
Dual Advantage Fields

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Offline goal-conditioned reinforcement learning requires both long-horizon reachability estimates and local action comparisons. Dual goal representations provide
2 value fields that capture global goal reachability, but they do not directly specify
3 which action should be preferred at a given state. We propose Dual Advantage
4 Fields, a policy-extraction method that turns a bilinear dual value model into a
5 local advantage signal. Under bilinear dual parameterization, the goal embedding
6 is the gradient of the value field with respect to the state representation. DAF
7 learns an action-effect model that predicts the discounted feature displacement induced
8 by an action and scores actions by the alignment between this displacement and
9 the goal direction. In the realizable case, this score equals the goal-conditioned
10 Bellman advantage, yielding a standard local policy-improvement guarantee. On
11 OGBench locomotion, manipulation, and puzzle tasks, DAF improves aggregate
12 RLiab metrics and performs strongly in settings where locally correct actions
13 differ from direct movement toward the final goal.
14

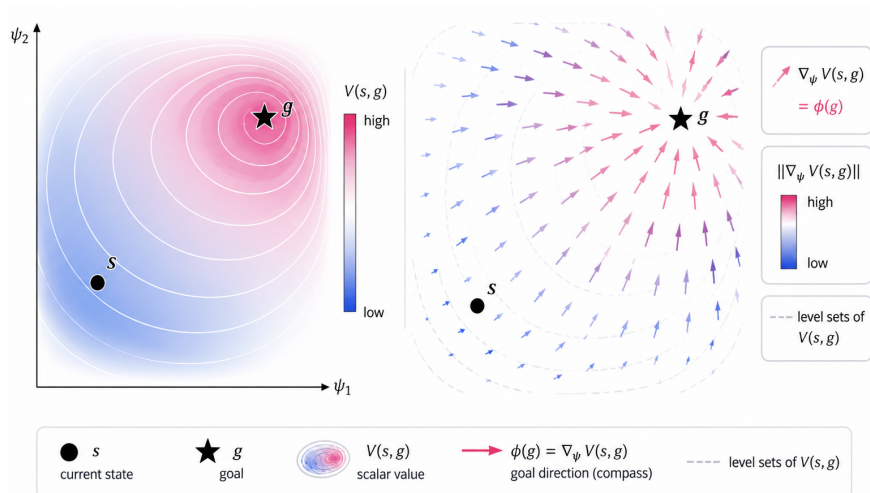


Figure 1: **Dual Advantage Fields.** A dual goal-conditioned value model defines a global value surface over state representations. DAF converts this global surface into a local action-comparative signal by predicting how each action moves the state representation and measuring whether this movement aligns with the goal direction.

15 1 Introduction

16 Goal-conditioned reinforcement learning (GCRL) aims to learn policies that reach arbitrary goals
17 from a fixed dataset of prior experience. This requires solving two different problems at once. First,

18 the agent must reason globally: it must infer how states are connected over long horizons so that
19 behavior observed in one part of the dataset can be stitched together with behavior observed elsewhere.
20 Second, the agent must act locally: at the current state, it must decide which available action makes
21 the most progress toward the requested goal. A good goal-conditioned agent therefore needs both a
22 global map of reachability and a local compass for action selection.

23 Recent dual goal representations provide a strong answer to the first problem. They parameterize
24 a goal-conditioned value function as a bilinear interaction between a state embedding and a goal
25 embedding $V_\theta(s, g) = \psi_\theta(s)^\top \phi_\theta(g)$. This structure induces a value surface for each goal, where
26 states that are more reachable or desirable for the goal receive higher values. Such value surfaces
27 are well suited for long-horizon reasoning: they encode temporal structure, support stitching across
28 offline trajectories, and generalize across state-goal pairs. However, a value surface alone does not
29 directly answer the local control question. It says how good the current state is for a goal, but not
30 which action should be preferred among the actions available at that state.

31 This distinction is central in offline GCRL. Policy extraction requires an action-comparative signal.
32 Two actions can start from the same state and therefore share the same value $V_\theta(s, g)$, while only
33 one of them may move the agent toward the goal. What is missing is not another global estimate of
34 reachability, but a local advantage-like quantity: a way to score whether an action changes the state
35 in a direction that improves goal-conditioned value.

36 Our key observation is that this local signal is already implicit in the geometry of dual representations.
37 Under the bilinear parameterization above, the goal embedding $\phi_\theta(g)$ is the direction in state-
38 representation space along which the goal-conditioned value increases:

$$\nabla_\psi V_\theta(s, g) = \phi_\theta(g).$$

39 Thus, if an action induces a displacement in the state representation, its usefulness for the goal
40 can be evaluated by a simple geometric test: does the predicted displacement align with the goal
41 direction? This turns goal-conditioned policy improvement into a local alignment problem in the
42 dual representation space.

43 We introduce *Dual Advantage Fields* (DAF), a policy-extraction method that makes this geometry
44 explicit; see Figure 1. DAF learns an action-effect model that predicts the discounted change in
45 the state representation caused by an action. It then scores actions by the inner product between
46 this predicted action effect and the goal embedding. The resulting score is local, goal-conditioned,
47 and action-comparative: it prefers actions whose predicted latent effect points in the direction of
48 increasing value for the goal.

49 This perspective leads to a simple principle for offline GCRL: *global value fields should be paired*
50 *with local advantage fields*. Dual representations provide the global map; DAF extracts from the same
51 representation space the local compass needed for policy improvement. This yields an efficient actor-
52 free mechanism for policy extraction: rather than learning a separate goal-conditioned action-value
53 function, DAF reuses the geometry of the dual critic to obtain an advantage-like score for weighting
54 offline actions.

55 Our contributions are:

- 56 • We show that, under the standard dual goal representation parameterization, the goal em-
57 bedding can be interpreted as the gradient direction of the goal-conditioned value field with
58 respect to the learned state representation.
- 59 • We introduce *Dual Advantage Fields*, which learn action-effect vectors and score actions
60 by their alignment with this goal direction, producing a local advantage-like signal for
61 goal-conditioned policy extraction.
- 62 • We use this signal to extract policies from offline data without training a separate goal-
63 conditioned action-value function, and evaluate the resulting method across challenging
64 offline GCRL benchmarks.

65 2 Preliminaries

66 **Goal-conditioned Reinforcement Learning.** We study *offline* goal-conditioned reinforcement
67 learning (GCRL) [5, 14, 15, 18]: the learner has access to a fixed offline dataset of transitions

68 but cannot collect new experience in the environment [20]. The objective is to infer an optimal
69 goal-conditioned policy even for unseen during training combinations of state-goal pairs.

70 Let \mathcal{S} and \mathcal{A} denote state and action spaces, and let $\mathcal{G} \subseteq \mathcal{S}$ (or an abstract goal space) denote goals.
71 At each step the environment emits a transition (s, a, s') according to an unknown Markov kernel
72 $P(s' | s, a)$. A goal $g \in \mathcal{G}$ induces a reward signal $r(s, a, g)$: in sparse goal-reaching problems this is
73 often zero until a success condition holds. A stochastic policy $\pi(a | s, g)$ induces the usual discounted
74 return with discount $\gamma \in (0, 1)$. The goal-conditioned value and action-value functions are

$$Q^\pi(s, a, g) := \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, g) \mid s_0 = s, a_0 = a \right], \quad V^\pi(s, g) := \mathbb{E}_{a \sim \pi(\cdot | s, g)} [Q^\pi(s, a, g)].$$

$$75 \quad Q^\pi(s, a, g) = \mathbb{E}_{s' \sim P(\cdot | s, a)} [r(s, a, g) + \gamma V^\pi(s', g)], \quad V^\pi(s, g) = \mathbb{E}_{a \sim \pi(\cdot | s, g)} [Q^\pi(s, a, g)]. \quad (1)$$

76 Recent GCRL methods combine several ideas, including representation learning, quasimetric objec-
77 tives [11, 18, 25], and hierarchical horizon reduction [7, 17, 19] over value functions, Q -functions,
78 and actors. These design choices are often complementary, but existing methods still show domain-
79 specific strengths: hierarchical methods tend to excel in long-horizon locomotion, while quasimetric
80 representations often work well for manipulation. In contrast, DAF emphasizes local policy improve-
81 ment during training while retaining long-horizon reasoning, leading to more consistent performance
82 across both domains.

83 **Hierarchical Implicit Q-Learning (HIQL).** In GCRL, accurately estimating the value function
84 for distant goals is the main challenge in solving complex long-horizon tasks [19]. To address this
85 issue, HIQL [19] proposed a hierarchical policy structure that utilizes a value function learned with
86 IQL [13]. This hierarchical design enables the agent to produce effective actions even when value
87 estimates for distant goals are noisy or unreliable. More specifically, HIQL trains a goal-conditioned
88 state-value function V with the following loss:

$$\mathcal{L}(V) = \mathbb{E}_{(s, s') \sim \mathcal{D}, g \sim p(g)} [L_2^\tau(r(s, g) + \gamma \bar{V}(s', g) - V(s, g))], \quad (2)$$

89 where the expectile loss is defined as $L_2^\tau(u) = |\tau - \mathbf{1}(u < 0)|u|^2$, with $\tau > 0.5$, and \bar{V} denotes the
90 target V network.¹ Following prior works [3, 19, 25], we adopt the sparse reward $r(s, g) = -\mathbf{1}\{s \neq$
91 $g\}$. Under this reward, the optimal value $|V^*(s, g)|$ corresponds to the *discounted temporal distance*,
92 *i.e.*, a discounted measure of the minimum number of environment steps required to reach the goal g
93 from state s . HIQL separates policy extraction² into two levels: a high-level policy $\pi^h(s_{t+k} | s_t, g)$
94 generates a k -step subgoal to guide progress toward the goal, while a low-level policy $\pi^\ell(a_t | s_t, s_{t+k})$
95 produces primitive actions to reach the subgoal. Both policies are extracted using advantage-weighted
96 regression (AWR) [23, 26] with the following objective:

$$\mathcal{J}(\pi^h) = \mathbb{E}_{(s_t, s_{t+k}, g) \sim \mathcal{D}} [\exp(\beta^h \cdot A^h(s_t, s_{t+k}, g)) \log \pi^h(s_{t+k} | s_t, g)], \quad (3)$$

$$\mathcal{J}(\pi^\ell) = \mathbb{E}_{(s_t, a_t, s_{t+1}, s_{t+k}) \sim \mathcal{D}} [\exp(\beta^\ell \cdot A^\ell(s_t, s_{t+1}, s_{t+k})) \log \pi^\ell(a_t | s_t, s_{t+k})], \quad (4)$$

97 where β^h and β^ℓ are inverse temperature parameters, $A^h(s_t, s_{t+k}, g) = V^h(s_{t+k}, g) - V^h(s_t, g)$
98 denotes the high-level policy advantage, and $A^\ell(s_t, s_{t+1}, s_{t+k}) = V^\ell(s_{t+1}, s_{t+k}) - V^\ell(s_t, s_{t+k})$
99 denotes the low-level policy advantage. HIQL uses a single goal-conditioned value function V , which
100 is shared between both π^h and π^ℓ (*i.e.*, $V^h = V^\ell = V$). However, despite this design, HIQL still
101 struggles with long-horizon, complex tasks, as shown in the GCRL benchmark, OGBench [20].

