
Abstraction for Offline Goal-Conditioned Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Markov Decision Processes (MDPs) often exhibit significant redundancy due to
2 symmetries and shared structure across state-goal pairs in real-world Goal-
3 Conditioned Reinforcement Learning (GCRL). While hierarchical policies have
4 been motivated for horizon reduction via *temporal* abstraction in offline GCRL,
5 we demonstrate that hierarchy also enables *absolute* abstraction. By introducing
6 *relativised* options as well as *distinct representations* for different levels of the hier-
7 archy, we demonstrate how an agent can reuse experience across similar contexts of
8 the state-space. Based on this framework, we introduce two simple algorithms for
9 learning relativised options and abstracting from the absolute frame of reference.
10 Our experiments show that such inductive biases significantly improve performance
11 in offline GCRL.

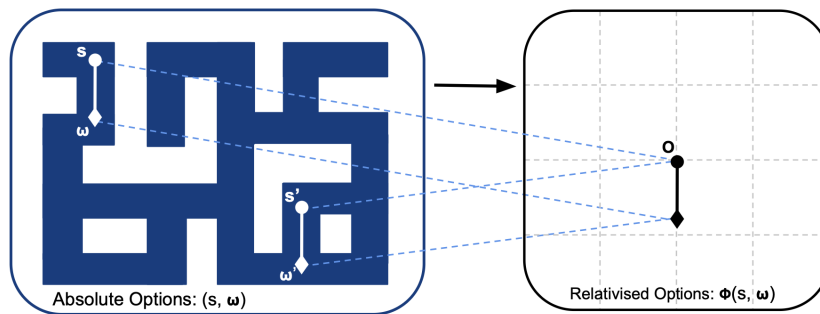


Figure 1: **Abstractive RL (ARL)**. By learning relativised options, ARL enables the reuse of experience across similar contexts of the state-space.

12 1 Introduction

13 Offline Goal-Conditioned Reinforcement Learning (GCRL) [1–5] provides a principled framework for
14 training a general-purpose agent to solve complex long-horizon tasks from static datasets. However,
15 in practice, existing methods have struggled to learn effective policies from offline data, partly due
16 to imperfect dataset coverage of the state-action space [3, 6]. Furthermore, since common offline
17 RL algorithms [7–9] regularise actions towards those close to the dataset distribution to mitigate
18 issues such as distribution shift [3, 6], an agent may fail to recover an optimal policy if a dataset only
19 contains low-quality actions in certain regions of the state space. This is a key challenge in offline RL
20 [10] and leads to difficulties in value estimation, policy extraction, and policy generalisation [11].

21 Recent work suggests that horizon reduction is essential for scaling offline RL [12], motivating the
22 use of hierarchical policies, or options [13], for temporal abstraction [14]. We extend this perspective
23 by arguing that hierarchy offers an additional advantage: *absolute* abstraction. By using *absolute*
24 abstraction, which we define as using *relativised* options and *distinct representations* at different
25 levels of the hierarchy, an agent abstracts away from the absolute frame of reference and can reuse
26 experience across similar contexts of the state-space.

27 To illustrate this point, consider an agent undertaking a locomotion task. While the dataset may
28 only contain a limited number of suboptimal demonstrations of the full task, many of the constituent
29 relativised options (subtasks such as simply moving forward or navigating a corner) might be well
30 represented across demonstrations with different goals. By defining options relative to the agent’s
31 local context (e.g. *navigate to the corner directly ahead* rather than *navigate to corner A*) and
32 decoupling low-level actions (i.e motor-actuation) from redundant high-level information, the agent
33 can leverage many more subtask examples to learn a policy.

34 In principle, relativised options [15] exploit redundancy and symmetry in MDPs to allow behaviours
35 to generalise across states. However, in practice, implementations remain limited to toy examples due
36 to the inherent difficulty of identifying MDP homomorphisms or learning such relative representations.
37 In this work, we introduce *Abstractive Reinforcement Learning* (ARL), a general framework that
38 learns relativised options via action similarity. Based on this, it defines high-level similarity and
39 low-level similarity respectively as state-goal pairs inducing similar options and state-option pairs
40 inducing similar immediate actions.

41 We propose two simple algorithms that comply with the ARL framework: the first can be applied
42 generally, simply using action similarity to implicitly learn relativised options; the second intro-
43 duces a representational inductive bias for the low-level MDP by explicitly enforcing translational
44 invariance, improving generalisation in certain high-dimensional manipulation tasks. Our experi-
45 ments demonstrate that such relativised options and inductive biases result in better performance in
46 high-dimensional offline GCRL.

47 **Contributions.** Concisely, this work addresses the following question:

48 *Can hierarchy enable more robust policy extraction in regions where the dataset suffers from*
49 *low-quality transitions?*

50 Our contributions are two-fold. We first motivate hierarchy in offline RL through abstraction from
51 the absolute frame of reference. We show how relativised options and distinct representations at
52 different levels of the hierarchy can enable data reuse, bounding the maximum error in expected
53 return compared to a flat policy. Secondly, we introduce two simple algorithms that learn relativised
54 options and abstract from the absolute frame of reference. These algorithms outperform both flat
55 policies and hierarchical ones that are anchored in the absolute state-space in high-dimensional tasks
56 — without introducing additional hyperparameters.

57 Explicitly, for an RL practitioner, our work demonstrates that: (i) options should be learned via
58 *action similarity* rather than value similarity i.e. jointly optimised with the low-level policy; (ii) we
59 necessitate two value functions to decouple the high-level from the low-level decision process; and
60 (iii) that imposing translation invariance on the low-level MDP can improve policy generalisation in
61 high-dimensional manipulation tasks. For an RL researcher, our work opens up new avenues, such
62 as methods to learn relativised options via action chunking [16, 17], or incorporating more flexible
63 inductive biases using ideas from Geometric Deep Learning (GDL) [18, 19].

64 2 Preliminaries

65 In Offline [3, 20] Goal-Conditioned Reinforcement Learning (GCRL) [1] an agent seeks to learn
66 a universal policy [2] from a fixed dataset, enabling it to reach arbitrary goal states in the smallest
67 number of timesteps [5].

68 2.1 Problem Setting

69 We consider a standard Markov Decision Process (MDP) [21] defined by the tuple $\mathcal{M} :=$
70 $(p_{s_0}, \mathcal{S}, \mathcal{G}, \mathcal{A}, \mathcal{T}, \beta_g, \gamma)$, where p_{s_0} is the initial state distribution, \mathcal{S}, \mathcal{G} and \mathcal{A} respectively denote the

71 state, goal and action space, \mathcal{T} is the transition function, β_g is a goal-conditioned pseudo-termination
 72 function, and $\gamma < 1$ is the discount factor. At the beginning of the episode, a state s_0 is sampled from
 73 p_{s_0} . A goal state g is uniformly sampled from the goal space \mathcal{G} , and is fixed for the entire episode.
 74 The goal space may be defined over all or a subset of the state dimensions. At each timestep $t \geq 0$, an
 75 agent takes an action a_t conditioned on its current state s_t and goal state g , and transitions to a new
 76 state $s_{t+1} = \mathcal{T}(\cdot | s_t, a_t)$. Episodes terminate according to the goal-conditioned pseudo-termination
 77 function [22] $\beta_g : \mathcal{S} \rightarrow \{0, 1\}$, where $\beta_g(s) = 1$ if and only if the goal has been reached. Following
 78 Andrychowicz et al. [23], we focus on the problem of sparse and binary rewards, which is motivated
 79 in robotics, for example. The agent receives a reward of -1 on all steps, and a reward of 0 upon
 80 reaching the goal $r_t = -\mathbb{1}\{\beta_g(s_t) = 0\}$. The aim of the agent is to learn a universal policy [2]
 81 conditioned on its state and the goal $\pi : \mathcal{S} \times \mathcal{G} \rightarrow \Delta(\mathcal{A})$ that maximises the sum of discounted
 82 returns $\mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r_t]$ from a fixed dataset \mathcal{D} . The dataset contains N state-action trajectories of
 83 length H , collected using an arbitrary policy in a task-agnostic manner. Following Park et al. [24],
 84 we adopt the fully observed deterministic MDP framework for simplicity, though we additionally
 85 evaluate in a stochastic environment.

86 2.2 Offline Reinforcement Learning

87 Offline RL is commonly formulated as a two-step procedure: first, value learning, and subsequently
 88 policy extraction from this value function. A major challenge in offline RL is extracting an optimal
 89 policy from suboptimal data. This difficulty is exacerbated by distribution shift [3, 10], where the
 90 learned policy induces state-action pairs that are poorly supported by the dataset, leading to unreliable
 91 value estimates.

92 To address this issue, model-free algorithms typically incorporate some form of conservatism during
 93 value learning, regularise policy extraction towards the dataset distribution via behaviour-cloning,
 94 or combine both strategies. The value function might be learned with a distribution-constrained
 95 objective [25, 7], to minimise uncertainty [26] or encourage structured representations [27]. Similarly
 96 for policy extraction, the agent is regularised by using an additional behaviour-cloning loss term [28],
 97 weighted behaviour-cloning [8], or rejection-sampling [29] of a behaviour-cloned policy. In all cases:

98 *Due to their underlying bias towards the dataset distribution, flat offline RL algorithms struggle in*
 99 *regions where the dataset suffers from low-quality transitions.*

100 2.3 Options

101 Options [13] provide a framework for decision making at different levels of temporal abstraction.
 102 In the original options framework, each option $\omega \sim \Omega$ represents temporally extended behaviour,
 103 executed until termination according to a learned or predefined termination condition [30]. In contrast,
 104 hierarchical approaches in Offline GCRL typically resample options at every timestep [24, 12, 31],
 105 effectively removing temporal commitment and avoiding the need to specify a termination function:

$$\pi(a | s, g) = \pi_l(a | s, \omega) \quad \omega \sim \pi_h(\cdot | s, g),$$

106 where π_h is the high-level policy that samples an option conditioned on the current state and goal,
 107 and π_l is a low-level policy that samples an action conditioned on the same state and that sampled
 108 option. Such approaches still differ from a flat policy, since the low-level policy is conditioned on
 109 the option rather than the goal. Since local gradient information can be uninformative or misleading
 110 when optimising for distant goals, a key benefit of this hierarchical inductive bias is that it reduces
 111 the effective horizon for both the high-level and low-level MDP [24, 12].

112 3 Related Work

113 Extensive literature motivates hierarchy for temporal abstraction and policy horizon reduction [13,
 114 32, 33], but apart from Nachum et al. [34], [35] and Levy et al. [36], very little work explicitly
 115 motivates hierarchy via representation learning and data reuse. While original work used a set of
 116 hardcoded options [13], since then options have been learned by identifying bottleneck states (e.g.
 117 [37, 38]), or by jointly learning the options with the low-level policy [30]. Most work defines options
 118 in the original state-space, but, for example, Vezhnevets et al. [14] learn options in an embedding
 119 space. Ravindran and Barto [15] introduce the concept of relativised options, but non-hardcoded

120 implementations of relativised options have remained scarce. We remark that options can be defined
 121 in a latent space, but still be anchored to the absolute frame of reference. All of these works study
 122 hierarchical RL in the online setting, while we focus on the offline setting.

123 Offline GCRL naturally lends itself to hierarchical formulations, where high-level policies specify
 124 intermediate goals, and low-level policies learn intermediate goal-reaching behaviours [5, 11, 12,
 125 24, 31, 39]. In hierarchical offline GCRL, options have again generally taken the form of absolute
 126 options in the original state-space (e.g. [12, 31]). Our approach is most related to HIQL [24], which
 127 learns latent options, but differs in two fundamental ways. First, unlike HIQL, which uses one value
 128 function, we use two distinct value functions, enabling different representations at different levels of
 129 the hierarchy (Section 4.3). Second, rather than basing option embeddings on value similarity, we
 130 base them on action similarity (Section 4.2). This changes the organisation of the latent space, since
 131 two state-waypoints might induce similar values but differ entirely in their low-level actions.

132 The theory of state abstraction identifies conditions under which information can be compressed
 133 while maintaining policy optimality at different levels of abstraction [40]. Prior work in State
 134 Representation Learning (SRL) has explored homomorphisms, bisimulation metrics, and contrastive
 135 objectives to map together semantically similar states (e.g. [33, 41, 42]). However, these methods
 136 often require auxiliary loss terms that necessitate hyperparameter tuning and can lead to training
 137 instability. We refer to Echchahed and Castro [43] for an overview. Furthermore, none of these
 138 methods explicitly address hierarchical goal-conditioned representation learning.

139 4 Abstractive Reinforcement Learning (ARL)

140 In this section, we introduce Abstractive Reinforcement Learning (ARL), a framework for learning
 141 abstractions in offline hierarchical GCRL to improve robustness in regions where the dataset suffers
 142 from low-quality transitions. Based on this framework, we introduce two simple algorithms: the first
 143 learns relativised options via action similarity, while the second additionally imposes translational
 144 invariance on the low-level MDP. Together, these algorithms demonstrate how (i) relativised options
 145 and (ii) representational inductive biases can improve generalisation in offline GCRL.

