Flattening Hierarchies with Policy Bootstrapping

Anonymous authors

Paper under double-blind review

Abstract

1	Offline goal-conditioned reinforcement learning (GCRL) is a promising approach for
2	pretraining generalist policies on large datasets of reward-free trajectories, akin to the
3	self-supervised objectives used to train foundation models for computer vision and natu-
4	ral language processing. However, scaling GCRL to longer horizons remains challenging
5	due to the combination of sparse rewards and discounting, which obscures the com-
6	parative advantages of primitive actions with respect to distant goals. Hierarchical RL
7	methods achieve strong empirical results on long-horizon goal-reaching tasks, but their
8	reliance on modular, timescale-specific policies and subgoal generation introduces signif-
9	icant additional complexity and hinders scaling to high-dimensional goal spaces. In this
10	work, we introduce an algorithm to train a flat (non-hierarchical) goal-conditioned policy
11	by bootstrapping on subgoal-conditioned policies with advantage-weighted importance
12	sampling. Our approach eliminates the need for a generative model over the (sub)goal
13	space, which we find is key for scaling to high-dimensional control in large state spaces.
14	We further show that existing hierarchical and bootstrapping-based approaches corre-
15	spond to specific design choices within our derivation. Across a comprehensive suite of
16	state- and pixel-based locomotion and manipulation benchmarks, our method matches or
17	surpasses state-of-the-art offline GCRL algorithms and scales to complex, long-horizon
18	tasks where prior approaches fail.

19 1 Introduction

20 Goal-conditioned reinforcement learning (GCRL) specifies tasks by desired outcomes, alleviating the 21 burden of defining reward functions over the state-space and enabling the training of general policies 22 capable of achieving a wide range of goals. Offline GCRL extends this paradigm to leverage existing 23 datasets of reward-free trajectories and has been likened to the simple self-supervised objectives that 24 have been successful in training foundation models for other areas of machine learning (Yang et al., 25 2023; Park et al., 2024a). However, the conceptual simplicity of GCRL belies practical challenges in 26 learning accurate value functions and, consequently, effective policies for goals requiring complex, 27 long-horizon behaviors. These limitations call into question its applicability as a general and scalable 28 objective for learning *foundation policies* (Black et al., 2024; Park et al., 2024d; Physical Intelligence 29 et al., 2025) that can be efficiently adapted to a diverse array of control tasks.

30 Hierarchical reinforcement learning (HRL) is commonly used to address these challenges and is 31 particularly well-suited to the recursive subgoal structure of goal-reaching tasks, where reaching 32 distant goals entails first passing through intermediate subgoal states. Goal-conditioned HRL exploits 33 this structure by learning a hierarchy composed of multiple levels: one or more high-level policies, 34 tasked with generating intermediate subgoals between the current state and the goal; and a low-35 level actor, which operates over the primitive action space to achieve the assigned subgoals. These 36 approaches have achieved state-of-the-art results in both online (Nachum et al., 2018; Levy et al., 37 2019) and offline GCRL (Park et al., 2024c), and are especially effective in long-horizon tasks. 38 However, despite the strong empirical performance of HRL, it suffers from major limitations as a 39 scalable pretraining strategy. In particular, the modularity of hierarchical policy architectures, fixed

40 to specific levels of temporal abstraction, precludes unified task representations and necessitates

41 learning a generative model over the subgoal space to interface between policy levels.

42 Learning to predict intermediate goals in a space that may be as high-dimensional as the raw 43 observations poses a difficult generative modeling problem. To ensure that subgoals are physically 44 realistic and reachable in the allotted time, previous work often implements additional processing 45 and verification of proposed subgoals (Zhang et al., 2020; Hatch et al., 2024; Czechowski et al., 46 2024; Zawalski et al., 2024). An alternative is to instead predict in a compact learned latent subgoal 47 space, but simultaneously optimizing subgoal representations and policies results in a nonstationary 48 input distribution to the low-level actor, which can slow and destabilize training (Vezhnevets et al., 49 2017; Levy et al., 2019). The choice of objective for learning such representations, ranging from 50 autoregressive prediction (Seo et al., 2022; Zeng et al., 2023) to metric learning (Tian et al., 2020; 51 Nair et al., 2023; Ma et al., 2023), remains an open question and adds significant complexity to the 52 design and tuning of hierarchical methods. 53 Following the tantalizing promise that flat, one-step policies can be optimal in fully observable, 54 Markovian settings (Puterman, 2005), this work aims to isolate the core advantages of hierarchies 55 for offline GCRL and distill them into a simpler training recipe for a single, unified policy. We 56 begin our empirical analysis by revisiting a state-of-the-art hierarchical method for offline GCRL 57 that significantly outperforms previous approaches on a range of long-horizon goal-reaching tasks. 58 Beyond the original explanation based on improved value function signal-to-noise ratio, we find that 59 separately training a low-level policy on nearby subgoals improves sampling efficiency. We reframe

60 this hierarchical approach as a form of implicit test-time bootstrapping on subgoal-conditioned 61 policies, revealing a conceptual connection to earlier methods that learn subgoal generators and

62 bootstrap directly from subgoal-conditioned policies to train a flat, unified goal-conditioned policy.

63 Building on these insights, we present an inference-based theoretical framework that unifies these 64 ideas and yields Subgoal Advantage-Weighted Policy Bootstrapping (SAW), a novel policy 65 extraction objective for offline GCRL. SAW uses advantage-weighted importance sampling to 66 bootstrap on subgoals sampled directly from data, capturing the long-horizon strengths of hierarchies in a single, flat policy without requiring a generative subgoal model. In evaluations across 20 67 68 state- and pixel-based offline GCRL datasets, our method matches or surpasses all baselines in 69 diverse locomotion and manipulation tasks and scales especially well to complex, long-horizon tasks, 70 being the only existing approach to achieve nontrivial success in the humanoidmaze-giant 71 environment.

72 2 Related Work

73 Our work builds on a rich body of literature encompassing goal-conditioned RL (Kaelbling, 1993), 74 offline RL (Lange et al., 2012; Levine et al., 2020), and hierarchical RL (Dayan & Hinton, 1992; 75 Sutton et al., 1999; Stone, 2008; Bacon et al., 2016; Vezhnevets et al., 2017). The generality of the 76 GCRL formulation enables powerful self-supervised training strategies such as hindsight relabeling 77 (Andrychowicz et al., 2018; Ghosh et al., 2020) and state occupancy matching (Ma et al., 2022). 78 These are often combined with approaches that exploit the recursive subgoal structure of GCRL: 79 either implicitly via quasimetric learning (Wang et al., 2023), probabilistic interpretations (Hoang 80 et al., 2021; Zhang et al., 2021b), and contrastive learning (Eysenbach et al., 2022; Zheng et al., 81 2024); or explicitly through hierarchical decomposition into subtasks (Nachum et al., 2018; Levy 82 et al., 2019; Gupta et al., 2020; Park et al., 2024c). Despite these advances, learning remains difficult 83 for distant goals due to sparse rewards and discounting over time. Many methods use the key insight 84 that actions which are effective for reaching an intermediate subgoal between the current state and the 85 goal are also effective for reaching the final goal. Such subgoals are typically selected via planning 86 (Huang et al., 2019; Zhang et al., 2021a; Hafner et al., 2022), searching within the replay buffer 87 (Eysenbach et al., 2019), or, most commonly, sampling from generative models. Hierarchical methods 88 in particular generate subgoals during inference and use them to query "subpolicies" trained on 89 shorter horizon goals, which are generally easier to learn (Strehl et al., 2009; Azar et al., 2017). Our



Figure 1: Learning with subgoals. Both HIQL and RIS "imagine" subgoals (thought bubbles) en route to the goal (red star) with generative models. However, HIQL samples actions directly from the subgoal-conditioned policy, while RIS regresses (black arrow) a flat goal-conditioned policy towards the subgoal-conditioned action distribution during training. SAW also performs regression but only uses "real" subgoals from the dataset D, weighting the regression more heavily towards distributions conditioned on good subgoals and less (gray arrow) towards bad ones.

90 method also leverages the ease of training subpolicies to effectively learn long-horizon behaviors, but 91 aims to learn a flat, unified policy while avoiding the complexity of training generative models to

92 synthesize new subgoals.

93 Policy bootstrapping. Our work is most closely related to Reinforcement learning with Imagined 94 Subgoals (Chane-Sane et al., 2021, RIS), which, to our knowledge, is the only prior work that performs 95 bootstrapping on *policies*, albeit in the online setting. Similar to goal-conditioned hierarchies, RIS 96 learns a generative model to synthesize "imagined" subgoals that lie between the current state and 97 the goal. Unlike HRL approaches, however, it regresses the full-goal-conditioned policy towards the subgoal-conditioned target, treating the latter as a prior to guide learning and exploration [Figure 1]. 98 99 While RIS yields a flat policy for inference, it still requires the full complexity of a hierarchical policy, 100 including a generative model over the goal space. In contrast, our work extends the core benefits 101 of subgoal-based bootstrapping to offline GCRL with an advantage-based importance weight on 102 subgoals sampled from dataset trajectories, eliminating the need for a subgoal generator altogether.

103 **3** Preliminaries

Problem setting: We consider the problem of offline goal-conditioned RL, described by a Markov decision process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P})$ where \mathcal{S} is the state space, \mathcal{A} the action space, \mathcal{R} : $\mathcal{S} \times \mathcal{S} \to \mathbb{R}$ the goal-conditioned reward function (where we assume that the goal space \mathcal{G} is equivalent to the state space \mathcal{S}), and $\mathcal{P} : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$ the transition function. In the offline setting, we are given a dataset \mathcal{D} of trajectories $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$ previously collected by some arbitrary policy (or multiple policies), and must learn a policy that can reach a specified goal state gfrom an initial state $s_0 \in \mathcal{S}$ without further interaction in the environment, maximizing the objective

$$J(\pi) = \mathbb{E}_{g \sim p(g), \tau \sim p^{\pi}(\tau)} \left[\sum_{t=0}^{\infty} \gamma^{t} r\left(s_{t}, g\right) \right],$$
(1)

where p(g) is the goal distribution and $p^{\pi}(\tau)$ is the distribution of trajectories generated by the policy π and the transition function \mathcal{P} during (online) evaluation.