102 **Dual Goal Representations [22].** In goal-conditioned RL, the goal representation determines
103 what information the policy and value function use about the target state. Rather than conditioning
104 directly on the raw goal observation, which may contain irrelevant or exogenous factors, dual goal
105 representations encode a goal by its reachability relation to other states. A goal g is represented by

$$\phi^V(g) : s \mapsto d^*(s, g),$$

106 where $d^*(s, g)$ denotes the optimal temporal distance from state s to goal g . In practice, we approxi-
107 mate this functional through a bilinear goal-conditioned potential [9]:

$$V_\theta(s, g) = \psi_\theta(s)^\top \phi_\theta(g), \quad (5)$$

108 where $\psi_\theta : \mathcal{S} \rightarrow \mathbb{R}^d$ and $\phi_\theta : \mathcal{G} \rightarrow \mathbb{R}^d$ are state and goal embeddings. The goal embedding $\phi_\theta(g)$
109 then serves as a finite-dimensional dual representation: when paired with $\psi_\theta(s)$, it predicts a value or
110 distance-like quantity that reflects the environment’s reachability structure.

¹Since the inherent over-estimation problem of IQL, we assume that the environment dynamics is deterministic.

²Policy extraction refers to learning a policy from a learned value function, emphasizing the separation between value learning and policy learning.

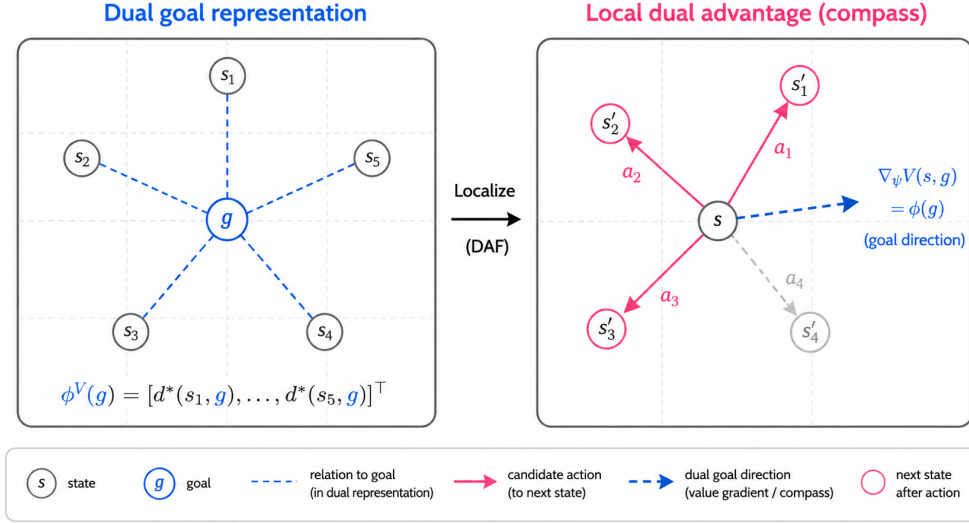


Figure 2: **Dual Advantage Fields.** Under a bilinear goal-conditioned value model, the goal embedding defines a direction in representation space. DAF scores an action by projecting its induced feature displacement onto this goal direction, yielding a local advantage-like signal for policy improvement.

111 3 Dual Advantage Fields

112 Our method is based on a simple insight from bilinear value decomposition in Eq. (5). Holding the
 113 goal fixed and viewing the value as a function of the state embedding ψ , we have

114 **Proposition 3.1.** *Under the bilinear goal-conditioned value model $V_\theta(s, g) = \psi_\theta(s)^\top \phi_\theta(g)$, the
 115 gradient of the value with respect to the state embedding is the goal embedding:*

$$\nabla_\psi V_\theta(s, g) = \nabla_\psi (\psi^\top \phi_\theta(g)) = \phi_\theta(g). \quad (6)$$

116 Thus, the goal embedding $\phi_\theta(g)$ is the value-gradient direction in representation space under the
 117 Euclidean geometry of the learned embedding. Please, see Figure 2 for intuition. For any transition
 118 from s to s' , the change in the bilinear value is exactly

$$V_\theta(s', g) - V_\theta(s, g) = \phi_\theta(g)^\top (\psi_\theta(s') - \psi_\theta(s)). \quad (7)$$

119 We use this identity to construct an advantage-like local policy improvement signal. For a policy π ,
 120 the standard goal-conditioned advantage is

$$A^\pi(s, a, g) = \mathbb{E}_{s' \sim p(\cdot | s, a)} [r(s, a, g) + \gamma V^\pi(s', g) - V^\pi(s, g)]. \quad (8)$$

121 Replacing V^π with the learned bilinear field V_θ gives the model-induced Bellman advantage

$$A_\theta(s, a, g) = \mathbb{E}_{s' \sim p(\cdot | s, a)} [r(s, a, g) + \phi_\theta(g)^\top (\gamma \psi_\theta(s') - \psi_\theta(s))]. \quad (9)$$

122 In offline learning, each dataset transition (s, a, s') provides a sample estimate of this quantity:

123 **Corollary 3.2.** *The sample-level Dual Advantage Field score is*

$$\boxed{\hat{A}_\theta(s, a, s', g) = r(s, a, g) + \phi_\theta(g)^\top (\gamma \psi_\theta(s') - \psi_\theta(s))}. \quad (10)$$

124 **Local policy improvement.** In the realizable case, the DAF score is exactly the goal-conditioned
 125 Bellman advantage. Specifically, if $V^\pi(s, g) = \psi(s)^\top \phi(g)$ and $u(s, a) = \mathbb{E}_{s' \sim P(\cdot | s, a)} [\gamma \psi(s') -$
 126 $\psi(s)]$, then

$$r(s, a, g) + u(s, a)^\top \phi(g) = A^\pi(s, a, g).$$

127 Thus, increasing the probability of actions (alignment) with positive DAF score is a standard policy-
 128 improvement step. Repeated exact DAF improvement therefore recovers an optimal primitive

129 goal-conditioned policy; in particular, its limiting policy is at least as good as any policy restricted to
 130 a fixed hierarchical class. We provide the formal statement and proof in Appendix ??.

131 Equation (10) defines the DAF score. The term $\gamma\psi_\theta(s') - \psi_\theta(s)$ is the discounted feature displacement
 132 caused by action a , and $\phi_\theta(g)$ is the value-gradient direction toward goal g . Their inner product
 133 measures the one-step increase in the bilinear value field, with the reward term completing the
 134 Bellman advantage. Thus, \hat{A}_θ provides a local, goal-conditioned action-ranking signal derived from
 135 the learned dual value geometry. This follows the comparative view of policy improvement, where
 136 actions are improved by relative advantages rather than absolute value estimates [4].

137 3.1 Motivational Example

138 We illustrate the local geometry captured by Dual Advantage Fields (DAF) on the
 139 cube-single-play-v0-task1 manipulation task from OGBench [20]. This task highlights a
 140 common failure mode in goal-conditioned control: before the cube can be placed at the final target,
 141 the agent must first move the gripper into a pre-grasp configuration. Thus, a direction that points
 142 directly toward the terminal object location may be globally plausible but locally unhelpful.

143 DAF addresses this by scoring actions according to their local improvement of the learned goal-
 144 conditioned potential. By Eq. (6), the goal embedding $\phi_\theta(g)$ is the representation-space gradient
 145 of the bilinear value field. We define an action-effect model $u_\xi(s, a)$ that estimates the discounted
 146 feature displacement induced by action a ,

$$u_\xi(s, a) \approx \mathbb{E}_{s' \sim p(\cdot | s, a)} [\gamma\psi_\theta(s') - \psi_\theta(s)].$$

147 Ignoring reward terms that are constant across actions in the pre-grasp region, DAF scores actions by

$$z_\theta(s, a, g) = u_\xi(s, a)^\top \phi_\theta(g). \quad (11)$$

148 This score favors actions whose predicted feature displacement is aligned with the local direction of
 149 value increase toward the goal.

150 Figure 3 visualizes this effect. We sample query
 151 states $\{\tilde{s}_i\}_{i=1}^N$ by perturbing only the gripper
 152 position around the cube, while keeping the ob-
 153 ject state and final goal fixed. For each method
 154 $m \in \{\text{OTA}, \text{DAF}\}$, we decode its high-level
 155 subgoal prediction into an X-Y coordinates via
 156 probing,

$$\hat{x}_i^m = D_m(h_m(\tilde{s}_i, g)), \quad (12)$$

157 where h_m is the method-specific latent output
 158 and D_m is a linear probe fitted on demonstration
 159 states. The plotted direction is

$$d_i^m = \frac{\hat{x}_i^m - x_{ee}(\tilde{s}_i)}{\|\hat{x}_i^m - x_{ee}(\tilde{s}_i)\|_2}, \quad (13)$$

160 drawn from the gripper position $x_{ee}(\tilde{s}_i)$. Near
 161 the cube, DAF produces directions that point
 162 toward the object, matching the immediate pre-
 163 grasp behavior required before transport. OTA
 164 instead points toward the terminal placement region in this example, which is appropriate only after
 165 grasping. The example shows why local advantage fields can be more useful than a high-level
 166 subgoals alone: they select actions by whether they locally improve the goal-conditioned potential.

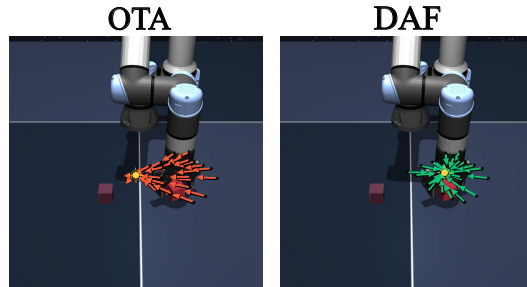


Figure 3: **Pre-grasp vector field in cube-single.** Arrows show decoded high-level directions from sampled gripper positions around the cube, with the cube and final goal fixed. DAF points locally toward the cube before grasping, while OTA points toward the terminal placement goal. The yellow marker denotes the mean decoded target.

167 4 Training and goal-conditioned policy extraction

168 Dayan and Singh [4] showed that policy improvement can be organized around *relative* measures of
 169 how actions compare at a state-*merits* that need not reduce to a fully trusted global value oracle. In
 170 the goal-conditioned setting, the Bellman advantage $A^\pi(s, a, g)$ in (8) is exactly such an object: it
 171 ranks a by the expected one-step gain in return, isolating the effect of the transition from the baseline

Algorithm 1 DAF training.