146 4.1 Objective

147 In principle, our aim is to learn abstractions that enable reuse of experience across similar contexts.
 148 By learning options that group together state-waypoint pairs, which, under an optimal policy, induce
 149 similar action sequences, options $\omega \in \Omega$ are relative rather than anchored to an absolute frame of
 150 reference. This naturally induces two hierarchical notions of similarity: a high-level similarity, where
 151 similar state-goal pairs induce similar options, and a low-level similarity, where similar state-option
 152 pairs induce similar immediate actions. Consequently, we can define high-level embeddings $\phi_h(s, g)$
 153 and low-level embeddings $\phi_l(s, \omega)$ that abstract away information irrelevant to their respective levels:

$$\pi(a | s, g) = \pi_l(a | \phi_l(s, \omega)) \quad \omega \sim \pi_h(\cdot | \phi_h(s, g)). \quad (1)$$

154 To allow the high-level and low-level decision processes to operate on different representations and at
 155 different temporal abstractions (i.e. different discount factors), such an approach necessitates two
 156 distinct value functions. In the following Motivation Box, we provide an intuition on how such
 157 abstractions can mitigate failures in regions where the dataset suffers from low-quality transitions by
 158 analysing the maximum error in offline RL for a finite-state, finite-action MDP.

Motivation Box: Bounding the Maximum Error

We consider learning an optimal policy in a finite-state, finite-action MDP from a fixed dataset \mathcal{D} of size N . We refer to Appendix B for a complete derivation and full definitions. We assume at most two possible next states for each state-action pair, which is realistic given our deterministic transition assumption (Section 2). For a flat goal-conditioned policy, the Probably Approximately Correct (PAC) Learning error ϵ in expected return is given by [44–46]:

$$\epsilon \propto \sqrt{\frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}| \cdot \kappa}{(1 - \gamma)^3 N}}, \quad \kappa = \sup_{s \in \mathcal{S}, a \in \mathcal{A}, g \in \mathcal{S}} \frac{d^{\pi^*}(s, a, g)}{d^{\pi^{BC}}(s, a, g)},$$

with constant probability of $1 - \delta$. Here, the proportionality constant depends on δ . κ is the concentrability coefficient that accounts for the distribution shift in data collected by the offline

159

behaviour cloning policy π^{BC} , and the data that would have been collected under the optimal policy π^* . $d^\pi(s, a, g)$ is the discounted occupancy measure (stationary distribution) of the policy π . Intuitively, if the dataset does not include the state-action pairs required to learn the optimal policy, κ (and hence the error ϵ) will be large.

We build on the work of Robert et al. [47] and Li et al. [40] to show that, by using a hierarchical policy with absolute abstraction, the maximum error is bounded by:

$$\epsilon^{\text{hierarchy, rep}} \propto \sqrt{\frac{|\mathcal{C}_h||\Omega| \cdot \kappa_h}{(1-\gamma^n)^3 N}} + \sqrt{\frac{|\mathcal{C}_l|\mathcal{A}|n^3 \cdot \kappa_l}{N}}.$$

This reparameterisation reduces the error bound in four ways:

1. **Horizon reduction:** the high-level policy faces an effective discount factor of γ^n (with $n \geq 1$) rather than γ , reducing the denominator’s sensitivity: $\frac{1}{(1-\gamma^n)^3} \leq \frac{1}{(1-\gamma)^3}^a$.
2. **Cardinality Reduction:** by mapping (s, g) and (s, ω) pairs into equivalence classes $c_h \in \mathcal{C}_h$ and $c_l \in \mathcal{C}_l$, the effective state-space is reduced: $|\mathcal{C}_h| \leq |\mathcal{S}||\mathcal{G}|$ and $|\mathcal{C}_l| \leq |\mathcal{S}||\Omega|$.
3. **Option Efficiency:** relativised options ensure that the option space is small and invariant to absolute position: $|\Omega^{\text{rel}}| \leq |\Omega^{\text{abs}}|$, where Ω^{abs} represents an option space anchored in an absolute frame of reference, and Ω^{rel} represents one in a relative frame of reference.
4. **Concentrability Improvement:** because the concentrability coefficients are now defined using reparameterised policies over the reparameterised latent spaces, the probability mass of the dataset is aggregated across similar contexts:

$$\kappa_h^{\text{rep}} = \sup_{c_h \in \mathcal{C}_h, \omega \in \Omega} \frac{d^{\pi_h^*}(c_h, \omega)}{d^{\pi_h^{\text{BC}}}(c_h, \omega)} \quad \text{and} \quad \kappa_l^{\text{rep}} = \sup_{c_l \in \mathcal{C}_l, a \in \mathcal{A}} \frac{d^{\pi_l^*}(c_l, a)}{d^{\pi_l^{\text{BC}}}(c_l, a)},$$

where

$$d^{\pi_h}(c_h, \omega) = \sum_{(s, g) \in \phi_h^{-1}(c_h)} d_h^\pi(s, g, \omega) \quad \text{and} \quad d^{\pi_l}(c_l, a) = \sum_{(s, \omega) \in \phi_l^{-1}(c_l)} d_l^\pi(s, \omega, a).$$

Even if the dataset has zero mass on a specific state-goal-option (s, g, ω) or state-option-action (s, ω, a) , it likely has mass on the abstract-context-option (c_h, ω) or abstract-context-action (c_l, a) . This reduces the likelihood of a support mismatch, where the optimal policy requires a state transition on which the dataset places zero mass. By aggregating similar contexts we now perform a ratio over sums such that

$$\kappa_h^{\text{rep}} \leq \kappa_h \quad \text{and} \quad \kappa_l^{\text{rep}} \leq \kappa_l.$$

Consequently, for a fixed N and unlike a flat policy, such a reparameterisation could enable learning more optimal behaviour in regions that suffer from low-quality data.

^aNote that while hierarchy introduces an additive error term for the low-level policy, this is typically dominated by the exponential reduction in the high-level error’s horizon-dependent constant, especially in tasks where $\gamma \rightarrow 1$ [12].

160

161 4.2 Abstractive RL Implicitly Learning Relativised Options

162 Based on this objective, we propose a minimal amendment to HIQL [24] to encourage learning
 163 relativised options, which directly address the third point in the Motivation Box (Section 4.1). Rather
 164 than learning option representations with the value function, we propose learning them via the
 165 low-level policy. Following HIQL, we also bound the option space to a hypersphere to introduce
 166 geometric regularisation. However, unlike HIQL, we learn representations via the low-level policy
 167 to push together state-waypoint pairs with similar low-level actions rather than state-waypoint pairs
 168 with similar values. Also unlike HIQL (and, as motivated in Section 4.1) we use two value functions
 169 rather than one.

170 **Low-Level Value.** Although ARL is agnostic to the choice of loss, we train the low-level value V_l
 171 and critic Q_l using Implicit Q Learning [7] to match our benchmark algorithms:

$$\mathcal{L}_{V_l} = \mathbb{E}_{(s,a) \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot|s,a)} \left[\ell_{\tau}^2 \left(V_l(s, g_s) - \tilde{Q}_l(s, g_s, a) \right) \right] \quad (2)$$

$$\mathcal{L}_{Q_l} = \mathbb{E}_{(s,a,s') \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot|s,a)} \left[\left(Q_l(s, g_s, a) - r(s, g_s) - \gamma_l V_l(s', g_s) \right)^2 \right] \quad (3)$$

172 where, \tilde{Q}_l denotes the target network, γ_l the low-level discount factor, ℓ_{τ}^2 the expectile loss, and
 173 $g_s \in \mathcal{S}$ a waypoint to the goal.

174 **Low-Level Policy.** We learn the option embeddings ϕ_{ω} jointly with the low-level policy. As with
 175 the value function, ARL is agnostic to the choice of policy and policy extraction algorithm. In our
 176 implementations we use Advantage-Weighted Regression (AWR) [8]:

$$\mathcal{L}_{\pi_l, \phi_{\omega}} = -\mathbb{E}_{(s,a,s' \dots g_s) \sim \mathcal{D}} \left[e^{\alpha_l A_l(s,a,s',g_s)} \log \pi_l \left(a \mid s, \hat{\phi}_{\omega}(s, g_s) \right) \right], \quad (4)$$

177 where $A_l(s, a, s', g_s) = Q_l(s, g_s, a) - V_l(s, g_s)$ represents the advantage associated with action a .

178 **High-Level Value.** To stabilise training and mitigate the issue of simultaneously learning the option
 179 representation, which can lead to training instability, we learn an action-free high-level value function,
 180 which can be learned directly from state trajectories without requiring explicit option labels. Note that,
 181 apart from being action-free, ARL is agnostic to the choice of high-level value learning and could be
 182 implemented with value horizon reduction such as TD- n or TRL [48]. In our implementations, we
 183 use one-step IVL [49, 50]:

$$\mathcal{L}_{V_h} = \mathbb{E}_{(s,a) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot|s,a)} \left[\ell_{\tau}^2 \left(V_h(s, g) - r(s, g) - \gamma \tilde{V}_h(s', g) \right) \right], \quad (5)$$

184 where \tilde{V}_h represents the high-level target network. Although this biases the high-level value function
 185 towards being optimistic in stochastic environments, future work could incorporate a notion of
 186 reachability from the low-level value function.

187 **High-Level Policy.** Again, ARL is agnostic to the choice of high-level policy extraction. We
 188 hypothesise the high-level policy to be multi-modal, corresponding to distinct and equally optimal
 189 options, but note that choice of high-level policy is orthogonal to this work. A high-level Q-function
 190 could also be fitted to the high-level value function, which would enable use of Behaviour-Cloned
 191 Deep Deterministic Policy Gradient (DDPGBC) for policy extraction [28], for example; we include
 192 details in Appendix C. In our implementations, we use a Gaussian high-level policy, which we
 193 generally (see Appendix E) train using AWR:

$$\mathcal{L}_{\pi_h} = -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p_{\gamma}^{\mathcal{D}}(\cdot|s)} \left[e^{\alpha_h A_h(s, \hat{\phi}_{\omega}(s, g_s), g_s, g)} \log \pi_h \left(\hat{\phi}_{\omega}(s, g_s) \mid s, g \right) \right], \quad (6)$$

194 where A_h represents the advantage associated with option $\hat{\phi}_{\omega}(s, g_s)$. We provide pseudocode in
 195 Algorithm 1 in Appendix C. Full experiment details, hyperparameters, sampling methods and seeds
 196 are found in Appendices C and E, and in our codebase.

197 4.3 Abstractive RL Explicitly Enforcing Translation Invariance

198 We now introduce a second algorithm, which explicitly imposes translation invariance on state-
 199 waypoint representations in the low-level decision process in order to learn from similar contexts
 200 across the state-space. We hypothesise that such an inductive bias could be useful in manipulation
 201 tasks, for example. We propose this algorithm as a proof of concept that using different representations
 202 at different levels of the hierarchy can improve generalisation in Offline RL.

203 We define *relativised states* as an unnormalised displacement vector:

$$v = g_s - s,$$

204 resulting in a single vector that simultaneously encodes both the state and waypoint.

205 Since hard-coding representations can impose representational constraints, our approach exploits
 206 the two-step procedure of Offline RL as a compromise. To enforce experience reuse, we define the
 207 low-level value function in terms of relativised states: $V_l(g_s - s)$. Although this introduces a repre-
 208 sentational constraint (where state-waypoint pairs with distinct values may map to the same relative

vector¹), it explicitly collapses the state-waypoint space into a manifold of relative displacements, addressing the second and fourth points in the Motivation Box (Section 4.1). We remark that the issue of representational constraints could also be mitigated by computing the difference of encoded representations such that $v = \phi_{l_{g_s}}(g_s) - \phi_{l_s}(s)$, although we did not find this to be necessary to achieve superior performance in our experiments.

As in Section 4.2, to relativise options and address the third point in the Motivation Box, option-embeddings are learned with the low-level policy. However, now options are also explicitly relativised by defining them in terms of relativised states. We use soft-normalisation rather than length-normalisation to avoid numerical instability while introduce geometric regularisation and allow re-normalising of samples from the high-level policy [24] upon deployment:

$$o := \hat{\phi}_\omega(s, g_s) = \frac{\phi_\omega(g_s - s) \cdot \tanh \|\phi_\omega(g_s - s)\|}{\|\phi_\omega(g_s - s)\|} \cdot \sqrt{d},$$

where d is the dimension of the embedding ϕ_ω , and $\|\cdot\|$ denotes the standard Euclidean norm. Soft normalisation means that the option space includes the space within the hypersphere and allows options to incorporate implicit temporal awareness as their magnitude scales linearly when displacement is small.

To satisfy local constraints, we still condition the low-level policy on the absolute state: $\pi_l(\cdot | s, \hat{\phi}_\omega(s, g_s))$. We provide pseudocode in Algorithm 2 in Appendix C. Full experiment details, hyperparameters, sampling methods and seeds are found in Appendices C and E and in our codebase.

5 Experiments

The goal of our experiments is simple: to test whether relativised options and distinct representations at different levels of the hierarchy can lead to better policy generalisation in offline GCRL.

Benchmark and Ablations. We perform all experiments on the standard OGBench datasets, focusing on the more challenging locomotion and manipulation environments (i.e. selecting *giant* over *medium* or *large*). For completeness, we also evaluate on a stochastic setting (*teleport*), despite its mismatch with our deterministic assumption (Section 2). We exclude visual environments as they introduce additional challenges related to high-dimensional perception that are orthogonal to this work. Unlike prior work [12, 48], we do not use oracle representations, which simplify option learning in locomotion, and discard proprioceptive information that might be useful in manipulation.