113 **Offline value learning**: We use a goal-conditioned, action-free variant of implicit Q-learning 114 (Kostrikov et al., 2021, IQL) referred to as goal-conditioned implicit value learning (Park et al., 115 2024c, GCIVL). The original IQL formulation modifies standard value iteration for offline RL by 116 replacing the max operator with an expectile regression, in order to avoid value overestimation for out-of-distribution actions. GCIVL replaces the state-action value function with a value-only estimator

$$\min_{\psi} \mathcal{L}_{\text{GCIVL}}(\psi) = \mathbb{E}_{s,a,g \sim p^{\mathcal{D}}} \left[\ell_{\tau}^2 \left(r(s,g) + \gamma \bar{V}(s',g) - V(s,g) \right) \right],$$
(2)

119 where $\ell_{\tau}^2(x) = |\tau - \mathbb{1}(x < 0)|x^2$ is the expectile loss parameterized by $\tau \in [0.5, 1)$ and $\bar{V}(\cdot)$ 120 denotes a target value function. Note that GCIVL is optimistically biased in stochastic environments, 121 since it directly regresses towards high-value transitions without using Q-values to marginalize over 122 action-independent stochasticity.

123 **Offline policy extraction**: To learn a target subpolicy, we use Advantage-Weighted Regression (Peng 124 et al., 2019, AWR) to extracts a policy from a learned value function. AWR reweights state-action 125 pairs according to their exponentiated advantage with an inverse temperature hyperparameter α , via 126 the objective

$$\max_{\pi} \mathcal{J}_{AWR}(\pi) = \mathbb{E}_{s,a,g\sim\mathcal{D}} \left[e^{\alpha(Q(s,a,g) - V(s,g))} \log \pi(a \mid s,g) \right],\tag{3}$$

127 thus remaining within the support of the data without requiring an additional behavior cloning penalty.

128 4 Understanding Hierarchies in Offline GCRL

In this section, we seek to identify the core reasons behind the empirical success of hierarchies in offline GCRL that can be used to guide the design of a simpler training objective for a flat policy. We first review previous explanations for the benefits of HRL and propose an initial algorithm that seeks to capture these benefits in a flat policy, but find that it still fails to close the performance gap to Hierarchical Implicit Q-Learning (Park et al., 2024c, HIQL), a state-of-the art method. We then identify an additional practical benefit of hierarchical training schemes and show how HIQL exploits this from a policy bootstrapping perspective.

136 4.1 Hierarchies in online and offline GCRL

Previous investigations into the benefits of hierarchical RL in the online setting attribute their success to improved exploration (Stone, 2008) and training value functions with multi-step rewards (Nachum et al., 2019). They demonstrate that augmenting non-hierarchical agents in this manner can largely close the performance gap to hierarchical policies. However, the superior performance of hierarchical methods in the *offline* GCRL setting, where there is no exploration, calls this conventional wisdom into question.

143 HIQL is a state-of-the-art hierarchical offline GCRL method that extracts a high-level policy over 144 subgoals and a low-level policy over primitive actions from a single goal-conditioned value function 145 trained using standard one-step temporal difference learning. HIQL achieves significant performance 146 gains across a number of complex, long-horizon navigation tasks purely through improvement on 147 the policy extraction side, without needing multi-step rewards to train the value function as done in 148 Nachum et al. (2019). While this does not preclude the potential benefits of multi-step rewards for 149 offline GCRL, it does demonstrate that the advantages of hierarchies are not limited to temporally 150 extended value learning, in line with previous claims that the primary bottleneck in offline RL is 151 policy extraction and not value learning (Park et al., 2024b).

152 4.2 Value signal-to-noise ratio in offline GCRL

Instead, HIQL addresses a separate "signal-to-noise ratio" (SNR) issue in value functions conditioned on distant goals, where a combination of sparse rewards and discounting makes it nearly impossible to accurately determine the advantage of one primitive action over another with respect to distant goals. By separating policy extraction into two levels, the low-level actor can instead evaluate the

- relative advantage of actions with respect to nearby subgoals and the high-level policy can utilize
 multi-step advantage estimates to get a clearer learning signal with respect to distant goals.
- To test whether improved SNR in advantage estimates with respect to distant goals is indeed the key to HIQL's superior performance, we propose to utilize subgoals to directly improve advantage estimates
- in a simple baseline method we term goal-conditioned waypoint advantage estimation (GCWAE).
- 162 Briefly, we use the advantage of actions with respect to subgoals generated by a high-level policy as
- 163 an estimator of the undiscounted advantage with respect to the true goal

$$\tilde{A}(s_t, a_t, g) \approx A(s_t, a_t, w), \qquad (4)$$

164 where $w \sim \pi_{sg}^h(w \mid s_t, g)$ is a subgoal sampled from a high-level policy π^h , and the sg subscript indi-165 cates a stop-gradient operator. Apart from using this advantage to directly train a flat policy with AWR, 166 we use the same architectures, sampling distributions, and training objective for π^h as HIQL. Despite 167 large gains over one-step policy learning objectives in several navigation tasks, GCWAE still underper-168 forms its hierarchical counterpart, achieving a 55% success rate on antmaze-large-navigate 169 compared to 90% for HIQL without subgoal representations and 16% for GCIVL with AWR.

170 4.3 It's easier to find good (dataset) actions for closer goals

171 While diagnosing this discrepancy, we observed that training statistics for the two methods were 172 largely identical except for a striking difference in the mean action advantage $A(s_t, a_t, w)$. The advantage was significantly lower for GCWAE, which samples "imagined" subgoals $w \sim \pi^h(\cdot \mid s, g)$ 173 from a high-level policy, than HIQL, which samples directly from the k-step future state distribution 174 $w \sim p^{\mathcal{D}}(s_{t+k} \mid s_t)$ of the dataset. This leads us to an obvious but important insight: in most cases, 175 176 dataset actions are simply better with respect to subgoals sampled from nearby future states in the trajectory than to distant goals or "imagined" subgoals generated by a high-level policy. The dataset is 177 178 far more likely to contain high-advantage actions for goals sampled at the ends of short subsequences, 179 whereas optimal state-action pairs for more distant goals along the trajectory are much rarer due to 180 the combinatorial explosion of possible goal states as the goal-sampling horizon increases.

The practical benefits of being able to easily sample high-advantage state-action-goal tuples are hinted at in Park et al. (2024a), who pose the question "*Why can't we use random goals when training policies*?" after finding that offline GCRL algorithms empirically perform better when only sampling (policy) goals from future states in the same trajectory as the initial state. While their comparison focuses on in-trajectory versus random goals instead of nearer versus farther in-trajectory goals, we hypothesize both observations are driven by similar explanations.

187 4.4 Hierarchies perform test-time policy bootstrapping

Our observations suggest that training policies on nearby goals benefits both from better value SNR in advantage estimates and the ease of sampling good state-action-goal combinations. For brevity, we will refer to such policies trained only on goals of a restricted horizon length as "subpolicies," denoted by π^{sub} and analogous to the low-level policies π^{ℓ} in hierarchies. Now we ask: how do hierarchical methods take advantage of the relative ease of training subpolicies to reach distant goals, and can we use similar strategies to train flat policies?

HIQL separately trains a low-level subpolicy on goals sampled from states at most k steps into the future, reaping all the benefits of policy training with nearby goals. Similar to other goal-conditioned hierarchical methods (Nachum et al., 2018; Levy et al., 2019), it then uses the high-level policy to predict optimal subgoals between the current state and the goal at *test* time, and "bootstraps" by using the subgoal-conditioned action distribution as an estimate for the full goal-conditioned policy.

199 5 Subgoal Advantage-Weighted Policy Bootstrapping

We now seek to unify the above insights into an objective to learn a single, flat goal-reaching policy *without* the additional complexity of HRL. Following the bootstrapping perspective, a direct analogue to hierarchies would use the subpolicy to construct *training* targets, regressing the full goal-conditioned policy towards a target subpolicy conditioned on the output of a subgoal generator $P(a \mid s, \pi^h(w \mid s, g))$. This approach, taken by RIS [Figure 1], still inherits the full complexity of hierarchical policies and then some: it requires learning a subgoal generator, a subpolicy, and an additional flat policy.

207 5.1 Hierarchical RL as inference

To eliminate this additional machinery, we adopt the view of GCRL as probabilistic inference (Levine, 2018). In this framing, the bilevel objectives for HIQL's hierarchical policy and the KL bootstrapping term for RIS's flat policy can be derived from the same inference problem with different choices of variational posterior. Our main insight is that the expectation over generated subgoals can be expressed as an expectation over the dataset distribution with an advantage-based importance weight, yielding our SAW objective. We present an abridged version below and leave the full derivation to Supplementary Section D.

Similar to previous work (Abdolmaleki et al., 2018), we cast the infinite-horizon, discounted GCRL formulation as an inference problem by constructing a probabilistic model via the likelihood function $p(U = 1 | \tau, \{w\}, g) \propto \exp(\beta \sum_{t=0}^{\infty} \gamma^t A(s_t, w_t, g))$, where the binary variable U can be interpreted as the event of reaching the goal g as quickly as possible from state s_t by passing through subgoal state w. The subgoal advantage is defined as $A(s_t, w, g) =$ $-V(s_t, g) + \gamma^k V(w, g) + \sum_{t'=t}^{k-1} r(s_{t'}, g)$, where $w \sim p^{\mathcal{D}}(s_{t+k} | s_t)$. In practice, we follow HIQL and simplify the advantage estimate to $V(w, g) - V(s_t, g)$, i.e., the progress towards the goal achieved by reaching w.