- 1: **Input:** offline dataset \mathcal{D} of (s, a, s', g) ;
 - 2: **Initialize:** ψ_θ, ϕ_θ , displacement map u_ξ , policy π_ω , target networks $(Q^{\text{tgt}}, V^{\text{tgt}})$.
 - 3: **while** not converged **do**
 - 4: Sample a minibatch from \mathcal{D} .
 - 5: **Critic:** update ψ_θ, ϕ_θ so $V_\theta(s, g) = \psi_\theta(s)^\top \phi_\theta(g)$ (5) using target networks.
 - 6: **AFU coupling:** minimize the actor-free loss coupling V_θ to z_θ (14) {Appendix E} and minimize \mathcal{L}_{ae} (15) for u_ξ .
 - 7: **Policy:** $w_\theta \leftarrow \min\{\exp(\alpha z_\theta(s, a, g)), W_{\max}\}$; minimize $-\mathbb{E}_{\mathcal{D}}[w_\theta \log \pi_\omega(a | s, c)]$ over ω .
 - 8: Update target networks.
 - 9: **end while**
-

172 $V^\pi(s, g)$. Our bilinear potential (5) turns this comparison into explicit geometry in ψ . Under the
173 model V_θ , the backup contribution $\gamma V_\theta(s', g) - V_\theta(s, g)$ equals $\phi_\theta(g)^\top (\gamma \psi_\theta(s') - \psi_\theta(s))$ by (7),
174 so the analogue of the advantage (8) with V^π replaced by V_θ is the closed form (9)-(10). The goal
175 embedding $\phi_\theta(g)$ acts as $\nabla_\psi V_\theta$ (Eq. (6)): the inner product in (10) measures whether the *local* feature
176 displacement induced by a is aligned with steepest increase of the learned potential toward g . Thus
177 Dayan’s comparative view of improvement is instantiated here as projection of one-step ψ -dynamics
178 onto the value-gradient direction.

179 In practice we estimate the discounted increment $\gamma \psi_\theta(s') - \psi_\theta(s)$ with a map $u_\xi(s, a)$ trained on
180 offline transitions (Sec. 4.1), and absorb r in the critic stack where noted. The raw dual score is

$$z_\theta(s, a, g) := u_\xi(s, a)^\top \phi_\theta(g), \quad (14)$$

181 which agrees with (10) when $u_\xi(s, a) \approx \gamma \psi_\theta(s') - \psi_\theta(s)$ and rewards are handled by the value
182 heads feeding the same Bellman targets.

183 4.1 Offline critic and feature dynamics

184 We learn $(\psi_\theta, \phi_\theta)$, and the displacement map u_ξ from offline tuples (s, a, s', g) [20]. For stability we
185 used a common approach in offline RL [13] that learns twin critics $Q_\theta^{(1)}, Q_\theta^{(2)}$, and the bilinear head
186 $V_\theta(s, g) = \psi_\theta(s)^\top \phi_\theta(g)$ is tied to pessimistic Q -estimates via expectile regression and to Bellman
187 backups on $Q_\theta^{(j)}$. To avoid brittle $\max_a Q$ operators in continuous control [16], we add an *actor-free*
188 coupling between V_θ and the scalar dual score z_θ from (14), following Perrin-Gilbert [24]; the explicit
189 construction is deferred to Appendix E. Finally, we ground u_ξ with the auxiliary loss

$$\mathcal{L}_{\text{ae}} = \mathbb{E} \left[\|u_\xi(s, a) - \text{sg}(\gamma \psi_\theta(s') - \psi_\theta(s))\|_2^2 \right], \quad (15)$$

190 with sg stopping gradients through the target, so u_ξ tracks one-step feature dynamics on \mathcal{D} .

191 4.2 Policy extraction

192 Let $\pi_\omega(a | s, c)$ denote the policy with conditioning c on g through $\phi_\theta(g)$ (and optionally s).
193 Advantage-weighted regression [23] uses weights

$$w_\theta(s, a, g) = \min \left\{ \exp(\alpha z_\theta(s, a, g)), W_{\max} \right\} \quad (16)$$

194 with temperature $\alpha > 0$ and cap W_{\max} , and minimizes $-\mathbb{E}_{\mathcal{D}}[w_\theta \log \pi_\omega(a | s, c)]$. Because z_θ
195 does not depend on ω , this is weighted behavior cloning that up-weights actions whose local ψ -
196 displacement aligns with the goal direction $\phi_\theta(g)$, i.e. actions that the bilinear model classifies as
197 improving the goal-conditioned potential in the sense of (10).

198 **Hierarchical goals.** For long horizons, a high-level policy over subgoals can be trained alongside
199 the low-level stack above, with value differences along options as in hierarchical offline GCRL [19];
200 option-aware temporally abstracted value learning offers a related hierarchical baseline [2].

201 *Remark 4.1* (Variants). Concrete instantiations differ by which of the optional terms above are active;
202 experimental details are summarized in Section 5 and Appendix A.

203 **5 Experiments**

204 In this section, we empirically validate the findings developed in the previous sections on the
 205 OGBench benchmark [20]. OGBench is designed to evaluate several core capabilities required by
 206 offline goal-conditioned reinforcement learning (GCRL), including long-horizon reasoning, trajectory
 207 stitching, generalization to unseen goals, robustness to suboptimal data, and control under imperfect
 208 offline coverage. We focus on the state-based locomotion and manipulation tasks used in prior
 209 work, which allows us to test whether DAF provides consistent improvements across domains with
 210 substantially different control structure.

211 All methods are trained purely offline on the provided datasets and are evaluated without additional
 212 environment interaction during training. We report success-based performance in $[0, 1]$, where higher
 213 values indicate better goal reaching. For each environment, we evaluate the corresponding OGBench
 214 dataset regimes. In maze-style locomotion, we use `navigate` and `stitch` datasets: `navigate` data is
 215 collected from noisy expert policies that traverse the environment, while `stitch` data contains shorter
 216 trajectory segments and therefore requires composing partial behaviors into longer goal-reaching
 217 solutions. In manipulation, we use `play` and `noisy` datasets: `play` data contains natural temporally
 218 correlated interactions generated by scripted policies, whereas `noisy` data increases state-action
 219 coverage through less structured exploration noise, making the offline data more suboptimal.

220 **Baselines.** We compare against a representative set of recent and relevant methods for offline
 221 GCRL, including HIQL [19], OTA [2], MQE [17], CRL [5], GCIQL [13], and GCIVL [10]. When
 222 applicable, we also include their corresponding variants that learn representations in the form of
 223 dual-goal representations [22]. These baselines cover the main families of methods used in offline
 224 GCRL, including horizon-reduction methods [21] and methods based on representation priors such
 225 as quasimetrics.

226 **What DAF does in each dataset.** Across all datasets, DAF uses the same policy-extraction
 227 principle: it scores offline actions by the alignment between their predicted local feature displacement
 228 and the goal direction induced by the dual value representation. Concretely, the action-effect model
 229 estimates $\gamma\psi_\theta(s') - \psi_\theta(s)$, and the dual score projects this displacement onto $\phi_\theta(g)$. Thus, DAF
 230 uses the learned value field not only as a global map of reachability, but also as a local compass for
 231 choosing among actions available in the offline dataset.

232 **Maze locomotion: long-horizon navigation and stitching.** We first evaluate on
 233 `humanoidmaze` and `antmaze`, shown in Table 1. These environments isolate the long-
 234 horizon navigation aspect of offline GCRL. The agent must reach target states from diverse initial
 235 states using only fixed offline data. The `antmaze` tasks require quadruped locomotion
 236 through maze layouts, while `humanoidmaze` is more challenging because it combines full-body
 237 humanoid control with long-horizon goal reaching. We include both `navigate` and `stitch`
 238 variants because they test complementary capabilities: `navigate` evaluates whether the
 239 method can exploit noisy expert trajectories, while `stitch` evaluates whether the method can
 240 compose shorter trajectory fragments into successful long-horizon behavior.

241 These tasks are important because many prior offline GCRL methods are designed around horizon
 242 reduction or hierarchical subgoal prediction. DAF is not primarily a hierarchical method: instead, it
 243 extracts local action preferences from a dual value field. Strong performance on these mazes therefore
 244 tests whether local advantage-field extraction can preserve the long-horizon structure needed for
 245 navigation. DAF is competitive with the strongest horizon-reduction baselines on `navigate` datasets
 246 and obtains the best results on the harder `stitch` cases where composing partial trajectories is
 247 essential.

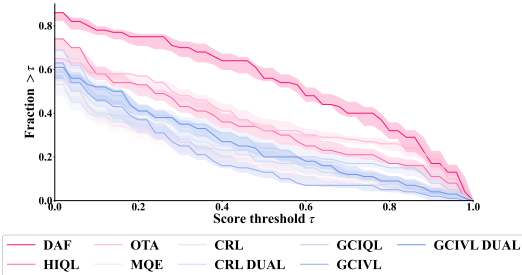


Figure 4: Performance profile across all tasks and environments. DAF achieves a better distribution of scores than the baselines across the OGBench evaluation suite.

Table 1: Maze locomotion results on humanoidmaze and antmaze. These tasks test long-horizon goal reaching from fixed offline data. The navigate datasets evaluate learning from noisy expert trajectories, while the stitch datasets evaluate whether a method can compose shorter trajectory segments into successful goal-reaching behavior. Best values are highlighted in blue.

ENV.	DATASET	DIMENSION	DAF	HIQL	OTA	MQE	CRL	CRL DUAL	GCIQL	GCIVL	GCIVL DUAL
HUMANOIDMAZE	NAVIGATE	MEDIUM	0.93 ± 0.03	0.91 ± 0.01	0.95 ± 0.01	0.49 ± 0.09	0.59 ± 0.03	0.62 ± 0.03	0.31 ± 0.04	0.31 ± 0.03	0.32 ± 0.03
		LARGE	0.66 ± 0.03	0.45 ± 0.04	0.83 ± 0.03	0.20 ± 0.07	0.26 ± 0.03	0.21 ± 0.05	0.04 ± 0.01	0.05 ± 0.01	0.04 ± 0.01
	STITCH	MEDIUM	0.90 ± 0.04	0.86 ± 0.03	0.92 ± 0.01	0.62 ± 0.09	0.53 ± 0.03	0.57 ± 0.01	0.15 ± 0.03	0.14 ± 0.01	0.20 ± 0.02
		LARGE	0.48 ± 0.06	0.32 ± 0.04	0.43 ± 0.04	0.18 ± 0.03	0.11 ± 0.02	0.06 ± 0.03	0.02 ± 0.00	0.02 ± 0.01	0.02 ± 0.01
ANTMAZE	NAVIGATE	TELEPORT	0.51 ± 0.08	0.46 ± 0.03	0.53 ± 0.03	0.49 ± 0.04	0.60 ± 0.01	0.57 ± 0.04	0.37 ± 0.02	0.44 ± 0.02	0.41 ± 0.02
		MEDIUM	0.98 ± 0.01	0.96 ± 0.01	0.97 ± 0.01	0.88 ± 0.05	0.96 ± 0.01	0.95 ± 0.02	0.76 ± 0.06	0.72 ± 0.06	0.79 ± 0.03
	STITCH	TELEPORT	0.50 ± 0.05	0.38 ± 0.02	0.38 ± 0.01	0.40 ± 0.03	0.23 ± 0.03	0.15 ± 0.04	0.23 ± 0.02	0.45 ± 0.05	0.39 ± 0.05
		MEDIUM	0.97 ± 0.02	0.96 ± 0.02	0.95 ± 0.02	0.96 ± 0.01	0.52 ± 0.04	0.52 ± 0.07	0.39 ± 0.06	0.48 ± 0.03	0.52 ± 0.03

Table 2: Object-manipulation results on cube and scene. These datasets test whether an offline GCRL method can extract precise local skills from play data and remain robust under the less structured noisy regime. Best values are highlighted in blue.