We benchmark ARL with implicitly learned relativised options (Section 4.2, **ARLi**), and explicitly enforced translation invariance (Section 4.3, **ARLe**), against the original version of HIQL [24] (**HIQL1vr**), which uses a single value function and learns option representations via this value function. To isolate the effect of the relativised representation rather than any differences arising due to structure of the value function (ARL employs two value functions), we compare against variants of HIQL that also use two value functions. We include a variant with two value functions that does not include option representation learning (**HIQL2vr**), and a variant with two value functions that learns option representations via the low-level value function (**HIQL2vr**), with the intention of mirroring HIQL1vr. Finally, we also compare against goal-conditioned IQL [7], the best-performing flat policy from Park et al. [5] (**IQL**). We refer to Appendix C and our codebase for full implementation details.

Since hyperparameter tuning is expensive and ARL, which is based on inductive biases, introduces no additional hyperparameters over HIQL, we simply adopt those tuned for HIQL [5, 12] (see Appendix E). The fact that ARL achieves strong performance under these hyperparameters attests to its efficacy. ARL is agnostic to the choice of policy class and value learning objective, so, to avoid related confounding factors, we use a Gaussian and one-step TD for all experiments.

Results. Our results (Table 1 and left of Figure 2) show that both variants of ARL outperform both the flat and absolute hierarchical policies, achieving a mean success rate that is at least 10 percentage points higher than the benchmarks when aggregating over all tasks. We highlight that this is without necessitating hyperparameter tuning.

¹For instance, in a maze, a state s and waypoint g_s separated by a wall may map to the same relative vector as a pair in open space, yet induce vastly different value estimates.

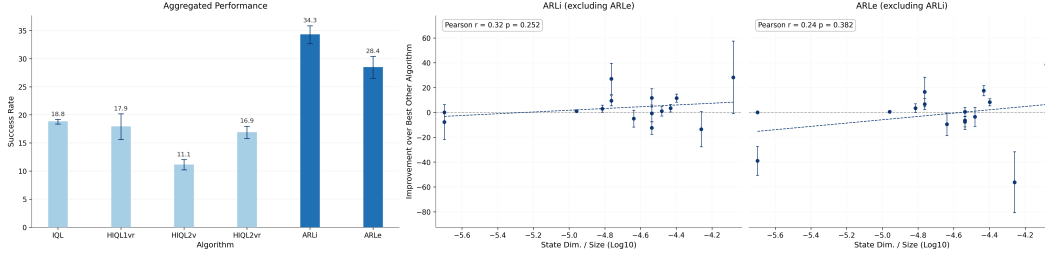


Figure 2: **Analysis.** Aggregate Performance across all tasks (**left**) and ARLi’s (**middle**) and ARLe’s (**right**) performance improvements over next-best performing algorithm against number of state dimensions per dataset sample. Bootstrapped 95% CI over 4 seeds and 20 evaluation runs.

255 Simply by learning relativised options based on action similarity rather than value similarity, ARLi
 256 consistently improves over the hierarchical benchmarks (excluding ARLe). In particular, it more
 257 than doubles the success rate in both the *humanoidmaze* environments, and almost doubles it for
 258 *scene-play-v0*.

259 ARLe performs especially well in the manipulation environments. Most notably, in *puzzle-4x4-play-*
 260 *v0*, an 83-dimensional manipulation task, ARLe achieves an 88% success rate, outperforming the
 261 next-best benchmark by 39 percentage points (excluding ARLi). We hypothesise that translational
 262 invariance is particularly beneficial in high-dimensional settings with underlying symmetries and
 263 sparse state-space coverage. To investigate this, we plot improvement over the next-best benchmark²
 264 against state dimensionality normalised by dataset size (right of Figure 2). Although correlation does
 265 not imply causation, and the observed correlations are weak and not statistically significant (amplified
 266 by a small number of environments), both ARLi and ARLe exhibit positive improvement trends with
 267 increasing dimensional sparsity. In comparison, HIQL1vr and HIQL2vr, for example, show stronger
 268 negative trends with increasing sparsity (Figure 6 in Appendix D) under the same methodology.
 269 Intuitively, it makes sense that relativised options and experience reuse become increasingly important
 270 as sparsity increases.

Table 1: **Results.** We report each method’s average (binary) success rate (%) across the five test-time goals. Bootstrapped 95% CI over 4 seeds and 20 evaluation runs. Blue bold indicates the highest mean; black bold overlapping confidence intervals.

Task	Size	Dim.	IQL	HIQL1vr	HIQL2v	HIQL2vr	ARLi	ARLe
pointmaze-giant-navigate-v0	1M	2	0±0	55±11	21±7	46±6	48±8	16±2
pointmaze-giant-stitch-v0	1M	2	0±0	0±0	0±0	0±0	0±0	0±0
antmaze-giant-navigate-v0	1M	29	0±0	38±3	40±3	48±6	48±3	42±4
antmaze-giant-stitch-v0	1M	29	0±0	7±7	15±3	20±3	32±7	21±1
antmaze-teleport-stitch-v0	1M	29	49±2	29±4	43±3	40±6	36±5	41±4
humanoidmaze-giant-navigate-v0	4M	69	1±1	22±11	12±5	11±10	49±6	38±4
humanoidmaze-giant-stitch-v0	4M	69	0±0	4±3	0±0	0±0	13±2	10±3
cube-double-play-v0	1M	37	50±3	2±0	0±0	3±0	53±2	67±3
cube-triple-play-v0	3M	46	11±2	7±3	0±0	1±1	14±2	15±3
cube-quadruple-play-v0	5M	55	0±0	0±0	0±0	0±0	1±0	0±0
puzzle-3x3-play-v0	1M	55	100±0	27±8	0±0	24±3	86±14	44±25
puzzle-4x4-play-v0	1M	83	30±3	49±27	0±0	17±6	78±11	88±6
puzzle-4x5-play-v0	3M	99	15±3	17±3	0±0	12±3	18±3	14±7
puzzle-4x6-play-v0	5M	115	13±1	0±0	0±0	18±2	14±7	9±9
scene-play-v0	1M	40	13±2	12±3	12±7	13±1	24±3	21±3

²We try to mitigate confounding factors such as horizon length, absolute dataset size and policy expressivity (e.g. whether unimodal or multimodal) by plotting performance gains over the next-best algorithm rather than absolute success rate.

271 To better understand the effect of imposing translational invariance, we visualise the
 272 low-level value functions for ARLe and HIQL2v in the *antmaze* locomotion environ-
 273 ment (Figure 3). Due to explicitly collapsing equivalent relative states, ARLe learns
 274 a substantially smoother low-level value function. We now address potential questions.
 275

276 **Why not tune hyperparameters for**
 277 **ARL?** Offline RL typically requires
 278 significant online tuning [51], which
 279 is expensive and would be unscal-
 280 able for training a billion-parameter
 281 general-purpose agent [52, 53]. By
 282 using inductive biases rather than rep-
 283 resentational losses, we avoid intro-
 284 ducing additional hyperparameters:
 285 ARL achieving significant perfor-
 286 mance gains under hyperparameters
 287 tuned for HIQL indicates robustness.

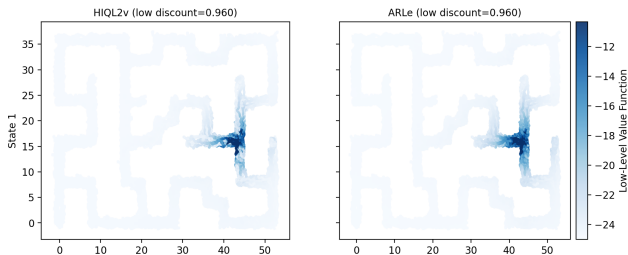


Figure 3: **Low-level** value function for **HIQL2v** (left) and **ARLe** (right) (task 4, *antmaze-giant-stitch-v0*).

288 **Why introduce ARLe if ARLi is so consistent?** ARLe demonstrates how hierarchical structures
 289 and relativised options enable level-specific representations, leveraging this inductive bias to improve
 290 performance by 10 percentage points over ARLi in two out of the four sparsest datasets (*puzzle-4x4-*
 291 *play-v0* and *cube-double-play-v0*).

292 **Why does ARLe perform poorly in certain tasks?** ARLe’s performance depends on the alignment
 293 between its inductive bias and the environment’s structure. Assuming local translational invariance
 294 can benefit high-dimensional manipulation through experience reuse but can be detrimental in dense,
 295 low-dimensional environments like *pointmaze*, since the representation is inherently lossy. We also
 296 hypothesise that mapping 2D relative displacements into a 10D latent hypersphere introduces repre-
 297 sentational noise. While future work could leverage GDL to learn more flexible symmetries, ARLe
 298 demonstrates that decoupling representations across the hierarchy enables a degree of experience
 299 reuse fundamentally inaccessible to flat or absolute-frame architectures. When the inductive bias is
 300 well-matched to the environment, it significantly enhances policy generalisation.

301 **Why do all algorithms perform poorly in certain tasks?** We hypothesise this to be due to other
 302 confounding factors such as: (i) poor high-level value learning and a lack of gradient in long horizon
 303 tasks (we include plots of the high-level value function in Figure 5 in Appendix D); (ii) the high-level
 304 policy being unimodal rather than multimodal; (iii) not training for enough gradient steps for large
 305 dataset sizes; (iv) not learning goal representations for the high-level policy. The aforementioned
 306 issues could be mitigated by combining ARL with value horizon reduction methods such as TD-
 307 *n* or TRL [48], learning a flow-policy [54] rather than a Gaussian (especially for the high-level
 308 policy), training for more steps (we run all experiments for 1M), and using, for example, Dual-Goal
 309 Representations [55] for the high-level decision process.

310 6 Conclusion

311 In this work we motivate hierarchy in offline RL through *absolute* abstraction. By learning relativised
 312 options and using distinct representations at different levels of the hierarchy, agents can reuse optimal
 313 experience across similar contexts of the state-space, enabling better performance in regions of the
 314 dataset only supported by low-quality data. Based on our framework, we introduce two simple
 315 algorithms for learning relativised options via action similarity and explicitly enforcing translational
 316 invariance on the low-level decision process. Our experiments demonstrate that such relativised
 317 options and inductive biases improve policy generalisation in high-dimensional offline GCRL. This
 318 proof of concept opens many avenues for future research, including imposing more flexible inductive
 319 biases, or leveraging action-chunking to learn relativised options over action sequences. We hope
 320 that this work motivates progress towards scalable offline RL.

References

- 321
- 322 [1] Leslie Pack Kaelbling. Learning to achieve goals. In *Proceedings of the Thirteenth International*
323 *Joint Conference on Artificial Intelligence (IJCAI)*, pages 1094–1099, 1993.
- 324 [2] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function approxi-
325 mators. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37
326 of *Proceedings of Machine Learning Research*, pages 1312–1320, 2015.
- 327 [3] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline Reinforcement Learning:
328 Tutorial, Review, and Perspectives on Open Problems, November 2020. URL [http://arxiv.](http://arxiv.org/abs/2005.01643)
329 [org/abs/2005.01643](http://arxiv.org/abs/2005.01643). arXiv:2005.01643 [cs].
- 330 [4] Sergey Levine. Understanding the World Through Action, October 2021. URL [http://arxiv.](http://arxiv.org/abs/2110.12543)
331 [org/abs/2110.12543](http://arxiv.org/abs/2110.12543). arXiv:2110.12543 [cs].
- 332 [5] Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. OGBench: Bench-
333 marking Offline Goal-Conditioned RL, February 2025. URL [http://arxiv.org/abs/2410.](http://arxiv.org/abs/2410.20092)
334 [20092](http://arxiv.org/abs/2410.20092). arXiv:2410.20092 [cs].
- 335 [6] Rafael Figueiredo Prudencio, Marcos R. O. A. Maximo, and Esther Luna Colombini. A
336 Survey on Offline Reinforcement Learning: Taxonomy, Review, and Open Problems. *IEEE*
337 *Transactions on Neural Networks and Learning Systems*, 35(8):10237–10257, August 2024.
338 ISSN 2162-237X, 2162-2388. doi: 10.1109/TNNLS.2023.3250269. URL [http://arxiv.](http://arxiv.org/abs/2203.01387)
339 [org/abs/2203.01387](http://arxiv.org/abs/2203.01387). arXiv:2203.01387 [cs].
- 340 [7] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline Reinforcement Learning with Implicit
341 Q-Learning, October 2021. URL <http://arxiv.org/abs/2110.06169>. arXiv:2110.06169
342 [cs].
- 343 [8] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-Weighted Re-
344 gression: Simple and Scalable Off-Policy Reinforcement Learning, October 2019. URL
345 <http://arxiv.org/abs/1910.00177>. arXiv:1910.00177 [cs].
- 346 [9] Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting
347 the Minimalist Approach to Offline Reinforcement Learning, October 2023. URL [http:](http://arxiv.org/abs/2305.09836)
348 [//arxiv.org/abs/2305.09836](http://arxiv.org/abs/2305.09836). arXiv:2305.09836 [cs].
- 349 [10] Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of Real-
350 World Reinforcement Learning, April 2019. URL <http://arxiv.org/abs/1904.12901>.
351 arXiv:1904.12901 [cs].
- 352 [11] Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is Value Learning Really the
353 Main Bottleneck in Offline RL?, October 2024. URL <http://arxiv.org/abs/2406.09329>.
354 arXiv:2406.09329 [cs].
- 355 [12] Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey
356 Levine. Horizon Reduction Makes RL Scalable, October 2025. URL [http://arxiv.org/](http://arxiv.org/abs/2506.04168)
357 [abs/2506.04168](http://arxiv.org/abs/2506.04168). arXiv:2506.04168 [cs].
- 358 [13] Richard S. Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A
359 framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112
360 (1-2):181–211, August 1999. ISSN 00043702. doi: 10.1016/S0004-3702(99)00052-1. URL
361 <https://linkinghub.elsevier.com/retrieve/pii/S0004370299000521>.
- 362 [14] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg,
363 David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning.
364 In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of
365 *Proceedings of Machine Learning Research*, pages 3540–3549, 2017.
- 366 [15] Balaraman Ravindran and Andrew G. Barto. Model minimization in hierarchical reinforcement
367 learning.