Without loss of generality (since we can represent any flat Markovian policy simply by setting $\pi^{h}(\cdot | s, g)$ to a point distribution on g), we use an inductive bias on the subgoal structure of GCRL to consider prior distributions π of a factored *hierarchical* form

$$p_{\pi}(\tau \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \pi^{\ell}(a_t \mid s_t, w_t) \pi^{h}(w_t \mid s_t, g).$$

The distinctions between hierarchical approaches like HIQL and non-hierarchical approaches such as RIS and SAW begin with our choice of variational posterior. For the former, we would consider

similarly factored distributions, whereas for the latter, we use a *flat* policy $\pi_{\theta}(\tau \mid g)$ that factors as

$$\pi_{\theta}(\tau \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \pi_{\theta}(a_t \mid s_t, g),$$

assuming that the dataset policies are Markovian. We also introduce a variational posterior $q(\{w\} \mid g)$ which factors over a sequence of waypoints $\{w\} = \{w_0, w_1, \ldots\}$ as

$$q(\{w\} \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \pi^{\text{sub}}(a_t \mid s_t, w_t) q(w_t \mid s_t, g).$$

where we treat the target subpolicy π^{sub} as fixed. Using these definitions, we define the evidence lower bound (ELBO) on the optimality likelihood $p_{\pi}(U=1)$ for policy π and goal distribution p(g)

$$\log p_{\pi}(U=1) = \log \int p(g) p_{\pi}(\tau, \{w\} \mid g) p(U=1 \mid \tau, \{w\}, g) d\{w\} d\tau dg$$
$$\geq \mathbb{E}_{\pi_{\theta}(\tau \mid g), q(\{w\} \mid g), p(g)} \log \left[\frac{p(U=1, \tau, \{w\} \mid g)}{\pi_{\theta}(\tau \mid g) q(\{w\} \mid g)} \right] = \mathcal{J}(\pi_{\theta}, q).$$

Algorithm 1 Subgoal Advantage-Weighted Policy Bootstrapping (SAW)

- 1: **Input**: offline dataset \mathcal{D} , goal distribution p(g).
- 2: Initialize value function V_{ϕ} , target subpolicy π_{ψ} , and policy π_{θ} .
- 3: while not converged do
- 4: Train value function: $\phi \leftarrow \phi \lambda \nabla_{\phi} \mathcal{L}_{\text{GCIVL}}(\phi)$ with $(s_t, s_{t+1}) \sim p^{\mathcal{D}}, g \sim p(g)$ [Equation 2]
- 5: end while
- 6: while not converged do
- 7: Train target subpolicy: $\psi \leftarrow \psi \lambda \nabla_{\psi} \mathcal{J}_{AWR}(\psi)$ with $(s_t, a, w) \sim p^{\mathcal{D}}$ [Equation 3]
- 8: end while
- 9: while not converged do
- 10: Train policy: $\theta \leftarrow \theta \lambda \nabla_{\theta} \mathcal{J}_{SAW}(\theta)$ with $(s_t, a, w) \sim p^{\mathcal{D}}, g \sim p(g)$ [Equation 7]
- 11: end while
- 233 Expanding distributions according to their factorizations, dropping terms that are independent of the
- variationals, and rewriting the discounted sum over time as an expectation over the (unnormalized)
- discounted stationary state distribution $\mu_{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t p(s_t = s \mid \pi)$ results in the final objective

$$\mathcal{J}(\theta,q) = \mathbb{E}_{\mu(s),p(g)} \left[\mathbb{E}_{q(w|s,g)} \left[A(s,w,g) \right] - \mathbb{E}_{\pi_{\theta}(a|s,g)} \left[D_{\mathrm{KL}}(q(w|s,g) \| \pi^{h}(w|s,g)) \right] - \mathbb{E}_{q(w|s,g)} \left[D_{\mathrm{KL}}(\pi_{\theta}(a|s,g) \| \pi^{\ell}(a|s,w)) \right] \right], \quad (5)$$

where we optimize an approximation of \mathcal{J} by sampling from the dataset distribution $\mu_{\mathcal{D}}(s)$ (Schulman et al., 2017; Abdolmaleki et al., 2018; Peng et al., 2019).

238 5.2 Eliminating the subgoal generator

Both RIS and HIQL directly parameterize q(w | s, g) with a subgoal generator and optimize according the first line of Equation 5, which can be solved analytically and results in an objective similar to

241 AWR [Equation 3] when projected into the space of parameterized policies

$$\mathcal{J}(q) = \mathbb{E}_{\mu(s), p(g)} \left[\mathbb{E}_{q(w|s,g)}[A(s,w,g)] - \mathbb{E}_{\pi_{\theta}(a|s,g)} \left[D_{\mathrm{KL}}(q(w \mid s,g) \| \pi^{h}(w \mid s,g)) \right] \right]$$

$$\propto \mathbb{E}_{\mu(s), p(g)} \left[\exp(A(s,w,g)) \log q(w \mid s,g) \right].$$

242 Then, RIS trains a flat policy π_{θ} using the third term in Equation 5, which is an expectation over

the subgoals generated by q. Our key insight is that, rather than learning a generative model to approximate p(w | s, g, U = 1), we can instead use a simple application of Bayes' rule:

$$p(w \mid s, g, U = 1) \propto p^{\mathcal{D}}(w \mid s)p(U = 1 \mid s, w, g)$$
$$\propto p^{\mathcal{D}}(w \mid s)\exp(A(s, w, g)),$$

which replaces the expectation over q to yield our subgoal advantage-weighted bootstrapping term

$$\mathbb{E}_{\mu(s),p^{\mathcal{D}}(w|s),p(g)}[\exp(A(s,w,g))D_{\mathrm{KL}}(\pi_{\theta}(a\mid s,g)\|\pi^{\ell}(a\mid s,w))].$$
(6)

We separately learn an approximation to $\pi^{\ell}(a \mid s, w)$, which we do in practice by training a target subpolicy with AWR in a similar fashion to HIQL (whereas RIS uses a exponential moving average of its online goal-conditioned policy as a target). While we omit this for clarity, a subpolicy AWR term can be incorporated into our objective by introducing another optimality variable $p(O_t = 1 \mid \tau, \{w\})$ and a posterior for π^{ℓ} , similar to our derivation for HIQL in Supplementary Section D.

251 While approximating $p(w \mid s, U = 1)$ with q directly and using our importance weight on the dataset 252 distribution are mathematically equivalent, the latter does introduce sampling-based limitations, 253 which we discuss in Appendix A. However, we show empirically that the benefits from lifting the 254 burden of learning a distribution over a high-dimensional subgoal space far outweigh these drawbacks, 255 especially in large state spaces with high intrinsic dimensionality.



Figure 2: **OGBench tasks**. We train SAW on 20 datasets collected from 7 different environments (pictured above) and perform evaluations across 5 state-goal pairs for each dataset.

256 5.3 The SAW objective

The importance weight in Equation 6 allows the policy to bootstrap from subgoals sampled directly from dataset trajectories by ensuring that only subpolicies conditioned on high-advantage subgoals influence the direction of the goal-conditioned policy. We combine our bootstrapping term with an additional learning signal from a (one-step) policy extraction objective utilizing the value function, which improves performance in stitching-heavy environments [Supplementary Section H]. Here, we use one-step AWR [Equation 3], yielding the full SAW objective:

$$\mathcal{J}(\theta) = \mathbb{E}_{p^{\mathcal{D}}(s,a,w),p(g)} \Big[e^{\alpha A(s,a,g)} \log \pi_{\theta}(a \mid s,g) - e^{\beta A(s,w,g)} D_{\mathrm{KL}} \left(\pi_{\theta}(\cdot \mid s,g) \| \pi_{\psi}^{\mathrm{sub}}(\cdot \mid s,w) \right) \Big]$$
(7)

263 where α and β are inverse temperature hyperparameters. This objective provides a convenient 264 dynamic balance between its two terms: as the goal horizon increases, the differences in action values 265 and therefore the contribution of one-step term decreases. This, in turn, downweights the noisier 266 value-based learning signal and shifts emphasis toward the policy bootstrapping term. Finally, we use 267 GCIVL to learn V, resulting in the full training scheme outlined in Algorithm 1.

268 6 Experiments

To assess SAW's ability to reason over long horizons and handle high-dimensional observations, we conduct experiments across 20 datasets corresponding to 7 locomotion and manipulation environments [Figure 2] with both state- and pixel-based observation spaces. We report performance averaged over 5 state-goal pairs for each dataset, yielding **100** total evaluation tasks. Implementation details and hyperparameter settings are discussed in Supplementary Sections F and G, respectively.

274 6.1 Experimental setup

275 We select several environments and their corresponding datasets from the recently released OGBench 276 suite (Park et al., 2024a), a comprehensive benchmark specifically designed for offline GCRL. 277 OGBench provides multiple state-goal pairs for evaluation and datasets tailored to evaluate desirable 278 properties of offline GCRL algorithms, such as the ability to reason over long horizons and stitch 279 across multiple trajectories or combinatorial goal sequences. We use the baselines from the original 280 OGBench paper, which include both one-step and hierarchical state-of-the-art offline GCRL methods. 281 We briefly describe each category of tasks below and baseline algorithms in Appendix C, and 282 encourage readers to refer to the Park et al. (2024a) for further details.

Locomotion: Locomotion tasks require the agent to control a simulated robot to navigate through a maze and reach a designated goal. The agent embodiment varies from a simple 2D point mass with two-dimensional action and observation spaces to a humanoid robot with 21 degrees of freedom and a 69-dimensional state space. In the visual variants, the agent receives a third-person, egocentric $64 \times 64 \times 3$ pixel-based observations, with its location within the maze indicated by the floor color. Maze layouts range from medium to giant, where tasks in the humanoidmaze version of the latter require up to 3000 environment steps to complete.