ENV.	DATASET	DIMENSION	DAF	HIQL	OTA	MQE	CRL	CRL DUAL	GCIQL	GCIVL	GCIVL DUAL
CUBE	PLAY	DOUBLE	0.41 ± 0.04	0.13 ± 0.01	0.05 ± 0.01	0.03 ± 0.00	0.16 ± 0.01	0.38 ± 0.06	0.35 ± 0.06	0.33 ± 0.05	0.58 ± 0.04
		TRIPLE	0.17 ± 0.03	0.05 ± 0.02	0.02 ± 0.00	0.01 ± 0.00	0.06 ± 0.02	0.05 ± 0.05	0.02 ± 0.01	0.01 ± 0.01	0.01 ± 0.00
		QUADRUPLE	0.03 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
	NOISY	DOUBLE	0.33 ± 0.05	0.03 ± 0.01	0.05 ± 0.03	0.07 ± 0.01	0.04 ± 0.02	0.08 ± 0.02	0.24 ± 0.06	0.17 ± 0.04	0.26 ± 0.02
		TRIPLE	0.23 ± 0.01	0.04 ± 0.01	0.01 ± 0.00	0.04 ± 0.02	0.03 ± 0.01	0.06 ± 0.02	0.05 ± 0.01	0.11 ± 0.02	0.09 ± 0.03
		QUADRUPLE	0.02 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.01 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
SCENE	PLAY	0.81 ± 0.04	0.55 ± 0.09	0.34 ± 0.04	0.20 ± 0.03	0.29 ± 0.02	0.56 ± 0.06	0.53 ± 0.02	0.51 ± 0.05	0.78 ± 0.07	
	NOISY	0.43 ± 0.03	0.27 ± 0.02	0.10 ± 0.02	0.07 ± 0.02	0.02 ± 0.01	0.06 ± 0.01	0.29 ± 0.02	0.31 ± 0.05	0.45 ± 0.02	

256 **Object manipulation: local control from imperfect demonstrations.** Next, we evaluate on cube
 257 and scene, shown in Table 2. Unlike maze navigation, these tasks require precise object-centric
 258 control. The cube tasks include pick-and-place, stacking, swapping, and multi-object rearrangement,
 259 while scene tasks require sequencing interactions with objects such as cubes, drawers, windows, and
 260 buttons.

261 These datasets test DAF’s central motivation: globally plausible behavior can be locally wrong. For
 262 example, moving toward a final object placement may be inappropriate before reaching a pre-grasp
 263 state. DAF addresses this by ranking dataset actions according to whether their predicted feature
 264 displacement aligns with the goal direction. This is especially useful in play and noisy datasets,
 265 where demonstrations contain useful local skills but also incomplete or suboptimal trajectories.

266 **Puzzle rearrangement: continuous control with combinatorial structure.** Finally, we evaluate
 267 on puzzle, shown in Table 3. These environments are robotic versions of Lights Out: pressing
 268 one button changes the state of neighboring buttons. They therefore combine continuous control
 269 with combinatorial generalization over discrete configurations. The 3x3 and 4x4 variants further
 270 increase the configuration space, testing whether goal representations generalize beyond simple object
 271 reaching.

272 Puzzle tasks stress a failure mode not captured by maze navigation or standard manipulation. Here,
 273 each local action can affect a larger configuration, so policy extraction must compare actions by their
 274 downstream effect on the goal. DAF is suited to this setting because it scores actions by whether their
 275 predicted local transition improves the goal-conditioned value field.

276 **Aggregate comparison.** Tables 1, 2, and 3 show that DAF performs strongly across different kinds
 277 of offline coverage and control structure. The aggregate comparison in Figure 5 further summarizes
 278 performance across all tasks using RLiable metrics [1]. We report Median, interquartile mean
 279 (IQM), Mean, and Optimality Gap with stratified-bootstrap confidence intervals. The IQM reduces
 280 sensitivity to outlier tasks, while the optimality gap measures the average remaining shortfall from
 281 perfect success. Overall, DAF improves the aggregate metrics while also achieving strong per-task
 282 performance, indicating that the gains are not driven by a single environment family.

Table 3: Puzzle rearrangement results on `puzzle`. These tasks test structured spatial reasoning: each local button press changes neighboring button states, so the policy must combine continuous control with combinatorial goal generalization. Best values are highlighted in blue.

ENV.	DATASET	DIMENSION	DAF	HIQL	OTA	MQE	CRL	CRL DUAL	GCIQL	GCIVL	GCIVL DUAL
PUZZLE	PLAY	3X3	0.74 \pm 0.04	0.17 \pm 0.03	0.64 \pm 0.06	0.11 \pm 0.00	0.07 \pm 0.01	0.09 \pm 0.03	0.98 \pm 0.02	0.08 \pm 0.02	0.08 \pm 0.01
		4X4	0.40 \pm 0.05	0.17 \pm 0.03	0.53 \pm 0.05	0.17 \pm 0.03	0.02 \pm 0.01	0.07 \pm 0.02	0.31 \pm 0.02	0.26 \pm 0.02	0.30 \pm 0.05
PUZZLE	NOISY	3X3	0.98 \pm 0.04	0.70 \pm 0.10	0.64 \pm 0.13	0.03 \pm 0.01	0.37 \pm 0.05	0.42 \pm 0.05	0.95 \pm 0.01	0.44 \pm 0.15	0.50 \pm 0.22
		4X4	0.47 \pm 0.03	0.31 \pm 0.07	0.01 \pm 0.00	0.02 \pm 0.01	0.00 \pm 0.00	0.00 \pm 0.00	0.33 \pm 0.11	0.24 \pm 0.02	0.25 \pm 0.02

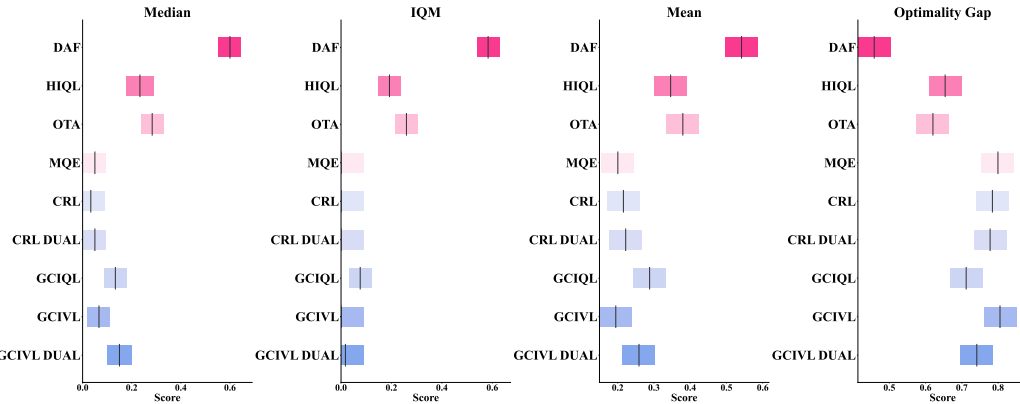


Figure 5: **Performance comparison.** Following the protocol proposed by Agarwal et al. [1], we report aggregate RLiable metrics, including Median, IQM, Mean, and Optimality Gap, with stratified-bootstrap confidence intervals across the offline GCRL environments. The colored horizontal segments denote confidence intervals, and the dark vertical markers denote point estimates.

283 6 Broader Impact and Limitations

284 DAF extracts goal-conditioned policies from offline data without additional environment interaction.
 285 Like other offline RL methods, it is reliable only when the dataset sufficiently covers the actions
 286 needed for improvement; poor coverage can produce incorrect action rankings. DAF also relies on the
 287 learned dual representation and action-effect model. Although $\nabla_{\psi} V_{\theta}(s, g) = \phi_{\theta}(g)$ holds exactly for
 288 the bilinear head, this direction is useful only if the representation encodes reachability. In stochastic
 289 or poorly covered regions, $u_{\xi}(s, a)$ may predict inaccurate feature displacements. Future work should
 290 study uncertainty-aware or distributional action-effect models and extend DAF to image-based and
 291 more stochastic goal-reaching settings.

292 7 Conclusion

293 We introduced Dual Advantage Fields (DAF), a method that turns dual goal representations into
 294 local policy-improvement signals. Under the bilinear value parameterization, the goal embedding
 295 acts as the gradient of the goal-conditioned value field with respect to the state representation. DAF
 296 uses this observation to score actions by the alignment between their predicted feature displacement
 297 and the goal direction. Empirically, DAF improves aggregate performance across offline GCRL
 298 benchmarks and is especially effective in manipulation tasks where local directional choices are
 299 important. Overall, the results suggest that dual representations should be used not only as global
 300 value maps, but also as local advantage fields for goal-conditioned policy extraction.

References

- 301
- 302 [1] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C. Courville, and Marc G.
303 Bellemare. Deep reinforcement learning at the edge of the statistical precipice. In *Advances in*
304 *Neural Information Processing Systems (NeurIPS)*, 2021.
- 305 [2] Hongjoon Ahn, Heewoong Choi, Jisu Han, and Taesup Moon. Option-aware temporally
306 abstracted value for offline goal-conditioned reinforcement learning, 2025. URL <https://arxiv.org/abs/2505.12737>.
307
- 308 [3] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,
309 Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience
310 replay. *Advances in neural information processing systems*, 30, 2017.
- 311 [4] Peter Dayan and Satinder P. Singh. Improving policies without measuring merits. In Gerald
312 Tesauro, David Touretzky, and Todd Leen, editors, *Advances in Neural Information Processing*
313 *Systems*, volume 8. MIT Press, 1995. URL [https://proceedings.neurips.cc/paper/](https://proceedings.neurips.cc/paper/1995/hash/208e43f0e45c4c78cafadb83d2888cb6-Abstract.html)
314 [1995/hash/208e43f0e45c4c78cafadb83d2888cb6-Abstract.html](https://proceedings.neurips.cc/paper/1995/hash/208e43f0e45c4c78cafadb83d2888cb6-Abstract.html).
- 315 [5] Benjamin Eysenbach, Ruslan Salakhutdinov, and Sergey Levine. Contrastive learning as goal-
316 conditioned reinforcement learning. In *Advances in Neural Information Processing Systems*,
317 volume 35, 2022.
- 318 [6] Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement learning from passive
319 data via latent intentions. In *Proceedings of the 40th International Conference on Machine*
320 *Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 11321–11339.
321 PMLR, 2023.
- 322 [7] Vittorio Giammarino and Ahmed H Qureshi. Goal reaching with eikonal-constrained hier-
323 archical quasimetric reinforcement learning. In *The Fourteenth International Conference on*
324 *Learning Representations*, 2026. URL <https://openreview.net/forum?id=5WhsCB0Vty>.
- 325 [8] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*
326 *arXiv:1606.08415*, 2016.
- 327 [9] Zhang-Wei Hong, Ge Yang, and Pulkit Agrawal. Bilinear value networks. In *International*
328 *Conference on Learning Representations*, 2022. arXiv:2204.13695.
- 329 [10] Kaiqiang Ke, Qian Lin, Zongkai Liu, Shenghong He, and Chao Yu. Conservative offline
330 goal-conditioned implicit v-learning. In *Forty-second International Conference on Machine*
331 *Learning*, 2025. URL <https://openreview.net/forum?id=5ryn8tYWHL>.
- 332 [11] Kaiqiang Ke, Zhonghai Ruan, Shengwen Tan, and Weixia Wu. Hierarchical quasimetric
333 reinforcement learning. In *Proceedings of the 2025 International Conference on Machine*
334 *Learning and Neural Networks*, pages 34–41, 2025.
- 335 [12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Interna-*
336 *tional Conference on Learning Representations (ICLR)*, 2015.
- 337 [13] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
338 Q-learning. In *International Conference on Learning Representations*, 2022.
- 339 [14] Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning:
340 Problems and solutions. *arXiv preprint arXiv:2201.08299*, 2022.
- 341 [15] Jason Yecheng Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline goal-conditioned
342 reinforcement learning via f -advantage regression. *Advances in neural information processing*
343 *systems*, 35:310–323, 2022.
- 344 [16] Gabriel Matheron, Nicolas Perrin, and Olivier Sigaud. Understanding the impact of the max
345 operation in value-based deep reinforcement learning. In *Advances in Neural Information*
346 *Processing Systems*, volume 33, 2020.

