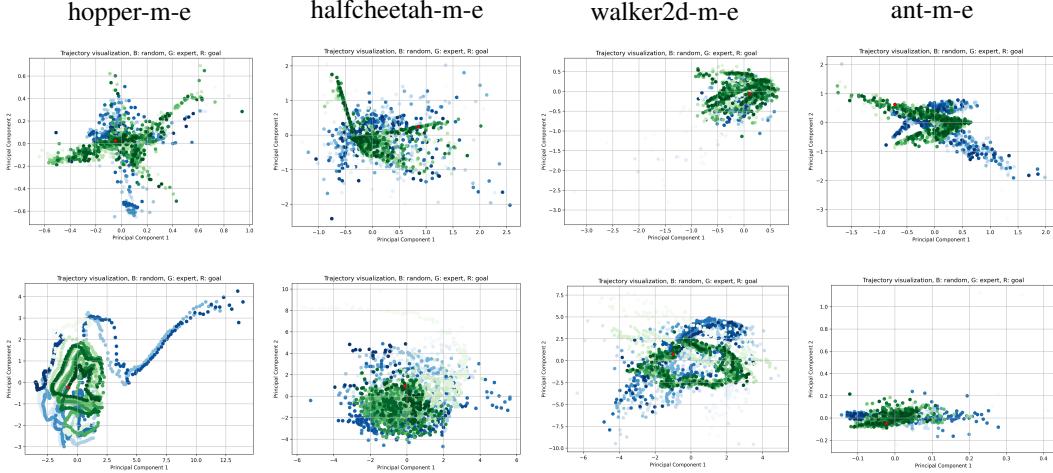


Figure 3: MDS plot of in-distribution representations with Multilinear bottleneck: We sample an expert trajectory that is not a part of the medium-expert data

221 In Figures 2, we observe VIP to collapse representations for most L2 tasks likely due to failing to
 222 implement a necessary embedding collision objective in its practical algorithm. Our representations
 223 for Hopper and Walker2d environments are most interpretable - The representation of expert trajectory
 224 forms a complete loop indicating the periodic pattern of states these environments encounter in an
 225 expert trajectory.

226 **6.2 Optimal Value Functions with Representation Bottleneck**

227 We evaluate the ability of VIP and to learn optimal value functions which is a direct indicator of the
 228 quality of representations learned. In this section, we sample an in-distribution expert trajectory - one
 229 expert trajectory that is in-distribution but does not exist in the medium-expert datasets and a random
 230 trajectory. We then fix the end state of the expert trajectory as the goal and plot the learned value
 231 function at every state of the trajectory.

232 In Figure 4 and Figure 5 we observe that 's value reflect the correct patterns we expect in the expert
 233 trajectory. For Hopper and Walker2d environments where the agent follows a periodic trajectory
 234 and observes similar state multiple times, the learned value function captures this behavior. For
 235 HalfCheetah and Ant, which are not periodic in nature as a result of choice of their observation space,
 236 the value function increases until the goal is reached. In contrast, VIP's value function is noisy and
 237 does not reflect the expected trends.

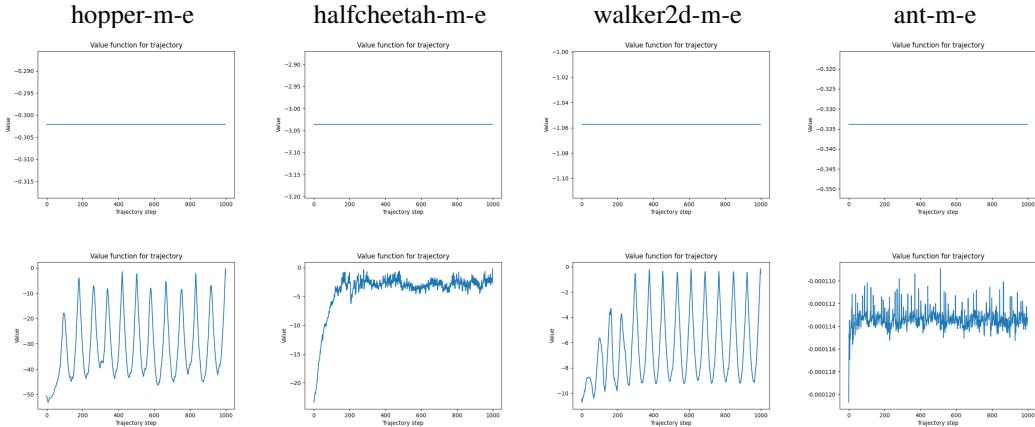


Figure 4: Optimal value function prediction on in-distribution trajectories under a bottlenecked L2 representation. VIP representations collapses causing the value function to collapse. learns meaningful value functions.