Environment	Dataset	GCBC	GCIVL	GCIQL	QRL	CRL	HIQL	SAW
	pointmaze-medium-navigate-v0	$9_{\pm 6}$	63 ± 6	$53 \pm s$	82 ± 5	29 ± 7	$79_{\pm 5}$	$97_{\pm 2}$
pointmaze	pointmaze-large-navigate-v0	29 ± 6	45 ± 5	34 ± 3	$86_{\pm 9}$	$39{\scriptstyle~\pm7}$	$58_{\pm 5}$	$85_{\pm 10}$
	pointmaze-giant-navigate-v0	$1_{\pm 2}$	0 ± 0	$0{\scriptstyle~\pm 0}$	$68{\scriptstyle~\pm7}$	$27{\scriptstyle~\pm10}$	$46{\scriptstyle~\pm 9}$	$68{\scriptstyle~\pm 8}$
	antmaze-medium-navigate-v0	$29_{\pm 4}$	72 ± 8	$71_{\pm 4}$	$88_{\pm 3}$	95 ±1	96 ±1	97 ±1
antmaze	antmaze-large-navigate-v0	$24_{\pm 2}$	16 ± 5	34 ± 4	$75_{\pm 6}$	$83_{\pm 4}$	$91_{\pm 2}$	90 ±3
	antmaze-giant-navigate-v0	$0_{\pm 0}$	0 ± 0	0 ± 0	$14{\scriptstyle~\pm3}$	$16{\scriptstyle~\pm3}$	$65{\scriptstyle~\pm5}$	$73{\scriptstyle~\pm 4}$
	humanoidmaze-medium-navigate-v0	8 ±2	$24_{\pm 2}$	$27_{\pm 2}$	$21 \pm s$	$60_{\pm 4}$	89 ± 2	88 ±3
humanoidmaze	humanoidmaze-large-navigate-v0	$1_{\pm 0}$	$2_{\pm 1}$	$2_{\pm 1}$	$5_{\pm 1}$	24 ± 4	$49{\scriptstyle~\pm4}$	$46{\scriptstyle~\pm 4}$
	humanoidmaze-giant-navigate-v0	$0_{\pm 0}$	0 ± 0	0 ± 0	$1_{\pm 0}$	$3{\scriptstyle~\pm 2}$	$12{\scriptstyle~\pm4}$	$35{}_{\pm 4}$
	cube-single-play-v0	6 ±2	53 ± 4	$68_{\pm 6}$	$5_{\pm 1}$	$19_{\pm 2}$	$15_{\pm 3}$	$72^{*}_{\pm 5}$
cube	cube-double-play-v0	$1_{\pm 1}$	36 ± 3	$40{\scriptstyle~\pm5}$	1 ± 0	$10_{\pm 2}$	$6_{\pm 2}$	$40{\scriptstyle~\pm7}$
	cube-triple-play-v0	$1_{\pm 1}$	$1_{\pm 0}$	$3{\scriptstyle~\pm1}$	$0{\scriptstyle~\pm 0}$	$4{}_{\pm 1}$	$3{\scriptstyle~\pm1}$	$4_{\pm 2}$
scene	scene-play-v0	$5_{\pm 1}$	$42{\scriptstyle~\pm4}$	$51_{\pm 4}$	$5_{\pm 1}$	$19{\scriptstyle~\pm 2}$	$38{\scriptstyle~\pm3}$	$63_{\pm 6}$
	visual-antmaze-medium-navigate-v0	11 ± 2	$22_{\pm 2}$	11 ±1	0 ± 0	$94_{\pm 1}$	$93{\scriptstyle~\pm4}$	$95_{\pm 0}$
visual-antmaze	visual-antmaze-large-navigate-v0	4 ± 0	$5_{\pm 1}$	$4_{\pm 1}$	0 ± 0	$84{}_{\pm1}$	$53_{\pm 9}$	$82{\scriptstyle~\pm4}$
	visual-antmaze-giant-navigate-v0	$0_{\pm 0}$	$1_{\pm 1}$	$0{\scriptstyle~\pm 0}$	$0{\scriptstyle~\pm 0}$	$47{\scriptstyle~\pm2}$	$6_{\pm 4}$	$10{\scriptstyle~\pm 2}$
	visual-cube-single-play-v0	$5_{\pm 1}$	60 ±5	$30_{\pm 5}$	$41{\scriptstyle~\pm 15}$	$31{\scriptstyle~\pm 15}$	$89_{\pm 0}$	$88_{\pm 3}$
visual-cube	visual-cube-double-play-v0	$1_{\pm 1}$	10 ± 2	$1_{\pm 1}$	$5_{\pm 0}$	2 ± 1	$39{\scriptstyle~\pm 2}$	$40{\scriptstyle~\pm3}$
	visual-cube-triple-play-v0	$15{\scriptstyle~\pm 2}$	14 ± 2	$15_{\pm 1}$	$16{\scriptstyle~\pm1}$	17 ± 2	$21{}_{\pm 0}$	$20{}_{\pm 1}$
visual-scene	visual-scene-play-v0	$12{\scriptstyle~\pm 2}$	$25{\scriptstyle~\pm3}$	$12{\scriptstyle~\pm 2}$	$10{\scriptstyle~\pm1}$	$11{\scriptstyle~\pm 2}$	$49{\scriptstyle~\pm4}$	47 ± 6

Table 1: Evaluating SAW on state- and pixel-based offline goal-conditioned RL tasks. We compare our method's average (binary) success rate (%) against the numbers reported in Park et al. (2024a) across the five test-time goals for each environment, averaged over 8 seeds (4 seeds for pixel-based visual tasks) with standard deviations after the \pm sign. Numbers within 5% of the best value in the row are in **bold**. Results with an asterisk (*) use different value learning hyperparameters and are discussed further in Section 6.3.

Manipulation: Manipulation tasks use a 6-DoF UR5e robot arm to manipulate object(s), including up to four cubes and a more diverse scene environment that includes buttons, windows, and drawers. The multi-cube and scene environments are designed to test an agent's ability to perform sequential, long-horizon goal stitching and compose together multiple atomic behaviors. The visual variants also provide $64 \times 64 \times 3$ pixel-based observations where certain parts of the environment and robot arm are made semitransparent to ease state estimation.

296 6.2 Locomotion results

297 State-based locomotion: As a method designed for long-horizon reasoning, SAW excels in all 298 variants of the state-based locomotion tasks. It scales particularly well to long horizons, exhibiting 299 the best performance of 73% across all tasks in antmaze-giant-navigate and is the first 300 method to achieve non-trivial success in humanoidmaze-giant-navigate, reaching 35% 301 success compared to the previous state-of-the-art of 12% (Park et al., 2024a). We demonstrate that 302 training subpolicies with subgoal representations scales poorly to the giant maze environments 303 [Figure 3] but are critical to HIQL's performance, emphasizing a fundamental tradeoff in hierarchical 304 methods: subgoal representations are essential for making high-level policy prediction tractable, but 305 those same representations can constrain policy expressiveness and limit overall performance. While 306 other subgoal representation learning objectives may perform better than those derived from the 307 value function, as in HIQL, this highlights the additional design complexity and tuning required for 308 HRL methods. We also implement an offline variant of RIS [Appendix C] and find that it performs 309 significantly worse than SAW with subgoal representations, which we suspect can be explained by 310 our insights in Section 4.3.

Pixel-based locomotion: SAW maintains strong performance when given $64 \times 64 \times 3$ visual observations and scales much better to visual-antmaze-large than does its hierarchical counterpart. However, we do see a significant performance drop in the giant variant relative to the results in the state-based observation space. As a possible explanation for this discrepancy, we observed that value function training diverged for HIQL and SAW in visual-antmaze-giant as well as all visual-humanoidmaze sizes (omitted since no method achieved non-trivial performance). This occurred even without shared policy gradients, suggesting that additional work is



Figure 3: **Subgoal representations scale poorly to high-dimensional control in large state spaces**. Using HIQL's subgoal representations (taken from an intermediate layer of the value function) for SAW's target subpolicy harms performance compared to training directly on observations. However, HIQL fails to learn meaningful behaviors when predicting subgoals directly in the raw observation space. RIS, which bootstraps on generated subgoals at every step, performs the worst of the three.

318 needed to scale offline value learning objectives to very long-horizon tasks with high-dimensional 319 visual observations.

320 6.3 Manipulation results

321 State-based manipulation: SAW consistently matches state-of-the-art performance in cube environ-322 ments and significantly outperforms existing methods in the 5 scene tasks, which require extended 323 compositional reasoning. Interestingly, we found that methods which use expectile regression-based 324 offline value learning methods (GCIVL, GCIQL, HIQL, and SAW) are highly sensitive to value learn-325 ing hyperparameters in the cube-single environment. Indeed, SAW performs more than twice as 326 well on cube-single-play-v0 with settings of $\tau = 0.9$ and $\beta = 0.3$, reaching state-of-the-art 327 performance (72% ± 5 vs. 32% ± 4 with $\tau = 0.7$). While SAW is agnostic to the choice of value 328 learning objective, we make special mention of these changes since they depart from the OGBench 329 convention of fixing value learning hyperparameters for each method across all datasets.

Pixel-based manipulation: In contrast to the state-based environments, SAW and HIQL achieve near-equivalent performance in visual manipulation. This suggests that representation learning, and not long-horizon reasoning or goal stitching, is the primary bottleneck in the visual manipulation environments. While we do not claim any representation learning innovations for this paper, our results nonetheless demonstrate that SAW is able to utilize similar encoder-sharing tricks as HIQL to scale to high-dimensional observation spaces.

336 7 Discussion

337 We presented Subgoal Advantage-Weighted Policy Bootstrapping (SAW), a simple yet effective 338 policy extraction objective that leverages the subgoal structure of goal-conditioned tasks to scale 339 to long-horizon tasks, without learning generative subgoal models. SAW consistently matches or 340 surpasses current state-of-the-art methods across a wide variety of locomotion and manipulation tasks 341 that require different timescales of control, whereas existing methods tend to specialize in particular 342 task categories. Our method especially distinguishes itself in long-horizon reasoning, excelling in the 343 most difficult locomotion tasks and scene-based manipulation. While the simplicity of our objective 344 does introduce some practical limitations related to subgoal sampling, which we discuss in Appendix 345 A, we find that avoiding explicit subgoal prediction is crucial for maintaining performance in large 346 state spaces. By demonstrating a scalable approach to train unified policies for offline GCRL, we 347 believe that SAW takes a step toward realizing the full potential of robotic foundation models in 348 addressing the long-horizon, high-dimensional challenges of real-world control.

349 A Limitations

350 A theoretical limitation of our approach, which is common to all hierarchical methods as well 351 as RIS, occurs in our assumption that the optimal policy can be represented in the factored form 352 $\pi^{\ell}(a \mid s, w)\pi^{h}(w \mid s, g)$. While this is true in theory (since we could trivially set $\pi^{h}(w \mid s, g)$ to a 353 point distribution at g), practical algorithms typically fix the distance of the subgoals to a shorter distance of k steps (or the midpoint in RIS), where subgoals s_{t+k} are sampled from the future 354 state distribution $p_{\text{trai}}^{\mathcal{D}}(s_{t+k} \mid s_t)$. Intuitively, if the dataset contains only suboptimal trajectories 355 356 towards waypoints occurring k steps later which are reachable in fewer than k steps with an optimal 357 low-level policy π^{ℓ} , then the space of sampled subgoals will not contain these further state-subgoal 358 pairs and any approximation of $\pi^h(w \mid s, g, U = 1)$ (whether an explicit subgoal generator or our 359 importance-weighted approach), will suffer additional approximation gaps.