- 347 [17] Vivek Myers, Chongyi Zheng, Anca Dragan, Sergey Levine, and Benjamin Eysenbach. Learning
348 temporal distances: Contrastive successor features can provide a metric structure for decision-
349 making. *arXiv preprint arXiv:2406.17098*, 2024.
- 350 [18] Vivek Myers, Bill Chunyuan Zheng, Benjamin Eysenbach, and Sergey Levine. Offline
351 goal-conditioned reinforcement learning with quasimetric representations. *arXiv preprint*
352 *arXiv:2509.20478*, 2025.
- 353 [19] Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. HIQL: Offline goal-
354 conditioned RL with latent states as actions. In *Advances in Neural Information Processing*
355 *Systems*, 2023. arXiv:2307.11949.
- 356 [20] Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. OGBench: Benchmark-
357 ing offline goal-conditioned RL. In *International Conference on Learning Representations*,
358 2025. arXiv:2410.20092.
- 359 [21] Seohong Park, Kevin Frans, Deepinder Mann, Benjamin Eysenbach, Aviral Kumar, and Sergey
360 Levine. Horizon reduction makes rl scalable. *arXiv preprint arXiv:2506.04168*, 2025.
- 361 [22] Seohong Park, Deepinder Mann, and Sergey Levine. Dual goal representations. *arXiv preprint*
362 *arXiv:2510.06714*, 2025.
- 363 [23] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:
364 Simple and scalable off-policy reinforcement learning, 2019. URL [https://arxiv.org/](https://arxiv.org/abs/1910.00177)
365 [abs/1910.00177](https://arxiv.org/abs/1910.00177).
- 366 [24] Nicolas Perrin-Gilbert. AFU: Actor-free critic updates in off-policy RL for continuous control,
367 2024. URL <https://arxiv.org/abs/2404.16159>.
- 368 [25] Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal goal-reaching
369 reinforcement learning via quasimetric learning. In *International Conference on Machine*
370 *Learning*, pages 36411–36430. PMLR, 2023.
- 371 [26] Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E
372 Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized
373 regression. *Advances in Neural Information Processing Systems*, 33:7768–7778, 2020.

Table 4: Network configuration for DAF on OGBench.

Configuration	Value
Gradient steps	10^6
Optimizer	Adam [12]
Nonlinearity	GELU [8]
Target network update rate	0.005
Goal representation dimensionality	256
Batch size	1024
Action-effect MLP dimensions	(512, 512, 512)
Policy MLP dimensions	(512, 512, 512)
Layer norm in MLPs	True
Discount (γ)	0.99 (0.995 for giant- environments)
Learning rate	0.0003

Table 5: Coefficient α for each environment

Environment	α
scene-*	10.0
antmaze	10.0
humanoidmaze	10.0
puzzle-3x3-*	3.0
puzzle-4x4-*	0.1
cube-double-*	3.0
cube-triple-*	3.0
cube-quadruple-*	10.0

374 A Implementation and Reproducibility

375 Our method and baselines are implemented on top of the implementations given in OGBench [20]
 376 and Dual Goal Representations [22] codebases. Our method is employed upon hierarchy of actors,
 377 with low actor being updated by dual score (Equation (14)) and high actor by AWR (Equation (3)).

378 Table 4 details the common hyperparameters for all methods on OGBench. Table 5 shows the α
 379 regularization hyperparameter that was found to be the best for performance of DAF. We also report
 380 the ablation studies on important architectural aspects of our proposed method: AFU Coupling
 381 [24], presence of action-effect module (Equation (15)), hierarchical actor [19] and integrating dual
 382 representations [22] upon the hierarchical backbone.

383 B Related Works

384 B.1 Goal-conditioned Implicit Q-Learning (GCIQL)

385 Implicit Q-Learning (IQL) [13] stabilizes offline RL by avoiding queries to out-of-distribution (OOD)
 386 actions through two key components: a state-value function $V_\psi(s)$ and an action-value function
 387 $Q_\theta(s, a)$. The value functions are trained via:

$$\mathcal{L}_Q(\theta) = \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[(r(s, a) + \gamma V_\psi(s') - Q_\theta(s, a))^2 \right], \quad (17)$$

$$\mathcal{L}_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^\tau(Q_{\bar{\theta}}(s, a) - V_\psi(s))], \quad (18)$$

388 where $L_2^\tau(x) = |\tau - \mathbf{1}(x < 0)|x^2$ and $\tau \in [0.5, 1)$ controls conservatism (higher τ prioritizes
 389 optimistic returns), and θ are the parameters of the target Q network. The policy $\pi_\phi(a|s)$ is then
 390 extracted via advantage-weighted regression (AWR) [23, 26]:

Table 6: Ablation on OGBench. We report full DAF and four requested ablations: removing AFU coupling, removing the action-effect model (using direct one-step value-difference scoring), removing hierarchy, and using a dual-representation hierarchical baseline. Best values are highlighted in blue.

ENV.	DATASET	DIMENSION	DAF	NO AFU COUPLING	NO ACTION-EFFECT	DAF W/O HIERARCHY	DUAL-REP BASELINE + HIERARCHY
HUMANOIDMAZE	NAVIGATE	MEDIUM	0.93 ± 0.03	0.18 ± 0.03	0.35 ± 0.03	0.38 ± 0.04	0.06 ± 0.01
		LARGE	0.66 ± 0.03	0.02 ± 0.01	0.04 ± 0.01	0.39 ± 0.05	0.01 ± 0.00
	STITCH	MEDIUM	0.90 ± 0.04	0.47 ± 0.06	0.50 ± 0.07	0.32 ± 0.06	0.05 ± 0.02
		LARGE	0.48 ± 0.06	0.06 ± 0.02	0.08 ± 0.01	0.65 ± 0.07	0.01 ± 0.01
ANTMAZE	NAVIGATE	TELEPORT	0.51 ± 0.08	0.33 ± 0.05	0.35 ± 0.05	0.39 ± 0.05	0.18 ± 0.07
		MEDIUM	0.98 ± 0.01	0.78 ± 0.05	0.93 ± 0.05	0.19 ± 0.04	0.89 ± 0.02
	STITCH	TELEPORT	0.50 ± 0.05	0.19 ± 0.04	0.17 ± 0.04	0.00 ± 0.01	0.09 ± 0.02
		MEDIUM	0.97 ± 0.02	0.42 ± 0.05	0.42 ± 0.10	0.00 ± 0.01	0.21 ± 0.06
CUBE	PLAY	DOUBLE	0.41 ± 0.04	0.36 ± 0.07	0.51 ± 0.05	0.39 ± 0.05	0.02 ± 0.01
		TRIPLE	0.17 ± 0.03	0.02 ± 0.01	0.04 ± 0.02	0.07 ± 0.02	0.01 ± 0.01
		QUADRUPLE	0.03 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	0.01 ± 0.01	0.00 ± 0.00
	NOISY	DOUBLE	0.33 ± 0.05	0.26 ± 0.05	0.39 ± 0.03	0.35 ± 0.04	0.03 ± 0.01
		TRIPLE	0.23 ± 0.01	0.01 ± 0.01	0.05 ± 0.02	0.02 ± 0.01	0.01 ± 0.00
		QUADRUPLE	0.02 ± 0.01	0.00 ± 0.00	0.00 ± 0.00	0.01 ± 0.00	0.00 ± 0.00
SCENE	PLAY	-	0.81 ± 0.04	0.49 ± 0.05	0.52 ± 0.01	0.45 ± 0.04	0.20 ± 0.07
	NOISY	-	0.43 ± 0.03	0.29 ± 0.03	0.40 ± 0.04	0.37 ± 0.05	0.05 ± 0.02
PUZZLE	PLAY	3X3	0.74 ± 0.04	0.10 ± 0.01	0.15 ± 0.01	0.02 ± 0.01	0.07 ± 0.02
		4X4	0.40 ± 0.05	0.18 ± 0.03	0.17 ± 0.03	0.75 ± 0.06	0.01 ± 0.00
	NOISY	3X3	0.98 ± 0.04	0.15 ± 0.03	0.17 ± 0.02	0.37 ± 0.06	0.05 ± 0.01
		4X4	0.47 ± 0.03	0.07 ± 0.01	0.09 ± 0.01	0.03 ± 0.01	0.00 ± 0.01

$$J_{\pi}(\phi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [\exp(\beta \cdot A(s,a)) \log \pi_{\phi}(a|s)], \quad (19)$$

391 with $A(s,a) = Q_{\theta}(s,a) - V_{\psi}(s)$, and β is the inverse temperature parameter.

392 For goal-conditioned RL, IQL is extended to learn a goal-conditioned state-value function $V_{\psi}(s,g)$,
 393 preserving IQL’s key advantage of stable value learning without requiring explicit Q-function evalua-
 394 tions on out-of-distribution actions [6].

395 B.2 Option-aware Temporally Abstracted Value (OTA)

396 HIQL [19] addresses long horizons by introducing a hierarchy over subgoals, but still relies on flat
 397 temporal-difference updates to a high-level value. OTA [2] instead bakes temporal abstraction directly
 398 into the Bellman operator by learning *option-aware* values: for an option o that lasts $k(o)$ steps, the
 399 high-level Bellman target becomes

$$V(s,g) \approx \mathbb{E}[r^{(o)}(s,g) + \gamma^{k(o)} V(s',g)], \quad (20)$$

400 where $r^{(o)}$ is the cumulative option reward and s' is the option-termination state. Each update
 401 contracts the effective horizon from $d^*(s,g)$ to roughly $d^*(s,g)/k(o)$, so value differences and the
 402 corresponding high-level advantages are computed over multi-step options rather than single primitive
 403 actions. This leads to more stable high-level signals and better long-horizon stitching on OGBench,
 404 at the cost of committing to a particular temporal abstraction schedule.

405 B.3 Quasimetric representations and MQE-style methods

406 Recent work views goal-conditioned value learning as estimating an asymmetric “distance” $d(s,g)$
 407 between states and goals. Quasimetric approaches [18, 25] directly fit such distances with multistep
 408 returns: instead of bootstrapping only from immediate successors, they regress

$$d_{\theta}(s,g) \approx \mathbb{E}\left[\sum_{t=0}^{K-1} c(s_t,g) + d_{\theta}(s_K,g) \mid s_0 = s\right] \quad (21)$$

409 for random horizons K , while encouraging triangle-like inequalities $d_{\theta}(s,g) \leq d_{\theta}(s,\tilde{g}) + d_{\theta}(\tilde{g},g)$
 410 for sampled pivots \tilde{g} . This multistep quasimetric estimation (MQE) improves horizon generalization—
 411 including long-horizon stitching in visual domains—but typically requires stronger structural assump-
 412 tions on the value landscape than local TD methods and can be sensitive to misspecification of the
 413 quasimetric prior.