- 368 [16] Kevin Black, Manuel Y. Galliker, and Sergey Levine. Real-Time Execution of Action
369 Chunking Flow Policies, December 2025. URL <http://arxiv.org/abs/2506.07339>.
370 arXiv:2506.07339 [cs].
- 371 [17] Kwanyoung Park, Seohong Park, Youngwoon Lee, and Sergey Levine. Scalable Offline Model-
372 Based RL with Action Chunks, December 2025. URL <http://arxiv.org/abs/2512.08108>.
373 arXiv:2512.08108 [cs].
- 374 [18] Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric deep learning:
375 Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- 376 [19] Arsh Tangri, Nichols Crawford Taylor, Haojie Huang, and Robert Platt. Equivariant goal
377 conditioned contrastive reinforcement learning, 2025. doi: 10.48550/arXiv.2507.16139.
- 378 [20] Sascha Lange, Thomas Gabel, and Martin A. Riedmiller. Batch reinforcement learning. Springer,
379 2012.
- 380 [21] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press,
381 2 edition, 2018.
- 382 [22] Adam White, Joseph Modayil, and Richard S. Sutton. Scaling life-long off-policy learning.
383 *CoRR*, abs/1206.6262, 2012. URL <http://arxiv.org/abs/1206.6262>.
- 384 [23] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,
385 Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience Replay,
386 February 2018. URL <http://arxiv.org/abs/1707.01495>. arXiv:1707.01495 [cs].
- 387 [24] Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline Goal-
388 Conditioned RL with Latent States as Actions, March 2024. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2307.11949)
389 [2307.11949](http://arxiv.org/abs/2307.11949). arXiv:2307.11949 [cs].
- 390 [25] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-Learning for
391 Offline Reinforcement Learning, August 2020. URL <http://arxiv.org/abs/2006.04779>.
392 arXiv:2006.04779 [cs].
- 393 [26] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-Based Offline
394 Reinforcement Learning with Diversified Q-Ensemble, October 2021. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2110.01548)
395 [2110.01548](http://arxiv.org/abs/2110.01548). arXiv:2110.01548 [cs].
- 396 [27] Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov, and Sergey Levine. Contrastive
397 Learning as Goal-Conditioned Reinforcement Learning, February 2023. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2206.07568)
398 [2206.07568](http://arxiv.org/abs/2206.07568). arXiv:2206.07568 [cs].
- 399 [28] Scott Fujimoto and Shixiang Shane Gu. A Minimalist Approach to Offline Reinforcement
400 Learning, December 2021. URL <http://arxiv.org/abs/2106.06860>. arXiv:2106.06860
401 [cs].
- 402 [29] Seyed Kamyar Seyed Ghasemipour, Dale Schuurmans, and Shixiang Shane Gu. Emaq:
403 Expected-max q-learning operator for simple yet effective offline and online RL. *CoRR*,
404 abs/2007.11091, 2020. URL <https://arxiv.org/abs/2007.11091>.
- 405 [30] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The Option-Critic Architecture, December
406 2016. URL <http://arxiv.org/abs/1609.05140>. arXiv:1609.05140 [cs].
- 407 [31] Seungho Baek, Taegeon Park, Jongchan Park, Seungjun Oh, and Yusung Kim. Graph-Assisted
408 Stitching for Offline Hierarchical Reinforcement Learning, June 2025. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2506.07744)
409 [2506.07744](http://arxiv.org/abs/2506.07744). arXiv:2506.07744 [cs] version: 1.
- 410 [32] Richard S Sutton, Doina Precup, and Satinder Singh. Intra-Option Learning about Temporally
411 Abstract Actions.
- 412 [33] Balaraman Ravindran and Andrew G. Barto. Smdp homomorphisms: An algebraic approach to
413 abstraction in semi-markov decision processes. In *Probabilistic Planning*, pages 1011–1016,
414 2003.

- 415 [34] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-Efficient Hierarchical Reinforcement Learning, October 2018. URL <http://arxiv.org/abs/1805.08296>.
 416 arXiv:1805.08296 [cs].
 417
- 418 [35] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-Optimal Representation
 419 Learning for Hierarchical Reinforcement Learning, January 2019. URL <http://arxiv.org/abs/1810.01257>.
 420 arXiv:1810.01257 [cs].
- 421 [36] Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning Multi-Level Hierarchies with Hindsight, September 2019. URL <http://arxiv.org/abs/1712.00948>.
 422 arXiv:1712.00948 [cs].
 423
- 424 [37] Amy McGovern and Andrew G. Barto. Automatic discovery of subgoals in reinforcement
 425 learning using diverse density. In *Proceedings of the 18th International Conference on Machine*
 426 *Learning (ICML)*, pages 361–368, 2001.
- 427 [38] Martin Stolle and Doina Precup. Learning options in reinforcement learning. In *Proceedings of*
 428 *the 5th International Symposium on Abstraction, Reformulation and Approximation (SARA)*,
 429 pages 212–223, 2002.
- 430 [39] Jinning Li, Chen Tang, Masayoshi Tomizuka, and Wei Zhan. Hierarchical planning through
 431 goal-conditioned offline reinforcement learning, 2022. URL [https://arxiv.org/abs/2205.](https://arxiv.org/abs/2205.11790)
 432 [11790](https://arxiv.org/abs/2205.11790).
- 433 [40] Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a Unified Theory of State
 434 Abstraction for MDPs.
- 435 [41] Norman Ferns, Prakash Panangaden, and Doina Precup. Metrics for Finite Markov Decision
 436 Processes, July 2012. URL <http://arxiv.org/abs/1207.4114>. arXiv:1207.4114 [cs].
- 437 [42] Jianda Chen and Sinno Jialin Pan. Learning Representations via a Robust Behavioral Metric for
 438 Deep Reinforcement Learning.
- 439 [43] Ayoub Echchahed and Pablo Samuel Castro. A Survey of State Representation Learning
 440 for Deep Reinforcement Learning, June 2025. URL <http://arxiv.org/abs/2506.17518>.
 441 arXiv:2506.17518 [cs].
- 442 [44] Sham Kakade. *On the Sample Complexity of Reinforcement Learning*. Phd thesis, University
 443 College London, 2003.
- 444 [45] Remi Munos, Remi Munos, and Csaba Szepesvari. Finite-Time Bounds for Fitted Value
 445 Iteration.
- 446 [46] Tor Lattimore and Marcus Hutter. PAC Bounds for Discounted MDPs, February 2012. URL
 447 <http://arxiv.org/abs/1202.3890>. arXiv:1202.3890 [cs].
- 448 [47] Arnaud Robert, Ciara Pike-Burke, and Aldo Faisal. Sample Complexity of Goal-Conditioned
 449 Hierarchical Reinforcement Learning.
- 450 [48] Seohong Park, Aditya Oberai, Pranav Atreya, and Sergey Levine. Transitive RL: Value Learning
 451 via Divide and Conquer, February 2026. URL <http://arxiv.org/abs/2510.22512>.
 452 arXiv:2510.22512 [cs].
- 453 [49] Dibya Ghosh, Chethan Bhateja, and Sergey Levine. Reinforcement Learning from Passive
 454 Data via Latent Intentions, April 2023. URL <http://arxiv.org/abs/2304.04782>.
 455 arXiv:2304.04782 [cs].
- 456 [50] Haoran Xu, Li Jiang, Jianxiong Li, and Xianyuan Zhan. A policy-guided imitation approach for
 457 offline reinforcement learning, 2023. URL <https://arxiv.org/abs/2210.08323>.
- 458 [51] Matthew Thomas Jackson, Uljad Berdica, Jarek Liesen, Shimon Whiteson, and Jakob Nicolaus
 459 Foerster. A Clean Slate for Offline Reinforcement Learning, April 2025. URL [http://arxiv.](http://arxiv.org/abs/2504.11453)
 460 [org/abs/2504.11453](http://arxiv.org/abs/2504.11453). arXiv:2504.11453 [cs].

- 461 [52] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
462 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
463 models. *CoRR*, abs/2001.08361, 2020. URL <https://arxiv.org/abs/2001.08361>.
- 464 [53] Kevin Wang, Ishaan Javali, Michał Borkiewicz, Tomasz Trzciński, and Benjamin Eysenbach.
465 1000 Layer Networks for Self-Supervised RL: Scaling Depth Can Enable New Goal-Reaching
466 Capabilities, February 2026. URL <http://arxiv.org/abs/2503.14858>. arXiv:2503.14858
467 [cs].
- 468 [54] Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky
469 T. Q. Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow Matching Guide and Code,
470 December 2024. URL <http://arxiv.org/abs/2412.06264>. arXiv:2412.06264 [cs].
- 471 [55] Seohong Park, Deepinder Mann, and Sergey Levine. Dual Goal Representations, February
472 2026. URL <http://arxiv.org/abs/2510.06714>. arXiv:2510.06714 [cs].
- 473 [56] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL
474 <https://arxiv.org/abs/1412.6980>.
- 475 [57] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL
476 <https://arxiv.org/abs/1607.06450>.
- 477 [58] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with
478 gaussian error linear units. *CoRR*, abs/1606.08415, 2016. URL <http://arxiv.org/abs/1606.08415>.
- 480 [59] Adam Villaflor, Zhe Huang, Swapnil Pande, John Dolan, and Jeff Schneider. Addressing
481 optimism bias in sequence modeling for reinforcement learning, 2022. URL <https://arxiv.org/abs/2207.10295>.
- 482

483 A Sample Complexity in Online Goal-Conditioned RL

484 We provide some intuition into the choice of policy learning and representation using sample
485 complexity in finite state and action space problems. We build on the works of [44–47, 40].

486 In online GCRL, the aim is to find the optimal policy with the smallest number of samples or online
487 interactions, N , such that the error in optimal return is smaller than a fixed constant ϵ .

488 A.1 GCRL Sample Complexity

489 Consider the case of a discrete, discounted horizon MDP with a finite state space \mathcal{S} , action space
490 \mathcal{A} and discount factor $\gamma \in [0, 1)$. The Probably-Approximately Correct (PAC) Learning [44] upper-
491 bound sample complexity to find an ϵ optimal policy reaching a unique goal-state optimally (assuming
492 at most two possible next-states for each state/action pair) with constant probability of $1 - \delta$ is given
493 by [46]:

$$N^{\text{infinite, single goal}} \propto \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2}.$$

494 Hence, the minimax sample complexity required to find an ϵ optimal policy reaching any given state
495 (such that we have $|\mathcal{G}|$ unique goal-states) is given by:

$$N^{\text{infinite}} \propto \frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2}. \quad (7)$$

496 Consider the case of a discrete finite horizon (of length H) MDP with a finite state space \mathcal{S} , action
497 space \mathcal{A} and discount factor $\gamma \in [0, 1)$. The minimax sample complexity to find an ϵ optimal policy
498 reaching a unique goal-state optimally is given by:

$$N^{\text{finite}} \propto \frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}|H^3}{\epsilon^2}.$$

499 A.2 GCRL Hierarchical Sample Complexity

500 Using a hierarchical policy, we can break the distant goal g from our current state s into options
501 defined over an option space Ω . Hence, the high-level policy becomes

$$N^{\text{high}} \propto \frac{|\mathcal{S}||\mathcal{G}||\Omega|}{(1-\gamma^n)^3 \epsilon^2},$$

502 where we have substituted the option-space to Equation 7, and use the environment’s discount factor
503 raised to a factor of n , assuming that the high-level policy acts, on average, every n steps.

504 The low-level policy has a sample complexity of

$$N^{\text{low}} \propto \frac{|\mathcal{S}||\Omega||\mathcal{A}|n^3}{\epsilon^2},$$

505 since $|\Omega|$ options can be executed. Hence, the overall sample complexity of the hierarchical policy
506 $\pi(a | s, g) := \pi_l(a | s, \omega)$, $\omega \sim \pi_h(\cdot | s, g)$ is [47]:

$$N^{\text{hierarchy}} \propto \frac{|\mathcal{S}||\mathcal{G}||\Omega|}{(1-\gamma^n)^3 \epsilon^2} + \frac{|\mathcal{S}||\Omega||\mathcal{A}|n^3}{\epsilon^2}. \quad (8)$$

507 Note, that, in the case of no temporal abstraction, when $n = 1$, we approximately recover the original
508 sample complexity of a flat policy. Since then the option just becomes a single primitive action,
509 Equation 8 becomes:

$$N^{\text{hierarchy}} \propto \frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2} + \frac{|\mathcal{S}||\mathcal{A}|^2}{\epsilon^2} \approx \frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2},$$

510 where the approximation follows due to the dominating $\frac{1}{(1-\gamma)^3}$ factor in the first term, assuming
511 long-horizon problems such that $\gamma \rightarrow 1$ and that the goal-space is larger or equal to the size of the
512 action space $|\mathcal{G}| \geq |\mathcal{A}|^3$.

³Even though this might not be the case, usually the goal-space is unknown, so we train the policy to reach any state within the state-space i.e. such that $\mathcal{G} = \mathcal{S}$. Assuming that $|\mathcal{S}| \gg |\mathcal{A}|$ is a standard assumption.