360 However, as discussed in Section 4.3, we find that subgoals sampled from the future state distribution 361 empirically work well with respect to goals also sampled from the future state distribution, which 362 is common in practice (Gupta et al., 2020; Ghosh et al., 2020; Yang et al., 2022; Eysenbach et al., 363 2022; Park et al., 2024c). However, we expect subgoal generator-based methods to have the edge 364 when we sample from a goal distribution for which in-trajectory state-subgoal pairs tend to be highly 365 suboptimal. While a generative subgoal model can synthesize "imagined" subgoals on which to 366 bootstrap, our approach may require alternative subgoal-sampling strategies to reach the same level 367 of performance.

368 **B** Planning Invariance

369 As an aside, we note that the discussions in this paper are closely related to the recently introduced 370 concept of *planning invariance* (Myers et al., 2025), which describes a policy that takes similar 371 actions when directed towards a goal as when directed towards an intermediate waypoint en route 372 to that goal. In fact, we can say that subgoal-conditioned HRL methods achieve a form of planning 373 invariance by construction, since they simply use the actions yielded by waypoint-conditioned policies 374 to reach further goals. By minimizing the divergence between the full goal-conditioned policy and an 375 associated subgoal-conditioned policy, both SAW and RIS can also be seen as implicitly enforcing 376 planning invariance.

377 C Offline GCRL Baseline Algorithms

In this section, we briefly review the baseline algorithms referenced in Table 1. For more thorough implementation details, as well as goal-sampling distributions, interested readers may refer to Appendix C of Park et al. (2024a) as well as the original works.

Goal-conditioned behavioral cloning (GCBC): GCBC is an imitation learning approach that clones
 behaviors using hindsight goal relabeling on future states in the same trajectory.

Goal-conditioned implicit {Q, V}-learning (GCIQL & GCIVL): GCIQL is a goal-conditioned variant of implicit Q-learning (Kostrikov et al., 2021), which performs policy iteration with an expectile regression to avoid querying the learned Q-value function for out-of-distribution actions. Park et al. (2024c) introduced a V-only variant that directly regresses towards high-value transitions [Equation 2], using $r(s, g) + \gamma \overline{V}(s', g)$ as an estimator of Q(s, a, g). Since it does not learn Q-values and therefore cannot marginalize over non-causal factors, it is optimistically biased in stochastic environments.

- Although both baselines are value learning methods that can be used with multiple policy extraction objectives (including our own, which uses GCIVL), the OGBench implementations are paired with
- following objectives: Deep Deterministic Policy Gradient with a behavior cloning penalty term
- 393 (Fujimoto & Gu, 2021, DDPG+BC) for GCIQL, and AWR [Equation 3] for GCIVL.

Quasimetric RL (QRL): QRL (Wang et al., 2023) is a non-traditional value learning algorithm that uses Interval Quasimetric Embeddings (Wang & Isola, 2024, IQE) to enforce quasimetric properties (namely, the triangle inequality and identity of indiscernibles) of the goal-conditioned value function on distances between representations. It uses a constrained "maximal spreading" objective to estimate the shortest paths between states, then learns a one-step dynamics model combined with DDPG+BC to extract a policy from the learned representations.

400 **Contrastive RL (CRL)**: CRL (Eysenbach et al., 2022) is a representation learning algorithm which 401 uses contrastive learning to enforce that the inner product between the learned representations of a 402 state-action pair and the goal state corresponds to the discounted future state occupancy measure of 403 the goal state, which is estimated directly from data using Monte Carlo sampling. CRL then performs 404 one-step policy improvement by choosing actions that maximize the future occupancy of the desired 405 goal state.

406 **Hierarchical implicit Q-learning (HIQL)**: HIQL (Park et al., 2024c) is a policy extraction method 407 that learns two levels of hierarchical policy from the same goal-conditioned value function. The 408 low-level policy π^{ℓ} is trained using standard AWR, and the high-level policy π^{h} is trained using an 409 action-free, multi-step variant of AWR that treats (latent) subgoal states as "actions."

410 Reinforcement learning with imagined subgoals (RIS): RIS (Chane-Sane et al., 2021) is a policy 411 extraction method originally designed for the online GCRL setting, which learns a subgoal generator and a flat, goal-conditioned policy. Unlike SAW, RIS uses a fixed coefficient on the KL term, instead 412 learning a subgoal generator and bootstrapping directly on a target policy (parameterized by an 413 414 exponential moving average of online policy parameters rather than a separately learned subpolicy) 415 conditioned on "imagined" subgoals. It also incorporates a value-based policy learning objective 416 similar to our approach, but learns a Q-function and differentiates directly through the policy with 417 DDPG. 418 To modify RIS for the offline setting in our implementation for Figure 3, we fixed the coefficient

on the KL term to $\beta = 3.0$, trained a subgoal generator identical to the one in HIQL, and replaced the dataset subgoals in SAW with "imagined" subgoals. Otherwise, for fairness of comparison, our offline RIS implementation used the same hyperparameters and architectures as SAW, including a separate target subpolicy network instead of a soft copy of the online policy, subgoals at a fixed distance instead of at midpoints, and AWR instead of DDPG+BC for the policy extraction objective, which we found to perform better in locomotion environments.

425 **References**

Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin
 Riedmiller. Maximum a Posteriori Policy Optimisation, June 2018. URL http://arxiv.org/
 abs/1806.06920. arXiv:1806.06920 [cs].

- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,
 Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience Replay,
 February 2018. URL http://arxiv.org/abs/1707.01495. arXiv:1707.01495.
- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax Regret Bounds for Reinforcement Learning. In *Proceedings of the 34th International Conference on Machine Learning*, pp.
 263–272. PMLR, July 2017. URL https://proceedings.mlr.press/v70/azar17a.
 html. ISSN: 2640-3498.
- 436 Pierre-Luc Bacon, Jean Harb, and Doina Precup. The Option-Critic Architecture, December 2016.
 437 URL http://arxiv.org/abs/1609.05140. arXiv:1609.05140 [cs].
- 438Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo439Fusai, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming440Ke, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang441Shi, James Tanner, Quan Vuong, Anna Walling, Haohuan Wang, and Ury Zhilinsky. π_0 : A442Vision-Language-Action Flow Model for General Robot Control, October 2024. URL http:443//arxiv.org/abs/2410.24164. arXiv:2410.24164 [cs] version: 1.
- Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. Goal-Conditioned Reinforcement Learning
 with Imagined Subgoals. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 1430–1440. PMLR, July 2021. URL https://proceedings.mlr.press/
 v139/chane-sane21a.html. ISSN: 2640-3498.
- Konrad Czechowski, Tomasz Odrzygozdz, Marek Zbysinski, Michal Zawalski, Krzysztof Olejnik,
 Yuhuai Wu, Lukasz Kucinski, and Piotr Milos. Subgoal Search For Complex Reasoning Tasks,
 April 2024. URL http://arxiv.org/abs/2108.11204. arXiv:2108.11204.
- 451 Peter Dayan and Geoffrey E Hinton. Feudal Reinforcement Learning. In Ad452 vances in Neural Information Processing Systems, volume 5. Morgan-Kaufmann,
 453 1992. URL https://proceedings.neurips.cc/paper/1992/hash/
 454 d14220ee66aeec73c49038385428ec4c-Abstract.html.
- Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. Search on the Replay Buffer:
 Bridging Planning and Reinforcement Learning, June 2019. URL http://arxiv.org/abs/
 1906.05253. arXiv:1906.05253 [cs].
- Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R. Salakhutdinov.
 Contrastive Learning as Goal-Conditioned Reinforcement Learning. Advances in Neural Information Processing Systems, 35:35603–35620, December 2022. URL
 https://proceedings.neurips.cc/paper_files/paper/2022/hash/
 e7663e974c4ee7a2b475a4775201ce1f-Abstract-Conference.html.
- 463 Scott Fujimoto and Shixiang Shane Gu. A Minimalist Approach to Offline Reinforcement Learning,
 464 December 2021. URL http://arxiv.org/abs/2106.06860. arXiv:2106.06860.

Dibya Ghosh, Abhishek Gupta, Ashwin Reddy, Justin Fu, Coline Devin, Benjamin Eysenbach, and
 Sergey Levine. Learning to Reach Goals via Iterated Supervised Learning, October 2020. URL
 http://arxiv.org/abs/1912.06088. arXiv:1912.06088.

Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay Policy
Learning: Solving Long-Horizon Tasks via Imitation and Reinforcement Learning. In *Proceedings*of the Conference on Robot Learning, pp. 1025–1037. PMLR, May 2020. URL https://
proceedings.mlr.press/v100/gupta20a.html. ISSN: 2640-3498.