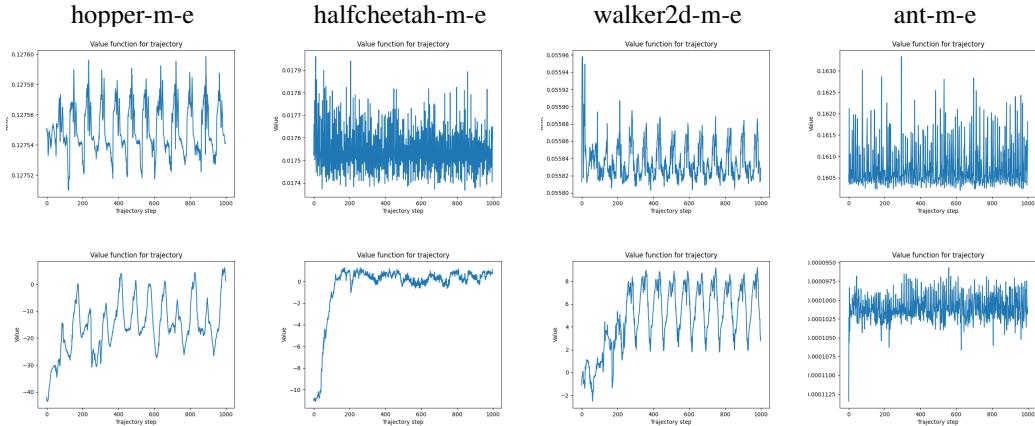


Figure 5: Optimal value function prediction on in-distribution trajectories under a bottlenecked Multilinear representation. VIP representations do not collapse but learn incorrect value functions for HalfCheetah and Ant environments. learns meaningful value functions.

238 **7 Conclusion**

239 Learning a general-purpose representation of the world for reinforcement learning has the potential
 240 to pave way to a foundational model for robotics. This work discusses the limitations of prior work,

241 Value Implicit Pretraining (VIP), that learns such representations by training optimal goal-reaching
242 value functions. Our core insight is assumptions made in VIP does not allow it to learn optimal
243 value functions and hence the right representations. Our work proposes a clean, simple and effective
244 alternative that overcomes these limitations and proposes an action-free objective for GCRL. To this
245 end, we support our claims by studying the learned representations on simulated tasks from the D4RL
246 benchmark.

247 **References**

- 248 Rishabh Agarwal, Marlos C Machado, Pablo Samuel Castro, and Marc G Bellemare. Contrastive
249 behavioral similarity embeddings for generalization in reinforcement learning.
- 250 Siddhant Agarwal, Aaron Courville, and Rishabh Agarwal. Behavior predictive representations
251 for generalization in reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021. URL
252 <https://openreview.net/forum?id=b5PJaxS6Jxg>.
- 253 Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos
254 Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling
255 egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*,
256 2018.
- 257 Ben Eysenbach, Sergey Levine, and Russ R Salakhutdinov. Replacing rewards with examples:
258 Example-based policy search via recursive classification. *Advances in Neural Information Processing
259 Systems*, 34:11541–11552, 2021.
- 260 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
261 data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 262 Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement learning from passive
263 data via latent intentions. In *International Conference on Machine Learning*, pp. 11321–11339.
264 PMLR, 2023.
- 265 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-
266 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan,
267 Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray,
268 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Car-
269 tillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erappalli, Christoph Feichtenhofer, Adriano
270 Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang,
271 Yifei Huang, Wenqi Jia, Leslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico
272 Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan
273 Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova,
274 Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo,
275 Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Crandall,
276 Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna
277 Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva,
278 Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba,
279 Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of
280 egocentric video, 2022.
- 281 Leslie Pack Kaelbling. Learning to achieve goals. In *IJCAI*, volume 2, pp. 1094–8. Citeseer, 1993.
- 282 Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy
283 Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training.
284 *arXiv preprint arXiv:2210.00030*, 2022.
- 285 Yecheng Jason Ma, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv:
286 Language-image representations and rewards for robotic control. In *International Conference on
287 Machine Learning*, pp. 23301–23320. PMLR, 2023.
- 288 Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha Silwal,
289 Aryan Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik, Dhruv Batra, Yixin Lin, Olek-
290 sandr Maksymets, Aravind Rajeswaran, and Franziska Meier. Where are we in the search for an
291 artificial visual cortex for embodied intelligence?, 2024.
- 292 Liyuan Mao, Haoran Xu, Weinan Zhang, and Xianyuan Zhan. Odice: Revealing the mystery of
293 distribution correction estimation via orthogonal-gradient update. *arXiv preprint arXiv:2402.00348*,
294 2024.
- 295 Ofir Nachum and Bo Dai. Reinforcement learning via fenchel-rockafellar duality. *arXiv preprint
296 arXiv:2001.01866*, 2020.

- 297 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal
298 visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- 299 Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander
300 Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic
301 learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.
- 302 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
303 Wiley & Sons, 2014.
- 304 Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bachman.
305 Data-efficient reinforcement learning with self-predictive representations, 2021.
- 306 Rutav M. Shah and Vikash Kumar. RRL: resnet as representation for reinforcement learning. *CoRR*,
307 abs/2107.03380, 2021. URL <https://arxiv.org/abs/2107.03380>.
- 308 Harshit Sikchi, Rohan Chitnis, Ahmed Touati, Alborz Geramifard, Amy Zhang, and Scott Niekum.
309 Score models for offline goal-conditioned reinforcement learning. *arXiv preprint arXiv:2311.02013*,
310 2023a.
- 311 Harshit Sikchi, Qinling Zheng, Amy Zhang, and Scott Niekum. Dual rl: Unification and new
312 methods for reinforcement and imitation learning. *arXiv preprint arXiv:2302.08560*, 2023b.
- 313 Aravind Srinivas, Michael Laskin, and Pieter Abbeel. Curl: Contrastive unsupervised representations
314 for reinforcement learning, 2020.
- 315 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 316 Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao,
317 Philippe Hansen-Estruch, Quan Vuong, Andre He, et al. Bridgedata v2: A dataset for robot
318 learning at scale. *arXiv preprint arXiv:2308.12952*, 2023.
- 319 Yi Zhao, Wenshuai Zhao, Rinu Boney, Juho Kannala, and Joni Pajarinen. Simplified temporal
320 consistency reinforcement learning, 2023.