513 Issues arise under a misspecified option horizon n . In this case, the sample complexity of the
 514 hierarchical policy becomes worse than that of a flat policy, as the skill space must then account for
 515 every sequence of primitive actions over n steps, such that $|\Omega| = |\mathcal{A}|$. The sample complexity of the
 516 hierarchical policy becomes

$$N^{\text{hierarchy}} \propto \frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}|^n}{(1-\gamma^n)^3\epsilon^2} + \frac{|\mathcal{S}||\Omega||\mathcal{A}|^n n^3}{\epsilon^2},$$

517 which blows up due to the exponential factor multiplying the action space.

518 A.3 GCRL State Representation Sample Complexity

519 State representation sample complexity exploits symmetry in the state-goal space to a mapping
 520 $\phi(s, g) \rightarrow c$ such that $\phi(s, g) = \phi(s', g')$ if $\pi^*(a | s, g) = \pi^*(a | s', g')$. The sample complexity of
 521 learning an optimal policy (Equation 7) becomes:

$$N^{\text{rep}} \propto \frac{|\mathcal{C}||\mathcal{A}|}{(1-\gamma)^3\epsilon^2}. \quad (9)$$

522 Note that $|\mathcal{C}| \leq |\mathcal{S}||\mathcal{G}|$, with $|\mathcal{C}| = |\mathcal{S}||\mathcal{G}|$ if each (s, g) has a distinct optimal action.

523 B Error in Offline Goal-Conditioned RL

524 In offline GCRL, the aim is to bound the error ϵ given an offline dataset \mathcal{D} of fixed size N .

525 Unlike in the online setting, where it is assumed that the agent can sample any state-action pair to
 526 learn the environment’s dynamics, in offline RL, a concentrability coefficient κ is incorporated to
 527 account for the distribution shift in data collected by the policy π^{BC} , and the data that would have
 528 been collected induced under the optimal policy π^* . Intuitively, if the dataset does not include the
 529 states required to learn the optimal policy, the algorithm may never learn that optimal policy. The
 530 concentrability coefficient κ is defined as:

$$\kappa = \sup_{s \in \mathcal{S}, a \in \mathcal{A}, g \in \mathcal{G}} \frac{d^{\pi^*}(s, a, g)}{d^{\pi^{\text{BC}}}(s, a, g)},$$

531 where $d^\pi(s, a, g)$ is the discounted occupancy measure (i.e. stationary distribution) of the policy π :

$$d^\pi(s, a, g) = (1-\gamma)\mathbb{E}_{\tau \sim \pi, s_0, g_0 \sim \text{Unif}(\mathcal{S})} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{1}(s_t = s, a_t = a, g_0 = g) \right],$$

532 and where the trajectory is induced by the MDP and following the policy $\pi(a | s, g)$. If the dataset
 533 is highly exploratory and covers the optimal paths well, κ will be small. If the dataset is narrow, or
 534 misses critical regions of the state-action space, κ will be large.

535 Rearranging Equation 7 and incorporating the concentrability coefficient, the offline bound for a
 536 goal-conditioned flat policy is given by:

$$\epsilon \propto \sqrt{\frac{|\mathcal{S}||\mathcal{G}||\mathcal{A}| \cdot \kappa}{(1-\gamma)^3 N}}.$$

537 B.1 GCRL Hierarchical Error

538 Since the size of the offline dataset is fixed, the only way to reduce the error ϵ is to use a hierarchical
 539 policy with distinct state-goal representations.

540 Using a hierarchical policy, the error term becomes:

$$\epsilon^{\text{hierarchy}} \propto \sqrt{\frac{|\mathcal{S}||\mathcal{G}||\Omega| \cdot \kappa_h}{(1-\gamma^n)^3 N}} + \sqrt{\frac{|\mathcal{S}||\Omega||\mathcal{A}| n^3 \cdot \kappa_l}{N}},$$

541 where the concentrability coefficients are defined as:

$$\kappa_h = \sup_{s \in \mathcal{S}, \omega \in \Omega, g \in \mathcal{S}} \frac{d^{\pi_h^*}(s, \omega, g)}{d^{\pi_h^{\text{BC}}}(s, \omega, g)} \quad \text{and} \quad \kappa_l = \sup_{s \in \mathcal{S}, a \in \mathcal{A}, \omega \in \Omega} \frac{d^{\pi_l^*}(s, a, \omega)}{d^{\pi_l^{\text{BC}}}(s, a, \omega)}.$$

542 Then, introducing two embeddings that group together state-goal pairs (s, g) requiring similar options
 543 ω such that $\phi_h(s, g) = \phi_h(s', g') = c_h$ if $\pi_h^*(\cdot | s, g) = \pi_h^*(\cdot | s', g')$ and $c_h \in \mathcal{C}_h$, and state-option
 544 (s, ω) requiring similar low-level actions $\phi_l(s, \omega) = \phi_l(s', \omega') = c_l$ if $\pi_l^*(\cdot | s, \omega) = \pi_l^*(\cdot | s', \omega')$ ⁴
 545 and $c_l \in \mathcal{C}_l$, the error term is bounded by:

$$\epsilon^{\text{hierarchy, rep}} \propto \sqrt{\frac{|\mathcal{C}_h||\Omega| \cdot \kappa_h}{(1 - \gamma^n)^3 N}} + \sqrt{\frac{|\mathcal{C}_l||\mathcal{A}|n^3 \cdot \kappa_l}{N}}.$$

546 By definition, $|\mathcal{C}_h||\Omega| \leq |\mathcal{S}||\mathcal{G}||\Omega|$ and $|\mathcal{C}_l||\mathcal{A}| \leq |\mathcal{S}||\Omega||\mathcal{A}|$.

547 Because the concentrability coefficients are now defined using the reparameterised policies over the
 548 reparameterised latent spaces \mathcal{C}_h and \mathcal{C}_l , i.e.

$$\kappa_h^{\text{rep}} = \sup_{c_h \in \mathcal{C}_h, \omega \in \Omega} \frac{d^{\pi_h^*}(c_h, \omega)}{d^{\pi_h^{\text{BC}}}(c_h, \omega)} \quad \text{and} \quad \kappa_l^{\text{rep}} = \sup_{c_l \in \mathcal{C}_l, a \in \mathcal{A}} \frac{d^{\pi_l^*}(c_l, a)}{d^{\pi_l^{\text{BC}}}(c_l, a)},$$

549 the probability mass of the offline dataset is aggregated across state-goal or state-option pairs that are
 550 equivalent under these embeddings:

$$d^{\pi_h}(c_h, \omega) = \sum_{(s, g) \in \phi_h^{-1}(c_h)} d^{\pi_h}(s, g, \omega) \quad \text{and} \quad d^{\pi_l}(c_l, a) = \sum_{(s, \omega) \in \phi_l^{-1}(c_l)} d^{\pi_l}(s, \omega, a).$$

551 This reduces the likelihood of a support mismatch, where the optimal policy requires a state transition
 552 on which the dataset places zero mass. Hence,

$$\kappa_h^{\text{rep}} \leq \kappa_h \quad \text{and} \quad \kappa_l^{\text{rep}} \leq \kappa_l.$$

⁴Note that we should not use these same embeddings for the value or critic functions, since even though $\pi_l^*(\cdot | s, \omega) = \pi_l^*(\cdot | s', \omega')$, this does not instantly imply that $Q_h(s, g, \omega) = Q_h(s', g', \omega')$: even though the ordering over actions might be the same, there might be an offset such that $Q_h(s, g, \omega) = Q_h(s', g', \omega) + c(s, g, s', g') \quad \forall \omega$. We refer the reader to Li et al. [40].

553 C Offline RL Algorithms

554 In the following section g_s is a waypoint to the goal, such that $g_s \in \mathcal{S}$. When sampled from the
 555 dataset \mathcal{D} , g_s is n steps ahead of the current state s . When sampling the waypoint g_s or goal g from
 556 $p^{\mathcal{D}}(\cdot | s, a)$, we sample either from a geometric distribution according to the specified discount factor,
 557 or a uniform distribution. Details are given in Appendix E.

558 C.1 Implicit Q-Learning (IQL)

559 The flat policy we benchmark is IQL [7], which trains a state-goal value function $V(s, g)$ and
 560 state-goal-action value function $Q(s, g, a)$ using the following losses:

$$\mathcal{L}_V = \mathbb{E}_{(s,a) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[\ell_{\tau}^2 \left(V(s, g) - \tilde{Q}(s, g, a) \right) \right],$$

561 where \tilde{Q} denotes the target network, and ℓ_{τ}^2 denotes the expectile loss $\ell_{\tau}^2(x) = |\tau - (x < 0)|x^2$, and

$$\mathcal{L}_Q = \mathbb{E}_{(s,a,s') \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[(Q(s, a, g) - r(s, g) - \gamma V(s', g))^2 \right].$$

562 A Gaussian policy is then extracted using the following DDPGBC [28] loss:

$$\mathcal{L}_{\pi}^{\text{DDPGBC}} = -\mathbb{E}_{(s,a,s') \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} [Q(s, g, \mu^{\pi}(s, g)) + \alpha \log \pi(a | s, g)], \quad (10)$$

563 which has been found to outperform AWR [11].

564 C.2 Hierarchical Implicit Q-Learning 1 Value Function with Representation Learning 565 (HIQL1vr)

566 We benchmark against HIQL [24], which trains a single, action-free state-goal value function $V(s, g)$
 567 using Implicit V-Learning (IVL) [5] and extracts hierarchical policies using AWR-like objectives.
 568 The parameterisation ϕ_{ω} for the low-level policy is learned with the value function $V(s, \hat{\phi}_{\omega}(s, g))$,
 569 where $\hat{\phi}_{\omega}(s, g) \in \mathbb{R}^d$, and normalised such that $\|\hat{\phi}_{\omega}(s, g)\|_2^2 = d$, where d is the dimension of the
 570 embedding. The IVL loss is given by:

$$\mathcal{L}_{V, \phi_{\omega}} = \mathbb{E}_{(s,a,s') \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[\ell_{\tau}^2 \left(V(s, \hat{\phi}_{\omega}(s, g)) - r(s, g) - \gamma \tilde{V}(s', \tilde{\phi}_{\omega}(s', g)) \right) \right],$$

571 where \tilde{V} and $\tilde{\phi}_{\omega}$ denote the target network and representations and ℓ_{τ}^2 the expectile loss, as before.
 572 The low-level and high-level policies are extracted as follows:

$$\mathcal{L}_{\pi_h} = -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[e^{\alpha_h (V(g_s, \hat{\phi}_{\omega}(g_s, g)) - V(s, \hat{\phi}_{\omega}(s, g)))} \log \pi_h \left(\hat{\phi}_{\omega}(s, g_s) | s, g \right) \right],$$

573

$$\mathcal{L}_{\pi_l} = -\mathbb{E}_{(s,a,s' \dots g_s) \sim \mathcal{D}} \left[e^{\alpha_l (V(s', \hat{\phi}_{\omega}(s', g_s)) - V(s, \hat{\phi}_{\omega}(s, g_s)))} \log \pi_l \left(a | s, \hat{\phi}_{\omega}(s, g_s) \right) \right].$$

574 Since Park et al. [11] found that DDPGBC is better at extracting a policy than AWR, and similarly to
 575 Park et al. [24]’s action-free value function, we then fit a sort of high-level, action-free Q function to
 576 the value function:

$$\mathcal{L}_{Q_h} = \mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[(Q_h(s, g, g_s) - V_h(g_s, g))^2 \right].$$

577 This is simply to allow some extrapolation during extraction of the high-level policy.

578 C.3 Hierarchical Implicit Q-Learning 2 Value Functions (HIQL2v)

579 Since ARL uses two value functions, we also train a second style of HIQL, where we use a low-level
 580 and high-level value function. The high-level value function is trained using IVL, while the low-level
 581 value function is trained using IQL:

$$\mathcal{L}_{V_h} = \mathbb{E}_{(s,a) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[\ell_{\tau}^2 \left(V_h(s, g) - r(s, g) - \gamma \tilde{V}_h(s', g) \right) \right],$$

582 and

$$\mathcal{L}_{V_l} = \mathbb{E}_{(s,a) \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\ell_{\tau}^2 \left(V_l(s, g_s) - \tilde{Q}_l(s, g_s, a) \right) \right],$$

$$\mathcal{L}_{Q_l} = \mathbb{E}_{(s,a,s') \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[(Q_l(s, g_s, a) - r(s, g_s) - \gamma_l V_l(s', g_s))^2 \right],$$

583 where \tilde{Q}_l and \tilde{V}_h denote the target networks, $\gamma_l = 1 - \frac{1}{n}$ denotes the low-level discount factor, and
 584 ℓ_τ^2 denotes the expectile loss $\ell_\tau^2(x) = |\tau - (x < 0)|x^2$.