- Danijar Hafner, Kuang-Huei Lee, Ian Fischer, and Pieter Abbeel. Deep Hierarchical Planning from 472 473 Pixels. May 2022. URL https://openreview.net/forum?id=wZk69kjy9_d.
- 474 Kyle B. Hatch, Ashwin Balakrishna, Oier Mees, Suraj Nair, Seohong Park, Blake Wulfe, Masha 475 Itkina, Benjamin Eysenbach, Sergey Levine, Thomas Kollar, and Benjamin Burchfiel. GHIL-Glue:
- 476 Hierarchical Control with Filtered Subgoal Images, October 2024. URL http://arxiv.org/ 477
- abs/2410.20018. arXiv:2410.20018 [cs].
- Christopher Hoang, Sungryull Sohn, Jongwook Choi, Wilka Carvalho, and Honglak Lee. Suc-478 479 cessor Feature Landmarks for Long-Horizon Goal-Conditioned Reinforcement Learning. In 480 Advances in Neural Information Processing Systems, volume 34, pp. 26963–26975. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper/2021/hash/ 481
- 482 e27c71957d1e6c223e0d48a165da2ee1-Abstract.html.
- Zhiao Huang, Fangchen Liu, and Hao Su. Mapping State Space using Landmarks for Universal 483 484 Goal Reaching. In Advances in Neural Information Processing Systems, volume 32. Curran Asso-485 ciates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/ 486 2019/hash/3b712de48137572f3849aabd5666a4e3-Abstract.html.
- 487 Leslie Pack Kaelbling. Learning to achieve goals. In IJCAI, volume 2, pp. 1094-8. Citeseer, 488 1993. URL https://citeseerx.ist.psu.edu/document?repid=rep1&type= 489 pdf&doi=6df43f70f383007a946448122b75918e3a9d6682.
- 490 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline Reinforcement Learning with Implicit Q-491 Learning. October 2021. URL https://openreview.net/forum?id=68n2s9ZJWF8.
- 492 Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch Reinforcement Learning. In Marco 493 Wiering and Martijn van Otterlo (eds.), Reinforcement Learning: State-of-the-Art, pp. 45–73. 494 Springer, Berlin, Heidelberg, 2012. ISBN 978-3-642-27645-3. DOI: 10.1007/978-3-642-27645-3_ 495 2. URL https://doi.org/10.1007/978-3-642-27645-3 2.
- 496 Sergey Levine. Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review, 497 May 2018. URL http://arxiv.org/abs/1805.00909. arXiv:1805.00909 [cs].
- 498 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline Reinforcement Learning: 499 Tutorial, Review, and Perspectives on Open Problems, November 2020. URL http://arxiv. 500 org/abs/2005.01643. arXiv:2005.01643.
- 501 Andrew Levy, George Konidaris, Robert Platt, and Kate Saenko. Learning Multi-Level Hier-502 archies with Hindsight, September 2019. URL http://arxiv.org/abs/1712.00948. arXiv:1712.00948 [cs]. 503
- 504 Jason Yecheng Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline 505 Goal-Conditioned Reinforcement Learning via f-Advantage Regression. Advances 506 in Neural Information Processing Systems, 35:310–323, December 2022. URL 507 https://proceedings.neurips.cc/paper files/paper/2022/hash/ 508 022a39052abf9ca467e268923057dfc0-Abstract-Conference.html.
- 509 Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy 510 Zhang, VIP: Towards Universal Visual Reward and Representation via Value-Implicit Pre-Training, 511 March 2023. URL http://arxiv.org/abs/2210.00030. arXiv:2210.00030 [cs].
- 512 Vivek Myers, Catherine Ji, and Benjamin Eysenbach. Horizon Generalization in Reinforcement 513 Learning, January 2025. URL http://arxiv.org/abs/2501.02709. arXiv:2501.02709 514 [cs].
- Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-Efficient Hierarchical 515
- Reinforcement Learning, October 2018. URL http://arxiv.org/abs/1805.08296. 516 517 arXiv:1805.08296.

- 518 Ofir Nachum, Haoran Tang, Xingyu Lu, Shixiang Gu, Honglak Lee, and Sergey Levine. Why
- Does Hierarchy (Sometimes) Work So Well in Reinforcement Learning?, December 2019. URL
 http://arxiv.org/abs/1909.10618. arXiv:1909.10618 [cs].

Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3M: A Universal
 Visual Representation for Robot Manipulation. In *Proceedings of The 6th Conference on Robot Learning*, pp. 892–909. PMLR, March 2023. URL https://proceedings.mlr.press/
 v205/nair23a.html. ISSN: 2640-3498.

- Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. OGBench: Benchmarking Of fline Goal-Conditioned RL, October 2024a. URL http://arxiv.org/abs/2410.20092.
 arXiv:2410.20092.
- Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is Value Learning Really the Main
 Bottleneck in Offline RL?, October 2024b. URL http://arxiv.org/abs/2406.09329.
 arXiv:2406.09329.

Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline GoalConditioned RL with Latent States as Actions, March 2024c. URL http://arxiv.org/abs/
2307.11949. arXiv:2307.11949 [cs].

- Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation Policies with Hilbert Representations,
 May 2024d. URL http://arxiv.org/abs/2402.15567. arXiv:2402.15567 [cs].
- Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-Weighted Regression:
 Simple and Scalable Off-Policy Reinforcement Learning, October 2019. URL http://arxiv.
 org/abs/1910.00177. arXiv:1910.00177.
- 539 Inc Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny 540 Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya 541 Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming 542 Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, 543 544 James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury 545 Zhilinsky. $\pi_{0.5}$: a Vision-Language-Action Model with Open-World Generalization, April 2025. 546 URL http://arxiv.org/abs/2504.16054. arXiv:2504.16054 [cs] version: 1.
- Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
 Wiley and Sons, 2005.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust
 Region Policy Optimization, April 2017. URL http://arxiv.org/abs/1502.05477.
 arXiv:1502.05477 [cs].
- Younggyo Seo, Kimin Lee, Stephen James, and Pieter Abbeel. Reinforcement Learning with ActionFree Pre-Training from Videos, June 2022. URL http://arxiv.org/abs/2203.13880.
 arXiv:2203.13880 [cs].
- Nicholas K. Jong and Todd Hester and Peter Stone. The Utility of Temporal Abstraction
 in Reinforcement Learning. 2008. URL https://www.cs.utexas.edu/~ai-lab/
 ?AAMAS08-jong.
- Alexander L. Strehl, Lihong Li, and Michael L. Littman. Reinforcement Learning in Finite MDPs:
 PAC Analysis. *Journal of Machine Learning Research*, 10(84):2413–2444, 2009. ISSN 1533-7928.
 URL http://jmlr.org/papers/v10/strehl09a.html.
- Richard S. Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112(1):181–
 211, August 1999. ISSN 0004-3702. DOI: 10.1016/S0004-3702(99)00052-1. URL https:
 //www.sciencedirect.com/science/article/pii/S0004370299000521.

Stephen Tian, Suraj Nair, Frederik Ebert, Sudeep Dasari, Benjamin Eysenbach, Chelsea Finn, and
 Sergey Levine. Model-Based Visual Planning with Self-Supervised Functional Distances. October
 2020. URL https://openreview.net/forum?id=UcoXdfrORC.

2020. OKL https://openreview.het/iorum/id=ocoxdriokc.

Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David
 Silver, and Koray Kavukcuoglu. FeUdal Networks for Hierarchical Reinforcement Learning,
 March 2017. URL http://arxiv.org/abs/1703.01161. arXiv:1703.01161 [cs].

Kevin Wang, Ishaan Javali, Michal Bortkiewicz, Tomasz Trzcinski, and Benjamin Eysenbach. 1000
Layer Networks for Self-Supervised RL: Scaling Depth Can Enable New Goal-Reaching Capabilities, March 2025. URL http://arxiv.org/abs/2503.14858. arXiv:2503.14858
[cs].

Tongzhou Wang and Phillip Isola. Understanding Contrastive Representation Learning through
 Alignment and Uniformity on the Hypersphere, August 2022. URL http://arxiv.org/
 abs/2005.10242. arXiv:2005.10242 [cs].

Tongzhou Wang and Phillip Isola. Improved Representation of Asymmetrical Distances with Interval
Quasimetric Embeddings, January 2024. URL http://arxiv.org/abs/2211.15120.
arXiv:2211.15120 [cs].

Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal Goal-Reaching Rein forcement Learning via Quasimetric Learning, November 2023. URL http://arxiv.org/
 abs/2304.01203. arXiv:2304.01203 [cs].

Rui Yang, Yiming Lu, Wenzhe Li, Hao Sun, Meng Fang, Yali Du, Xiu Li, Lei Han, and Chongjie
Zhang. Rethinking Goal-conditioned Supervised Learning and Its Connection to Offline RL,
February 2022. URL http://arxiv.org/abs/2202.04478. arXiv:2202.04478 [cs].

Rui Yang, Lin Yong, Xiaoteng Ma, Hao Hu, Chongjie Zhang, and Tong Zhang. What is Essential
for Unseen Goal Generalization of Offline Goal-conditioned RL? In *Proceedings of the 40th International Conference on Machine Learning*, pp. 39543–39571. PMLR, July 2023. URL
https://proceedings.mlr.press/v202/yang23q.html. ISSN: 2640-3498.

Michal Zawalski, Michal Tyrolski, Konrad Czechowski, Tomasz Odrzygozdz, Damian Stachura,
Piotr Piekos, Yuhuai Wu, Lukasz Kucinski, and Piotr Milos. Fast and Precise: Adjusting Planning
Horizon with Adaptive Subgoal Search, May 2024. URL http://arxiv.org/abs/2206.
00702. arXiv:2206.00702.

Zilai Zeng, Ce Zhang, Shijie Wang, and Chen Sun. Goal-Conditioned Predictive Coding for Offline
 Reinforcement Learning. Advances in Neural Information Processing Systems, 36:25528–25548,
 December 2023. URL https://proceedings.neurips.cc/paper_files/paper/
 2023/hash/51053d7b8473df7d5a2165b2a8ee9629-Abstract-Conference.

- 599 html.
- Lunjun Zhang, Ge Yang, and Bradly C. Stadie. World Model as a Graph: Learning Latent Landmarks
 for Planning. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 12611–
 12620. PMLR, July 2021a. URL https://proceedings.mlr.press/v139/zhang21x.
 html. ISSN: 2640-3498.
- Tianjun Zhang, Benjamin Eysenbach, Ruslan Salakhutdinov, Sergey Levine, and Joseph E. Gonzalez.
 C-Planning: An Automatic Curriculum for Learning Goal-Reaching Tasks, October 2021b. URL
 http://arxiv.org/abs/2110.12080. arXiv:2110.12080 [cs].
- Tianren Zhang, Shangqi Guo, Tian Tan, Xiaolin Hu, and Feng Chen. Generating
 Adjacency-Constrained Subgoals in Hierarchical Reinforcement Learning. In Advances *in Neural Information Processing Systems*, volume 33, pp. 21579–21590. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/hash/
 f5f3b8d720f34ebebceb7765e447268b-Abstract.html.

- 612 Chongyi Zheng, Ruslan Salakhutdinov, and Benjamin Eysenbach. Contrastive Difference Predictive
- 613 Coding, February 2024. URL http://arxiv.org/abs/2310.20141. arXiv:2310.20141.

