# 000 UNSUPERVISED-TO-ONLINE 002 REINFORCEMENT LEARNING

Anonymous authors

003

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

025

026 027 028

029

Paper under double-blind review

# ABSTRACT

Offline-to-online reinforcement learning (RL), a framework that trains a policy with offline RL and then further fine-tunes it with online RL, has been considered a promising recipe for data-driven decision-making. While sensible, this framework has drawbacks: it requires domain-specific offline RL pre-training for each task, and is often brittle in practice. In this work, we propose **unsupervised-to-online** RL (U20 RL), which replaces domain-specific supervised offline RL with unsu*pervised* offline RL, as a potentially better alternative to offline-to-online RL. U2O RL not only enables reusing a single pre-trained model for multiple downstream tasks, but also learns better representations, which often result in *even better* performance and stability than *supervised* offline-to-online RL. To instantiate U2O RL in practice, we propose a general recipe for U2O RL to bridge task-agnostic unsupervised offline skill-based policy pre-training and supervised online fine-tuning. Throughout our experiments in eleven state-based and pixel-based environments, we empirically demonstrate that U2O RL often achieves strong performance that matches or even outperforms previous offline-to-online RL approaches when the dataset consists of diverse trajectories, while being able to reuse a single pre-trained model for a number of different downstream tasks.

# 1 INTRODUCTION

Across natural language processing (NLP), computer vision (CV), and speech processing, ubiquitous in the recent successes of machine learning is the idea of adapting an expressive model pre-trained on large-scale data to domain-specific tasks via fine-tuning. In the domain of reinforcement learning (RL), offline-to-online RL has been considered an example of such a recipe for leveraging offline data for efficient online fine-tuning. Offline-to-online RL first trains a task-specific policy on a previously collected dataset with offline RL, and then continues training the policy with additional environment interactions to further improve performance.

But, is offline-to-online RL really the most effective way to leverage offline data for online RL? Offline-to-online RL indeed has several limitations. First, it pre-trains a policy with a *domain-specific* task reward, which precludes sharing a single pre-trained model for multiple downstream tasks. This 040 is in contrast to predominant pre-training recipes in large language models or visual representation 041 learning, where they pre-train large models with self-supervised or *unsupervised* objectives to learn 042 useful representations, which can facilitate learning a wide array of downstream tasks. Second, 043 naïve offline-to-online RL is often brittle in practice (Lee et al., 2022; Nakamoto et al., 2023). This 044 is mainly because pre-trained offline RL agents suffer the distributional shift between the offline and online interaction data (Lee et al., 2022; Nakamoto et al., 2023) or experience feature collapse (Section 5), which necessitates specialized, potentially complicated techniques. 046

In this work, our central hypothesis is that *unsupervised pre-training of diverse policies* from offline data can serve as an effective data-driven recipe for *online* RL, and can be more effective than even domain-specific ("supervised") offline pre-training. We call this recipe unsupervised-to-online RL (U2O RL). U2O RL has two appealing properties. First, unlike offline-to-online RL, a single pre-trained model can be fine-tuned for different downstream tasks. Since offline unsupervised RL does not require task information, we can pre-train diverse policies on unlabeled data before knowing downstream tasks. Second, by pre-training multi-task policies with diverse intrinsic rewards, the agent extracts rich representations from data, which often helps achieve *even better* final performance



Figure 1: **Illustration of U2O RL.** In this work, we propose to replace supervised offline RL with *unsupervised offline RL* in the offline-to-online RL framework. We call this scheme **unsupervised-to-online RL (U2O RL)**. U2O RL consists of three stages: (1) unsupervised offline RL pre-training, (2) bridging, and (3) online RL fine-tuning. In unsupervised offline RL pre-training, we train a multi-task skill policy  $\pi_{\theta}(a \mid s, z)$  instead of a single-task policy  $\pi_{\theta}(a \mid s)$ . Then, we convert the multi-task policy into a task-specific policy in the bridging phase. Finally, we fine-tune the skill policy with online environment interactions.

and stability than *supervised* offline-to-online RL. This resembles how general-purpose unsupervised pre-training in other domains, such as with LLMs or self-supervised representations (Brown et al., 2020; Devlin et al., 2019; Radford et al., 2019; He et al., 2021; 2020; Hénaff et al., 2020), improves over the performance of domain-specific specialist pre-training.

078 U2O RL consists of three stages: unsupervised offline policy pre-training, bridging, and online 079 fine-tuning (Figure 1). In the first unsupervised offline pre-training phase, we employ a skill-based 080 offline unsupervised RL or offline goal-conditioned RL method, which trains diverse behaviors 081 (or skills) with intrinsic rewards and provides an efficient mechanism to identify the best skill for 082 a given task reward. In the subsequent bridging and fine-tuning phases, we adapt the best skill 083 among the learned policies to the given downstream task reward with online RL. Here, to prevent a 084 potential mismatch between the intrinsic and task rewards, we propose a simple yet effective reward 085 scale matching technique that bridges the gap between the two training schemes and thus improves 086 performance and stability.