414 B.4 Conservative goal-conditioned implicit V-learning (GCIVL)

415 GCIVL-style methods extend IQL to goal-conditioned settings but can overestimate values for *uncon-*
416 *nected* state-goal pairs produced by cross-trajectory pairing. GCIVL [10] introduces conservative
417 penalties on such pairs together with a quasimetric formulation. Concretely, for a learned value or
418 distance $v_\theta(s, g)$ and a connectivity indicator $c(s, g) \in \{0, 1\}$ (reachable from \mathcal{D}), the GCIVL loss
419 augments Bellman terms with

$$\mathcal{L}_{\text{cons}}(\theta) = \lambda \mathbb{E}_{(s,g) \sim p_{\text{pair}}} \left[(1 - c(s, g)) (\max\{0, v_\theta(s, g) - \delta\})^2 \right], \quad (22)$$

420 penalizing large estimates on likely-unreachable pairs. This improves robustness on goal-stitching
421 tasks in OGBench, but depends on correctly identifying or regularizing unreachable pairs and still
422 operates on scalar values rather than local action-effect structure.

423 B.5 Contrastive representation learning (CRL)

424 Contrastive RL methods treat goal-conditioned control as a representation learning problem: they
425 learn embeddings so that inner products between state(-action) and goal features approximate a
426 goal-conditioned value or reachability score [5]. A typical loss takes the form

$$\mathcal{L}_{\text{CRL}} = -\mathbb{E} \left[\log \frac{\exp(\phi(s, a)^\top \psi(g^+)/\tau)}{\sum_{g' \in \mathcal{N}} \exp(\phi(s, a)^\top \psi(g')/\tau)} \right], \quad (23)$$

427 where (s, a, g^+) is a positive triple and \mathcal{N} is a set of negatives. Policies then act by choosing actions
428 whose embeddings are closest to the goal embedding. These approaches can learn powerful, task-
429 agnostic representations from unlabeled trajectories, but the contrastive loss is global rather than
430 local in the sense of our dual advantage field: it encourages correct ordering over large batches of
431 positive and negative pairs without explicitly privileging one-step action-induced displacements in
432 representation space.

433 C Additional Environment and Evaluation Details

434 **OGBench environments.** We evaluate on goal-conditioned offline reinforcement learning tasks
435 from OGBench [20]. OGBench is designed to test several capabilities that are central to offline GCRL,
436 including long-horizon reasoning, trajectory stitching, generalization to unseen goals, robustness to
437 suboptimal data, and control under stochasticity. In our main experiments, we focus on the state-based
438 locomotion and manipulation tasks used in prior work.

439 The locomotion tasks include maze-style navigation domains such as `pointmaze`, `antmaze`, and
440 `humanoidmaze`, as well as `antsoccer`. These tasks require the agent to reach target goal states
441 from diverse initial configurations using only offline data. The difficulty varies with maze size, agent
442 morphology, and dataset coverage. In particular, `humanoidmaze` requires full-body control and there-
443 fore combines low-level locomotion with long-horizon navigation, while `antsoccer` additionally
444 requires controlling a ball while navigating.

445 The manipulation tasks include `cube`, `scene`, and `puzzle`. The `cube` environments test basic object
446 manipulation through pick-and-place, stacking, swapping, and rearrangement of colored cubes. The
447 `scene` environment contains multiple interacting objects, such as a cube, drawer, window, and buttons,
448 and therefore requires sequencing several atomic behaviors to achieve the desired goal configuration.
449 The `puzzle` environments instantiate a robotic version of the Lights Out puzzle, where pressing one
450 button changes the state of neighboring buttons. These tasks are particularly challenging because the
451 agent must combine continuous robotic control with combinatorial generalization over many possible
452 configurations.

453 **Dataset variants.** For each environment, OGBench provides multiple dataset variants that differ in
454 coverage, trajectory quality, and the extent to which successful behavior can be recovered directly
455 from the dataset. In maze-style locomotion tasks, `navigate` datasets are collected from noisy
456 expert policies that traverse the environment, while `stitch` datasets contain shorter trajectory
457 segments and require the policy to compose partial behaviors into longer goal-reaching trajectories.
458 Some locomotion domains also provide `explore` datasets, which contain highly exploratory and
459 substantially suboptimal trajectories.

460 For manipulation tasks, OGBench provides play and noisy datasets. The play datasets contain
 461 natural interaction trajectories generated by scripted policies with temporally correlated behavior.
 462 These datasets often contain useful local skills but do not necessarily demonstrate each evaluation
 463 task end-to-end. The noisy datasets are collected with larger, less structured exploration noise,
 464 which increases state-action coverage but also makes the data more suboptimal. Together, these
 465 dataset variants test whether an offline GCRL method can learn useful local behaviors, compose them
 466 over long horizons, and remain robust when the data are imperfect or only partially aligned with the
 467 evaluation goals.

468 **Evaluation protocol.** We follow the standard OGBench protocol and report success-based perfor-
 469 mance on each task. For a method m , environment e , and random seed r , let $s_{m,e,r} \in [0, 1]$ denote
 470 the resulting success rate, averaged over the evaluation episodes and goals for that environment.
 471 Higher values indicate better goal-reaching performance. Unless otherwise stated, all methods are
 472 trained purely offline on the provided datasets and are evaluated without additional environment
 473 interaction during training.

474 **Aggregate metrics with RLiable.** In addition to per-environment results, we report aggregate
 475 statistics using the RLiable evaluation framework [1]. RLiable is useful in the few-seed regime
 476 because it summarizes performance across tasks while also quantifying uncertainty with stratified-
 477 bootstrap confidence intervals. Importantly, RLiable does not discard “noisy” runs or remove
 478 experiments. Instead, it estimates how sensitive aggregate conclusions are to the finite set of tasks
 479 and random seeds.

480 Let $S_m = \{s_{m,e,r}\}_{e,r}$ denote the collection of scores for method m across environments and seeds.
 481 We report the following aggregate metrics:

$$\text{Mean}(m) = \frac{1}{|\mathcal{E}||\mathcal{R}|} \sum_{e \in \mathcal{E}} \sum_{r \in \mathcal{R}} s_{m,e,r}, \quad (24)$$

$$\text{Median}(m) = \text{median}(\{s_{m,e,r}\}_{e,r}), \quad (25)$$

$$\text{IQM}(m) = \text{mean}(\{s_{m,e,r} : s_{m,e,r} \text{ lies between the 25th and 75th percentiles}\}), \quad (26)$$

$$\text{OptimalityGap}(m) = \frac{1}{|\mathcal{E}||\mathcal{R}|} \sum_{e \in \mathcal{E}} \sum_{r \in \mathcal{R}} \max(0, 1 - s_{m,e,r}). \quad (27)$$

482 The interquartile mean (IQM) averages the middle 50% of scores and is therefore less sensitive to
 483 extreme outlier tasks than the mean, while being more statistically efficient than the median. The
 484 optimality gap measures the average shortfall from the maximum normalized score of 1; thus, lower
 485 values are better. Since our scores are success rates in $[0, 1]$, the optimality gap is directly interpretable
 486 as the average remaining failure mass. If scores are reported as percentages, they are first divided by
 487 100 before computing the RLiable metrics.

488 For confidence intervals, we use stratified bootstrap resampling over tasks and seeds. Each bootstrap
 489 replicate preserves the task structure: for every environment, we resample seeds with replacement
 490 and then recompute the aggregate metric on the resampled score matrix. The reported intervals
 491 correspond to the empirical percentiles of the bootstrap distribution. This procedure avoids treating
 492 all scores as exchangeable independent samples and prevents environments with more runs from
 493 dominating the uncertainty estimate.

494 D Additional Results

495 We include the additional RLiable [1] plots in Figure 6.

496 E AFU-style coupling of the bilinear value and dual score

497 This section provides the actor-free coupling we use between the bilinear value $V_\theta(s, g) =$
 498 $\psi_\theta(s)^\top \phi_\theta(g)$ and the dual score, following the separation of roles emphasized by Perrin-Gilbert [24].
 499 The policy parameters do not receive gradients through this objective; policy learning uses only the
 500 weighted regression step.

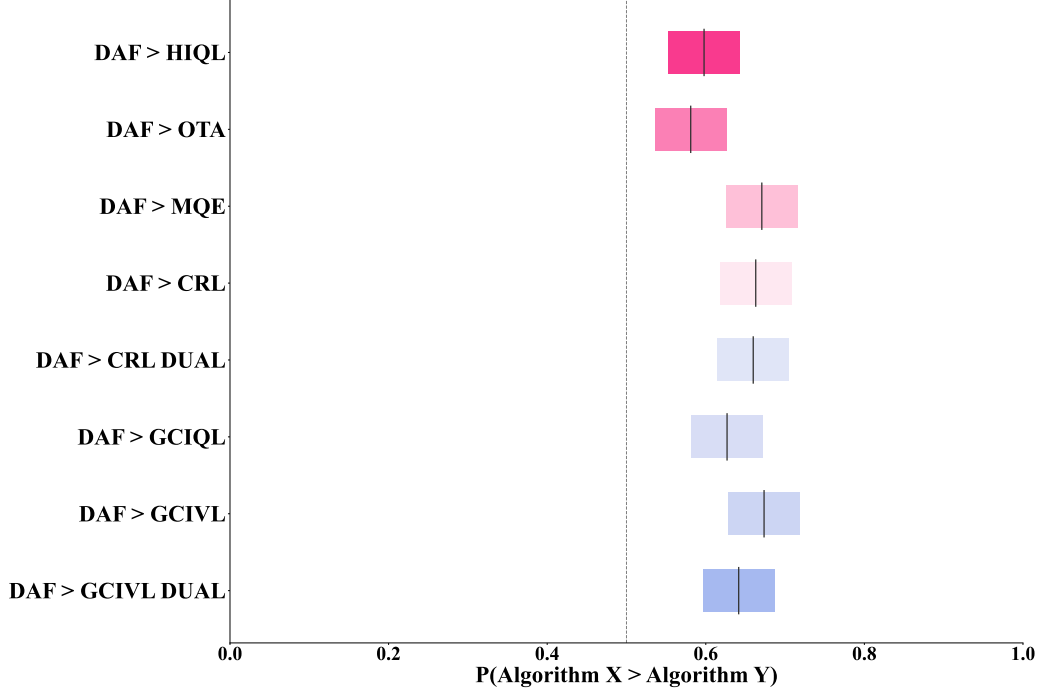


Figure 6: Reliable Probability of Improvement.