585 As for HIQL1vr, a high-level Q function is then learned only to allow DDPGBC high-level policy
 586 extraction:

$$\mathcal{L}_{Q_h} = \mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[(Q_h(s, g, g_s) - V_h(g_s, g))^2 \right].$$

587 Policies are then extracted using one of the following two loss functions for the low-level policy:

$$\begin{aligned} \mathcal{L}_{\pi_l}^{\text{AWR}} &= -\mathbb{E}_{(s, a, s' \dots g_s) \sim \mathcal{D}} \left[e^{\alpha_l (Q_l(s, g_s, a) - V_l(s, g_s))} \log \pi_l(a | s, g_s) \right], \\ \mathcal{L}_{\pi_l}^{\text{DDPGBC}} &= -\mathbb{E}_{(s, a, s', \dots, g_s) \sim \mathcal{D}} \left[Q_l(s, g_s, \mu^{\pi_l}(s, g_s)) + \alpha_l \log \pi_l(a | s, g_s) \right], \end{aligned}$$

588 and one of the following two loss functions for the high-level policy:

$$\begin{aligned} \mathcal{L}_{\pi_h}^{\text{AWR}} &= -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[e^{\alpha_h (Q_h(s, g, g_s) - V_h(s, g))} \log \pi_h(g_s | s, g) \right], \\ \mathcal{L}_{\pi_h}^{\text{DDPGBC}} &= -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[Q_h(s, g, \mu^{\pi_h}(s, g)) + \alpha_h \log \pi_h(g_s | s, g) \right]. \end{aligned}$$

589 **C.4 Hierarchical Implicit Q-Learning 2 Value Functions with Representation Learning** 590 **(HIQL2vr)**

591 Since HIQL1vr uses representation learning for the options, HIQL2vr learns option representations
 592 with the low-level value function. As in HIQL1vr (Section C.2), the representation is length-
 593 normalised. The high-level value function is trained as before, but the low-level value function,
 594 low-level Q function and high-level Q-function are trained using the following losses:

$$\begin{aligned} \mathcal{L}_{V_l, \hat{\phi}_\omega} &= \mathbb{E}_{(s, a) \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\ell_\tau^2 \left(V_l(s, \hat{\phi}_\omega(s, g_s)) - \tilde{Q}_l(s, g_s, a) \right) \right], \\ \mathcal{L}_{Q_l} &= \mathbb{E}_{(s, a, s') \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\left(Q_l(s, g_s, a) - r(s, g_s) - \gamma_l V_l(s', \hat{\phi}_\omega(s', g_s)) \right)^2 \right], \\ \mathcal{L}_{Q_h} &= \mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[\left(Q_h(s, g, \hat{\phi}_\omega(s, g_s)) - V_h(g_s, g) \right)^2 \right]. \end{aligned}$$

595 The policies are extracted using one of the following two loss functions for the low-level policy,

$$\begin{aligned} \mathcal{L}_{\pi_l}^{\text{AWR}} &= -\mathbb{E}_{(s, a, s' \dots g_s) \sim \mathcal{D}} \left[e^{\alpha_l (Q_l(s, g_s, a) - V_l(s, \hat{\phi}_\omega(s, g_s)))} \log \pi_l(a | s, \hat{\phi}_\omega(s, g_s)) \right], \\ \mathcal{L}_{\pi_l}^{\text{DDPGBC}} &= -\mathbb{E}_{(s, a, s', \dots, g_s) \sim \mathcal{D}} \left[Q_l(s, g_s, \mu^{\pi_l}(s, \hat{\phi}_\omega(s, g_s))) + \alpha_l \log \pi_l(a | s, \hat{\phi}_\omega(s, g_s)) \right], \end{aligned}$$

596 and one of the following two loss functions for the high-level policy:

$$\begin{aligned} \mathcal{L}_{\pi_h}^{\text{AWR}} &= -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[e^{\alpha_h (Q_h(s, g, \hat{\phi}_\omega(s, g_s)) - V_h(s, g))} \log \pi_h(\hat{\phi}_\omega(s, g_s) | s, g) \right], \\ \mathcal{L}_{\pi_h}^{\text{DDPGBC}} &= -\mathbb{E}_{(s \dots g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[Q_h(s, g, \mu^{\pi_h}(s, g)) + \alpha_h \log \pi_h(\hat{\phi}_\omega(s, g_s) | s, g) \right]. \end{aligned}$$

597 **C.5 Abstractive Reinforcement Learning Implicitly Learning Relativised Options (ARLi)**

598 As presented in the main body of the paper, this is a minimal amendment to HIQL2vr, but, to learn
 599 relativised options via action similarity, option representations are now learned with the low-level
 600 policy. Equations are identical to HIQL2vr, except for the following low-level value functions:

$$\begin{aligned} \mathcal{L}_{V_l} &= \mathbb{E}_{(s, a) \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\ell_\tau^2 \left(V_l(s, g_s) - \tilde{Q}_l(s, g_s, a) \right) \right], \\ \mathcal{L}_{Q_l} &= \mathbb{E}_{(s, a, s') \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\left(Q_l(s, g_s, a) - r(s, g_s) - \gamma_l V_l(s', g_s) \right)^2 \right], \end{aligned}$$

601 and low-level policy functions:

$$\begin{aligned} \mathcal{L}_{\pi_l, \hat{\phi}_\omega}^{\text{AWR}} &= -\mathbb{E}_{(s, a, s' \dots g_s) \sim \mathcal{D}} \left[e^{\alpha_l (Q_l(s, g_s, a) - V_l(s, g_s))} \log \pi_l(a | s, \hat{\phi}_\omega(s, g_s)) \right], \\ \mathcal{L}_{\pi_l, \hat{\phi}_\omega}^{\text{DDPGBC}} &= -\mathbb{E}_{(s, a, s', \dots, g_s) \sim \mathcal{D}} \left[Q_l(s, g_s, \mu^{\pi_l}(s, \hat{\phi}_\omega(s, g_s))) + \alpha_l \log \pi_l(a | s, \hat{\phi}_\omega(s, g_s)) \right]. \end{aligned}$$

Algorithm 1 Abstractive RL implicitly learning relativised options (ARLi)

Training

 Initialise low-level policy $\pi_l(a | s, \omega)$, high-level policy $\pi_h(\omega | s, g)$, and representation ϕ_ω .

while not converged **do**

 Sample batch \mathcal{D}

 ▶ **Hierarchical Policy**

 Update low-level policy π_l and representation ϕ_ω using $\pi_l(a | s, \hat{\phi}_\omega(s, g_s))$ (Equation 4)

 $\omega \leftarrow \text{stopgrad}(\hat{\phi}_\omega(s, g_s))$

 Update high-level policy $\pi_h(\omega | s, g)$ (Equation 6)

 ▶ **Hierarchical Value**

 Update low-level value function $V_l(s, g_s)$ (Equation 2)

 Update high-level value function $V_h(s, g)$ (Equation 5)

 Update critic $Q_l(s, g_s, a)$ (Equation 3)

end while
return $\pi_l, \pi_h, \phi_\omega$
Deployment (state s , goal g)

 Sample option from high-level policy $\omega \sim \pi_h(\cdot | s, g)$

 Length-Normalise $\omega \leftarrow \frac{\omega}{\|\omega\|} \cdot \sqrt{d}$

 Sample action from low-level policy $a \sim \pi_l(\cdot | s, \omega)$

 602 **C.6 Abstractive Reinforcement Learning Explicitly Enforcing Translational Invariance**
 603 **(ARLe)**

 604 ARL uses relativised options and relativised states for the low-level value function. Unlike ARLi, the
 605 low-level value and low-level critic now use *relativised* goals, and options are explicitly relativised.
 606 As for ARLi, the option representations are learned with the low-level policy. The high-level value
 607 function is learned identically to HIQL2v and ARLi, but the low-level value function, low-level Q
 608 function and high-level Q function are learned as follows:

$$\begin{aligned} \mathcal{L}_{V_l} &= \mathbb{E}_{(s,a) \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\ell_\tau^2 \left(V_l(g_s - s) - \tilde{Q}_l(s, g_s - s, a) \right) \right], \\ \mathcal{L}_{Q_l} &= \mathbb{E}_{(s,a,s') \sim \mathcal{D}, g_s \sim p_{\gamma_l}^{\mathcal{D}}(\cdot | s, a)} \left[\left(Q_l(s, g_s - s, a) - r(s, g_s) - \gamma_l V_l(g_s - s') \right)^2 \right], \\ \mathcal{L}_{Q_h} &= \mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p^{\mathcal{D}}(\cdot | s, a)} \left[\left(Q_h(s, g, \hat{\phi}_\omega(g_s - s)) - V_h(g_s, g) \right)^2 \right]. \end{aligned}$$

 609 where, as before, \tilde{Q}_l and \tilde{V}_h denote the target networks, $\gamma_l = 1 - \frac{1}{n}$ denotes the low-level discount
 610 factor, and ℓ_τ^2 denotes the expectile loss $\ell_\tau^2(x) = |\tau - (x < 0)|x^2$. The policies are extracted using
 611 one of the following two loss functions for the low-level policy,

$$\begin{aligned} \mathcal{L}_{\pi_l, \phi_\omega}^{\text{AWR}} &= -\mathbb{E}_{(s,a,s', \dots, g_s) \sim \mathcal{D}} \left[e^{\alpha_l (Q_l(s, g_s - s, a) - V_l(g_s - s))} \log \pi_l \left(a | s, \hat{\phi}_\omega(s, g_s) \right) \right], \\ \mathcal{L}_{\pi_l, \phi_\omega}^{\text{DDPGBC}} &= -\mathbb{E}_{(s,a,s', \dots, g_s) \sim \mathcal{D}} \left[Q_l(s, g_s, \mu^{\pi_l}(s, \hat{\phi}_\omega(s, g_s))) + \alpha_l \log \pi_l(a | s, \hat{\phi}_\omega(s, g_s)) \right], \end{aligned}$$

612 and one of the following two loss functions for the high-level policy:

$$\begin{aligned} \mathcal{L}_{\pi_h}^{\text{AWR}} &= -\mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[e^{\alpha_h (Q_h(s, g, \hat{\phi}_\omega(g_s - s)) - V_h(s, g))} \log \pi_h \left(\hat{\phi}_\omega(s, g_s) | s, g \right) \right], \\ \mathcal{L}_{\pi_h}^{\text{DDPGBC}} &= -\mathbb{E}_{(s, \dots, g_s) \sim \mathcal{D}, g \sim p_\gamma^{\mathcal{D}}(\cdot | s)} \left[Q_h(s, g, \mu^{\pi_h}(s, g)) + \alpha_h \log \pi_g(\hat{\phi}_\omega(s, g_s) | s, g) \right]. \end{aligned}$$

Algorithm 2 Abstractive RL explicitly enforcing translation invariance (ARLe)

TrainingInitialise low-level policy $\pi_l(a | s, \omega)$, high-level policy $\pi_h(\omega | s, g)$, and representation ϕ_ω .**while** not converged **do**Sample batch from \mathcal{D} ▷ **Hierarchical Policy**Update low-level policy π_l and representation ϕ_ω using $\pi_l(a | s, \hat{\phi}_\omega(s, g_s))$ (Equation 4) $\omega \leftarrow \text{stopgrad}(\hat{\phi}_\omega(s, g_s))$ Update high-level policy $\pi_h(\omega | s, g)$ (Equation 6)▷ **Hierarchical Value**Update low-level value function $V_l(g_s - s)$ (Equation 2)Update high-level value function $V_h(s, g)$ (Equation 5)Update critic $Q_l(s, g_s - s, a)$ (Equation 3)**end while****return** $\pi_l, \pi_h, \phi_\omega$ **Deployment (state s , goal g)**Sample option from high-level policy $\omega \sim \pi_h(\cdot | s, g)$ Soft-Normalise $\omega \leftarrow \frac{\omega \cdot \tanh(\|\omega\|)}{\|\omega\|} \cdot \sqrt{d}$ Sample action from low-level policy $a \sim \pi_l(\cdot | s, \omega)$

Table 2: **Full Results 1.** We report each method’s average (binary) success rate (%) across the five test-time goals on each task. Bootstrapped 95% CI over 4 seeds and 20 evaluation runs. Blue bold indicates the highest mean; black bold overlapping confidence intervals.

Environment	Task	IQL	HIQL1vr	HIQL2v	HIQL2vr	ARLi	ARLe
pointmaze-giant-navigate-v0	1	0±0	1±2	6±9	2±4	38±38	4±4
	2	0±0	66±22	64±24	75±19	86±7	31±14
	3	0±0	49±21	1±2	24±16	8±8	2±4
	4	0±0	71±33	4±6	60±20	10±10	22±21
	5	0±0	89±7	31±19	69±18	96±4	21±14
	Overall	0±0	55±11	21±7	46±6	48±8	16±2
pointmaze-giant-stitch-v0	1	0±0	0±0	0±0	0±0	0±0	0±0
	2	0±0	0±0	0±0	0±0	0±0	0±0
	3	0±0	0±0	0±0	0±0	0±0	0±0
	4	0±0	0±0	0±0	0±0	0±0	0±0
	5	0±0	0±0	0±0	0±0	0±0	0±0
	Overall	0±0	0±0	0±0	0±0	0±0	0±0
antmaze-giant-navigate-v0	1	0±0	18±8	19±11	29±9	18±5	25±4
	2	1±2	44±19	56±6	50±10	48±8	42±8
	3	0±0	34±7	25±8	39±11	51±18	35±10
	4	0±0	30±6	44±7	54±14	56±9	44±6
	5	0±0	64±6	57±9	71±16	66±6	64±11
	Overall	0±0	38±3	40±3	48±6	48±3	42±4
antmaze-giant-stitch-v0	1	0±0	8±8	9±4	30±5	44±14	26±7
	2	0±0	6±4	30±12	26±11	22±5	14±4
	3	0±0	2±2	5±8	4±4	15±6	10±13
	4	0±0	16±16	16±11	36±9	54±14	45±11
	5	0±0	4±4	15±6	6±2	26±21	10±0
	Overall	0±0	7±7	15±3	20±3	32±7	21±1
antmaze-teleport-stitch-v0	1	38±8	41±7	55±8	45±9	42±9	45±4
	2	55±4	36±11	61±6	46±7	45±13	55±11
	3	55±9	16±4	31±12	34±4	35±18	29±4
	4	45±9	24±9	22±7	39±11	27±4	36±14
	5	51±6	26±9	46±11	34±11	31±6	41±6
	Overall	49±2	29±4	43±3	40±6	36±5	41±4
humanoidmaze-giant-navigate-v0	1	1±2	16±9	8±8	2±4	40±9	29±11
	2	2±2	36±19	29±14	16±14	50±6	46±4
	3	0±0	12±5	11±6	10±9	41±9	32±5
	4	0±0	19±18	11±6	15±12	57±10	40±10
	5	0±0	26±9	1±2	10±13	56±14	45±4
	Overall	1±1	22±11	12±5	11±10	49±6	38±4
humanoidmaze-giant-stitch-v0	1	0±0	5±8	0±0	0±0	12±5	11±6
	2	2±2	11±8	2±2	0±0	31±16	21±9
	3	0±0	2±2	0±0	0±0	11±2	8±4
	4	0±0	0±0	0±0	0±0	9±9	9±7
	5	0±0	0±0	0±0	0±0	2±2	4±2
	Overall	0±0	4±3	0±0	0±0	13±2	10±3

Table 3: **Full Results 2.** We report each method’s average (binary) success rate (%) across the five test-time goals on each task. Bootstrapped 95% CI over 4 seeds and 20 evaluation runs. Blue bold indicates the highest mean; black bold overlapping confidence intervals.