614 **Supplementary Materials** 615 *The following content was not necessarily subject to peer review.*

617 D Derivations of HIQL, RIS, and SAW Objectives

618 We cast the infinite-horizon, discounted GCRL formulation as an inference problem by constructing

a probabilistic model via the likelihood function

$$p(U_t = 1 \mid \tau, \{w\}, g) \propto \exp\left(\beta \sum_{t=0}^{\infty} \gamma^t A(s_t, w_t, g)\right),$$

where β is an inverse temperature parameter and the binary variable U can be intuitively understood as the event of reaching the goal g as quickly as possible by passing through a subgoal w, or passing

622 through a subgoal w which is on the shortest path between s_t and g.

623 We consider prior distributions π of a factored *hierarchical* form

$$p_{\pi}(\tau \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \, \pi^{\ell}(a \mid s_t, w_t) \, \pi^h(w_t \mid s_t, g).$$

624 D.1 HIQL derivation

625 Since HIQL learns two levels of a policy, we can use a variational posterior of the same form

$$\pi_{\theta,\psi}(\tau \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \, \pi_{\theta}^{\ell}(a_t \mid s_t, w_t) \, \pi_{\psi}^{h}(w_t \mid s_t, g).$$

To incorporate training of the low-level policy, we also construct an additional probabilistic model for the optimality of primitive actions towards a waypoint w (note that this can also be done to

628 incorporate target policy training into the SAW objective, but we leave it out for brevity)

$$p(O_t = 1 \mid \tau, \{w\}) \propto \exp\left(\alpha \sum_{t=0}^{\infty} \gamma^t A(s_t, a_t, w_t)\right).$$

629 With these definitions, we define the evidence lower bound (ELBO) on the joint optimality likelihood 630 $p_{\pi}(O = 1, U = 1)$ for policy π

$$\log p_{\pi}(O = 1, U = 1) = \log \int p(g) p(O = 1, U = 1, \tau, \{w\} \mid g) d\{w\} d\tau dg$$

= $\log \int p(g) \pi_{\theta}(\tau \mid g) \pi_{\psi}(\{w\} \mid g) \frac{p(O = 1, U = 1, \tau, \{w\} \mid g)}{\pi_{\theta}(\tau \mid g) \pi_{\psi}(\{w\} \mid g)} d\{w\} d\tau dg$
= $\log \mathbb{E}_{\pi_{\theta}(\tau \mid g), \pi_{\psi}(\{w\} \mid g), p(g)} \left[\frac{p(O = 1, U = 1, \tau, \{w\} \mid g)}{\pi_{\theta}(\tau \mid g) \pi_{\psi}(\{w\} \mid g)} \right]$
 $\geq \mathbb{E}_{\pi_{\theta}(\tau \mid g), \pi_{\psi}(\{w\} \mid g), p(g)} \log \left[\frac{p(O = 1, U = 1, \tau, \{w\} \mid g)}{\pi_{\theta}(\tau \mid g) \pi_{\psi}(\{w\} \mid g)} \right] = \mathcal{J}(\theta, \psi).$

Expanding the fraction, moving the log inside, and dropping the start state distribution $p(s_0)$ and transition distributions $p(s_{t+1} | s_t, a_t)$, which are fixed with respect to θ and ψ , gives us

$$\mathbb{E}_{\pi_{\theta}(\tau|g), \pi_{\psi}(\{w\}|g), p(g)} \left[\alpha \sum_{t=0}^{\infty} \gamma^{t} A(s_{t}, a_{t}, w_{t}) + \beta \sum_{t=0}^{\infty} \gamma^{t} A(s_{t}, w_{t}, g) + \sum_{t=0}^{\infty} \log \left(\frac{\pi^{\ell}(a_{t} \mid s_{t}, w_{t}) \pi^{h}(w_{t} \mid s_{t}, g)}{\pi_{\theta}(a_{t} \mid s_{t}, g) \pi_{\psi}(w_{t} \mid s_{t}, g)} \right) \right]$$

633 We rewrite the discounted sum over time as an expectation over the (unnormalized) discounted 634 stationary state distribution $\mu_{\pi}(s) = \sum_{t=0}^{\infty} \gamma^t p(s_t = s \mid \pi)$ induced by policy π . In practice, however, 635 we optimize an approximation of $\mathcal{J}(\theta, \psi)$ by sampling from the dataset distribution over states $\mu_{\mathcal{D}}$. 636 For brevity, we omit the conditionals in the expectations below, defining $\pi_{\theta}(a) \coloneqq \pi_{\theta}(a \mid s, g)$ and 637 $\pi_{\psi}(w) \coloneqq \pi_{\psi}(w \mid s, g)$ in the expectations below

$$\mathbb{E}_{\mu(s),\pi_{\theta}(a),\pi_{\psi}(w),p(g)} \left[\alpha A(s,a,w) + \log \left[\frac{\pi^{\ell}(a \mid s,w)}{\pi_{\theta}(a \mid s,w)} \right], \right] \\
+ \mathbb{E}_{\mu(s),\pi_{\psi}(w),p(g)} \left[\beta A(s,w,g) + \log \left[\frac{\pi^{h}(w \mid s,g)}{\pi_{\psi}(w \mid s,g)} \right] \right] \\
= \mathbb{E}_{\mu(s),\pi_{\theta}(a),\pi_{\psi}(w),p(g)} \left[\alpha A(s,a,w) \right] - \mathbb{E}_{\mu(s),\pi_{\psi}(w),p(g)} \left[D_{\mathrm{KL}} \left[\pi_{\theta}(a \mid s,w) \| \pi^{\ell}(a \mid s,w) \right] \right] \\
+ \mathbb{E}_{\mu(s),\pi_{\psi}(w),p(g)} \left[\beta A(s,w,g) \right] - \mathbb{E}_{\mu(s),p(g)} \left[D_{\mathrm{KL}} \left[\pi_{\psi}(w \mid s,g) \| \pi^{h}(w \mid s,g) \right] \right]. \tag{8}$$

HIQL separately optimizes the two summation terms, which correspond to the low- and high-level policies, respectively. Forming the Lagrangian with the normalization condition and solving for the

optimal low- and high-level policies, as done in Abdolmaleki et al. (2018) and Peng et al. (2019),

641 yields the HIQL AWR objectives:

$$\begin{aligned} \mathcal{J}_{\pi^{\ell}}(\theta) &= \mathbb{E}_{\mu(s),\pi_{\theta}(a),\pi_{\psi}(w),p(g)} \left[\exp(\alpha A(s,a,w)) \log \pi_{\theta}^{\ell}(a \mid s,w) \right] \\ \mathcal{J}_{\pi^{h}}(\psi) &= \mathbb{E}_{\mu(s),\pi_{\psi}(w),p(g)} \left[\exp(\beta A(s,w,g)) \log \pi_{\psi}^{h}(w \mid s,g) \right]. \end{aligned}$$

While our derivation produces an on-policy expectation over actions and subgoals, the sampling distribution over states, actions, subgoals, and goals in the offline setting varies in practice (see Appendix C of Park et al. (2024a) for commonly used goal distributions).

645 D.2 RIS and SAW derivations

646 Unlike HIQL, both RIS and SAW seek to learn a unified flat policy, and therefore we choose a policy 647 posterior $\pi_{\theta}(\tau)$ that factors as

$$\pi_{\theta}(\tau \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \pi_{\theta}(a_t \mid s_t, g).$$

648 We also introduce a variational posterior q which factors over a sequence of waypoints $\{w\} = \{w_0, w_1, \ldots\}$ as

$$q(\{w\} \mid g) = p(s_0) \prod_{t=0}^{\infty} p(s_{t+1} \mid s_t, a_t) \pi^{\text{sub}}(a_t \mid s_t, w_t) q(w_t \mid s_t, g),$$

- 650 where π^{sub} is a target subpolicy and is treated as fixed with respect to the parameters of the posteriors.
- Using these definitions, we define the evidence lower bound (ELBO) on the likelihood of subgoal

652 optimality $p_{\pi}(U=1)$ for policy π

$$\log p_{\pi}(U=1) = \log \int p(g)p(U=1,\tau,\{w\} \mid g)d\tau d\{w\}dg$$

= $\log \int p(g)\pi_{\theta}(\tau \mid g)q(\{w\} \mid g)\frac{p(U=1,\tau,\{w\} \mid g)}{\pi_{\theta}(\tau \mid g)q(\{w\} \mid g)}d\{w\}d\tau dg$
= $\log \mathbb{E}_{\pi_{\theta}(\tau \mid g),q(\{w\} \mid g),p(g)} \left[\frac{p(U=1,\tau,\{w\} \mid g)}{\pi_{\theta}(\tau \mid g)q(\{w\} \mid g)}\right]$
 $\geq \mathbb{E}_{\pi_{\theta}(\tau \mid g),q(\{w\} \mid g),p(g)}\log \left[\frac{p(U=1,\tau,\{w\} \mid g)}{\pi_{\theta}(\tau \mid g)q(\{w\} \mid g)}\right] = \mathcal{J}(\pi_{\theta},q).$

Expanding the fraction, moving the log inside, and dropping the start state distribution $p(s_0)$, transition distributions $p(s_{t+1} | s_t, a_t)$, and target subpolicy $\pi^{\text{sub}}(a_t | s_t, w_t)$, which are fixed with respect to the variationals, leaves us with

$$\mathbb{E}_{\pi_{\theta}(\tau|g), q(\{w\}|g), p(g)} \left[\beta \sum_{t=0}^{\infty} \gamma^{t} A(s_{t}, w_{t}, g) + \sum_{t=0}^{\infty} \log \left[\frac{\pi^{\ell}(a_{t} \mid s_{t}, w_{t}) \pi^{h}(w_{t} \mid s_{t}, g)}{\pi_{\theta}(a_{t} \mid s_{t}, g) q(w_{t} \mid s_{t}, g)} \right] \right]$$