087 Our main contributions in this work are twofold. First, to the best of our knowledge, this is the first 880 work that makes the (potentially surprising) observation that it is often better to replace supervised 089 offline RL with unsupervised offline RL in the offline-to-online RL setting. We also identify the 090 reason behind this phenomenon: this is mainly because offline unsupervised pre-training learns better 091 representations than task-specific supervised offline RL. Second, we propose a general recipe to 092 bridge skill-based unsupervised offline RL pre-training and online RL fine-tuning. Through our experiments on eleven state-based and pixel-based environments, we demonstrate that U2O RL often 093 outperforms standard offline-to-online RL both in terms of sample efficiency and final performance, 094 while being able to reuse a single pre-trained model for multiple downstream tasks. 095

096

069

070

071

072

073

# 2 Related work

098

099 **Online RL from prior data**. Prior works have proposed several ways to leverage a previously 100 collected offline dataset to accelerate online RL training. They can be categorized into two main 101 groups: offline-to-online RL and off-policy online RL with offline data. Offline-to-online RL first 102 pre-trains a policy and a value function with offline RL (Lange et al., 2012; Levine et al., 2020; 103 Fujimoto & Gu, 2021; Fujimoto et al., 2019; Kumar et al., 2019; Tarasov et al., 2023a; Wu et al., 104 2019a; Kostrikov et al., 2021; Kumar et al., 2020; Hansen-Estruch et al., 2023; Kostrikov et al., 105 2022; Nair et al., 2020; Peng et al., 2019; Wang et al., 2020), and then continues to fine-tune them with additional online interactions (Lee et al., 2022; Nair et al., 2020; Nakamoto et al., 2023; Yu & 106 Zhang, 2023; Lei et al., 2023; Zheng et al., 2022; Mark et al., 2022; Zhao et al., 2023). Since naïve 107 offline-to-online RL is often unstable in practice due to the distributional shift between the dataset and

108 online interactions, prior works have proposed several techniques, such as balanced sampling (Lee 109 et al., 2022), actor-critic alignment (Yu & Zhang, 2023), adaptive conservatism (Wang et al., 2023a), 110 and return lower-bounding (Nakamoto et al., 2023). In this work, unlike offline-to-online RL, which 111 trains a policy with the target task reward, we offline pre-train a multi-task policy with unsupervised 112 (intrinsic) reward functions. This makes our single pre-trained policy reusable for any downstream task and learn richer representations. The other line of research, off-policy online RL, trains an online 113 RL agent from scratch on top of a replay buffer filled with offline data, without any pre-training (Ball 114 et al., 2023; Li et al., 2023; Luo et al., 2024; Song et al., 2023). While this simple approach often 115 leads to improved stability and performance (Ball et al., 2023), it does not leverage the benefits of pre-116 training; in contrast, we do leverage pre-training by learning useful features via offline unsupervised 117 RL, which we show leads to better fine-tuning performance in our experiments. 118

Unsupervised RL. The goal of unsupervised RL is to leverage unsupervised pre-training to facilitate 119 downstream reinforcement learning. Prior works have mainly focused on unsupervised representation 120 learning and unsupervised behavior learning. Unsupervised representation learning methods (Ser-121 manet et al., 2018; Shah & Kumar, 2021; Parisi et al., 2022; Xiao et al., 2022; Nair et al., 2022; 122 Ma et al., 2023b;a; Ghosh et al., 2023; Seo et al., 2022b;a; 2023) aim to extract useful (visual) 123 representations from data. These representations are then fed into the policy to accelerate task 124 learning. In this work, we focus on unsupervised behavior learning, which aims to pre-train policies 125 that can be directly adapted to downstream tasks. Among unsupervised behavior learning methods, 126 online unsupervised RL pre-trains useful policies by either maximizing state coverage (Pathak et al., 127 2017; 2019; Mendonca et al., 2021; Liu & Abbeel, 2021) or learning distinct skills (Gregor et al., 128 2016; Eysenbach et al., 2019b; Sharma et al., 2020; Park et al., 2024d) via reward-free interactions 129 with the environment. In this work, we consider offline unsupervised RL, which does not allow any environment interactions during the pre-training stage. 130