Environment	Task	IQL	HIQL1vr	HIQL2v	HIQL2vr	ARLi	ARLe
cube-double-play-v0	1	94±6	8±2	0±0	14±2	96±4	94±4
	2	46±11	0±0	0±0	0±0	59±6	78±11
	3	59±8	0±0	0±0	0±0	52±5	68±4
	4	15±6	0±0	0±0	0±0	14±7	35±6
	5	34±9	0±0	0±0	0±0	44±6	61±9
	Overall	50±3	2±0	0±0	3±0	53±2	67±3
cube-triple-play-v0	1	51±7	35±16	0±0	4±4	70±10	69±9
	2	1±2	0±0	0±0	0±0	1±2	2±2
	3	4±2	0±0	0±0	0±0	0±0	2±4
	4	0±0	0±0	0±0	0±0	0±0	0±0
	5	0±0	0±0	0±0	0±0	0±0	0±0
	Overall	11±2	7±3	0±0	1±1	14±2	15±3
cube-quadruple-play-v0	1	0±0	0±0	0±0	0±0	4±2	2±2
	2	0±0	0±0	0±0	0±0	0±0	0±0
	3	0±0	0±0	0±0	0±0	1±2	0±0
	4	0±0	0±0	0±0	0±0	0±0	0±0
	5	0±0	0±0	0±0	0±0	0±0	0±0
	Overall	0±0	0±0	0±0	0±0	1±0	0±0
puzzle-3x3-play-v0	1	100±0	100±0	0±0	100±0	100±0	100±0
	2	100±0	14±21	0±0	15±12	88±12	45±34
	3	99±2	4±6	0±0	0±0	78±19	22±30
	4	100±0	2±4	0±0	2±2	77±28	21±28
	5	100±0	15±12	0±0	1±2	89±11	29±32
	Overall	100±0	27±8	0±0	24±3	86±14	44±25
puzzle-4x4-play-v0	1	50±10	66±36	0±0	29±12	91±9	100±0
	2	6±4	32±24	0±0	15±6	70±15	75±18
	3	38±3	59±32	0±0	18±8	81±7	90±9
	4	29±13	48±29	0±0	10±4	72±18	89±6
	5	29±7	41±25	0±0	12±5	72±13	85±5
	Overall	30±3	49±27	0±0	17±6	78±11	88±6
puzzle-4x5-play-v0	1	72±12	82±11	0±0	62±15	90±13	68±36
	2	1±2	2±4	0±0	0±0	0±0	0±0
	3	0±0	0±0	0±0	0±0	0±0	0±0
	4	0±0	0±0	0±0	0±0	0±0	0±0
	5	0±0	0±0	0±0	0±0	0±0	0±0
	Overall	15±3	17±3	0±0	12±3	18±3	14±7
puzzle-4x6-play-v0	1	51±9	0±0	0±0	76±6	65±32	45±45
	2	15±6	0±0	0±0	16±9	2±2	0±0
	3	0±0	0±0	0±0	0±0	0±0	0±0
	4	0±0	0±0	0±0	0±0	0±0	0±0
	5	0±0	0±0	0±0	0±0	0±0	0±0
	Overall	13±1	0±0	0±0	18±2	14±7	9±9
scene-play-v0	1	31±6	22±5	26±23	26±7	50±6	44±13
	2	16±11	5±4	11±4	6±4	16±11	18±3
	3	6±8	8±8	5±8	10±4	19±4	12±2
	4	9±6	16±6	11±4	19±9	30±8	26±9
	5	1±2	9±6	8±5	4±4	8±5	6±4
	Overall	13±2	12±3	12±7	13±1	24±3	21±3

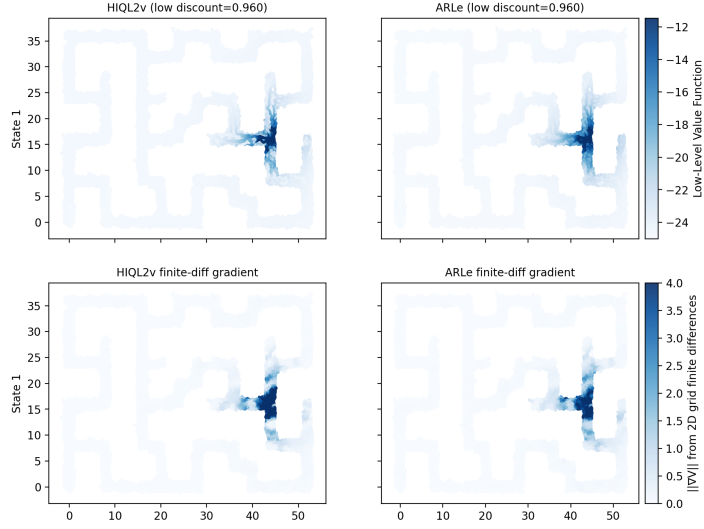


Figure 4: **Low-level** value functions (**top**) and **gradient** of low-level value function (**bottom**). IQL and HIQL1vr are excluded as they have a single value function.

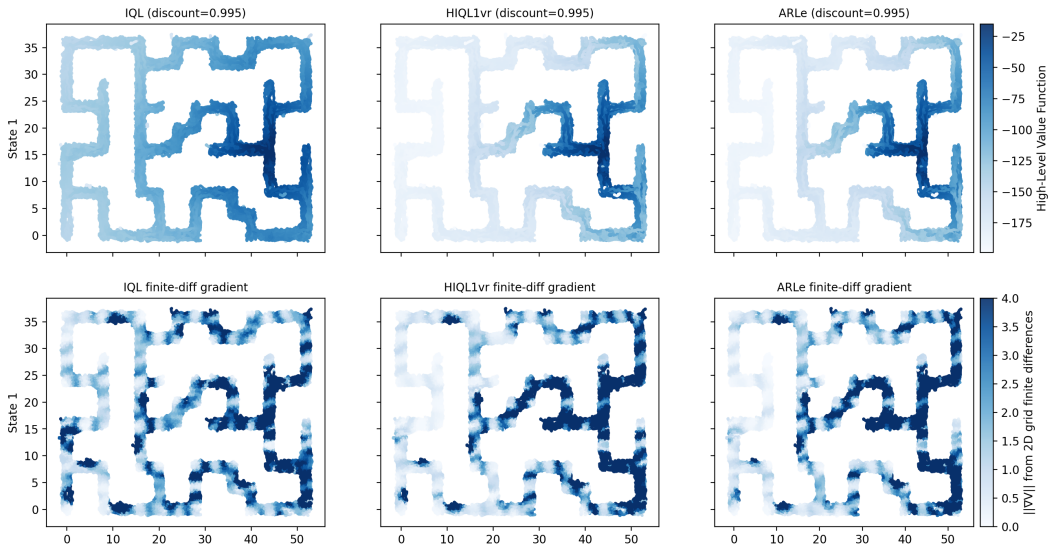


Figure 5: **High-level** value functions (**top**) and **gradient** of high-level value function (**bottom**). HIQL2v and HIQL2vr are excluded, as they have identical ones to ARL.

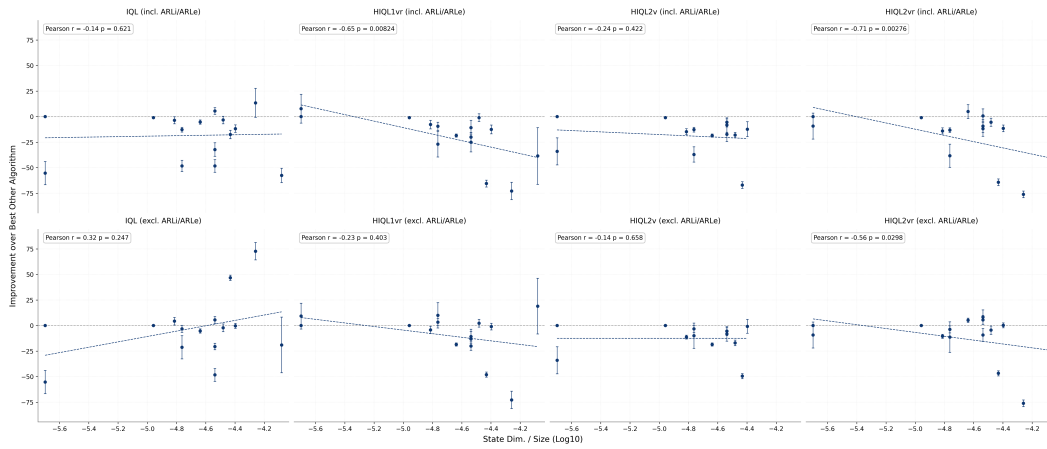


Figure 6: **Performance improvements** over next best performing algorithm against number of state dimensions per dataset sample: **IQL** (left), **HIQL1vr** (centre left), **HIQL2v** (centre right) and **HIQL2vr** (right). Including ARLi and ARLe (**top**), and excluding ARLi and ARLe (**bottom**). Bootstrapped 95% CI over 4 seeds and 20 evaluation runs.

614 **E Experimental Details**

615 We release all code and hyperparameters in a repository with this work.

616 **Datasets.** We use the standard OGBench datasets. States are randomly and uniformly sampled
 617 from the dataset. Goals for the value function learning and policy extraction are sampled using a
 618 certain probability of sampling the current state $p_{cur}^{\mathcal{D}}$, from the current trajectory $p_{traj}^{\mathcal{D}}$ (geometrically,
 619 according to the discount factor, or uniformly), or randomly from the dataset $p_{rand}^{\mathcal{D}}$. waypoints g_s
 620 are taken as the states n steps ahead of the current state.

Table 4: **Environment Characteristics.** Environment properties to provide intuition for interpreting results.

Task	State Dim.	Action Dim.	Max. Episode Length	Dataset Size
pointmaze-giant	2	2	1000	1M
antmaze-giant	29	8	1000	1M
antmaze-teleport	29	8	1000	1M
humanoidmaze-large	69	21	1000	4M
humanoidmaze-giant	69	21	4000	4M
cube-double	37	5	500	1M
cube-triple	46	5	1000	3M
cube-quadruple	55	5	1000	5M
puzzle-3x3	55	5	500	1M
puzzle-4x4	83	5	500	1M
puzzle-4x5	99	5	1000	3M
puzzle-4x6	115	5	1000	5M
scene-play	40	5	750	1M

621 **Reward Relabelling.** Unlike prior work on OGBench [5], we use the original environment reward
 622 functions for relabeling rewards rather than using a binary indicator of the state-index to ensure
 623 that the agent remains focused on the primary task objectives and prevents the value function from
 624 becoming overly specific to irrelevant dimensions. To ensure fair comparison, this relabeling strategy
 625 is applied consistently across all algorithms and tasks. Note that assuming access to the environment
 626 reward function in robotic tasks is an entirely valid assumption.

627 **Hyperparameters.** We provide the full list of hyperparameters. We follow those from Park et al.
 628 [5] and Park et al. [12]. Notably, while these parameters were specifically tuned for HIQL, we apply
 629 them to ARL without further adjustment. The fact that ARL achieves strong performance using
 630 parameters optimised for a different algorithm demonstrates its robustness. We use DDPGBC with a
 631 behaviour cloning strength of 0.1 to extract the high-level policy in manipulation tasks, which allows
 632 for more extrapolation [11]. This was generally found to outperform AWR.