656 Once again, we express the discounted sum over time as an expectation over the discounted stationary 657 state distribution $\mu(s)$ and omit the conditionals in the expectation over $\pi_{\theta}(a)$ and q(w) for brevity. 658 Simplifying gives us

$$\mathbb{E}_{\mu(s),\pi_{\theta}(a),q(w),p(g)} \left[\beta A(s,w,g) + \log \left[\frac{\pi^{h}(w \mid s,g)}{q(w \mid s,g)} \right] + \log \left[\frac{\pi^{\ell}(a \mid s,w)}{\pi_{\theta}(a \mid s,g)} \right] \right] \\
= \mathbb{E}_{\mu(s),q(w),p(g)} \left[\beta A(s,w,g) \right] - \mathbb{E}_{\mu(s),p(g)} \left[D_{\mathrm{KL}} \left(q \left(w \mid s,g \right) \| \pi^{h} \left(w \mid s,g \right) \right) \right] \\
+ \mathbb{E}_{\mu(s),q(w),p(g)} \left[D_{\mathrm{KL}} \left(\pi_{\theta} \left(a \mid s,g \right) \| \pi^{\ell} \left(a \mid s,w \right) \right) \right]$$

RIS: RIS partitions this objective into two parts and optimizes them separately, where the subgoal generator is trained according to the loss

$$\mathcal{J}^{h}(q) = \mathbb{E}_{\mu(s), q(w), p(g)}[A(s, w, g)] - \mathbb{E}_{\mu(s), p(g)}\left[D_{\mathrm{KL}}(q(w \mid s, g) \| \pi^{h}(w \mid s, g))\right],$$

661 which is identical to the objective for the HIQL high-level policy π_{ψ} and yields the same AWR-like 662 objective over subgoals $\exp(A(s, w, g)) \log q(w \mid s, g)$. We then incorporate the remaining KL 663 divergence term into the objective for the flat policy posterior π_{θ}

$$\mathcal{J}(\theta) = \mathbb{E}_{\mu(s), q(w), p(g)}[D_{\mathrm{KL}}(\pi_{\theta}(a \mid s, g) \| \pi^{\ell}(a \mid s, w))],$$

664 which is an expectation over subgoals drawn from the (simultaneously learned) subgoal generator.

665 SAW: Instead of directly learning q, we use Bayes' rule and our earlier definition of p(U = 1 | s, w, g)666 to directly approximate the posterior distribution over subgoals p(w | s, g, U = 1), where

$$p(w \mid s, g, U = 1) \propto p^{\mathcal{D}}(w \mid s)p(U = 1 \mid s, w, g)$$
$$\propto p^{\mathcal{D}}(w \mid s) \exp(A(s, w, g)).$$

667 Although the proportionality constant in the first line is the $p_{\pi}(U = 1)$, which is the subject of 668 our optimization, we note that approximating the expectation over subgoals w corresponds to the 669 expectation (E) step in a standard expectation-maximization (EM) procedure (Abdolmaleki et al., 670 2018). Because we are only seeking to fit the shape of the optimal variational posterior over subgoals

for the purposes of approximating the expectation over q, and not maximizing $p_{\pi}(U=1)$ (the M star), we can treat $p_{\pi}(U=1)$ as constant with respect to g to get

step), we can treat $p_{\pi}(U = 1)$ as constant with respect to q to get

$$\begin{aligned} \mathcal{J}(\theta) &= \mathbb{E}_{\mu(s),q(w),\,p(g)}[D_{\mathrm{KL}}(\pi_{\theta}(a \mid s, g) \| \pi^{\ell}(a \mid s, w))] \\ &= \int \mu(s)\,p(g)\,q(w \mid s)[D_{\mathrm{KL}}(\pi_{\theta}(a \mid s, g) \| \pi^{\ell}(a \mid s, w))]dw\,dg\,ds \\ &\propto \int \mu(s)\,p(g)\,p^{\mathcal{D}}(w \mid s)\exp(A(s, w, g))[D_{\mathrm{KL}}(\pi_{\theta}(a \mid s, g) \| \pi^{\ell}(a \mid s, w))]dw\,dg\,ds \\ &= \mathbb{E}_{\mu(s),p^{\mathcal{D}}(w \mid s),p(g)}\exp(A(s, w, g))[D_{\mathrm{KL}}(\pi_{\theta}(a \mid s, g) \| \pi^{\ell}(a \mid s, w))], \end{aligned}$$

673 which yields our subgoal advantage-weighted bootstrapping term in Equation 6.

674 E Computational Resources

All experiments were conducted on a cluster consisting of Nvidia GeForce RTX 3090 GPUs with 24
GB of VRAM and Nvidia GeForce RTX 3070 GPUs with 8 GB of VRAM. State-based experiments
take around 4 hours to run for the largest environments (humanoidmaze-giant-navigate)
and visual experiments up to 12 hours.

679 F Implementation Details

Target policy: While Chane-Sane et al. (2021) use a exponential moving average (EMA) of the online policy parameters θ as the target policy prior $\pi_{\overline{\theta}}$, we instead simply train a smaller policy network parameterized separately by ψ on (sub)goals sampled from k steps into the future, where k is a hyperparameter. We find that this leads to faster training and convergence, albeit with a small increase in computational complexity.

685 Architecture: During our experiments, we observed that the choice of network architecture for 686 both the value function and policy networks had a significant impact on performance in several 687 environments. Instead of taking in the raw concatenated state and goal inputs, HIQL prepends a 688 subgoal representation module consisting of an additional three-layer MLP followed by a bottleneck 689 layer of dimension 10 and a length-normalizing layer that projects state-(sub)goal representations to 690 the surface of a hypersphere with radius equal to the dimension of the input vector. The value and 691 low-level policy networks receive this representation in place of the goal information, as well as the 692 raw (in state-based environments) or encoded (in pixel-based environments) state information. We 693 found that simply adding these additional layers (separately) to the value and actor network encoders 694 significantly boosted performance in state-based locomotion tasks, with modifications to the former 695 improving training stability in pointmaze and modifications to the latter being critical for good 696 performance in the antmaze and humanoidmaze environments.

697 While we did not perform comprehensive architectural ablations due to computational limitations, we 698 note that the desirable properties of the unit hypersphere as a representation space are well-studied in 699 contrastive learning (Wang & Isola, 2022) and preliminary work by Wang et al. (2025) has explored 690 the benefits of scaling network depth for GCRL (albeit with negative results for the offline setting). 691 Further studying the properties of representations emerging from these architectural choices may 692 inform future work in representation learning for offline GCRL.

703 G Hyperparameters

We find that our method is robust to hyperparameter selection for different horizon lengths and environment types in locomotion tasks, but is more sensitive to choices of the value learning expectile parameter τ and the temperature parameter β of the divergence term in manipulation tasks (see Supplementary Section H for training curves of different β settings). Unless otherwise stated in Table 2, all common hyperparameters are the same as specified in Park et al. (2024a) and state, subgoal, and goal-sampling distributions are identical to those for HIQL.

Environment Type	Dataset	Expectile τ	AWR α	KLD β	Subgoal steps \boldsymbol{k}
	pointmaze-medium-navigate-v0	0.7	3.0	3.0	25
pointmaze	pointmaze-large-navigate-v0	0.7	3.0	3.0	25
	pointmaze-giant-navigate-v0	0.7	3.0	3.0	25
	antmaze-medium-navigate-v0	0.7	3.0	3.0	25
antmaze	antmaze-large-navigate-v0	0.7	3.0	3.0	25
	antmaze-giant-navigate-v0	0.7	3.0	3.0	25
	humanoidmaze-medium-navigate-v0	0.7	3.0	3.0	25
humanoidmaze	humanoidmaze-large-navigate-v0	0.7	3.0	3.0	25
	humanoidmaze-giant-navigate-v0	0.7	3.0	3.0	25
	visual-antmaze-medium-navigate-v0	0.7	3.0	3.0	25
visual-antmaze	visual-antmaze-large-navigate-v0	0.7	3.0	3.0	25
	visual-antmaze-giant-navigate-v0	0.7	3.0	3.0	25
	cube-single-play-v0	0.9	3.0	0.3	10
cube	cube-double-play-v0	0.7	3.0	1.0	10
	cube-triple-play-v0	0.7	3.0	1.0	10
scene	scene-play-v0	0.7	3.0	1.0	10
	visual-cube-single-play-v0	0.7	3.0	3.0	10
visual-cube	visual-cube-double-play-v0	0.7	3.0	3.0	10
	visual-cube-triple-play-v0	0.7	3.0	3.0	10
visual-scene	visual-scene-play-v0	0.7	3.0	3.0	10

Table 2: SAW hyperparameters. Each cell indicates the hyperparameters for the corresponding environment and dataset. From left to right, these hyperparameters are: the expectile parameter τ for GCIVL, the one-step AWR temperature α (used for training both the target and policy networks), the temperature on the KL divergence term β , and the number of subgoal steps k.

710 H Ablations

In this section, we ablate various components of our objective and assess its sensitivity to varioushyperparameters.

713 H.1 One-step AWR ablation

We ablate the one-step AWR term in our objective, which is akin to training purely on bootstrapped policies from a target policy (which itself is trained with AWR). Note that ablating the bootstrapping term simply recovers the **GCIVL** baseline. We observe that ablations to the one-step term primarily affect performance in short-horizon, stitching-heavy tasks such as the simpler manipulation environments. On the other hand, performance is largely unaffected in longer-horizon manipulation and locomotion tasks, confirming our initial hypotheses that the bulk of SAW's performance in more complex tasks is due to policy bootstrapping rather than one-step policy extraction.

721 H.2 Hyperparameter sensitivity

Here, we investigate SAW's sensitivity to the β inverse temperature hyperparameter and run different settings of $\beta \in \{0.3, 1.0, 3.0, 10.0\}$ across selected state-based environments. We observe a similar pattern to the one-step AWR ablation experiments, where the simpler manipulation environments are much more sensitive to hyperparameter settings compared to more complex, long-horizon tasks.



Figure 4: Training curves for cube-single-play and cube-double-play with one-step ablations.



Figure 5: Training curves for scene-play and antmaze-large-navigate with one-step ablations.



Figure 6: Training curves for cube-single-play and cube-double-play with different values of β (where the default hyperparameters settings are $\beta = 0.3$ and $\beta = 1$, respectively).



Figure 7: Training curves for scene-play and antmaze-large-navigate with different values of β (where the default hyperparameters settings are $\beta = 1$ and $\beta = 3$, respectively).