131 Offline unsupervised RL. Offline unsupervised RL methods focus on learning diverse policies (*i.e.*, 132 skills) from the dataset, rather than exploration, as online interactions are not permitted in this problem 133 setting. There exist three main approaches to offline unsupervised RL. Behavioral cloning-based 134 approaches extract skills from an offline dataset by training a generative model (e.g., variational 135 autoencoders (Kingma & Welling, 2014), Transformers (Vaswani et al., 2017), etc.) (Ajay et al., 2021; Pertsch et al., 2021; Singh et al., 2021). Offline goal-conditioned RL methods learn diverse 136 goal-reaching behaviors with goal-conditioned reward functions (Chebotar et al., 2021; Eysenbach 137 et al., 2022; Ma et al., 2022; Park et al., 2024b; Wang et al., 2023b; Yang et al., 2023; Fang et al., 138 2022; 2023). Offline unsupervised skill learning approaches learn diverse skills based on intrinsically 139 defined reward functions (Park et al., 2024c; Touati et al., 2022; Hu et al., 2023). Among these 140 approaches, we use methods in the second and third categories (i.e., goal- or skill-based unsupervised 141 offline RL) as part of our method. 142

Our goal in this work is to study how unsupervised offline RL, as opposed to supervised task-specific 143 offline RL, can be employed to facilitate online RL fine-tuning. While somewhat similar unsupervised 144 pre-training schemes have been explored in prior works, they either consider hierarchical RL (or 145 zero-shot RL) with frozen learned skills without fine-tuning (Ajay et al., 2021; Pertsch et al., 2021; 146 Touati et al., 2022; Park et al., 2024c; Hu et al., 2023), assume online-only RL (Laskin et al., 2021), 147 or are limited to the specific setting of goal-conditioned RL (Fang et al., 2022; 2023; Eysenbach et al., 148 2019a; Nasiriany et al., 2019). To the best of our knowledge, this is the first work that considers 149 the *fine-tuning* of skill policies pre-trained with unsupervised offline RL in the context of offline-to-150 online RL. Through our experiments, we show that our fine-tuning framework leads to significantly 151 better performance than previous approaches based on hierarchical RL, zero-shot RL, and standard 152 offline-to-online RL.

153 154

# 3 PRELIMINARIES

155 156

157 We formulate a decision making problem as a Markov decision process (MDP) (Sutton & Barto, 158 2018), which is defined by a tuple of  $(S, A, P, r, \rho, \gamma)$ , where S is the state space, A is the action 159 space,  $P: S \times A \to \Delta(S)$  is the transition dynamics,  $r: S \times A \times S \to \mathbb{R}$  is the task reward 160 function,  $\rho \in \Delta(S)$  is the initial state distribution, and  $\gamma \in (0, 1)$  is the discount factor. Our aim 161 is to learn a policy  $\pi: S \to \Delta(A)$  that maximizes the expectation of cumulative task rewards,  $\mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}, a_{t}, s_{t+1})].$ 

$$\mathcal{L}_{\mathsf{TD}}(\phi) = \mathbb{E}_{(s,a,s',r)\sim\mathcal{D}_{\mathsf{off}}}\left[ \left( r + \gamma \max_{a'} Q_{\bar{\phi}}(s',a') - Q_{\phi}(s,a) \right)^2 \right],\tag{1}$$

where  $Q_{\phi}$  denotes the parameterized action-value function, and  $Q_{\bar{\phi}}$  represents the target action-value function (Mnih et al., 2013), whose parameter  $\bar{\phi}$  is updated via Polyak averaging (Polyak & Juditsky, 1992) using  $\phi$ . We can then train a policy  $\pi$  to maximize  $\mathbb{E}_{a \sim \pi} [Q_{\phi}(s, a)]$ .

173 While this simple off-policy TD learning can be enough when the dataset has sufficiently large 174 state-action coverage, offline datasets in practice often have limited coverage, which makes the agent 175 susceptible to value overestimation and exploitation, as the agent cannot get corrective feedback from the environment (Levine et al., 2020). To address this issue, Kostrikov et al. (2022) proposed implicit 176 Q-learning (IQL), which fits an optimal action-value function without querying out-of-distribution 177 actions: IQL replaces the arg max operator, which potentially allows the agent to exploit Q-values 178 from out-of-distribution actions, with an expectile loss that implicitly approximates the maximum 179 value. Specifically, IQL minimizes the following losses: 180

$$\mathcal{L}^{Q}_{\text{IQL}}(\phi) = \mathbb{E}_{(s,a,s',r)\sim\mathcal{D}_{\text{off}}}\left[ (r + \gamma V_{\psi}(s') - Q_{\phi}(s,a))^{2} \right],$$
(2)

$$\mathcal{L}_{IQL}^{V}(\psi) = \mathbb{E}_{(s,a)\sim\mathcal{D}_{off}} \left[ \ell_{\tau}^{2}(Q_{\bar{\phi}}(s,a) - V_{\psi}(s)) \right], \tag{3}$$

where  $Q_{\phi}$  and  $Q_{\bar{\phi}}$  respectively denote the action-value and target action-value functions,  $V_{\psi}$  denotes the value function,  $\ell_{\tau}^2(x) = |\tau - \mathbb{1}(x < 0)|x^2$  denotes the expectile loss (Newey & Powell, 1987) and  $\tau$  denotes the expectile parameter. Intuitively, the asymmetric expectile loss in Equation 3 makes  $V_{\psi}$  implicitly approximate  $\max_a Q_{\bar{\phi}}(s, a)$  by penalizing positive errors more than negative errors.