Table 5: Hyperparameters

Hyperparameter	Value
Gradient steps	10^6
Optimiser	Adam [56]
Learning rate	0.0003
Batch size	1024
Layer Normalisation [57]	True
Nonlinearity	GELU [58]
Value MLP	[1024, 1024, 1024, 1024]
Actor MLP	[1024, 1024, 1024, 1024]
Representation MLP	[512, 512, 512]
Representation Dimension	10
Target network update rate	0.005
IQL Expectile τ	0.9 (IQL), 0.7 (HIQL1vr, HIQL2v, HIQL2vr, ARLi, ARLe)
Value ratio $(p_{\text{cur}}^{\mathcal{D}}, p_{\text{traj}}^{\mathcal{D}}, p_{\text{rand}}^{\mathcal{D}}, p_{\text{geom}}^{\mathcal{D}})$	(0.2, 0.5, 0.3, 0)
Low Value ratio $(p_{\text{cur}}^{\mathcal{D}}, p_{\text{traj}}^{\mathcal{D}}, p_{\text{rand}}^{\mathcal{D}}, p_{\text{geom}}^{\mathcal{D}})$	(0.10, 0.85, 0.05, 1)
High Value ratio $(p_{\text{cur}}^{\mathcal{D}}, p_{\text{traj}}^{\mathcal{D}}, p_{\text{rand}}^{\mathcal{D}}, p_{\text{geom}}^{\mathcal{D}})$	(0.2, 0.5, 0.3, 0)
Policy ratio $(p_{\text{cur}}^{\mathcal{D}}, p_{\text{traj}}^{\mathcal{D}}, p_{\text{rand}}^{\mathcal{D}}, p_{\text{geom}}^{\mathcal{D}})$	(0.0, 0.5, 0.5, 1)

Table 6: Task-Specific Hyperparameters

Task	n	γ	Loss π	α	Loss π_l	α_l	Loss π_h	α_h
pointmaze-giant-navigate-v0	25	0.995	DDPGBC	0.1	AWR	3.0	AWR	3.0
pointmaze-giant-stitch-v0	25	0.995	DDPGBC	0.1	AWR	3.0	AWR	3.0
antmaze-giant-navigate-v0	25	0.995	DDPGBC	0.1	AWR	3.0	AWR	3.0
antmaze-giant-stitch-v0	25	0.995	DDPGBC	0.1	AWR	3.0	AWR	3.0
antmaze-teleport-stitch-v0	25	0.990	DDPGBC	0.1	AWR	3.0	AWR	3.0
humanoidmaze-giant-navigate-v0	100	0.999	DDPGBC	0.1	AWR	3.0	AWR	3.0
humanoidmaze-giant-stitch-v0	100	0.999	DDPGBC	0.1	AWR	3.0	AWR	3.0
cube-double-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
cube-triple-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
cube-quadruple-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
puzzle-3x3-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
puzzle-4x4-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
puzzle-4x5-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
puzzle-4x6-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1
scene-play-v0	25	0.99	DDPGBC	1.0	AWR	3.0	DDPGBC	0.1

633 The low-level discount factor is computed as $\gamma_l = 1 - \frac{1}{n}$. The high-level discount factor is computed
634 as $\gamma_h = \gamma^n$.

635 **F Limitations**

636 ARLe’s performance depends on the alignment between its inductive bias (assuming translational
637 invariance) and the environment’s structure. We discuss these limitations with the experiments
638 (Section 5) and in the conclusion (Section 6).

639 Like other prior methods learning action-free value functions, learning an action-free high-level value
640 function biases our instantiations of ARL towards being optimistic in stochastic environments [24].
641 Such optimism bias could be addressed by disentangling controllable parts of the state [59], but we
642 leave this to future work. We also note that, since only the high-level value function is action-free,
643 performance degradation compared to IQL for both ARLe and ARLi in the stochastic environment
644 (*antmaze-teleport-stitch-v0*) is less significant than for HIQL1vr, for example.

645 Including more than 15 environments would have helped to strengthen our hypothesis in Section 6,
646 but we were limited by compute resources.

647 **G Compute**

648 All experiments were conducted on NVIDIA L40 GPUs, lasting 3 hours per run, including evaluation.

649 **Impact Statement**

650 This paper presents work whose goal is to advance the field of machine learning. There are many
651 potential societal consequences of our work, none of which we feel must be specifically highlighted
652 here.

653 **NeurIPS Paper Checklist**

654 **1. Claims**

655 Question: Do the main claims made in the abstract and introduction accurately reflect the
656 paper’s contributions and scope?

657 Answer: [\[Yes\]](#)

658 Justification: Methods Section 4 and Experiments Section 5.

659 Guidelines:

- 660 • The answer [\[N/A\]](#) means that the abstract and introduction do not include the claims
661 made in the paper.
- 662 • The abstract and/or introduction should clearly state the claims made, including the
663 contributions made in the paper and important assumptions and limitations. A [\[No\]](#) or
664 [\[N/A\]](#) answer to this question will not be perceived well by the reviewers.
- 665 • The claims made should match theoretical and experimental results, and reflect how
666 much the results can be expected to generalize to other settings.
- 667 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
668 are not attained by the paper.

669 **2. Limitations**

670 Question: Does the paper discuss the limitations of the work performed by the authors?

671 Answer: [\[Yes\]](#)

672 Justification: Appendix F and discussion in Sections 5 and 6.

673 Guidelines:

- 674 • The answer [\[N/A\]](#) means that the paper has no limitation while the answer [\[No\]](#) means
675 that the paper has limitations, but those are not discussed in the paper.
- 676 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 677 • The paper should point out any strong assumptions and how robust the results are to
678 violations of these assumptions (e.g., independence assumptions, noiseless settings,
679 model well-specification, asymptotic approximations only holding locally). The authors
680 should reflect on how these assumptions might be violated in practice and what the
681 implications would be.
- 682 • The authors should reflect on the scope of the claims made, e.g., if the approach was
683 only tested on a few datasets or with a few runs. In general, empirical results often
684 depend on implicit assumptions, which should be articulated.
- 685 • The authors should reflect on the factors that influence the performance of the approach.
686 For example, a facial recognition algorithm may perform poorly when image resolution
687 is low or images are taken in low lighting. Or a speech-to-text system might not be
688 used reliably to provide closed captions for online lectures because it fails to handle
689 technical jargon.
- 690 • The authors should discuss the computational efficiency of the proposed algorithms
691 and how they scale with dataset size.
- 692 • If applicable, the authors should discuss possible limitations of their approach to
693 address problems of privacy and fairness.
- 694 • While the authors might fear that complete honesty about limitations might be used by
695 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
696 limitations that aren’t acknowledged in the paper. The authors should use their best
697 judgment and recognize that individual actions in favor of transparency play an impor-
698 tant role in developing norms that preserve the integrity of the community. Reviewers
699 will be specifically instructed to not penalize honesty concerning limitations.

700 **3. Theory assumptions and proofs**

701 Question: For each theoretical result, does the paper provide the full set of assumptions and
702 a complete (and correct) proof?

703 Answer: [\[Yes\]](#)

704 Justification: Motivation Box in Section 4 and Appendix B.

705 Guidelines:

- 706 • The answer [N/A] means that the paper does not include theoretical results.
- 707 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 708 referenced.
- 709 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 710 • The proofs can either appear in the main paper or the supplemental material, but if
- 711 they appear in the supplemental material, the authors are encouraged to provide a short
- 712 proof sketch to provide intuition.
- 713 • Inversely, any informal proof provided in the core of the paper should be complemented
- 714 by formal proofs provided in appendix or supplemental material.
- 715 • Theorems and Lemmas that the proof relies upon should be properly referenced.

716 4. Experimental result reproducibility

717 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

718 perimental results of the paper to the extent that it affects the main claims and/or conclusions

719 of the paper (regardless of whether the code and data are provided or not)?

720 Answer: [Yes]

721 Justification: All loss functions and hyperparameters are detailed in Appendices C and E.

722 Code is provided, with the exact seeds and command lines required to generate results.

723 Guidelines:

- 724 • The answer [N/A] means that the paper does not include experiments.
- 725 • If the paper includes experiments, a [No] answer to this question will not be perceived
- 726 well by the reviewers: Making the paper reproducible is important, regardless of
- 727 whether the code and data are provided or not.
- 728 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 729 to make their results reproducible or verifiable.
- 730 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 731 For example, if the contribution is a novel architecture, describing the architecture fully
- 732 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 733 be necessary to either make it possible for others to replicate the model with the same
- 734 dataset, or provide access to the model. In general, releasing code and data is often
- 735 one good way to accomplish this, but reproducibility can also be provided via detailed
- 736 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 737 of a large language model), releasing of a model checkpoint, or other means that are
- 738 appropriate to the research performed.
- 739 • While NeurIPS does not require releasing code, the conference does require all submis-
- 740 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 741 nature of the contribution. For example
- 742 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 743 to reproduce that algorithm.
- 744 (b) If the contribution is primarily a new model architecture, the paper should describe
- 745 the architecture clearly and fully.
- 746 (c) If the contribution is a new model (e.g., a large language model), then there should
- 747 either be a way to access this model for reproducing the results or a way to reproduce
- 748 the model (e.g., with an open-source dataset or instructions for how to construct
- 749 the dataset).
- 750 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 751 authors are welcome to describe the particular way they provide for reproducibility.
- 752 In the case of closed-source models, it may be that access to the model is limited in
- 753 some way (e.g., to registered users), but it should be possible for other researchers
- 754 to have some path to reproducing or verifying the results.

755 5. Open access to data and code

756 Question: Does the paper provide open access to the data and code, with sufficient instruc-

757 tions to faithfully reproduce the main experimental results, as described in supplemental

758 material?

759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Answer: [Yes]

Justification: Code is provided in the supplementary material, with the exact seeds and command lines required to generate results. The data is referenced and open-sourced.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: Details in the Experiments Section 5 and Appendices E.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Bootstrapped 95% confidence intervals and Pearson Correlation for statistical significance tests.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 810
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
 - 811
 - 812 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - 813
 - 814
 - 815 • For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
 - 816
 - 817
 - 818 • If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
 - 819

8. Experiments compute resources

820

821 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

822

823

824 Answer: [Yes]

825 Justification: Appendix G.

826 Guidelines:

- 827 • The answer [N/A] means that the paper does not include experiments.
- 828 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 829
- 830 • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 831
- 832 • The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 833
- 834

835 9. Code of ethics

836 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

837

838 Answer: [Yes]

839 Justification: We confirm that the research conforms with the NeurIPS Code of Ethics in every respect.

840

841 Guidelines:

- 842 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 843
- 844 • If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- 845
- 846 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 847

848 10. Broader impacts

849 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

850

851 Answer: [Yes]

852 Justification: Appendix G.

853 Guidelines:

- 854 • The answer [N/A] means that there is no societal impact of the work performed.
- 855 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- 856
- 857 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- 858
- 859
- 860

- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 875
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

876 11. Safeguards

877 Question: Does the paper describe safeguards that have been put in place for responsible
878 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
879 image generators, or scraped datasets)?

880 Answer: [N/A]

881 Justification: The paper poses no such risks.

882 Guidelines:

- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- 891
- 892
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

893 12. Licenses for existing assets

894 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
895 the paper, properly credited and are the license and terms of use explicitly mentioned and
896 properly respected?

897 Answer: [Yes]

898 Justification: We cite the original papers that produced the code packages and datasets.

899 Guidelines:

- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- 911
- 912
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

913 • If this information is not available online, the authors are encouraged to reach out to
914 the asset’s creators.

915 **13. New assets**

916 Question: Are new assets introduced in the paper well documented and is the documentation
917 provided alongside the assets?

918 Answer: [N/A]

919 Justification: The paper does not release new assets.

920 Guidelines:

- 921 • The answer [N/A] means that the paper does not release new assets.
- 922 • Researchers should communicate the details of the dataset/code/model as part of their
923 submissions via structured templates. This includes details about training, license,
924 limitations, etc.
- 925 • The paper should discuss whether and how consent was obtained from people whose
926 asset is used.
- 927 • At submission time, remember to anonymize your assets (if applicable). You can either
928 create an anonymized URL or include an anonymized zip file.

929 **14. Crowdsourcing and research with human subjects**

930 Question: For crowdsourcing experiments and research with human subjects, does the paper
931 include the full text of instructions given to participants and screenshots, if applicable, as
932 well as details about compensation (if any)?

933 Answer: [N/A]

934 Justification: The paper does not involve crowdsourcing nor research with human subjects.

935 Guidelines:

- 936 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
937 with human subjects.
- 938 • Including this information in the supplemental material is fine, but if the main contribu-
939 tion of the paper involves human subjects, then as much detail as possible should be
940 included in the main paper.
- 941 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
942 or other labor should be paid at least the minimum wage in the country of the data
943 collector.

944 **15. Institutional review board (IRB) approvals or equivalent for research with human
945 subjects**

946 Question: Does the paper describe potential risks incurred by study participants, whether
947 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
948 approvals (or an equivalent approval/review based on the requirements of your country or
949 institution) were obtained?

950 Answer: [N/A]

951 Justification: The paper does not involve crowdsourcing nor research with human subjects.

952 Guidelines:

- 953 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
954 with human subjects.
- 955 • Depending on the country in which research is conducted, IRB approval (or equivalent)
956 may be required for any human subjects research. If you obtained IRB approval, you
957 should clearly state this in the paper.
- 958 • We recognize that the procedures for this may vary significantly between institutions
959 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
960 guidelines for their institution.
- 961 • For initial submissions, do not include any information that would break anonymity (if
962 applicable), such as the institution conducting the review.

963 **16. Declaration of LLM usage**

964 Question: Does the paper describe the usage of LLMs if it is an important, original, or
965 non-standard component of the core methods in this research? Note that if the LLM is used
966 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
967 scientific rigor, or originality of the research, declaration is not required.

968 Answer: [N/A]

969 Justification: The core method development in this research does not involve LLMs as any
970 important, original, or non-standard components.

971 Guidelines:

- 972 • The answer [N/A] means that the core method development in this research does not
973 involve LLMs as any important, original, or non-standard components.
- 974 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
975 be described.