501 **Surrogate dual score for the coupling.** The main text defines the raw dual score z_θ in (14). In the
 502 AFU objective below it is convenient to use a non-positive surrogate

$$\tilde{A}_\theta(s, a, g) := h(z_\theta(s, a, g)), \quad h : \mathbb{R} \rightarrow (-\infty, 0], \quad (28)$$

503 where h is any monotone transformation used in implementation to keep the coupling term bounded
 504 on the optimistic side while preserving action ordering. In our experiments we use `softplus`
 505 function. The same z_θ can still be used directly in advantage-weighted regression, as in the main text;
 506 (28) is only required for the piecewise coupling with V_θ .

507 **Scalar Bellman target.** Let

$$T(s, a, g) := r(s, g) + \gamma V_\theta^{\text{tgt}}(s', g), \quad (29)$$

508 with V_θ^{tgt} a slowly updated target network for the bilinear head.

509 **Conditional scaling of V_θ .** Let $U = \mathbb{1}[V_\theta + \tilde{A}_\theta < T]$ and $\rho \in (0, 1)$. Define

$$\tilde{V} := (1 - \rho U) V_\theta + \rho U \text{stopgrad}(V_\theta). \quad (30)$$

510 When the optimistic sum $V_\theta + \tilde{A}_\theta$ falls short of the Bellman target T , the mask down-weights direct
 511 updates to V_θ so that \tilde{A}_θ can absorb slack in the near-optimistic regime.

512 **Piecewise coupling loss.** With $x = \tilde{V} - T$ and $y = \tilde{A}_\theta$, set

$$Z(x, y) = \begin{cases} (x + y)^2, & x \geq 0, \\ x^2 + y^2, & x < 0. \end{cases} \quad (31)$$

513 Training minimizes $\mathbb{E}_{\mathcal{D}}[Z(x, y)]$ jointly over the parameters of V_θ (equivalently ψ_θ and, where tied,
 514 ϕ_θ) and of the heads that define \tilde{A}_θ (including u_ξ and ϕ_θ as used in z_θ). The asymmetric split between
 515 $x \geq 0$ and $x < 0$ mirrors the AFU construction: pessimistic errors on V and the dual score are not
 516 forced to cancel spuriously when the backup is optimistic.

517 **Feature dynamics auxiliary loss.** The loss \mathcal{L}_{ae} in (15) complements the coupling above: the
 518 AFU-style term enforces Bellman consistency between V_θ and \tilde{A}_θ , while \mathcal{L}_{ae} grounds u_ξ in explicit
 519 one-step feature dynamics on the offline dataset.

520 F Theoretical Analysis

521 This section establishes two complementary properties of Dual Advantage Fields (DAF). First, we
 522 show that under exact representability DAF recovers the true Bellman advantage and therefore
 523 constitutes a valid policy-improvement operator (Section F.1). Second, we analyse a didactic 1-D
 524 example and prove that, even when the learned goal embedding is corrupted by noise in irrelevant
 525 directions, DAF’s local advantage remains significantly more robust than both flat and hierarchical
 526 value-difference extraction (Section F.2).

527 F.1 DAF as exact policy-improvement signal

528 Fix a goal g and consider the goal-conditioned MDP with reward $r_g(s, a) := r(s, a, g)$. For a policy
 529 π , define the usual Bellman advantage

$$A^\pi(s, a, g) := \mathbb{E}_{s' \sim P(\cdot | s, a)} [r(s, a, g) + \gamma V^\pi(s', g) - V^\pi(s, g)].$$

530 This is the relative quantity that drives policy improvement: only the ordering of actions at a given
 531 state matters, not the absolute level of V^π .

532 Assume that the policy value is realisable by the bilinear dual field,

$$V^\pi(s, g) = \psi(s)^\top \phi(g),$$

533 and that the action-effect model is exact,

$$u(s, a, g) = \mathbb{E}_{s' \sim P(\cdot | s, a)} [\gamma \psi(s') - \psi(s)].$$

534 Then the DAF score

$$D^\pi(s, a, g) := r(s, a, g) + u(s, a, g)^\top \phi(g)$$

535 equals the true goal-conditioned advantage:

$$D^\pi(s, a, g) = A^\pi(s, a, g).$$

536 **Proposition F.1** (DAF local policy improvement). *Let π^+ be any goal-conditioned policy satisfying*

$$\mathbb{E}_{a \sim \pi^+(\cdot | s, g)} [D^\pi(s, a, g)] \geq 0 \quad \text{for all } s, g.$$

537 *Under the realizability and exact action-effect assumptions above,*

$$V^{\pi^+}(s, g) \geq V^\pi(s, g) \quad \text{for all } s, g.$$

538 *Proof.* Since $D^\pi = A^\pi$, the assumption gives $\mathbb{E}_{a \sim \pi^+(\cdot | s, g)} [A^\pi(s, a, g)] \geq 0$. This is exactly

$$(T_{\pi^+} V^\pi)(s, g) - V^\pi(s, g) = \mathbb{E}_{a \sim \pi^+(\cdot | s, g)} [A^\pi(s, a, g)] \geq 0,$$

539 where T_{π^+} is the Bellman operator for policy π^+ . Hence $T_{\pi^+} V^\pi \geq V^\pi$ pointwise, and by mono-
 540 tonicity of the Bellman operator, $T_{\pi^+}^k V^\pi \geq V^\pi$ for every $k \geq 1$. Taking $k \rightarrow \infty$ and using the
 541 contraction property of T_{π^+} yields $V^{\pi^+} \geq V^\pi$. \square

542 The advantage-weighted regression (AWR) update used by DAF is one such improvement in the
 543 exact on-policy case. If

$$\pi_\alpha^+(a | s, g) = \frac{\pi(a | s, g) \exp(\alpha D^\pi(s, a, g))}{\sum_b \pi(b | s, g) \exp(\alpha D^\pi(s, b, g))}, \quad \alpha \geq 0,$$

544 then a standard argument shows $\mathbb{E}_{a \sim \pi_\alpha^+} [D^\pi(s, a, g)] \geq \mathbb{E}_{a \sim \pi} [D^\pi(s, a, g)] = 0$, so the AWR policy
 545 satisfies the condition of Proposition F.1.

546 **Corollary F.2** (Exact DAF policy iteration). *In a finite discounted goal-conditioned MDP, suppose*
 547 *each iteration k uses exact representations for V^{π_k} and an exact action-effect model, and define*

$$\pi_{k+1}(\cdot \mid s, g) \in \arg \max_{\pi'} \mathbb{E}_{a \sim \pi'(\cdot \mid s, g)} [D^{\pi_k}(s, a, g)].$$

548 *Then π_{k+1} is the standard greedy policy-improvement step with respect to Q^{π_k} . Consequently,*
 549 *repeated exact DAF improvement is policy iteration and converges to an optimal goal-conditioned*
 550 *policy.*

551 *Proof.* Because $D^{\pi_k} = A^{\pi_k} = Q^{\pi_k} - V^{\pi_k}$, maximising D^{π_k} over actions is equivalent to maximising
 552 Q^{π_k} . The result follows from classical policy iteration for finite discounted MDPs, applied separately
 553 for each goal g . \square

554 **Relation to hierarchical policies.** Let Π denote the class of all stationary goal-conditioned
 555 primitive-action policies, and let $\Pi_{\text{hier}} \subseteq \Pi$ be any hierarchically constrained class (e.g. subgoal or
 556 option policies). The optimal primitive-action policy $\pi^* \in \arg \max_{\pi \in \Pi} V^{\pi}$ satisfies

$$V^{\pi^*}(s, g) \geq \sup_{\pi \in \Pi_{\text{hier}}} V^{\pi}(s, g) \quad \text{for all } s, g.$$

557 Thus, in the exact realisable limit, DAF policy iteration reaches a policy that is at least as good as the
 558 best policy in any fixed hierarchical class.

559 This comparison is a representational statement: hierarchy may improve learning by reducing the ef-
 560 fective horizon, but a fixed hierarchy can also introduce subgoal-level constraints that exclude the true
 561 optimal primitive-action policy. DAF instead performs improvement directly at the primitive-action
 562 level using the local dual advantage, while preserving the long-horizon reachability information
 563 encoded in the dual value field.

564 F.2 Robustness to learned embedding noise: a didactic example

565 We now turn to a more practical regime where the representation is learned from finite data and
 566 inevitably contains noise. The following analysis uses only the bilinear parameterisation and DAF
 567 scoring rule; no quasimetric or Eikonal assumptions are required.

568 F.2.1 Environment and representation model

569 **Line-world dynamics.** Consider deterministic states $s \in \{0, 1, \dots, T\}$ with a fixed goal $g = T >$
 570 0 . Two actions are available: right ($a = +1, s \rightarrow s + 1$) and left ($a = -1, s \rightarrow s - 1$). The episode
 571 terminates upon reaching g ; the reward is 0 at the goal and -1 otherwise. Hence the optimal policy
 572 always moves right for $s < T$, and the optimal (negative) value function is

$$V^*(s, g) = s - T, \quad s \leq T.$$

573 **Fixed state embedding.** The environment provides a feature map $\psi : \mathbb{Z} \rightarrow \mathbb{R}^d$ with $d = m +$
 574 2 ($m \geq 0$):

$$\psi(s) = [s, 1, f_1(s), \dots, f_m(s)]^\top,$$

575 where $\{f_i\}_{i=1}^m$ are bounded C^2 functions (or, in the discrete case, functions with well-defined first
 576 and second differences). The first two coordinates are “essential” for representing the linear optimal
 577 value; the remaining ones are *nuisance* dimensions that are irrelevant for the control task (e.g., visual
 578 textures, lighting gradients).

579 **True goal embedding.** The optimal value can be expressed via an inner product:

$$\phi^*(g) = [1, -T, 0, \dots, 0]^\top \implies \psi(s)^\top \phi^*(g) = s - T = V^*(s, g).$$

580 **F.2.2 Noise model for the learned goal embedding**

581 In offline training, the goal embedding $\phi(g)$ is estimated from a finite dataset. Because the
 582 temporal-difference loss only weakly constrains the coefficients of the nuisance coordinates (espe-
 583 cially if those coordinates vary slowly), the learned embedding can accumulate significant noise along
 584 those directions. We model this by an additive perturbation confined to the nuisance components:

$$\phi(g) = \phi^*(g) + \varepsilon, \quad \varepsilon = [0, 0, \eta_1, \dots, \eta_m]^\top,$$

585 where $\eta_i \sim \mathcal{N}(0, \sigma_i^2)$ are independent. The essential coordinates are assumed to be learned accurately
 586 for simplicity; allowing noise there would not change the qualitative conclusions.

587 Consequently the noisy value estimate at any state s is

$$\widehat{V}(s, g) = \psi(s)^\top \phi(g) = s - T + \sum_{i=1}^m \eta_i f_i(s).$$

588 For the subgoal s_{sub} we assume the same embedding function $\phi(\cdot)$ is applied and that its noise is
 589 independent of $\phi(g)$:

$$\phi(s_{\text{sub}}) = \phi^*(s_{\text{sub}}) + \tilde{\varepsilon}, \quad \tilde{\eta}_i \sim \mathcal{N}(0, \sigma_i^2) \text{ independent of } \eta_i.$$

590 **F.2.3 Action-effect model and policy extraction rules**

591 We assume that a separate action-effect model $u(s, a)$ has been trained to regress to the true one-step
 592 feature change $\psi(s+a) - \psi(s)$ and has converged to the exact quantity (realistic because the model
 593 sees abundant transitions and the dynamics are deterministic).

594 Thus

$$u(s, +1) = \psi(s+1) - \psi(s), \quad u(s, -1) = \psi(s-1) - \psi(s).$$

595 We compare three policy extraction methods, all built upon the same learned bilinear value \widehat{V} and the
 596 same u .

597 1. **Flat value-difference.** Choose the action that leads to the highest estimated next-state value:

$$a_V(s) = \arg \max_{a \in \{-1, +1\}} \widehat{V}(s+a, g).$$

598 This corresponds to the implicit advantage used in HIQL’s flat baseline (comparing $V(s+1, g)$
 599 and $V(s-1, g)$).

600 2. **DAF local advantage.** Score each action by the inner product of its predicted feature displacement
 601 and the goal embedding (Eq. 14 in the main paper):

$$a_{\text{DAF}}(s) = \arg \max_{a \in \{-1, +1\}} u(s, a)^\top \phi(g).$$

602 (The sparse reward, identical for both actions, is omitted from the comparison.)

603 3. **Hierarchical HIQL.** The hierarchical policy first selects a subgoal at distance $k \geq 2$ (to the right,
 604 $s_{\text{sub}} = s+k$) by comparing values of the candidate subgoals:

$$s_{\text{sub}} = \arg \max_{x \in \{s+k, s-k\}} [\widehat{V}(x, g) - \widehat{V}(s, g)].$$

605 Subsequently a low-level controller attempts to reach that subgoal, using the subgoal’s own
 606 embedding $\phi(s_{\text{sub}})$ and the same flat value-difference rule:

$$a_\ell(s) = \arg \max_{a \in \{-1, +1\}} \widehat{V}(s+a, s_{\text{sub}}).$$

607 An error occurs if either the subgoal choice is wrong or the low-level action is wrong; we bound
 608 this with a union argument as in Park et al. [19, Proposition 4.1].

609 **F.2.4 Error probabilities**

610 For any nuisance function f , define the first and second discrete differences at state s :

$$\Delta f(s) := f(s+1) - f(s-1), \quad \Delta^2 f(s) := f(s+1) + f(s-1) - 2f(s).$$

611 **Flat value-difference.**

$$\Delta_V(s) = \widehat{V}(s+1, g) - \widehat{V}(s-1, g) = 2 + \sum_{i=1}^m \eta_i \Delta f_i(s).$$

612 **DAF.**

$$\Delta_{\text{DAF}}(s) = (u(s, +1) - u(s, -1))^\top \phi(g) = 2 + \sum_{i=1}^m \eta_i \Delta^2 f_i(s).$$

613 **Hierarchical high-level.**

$$\Delta_{\text{high}}(s) = \widehat{V}(s+k, g) - \widehat{V}(s-k, g) = 2k + \sum_{i=1}^m \eta_i (f_i(s+k) - f_i(s-k)).$$

614 **Hierarchical low-level.** Conditioned on the subgoal $s+k$ being selected,

$$\Delta_{\text{low}}(s) = \widehat{V}(s+1, s+k) - \widehat{V}(s-1, s+k) = 2 + \sum_{i=1}^m \tilde{\eta}_i \Delta f_i(s).$$

615 All decision statistics are Gaussian. Let Φ be the standard normal c.d.f.

616 **Proposition F.3** (Error probabilities). *For any state $s \in \{1, \dots, T-1\}$ and subgoal step k ,*

$$\begin{aligned} \varepsilon_{\text{flat}}(s) &= \Phi \left(-\frac{2}{\sqrt{\sum_i \sigma_i^2 (\Delta f_i(s))^2}} \right), \\ \varepsilon_{\text{DAF}}(s) &= \Phi \left(-\frac{2}{\sqrt{\sum_i \sigma_i^2 (\Delta^2 f_i(s))^2}} \right), \\ \varepsilon_{\text{high}}(s) &= \Phi \left(-\frac{2k}{\sqrt{\sum_i \sigma_i^2 (f_i(s+k) - f_i(s-k))^2}} \right), \\ \varepsilon_{\text{low}}(s) &= \Phi \left(-\frac{2}{\sqrt{\sum_i \sigma_i^2 (\Delta f_i(s))^2}} \right). \end{aligned}$$

617 *The overall hierarchical error is bounded by*

$$\varepsilon_{\text{hier}}(s) \leq \varepsilon_{\text{high}}(s) + \varepsilon_{\text{low}}(s).$$

618 *Proof.* Each decision margin is a normal random variable with the stated mean and variance; mis-
 619 classification is the event “margin < 0 ”. The hierarchical bound follows from a union bound over the
 620 two decision stages, exactly as in Park et al. [19, Proposition 4.1]. \square

621 **F.2.5 Why DAF can be more robust**

622 The formulas in Proposition F.3 show that DAF’s noise enters through the *second differences* $\Delta^2 f_i(s)$,
 623 whereas all value-difference methods (flat and low-level) involve the *first differences* $\Delta f_i(s)$. The
 624 high-level comparison involves the even larger span $f_i(s+k) - f_i(s-k)$.

625 For many realistic nuisance functions, the second difference is much smaller than the first difference.
 626 Two concrete regimes make this quantitative.

627 **Corollary F.4** (Affine nuisance coordinates are eliminated by DAF). *If $f_i(s) = \alpha_i s + \beta_i$ for all i ,*
 628 *then $\Delta^2 f_i(s) = 0$ for every s ; hence $\Delta_{\text{DAF}}(s) \equiv 2$ and $\varepsilon_{\text{DAF}}(s) = 0$. In contrast,*

$$\varepsilon_{\text{flat}}(s) = \varepsilon_{\text{low}}(s) = \Phi\left(-\frac{1}{\sqrt{\sum_i \sigma_i^2 \alpha_i^2}}\right), \quad \varepsilon_{\text{high}}(s) = \Phi\left(-\frac{1}{\sqrt{\sum_i \sigma_i^2 \alpha_i^2}}\right).$$

629 *Thus DAF makes zero mistakes regardless of the horizon, while the flat and hierarchical baselines*
 630 *can suffer significant error whenever $\sum_i \sigma_i^2 \alpha_i^2$ is large.*

631 **Corollary F.5** (Low-curvature nuisance coordinates). *Suppose each f_i is twice differentiable with*
 632 *$|f_i''(s)| \leq C$ and that over a short interval the first difference can be expressed as $\Delta f_i(s) =$*
 633 *$2f_i'(s) + O(C)$, $\Delta^2 f_i(s) = 2f_i''(s) + O(C)$. If the local slope $f_i'(s)$ is large (e.g., a strong linear trend)*
 634 *while the curvature remains bounded, then $\sigma_{\text{DAF}}^2(s) = O(C^2)$ whereas $\sigma_{\text{flat}}^2(s) = 4 \sum_i \sigma_i^2 f_i'(s)^2$*
 635 *can be arbitrarily large. Consequently $\varepsilon_{\text{DAF}}(s)$ stays close to zero while $\varepsilon_{\text{flat}}(s)$ and $\varepsilon_{\text{low}}(s)$ may*
 636 *approach $\frac{1}{2}$.*

637 **Comparison with the hierarchical baseline.** Even with a well-chosen subgoal step k , the low-level
 638 controller still relies on first differences (Proposition F.3), inheriting the same vulnerability as the flat
 639 extraction. Moreover, the high-level stage introduces an additional source of error that scales with the
 640 span of the nuisance functions. As a result, a *single* DAF flat policy can achieve a lower error rate
 641 than a hierarchical policy that employs two value-difference decisions.

642 **Illustrative quantitative example.** Let $m = 1$ and $f_1(s) = s^2$. Then $\Delta f_1(s) = 4s$, $\Delta^2 f_1(s) = 2$,
 643 and

$$\varepsilon_{\text{DAF}}(s) = \Phi\left(-\frac{2}{\sigma_1 \cdot 2}\right), \quad \varepsilon_{\text{flat}}(s) = \varepsilon_{\text{low}}(s) = \Phi\left(-\frac{2}{\sigma_1 \cdot 4s}\right).$$

644 For a state far from the goal ($s = T - 1 \gg 1$), $\varepsilon_{\text{flat}}$ and ε_{low} are close to 0.5 if σ_1 is large, while
 645 ε_{DAF} remains bounded by a constant that does not grow with T .

646 G Compute Resources

647 All experiments were performed on servers with a single H100 GPU with 80 GB of GPU memory, 12
 648 CPU cores, and 244 GB of RAM. All metrics for the experiments were logged using the Weights &
 649 Biases platform. Overall, the Weights & Biases project of the paper had 17,359 tracked experiments
 650 at the time of submission and used an estimated ~ 407 days of GPU compute in total.

651 **NeurIPS Paper Checklist**

652 **1. Claims**

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

655 Answer: [Yes]

656 Justification: Claims are supported by the main experiments, see Table 1, Table 2, Table 3
657 and Figure 5.

658 Guidelines:

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

668 **2. Limitations**

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

670 Answer: [Yes]

671 Justification: We discuss the limitations in Section 6.

672 Guidelines:

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

699 **3. Theory assumptions and proofs**

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

702 Answer: [Yes]

703 Justification: See Section F

704 Guidelines:

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

715 4. Experimental result reproducibility

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

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

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

719 Answer: [Yes]

720 Justification: We describe the method in Algorithm 1, evaluation protocols in Section C and

721 provide hyperparameters in Section A.

722 Guidelines:

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

754 5. Open access to data and code

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

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

757 material?

758
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: We provide the anonymized source code along with the submission, which include instructions on how to reproduce experiments on a publicly available OGBench [20].

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: See Section A, Section 4.

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: We use five random seeds and report IQM and 95%-CI based on stratified bootstrapping, following the Agarwal et al. [1].

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

820 8. Experiments compute resources

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 Answer: [Yes]

823 Justification: We report information about compute resources in Section G.

824 Guidelines:

- 825 • The answer [N/A] means that the paper does not include experiments.
- 826
- 827 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- 828
- 829 • The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 830
- 831 • 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).
- 832
- 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 Answer: [Yes]

838 Justification: [Yes]

839 Guidelines:

- 840 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 841
- 842 • If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- 843
- 844 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 845
- 846

847 10. Broader impacts

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

849 Answer: [N/A]

850 Justification: [N/A]

851 Guidelines:

- 852 • The answer [N/A] means that there is no societal impact of the work performed.
- 853
- 854 • 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.
- 855
- 856 • 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.
- 857
- 858
- 859

- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 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).

875 11. Safeguards

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

879 Answer: [N/A]

880 Justification: [N/A]

881 Guidelines:

- 882
- 883
- 884
- 885
- 886
- 887
- 888
- 889
- 890
- 891
- 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.

892 12. Licenses for existing assets

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

896 Answer: [N/A]

897 Justification: [N/A]

898 Guidelines:

- 899
- 900
- 901
- 902
- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- 911
- 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.

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

914 **13. New assets**

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

917 Answer: [N/A]

918 Justification: [N/A]

919 Guidelines:

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

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

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

932 Answer: [N/A]

933 Justification: [N/A]

934 Guidelines:

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

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

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

949 Answer: [N/A]

950 Justification: [N/A]

951 Guidelines:

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

962 **16. Declaration of LLM usage**

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

967 Answer: [N/A]

968 Justification: [N/A]

969 Guidelines:

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