Hilbert foundation policy (HILP). Our unsupervised-to-online recipe requires an offline unsupervised RL algorithm that trains a skill policy  $\pi_{\theta}(a \mid s, z)$  from an unlabeled dataset, and we mainly use HILP (Park et al., 2024c) in our experiments. HILP consists of two phases. In the first phase, HILP trains a feature network  $\xi : S \to Z$  that embeds temporal distances (*i.e.*, shortest path lengths) between states into the latent space by enforcing the following equality:

$$d^*(s,g) = \|\xi(s) - \xi(g)\|_2 \tag{4}$$

for all  $s, g \in S$ , where  $d^*(s, g)$  denotes the temporal distance (the minimum number of steps required to reach g from s) between s and g. In practice, given the equivalence between goal-conditioned values and temporal distances,  $\xi$  can be trained with any offline goal-conditioned RL algorithm (Park et al., 2024b) (see Park et al. (2024c) for further details). After training  $\xi$ , HILP trains a skill policy  $\pi_{\theta}(a \mid s, z)$  with the following intrinsic reward using an off-the-shelf offline RL algorithm (Kostrikov et al., 2022; Fujimoto et al., 2018):

$$^{int}(s, a, s', z) = (\xi(s') - \xi(s))^{\top} z,$$
(5)

where z is sampled from the unit ball,  $\{z \in \mathbb{Z} : ||z|| = 1\}$ . Intuitively, Equation 5 encourages the agent to learn behaviors that move in every possible latent direction, resulting in diverse statespanning skills (Park et al., 2024c). Note that Equation 5 can be interpreted as the inner product between the task vector z and the feature vector  $f(s, a, s') := \xi(s') - \xi(s)$  in the successor feature framework (Dayan, 1993; Barreto et al., 2017).

208 209

210

201

202

167

168 169

181 182 183

194

# 4 UNSUPERVISED-TO-ONLINE RL (U2O RL)

η

Our main hypothesis in this work is that task-agnostic offline RL pre-training of *unsupervised* skills
 can be more effective than task-specific, supervised offline RL for online RL fine-tuning. We call this
 recipe unsupervised-to-online RL (U2O RL). In this section, we first describe the three stages of
 U2O RL (Figure 1): unsupervised offline policy pre-training (Section 4.1), bridging (Section 4.2),
 and online fine-tuning (Section 4.3). We then explain why *unsupervised*-to-online RL is potentially
 better than standard *supervised* offline-to-online RL (Section 4.4).

# 4.1 UNSUPERVISED OFFLINE POLICY PRE-TRAINING

In the first unsupervised offline policy pre-training phase (Figure 1 (bottom left)), we train diverse policies (or *skills*) with intrinsic rewards to extract a variety of useful behaviors as well as rich features from the offline dataset  $\mathcal{D}_{off}$ . In other words, instead of training a single-task policy  $\pi_{\theta}(a \mid s)$  with task rewards r(s, a, s') as in standard offline-to-online RL, we train a *multi-task* skill policy  $\pi_{\theta}(a \mid s, z)$  with a family of unsupervised, *intrinsic* rewards  $r^{int}(s, a, s', z)$ , where z is a skill latent vector sampled from a latent space  $\mathcal{Z} = \mathbb{R}^d$ . Even if  $\mathcal{D}_{off}$  contains reward labels, we do not use any reward information in this phase.

225 Among existing unsupervised offline policy pre-training methods (Section 2), we opt to employ successor feature-based methods (Dayan, 1993; Barreto et al., 2017; Wu et al., 2019b; Touati et al., 226 2022; Park et al., 2024c) or offline goal-conditioned RL methods (Chebotar et al., 2021; Park et al., 227 2024b) for our unsupervised pre-training, since they provide a convenient mechanism to identify 228 the best skill latent vector given a downstream task, which we will utilize in the next phase. More 229 concretely, we mainly choose to employ HILP (Park et al., 2024c) (Section 3) as an unsupervised 230 offline policy pre-training method in our experiments for its state-of-the-art performance in previous 231 benchmarks (Park et al., 2024c). We note, however, that any other unsupervised offline successor 232 feature-based skill learning methods (Touati et al., 2022) or offline goal-conditioned RL methods (Park 233 et al., 2024b) can also be used in place of HILP (see Appendix A.2).

- 234
- 235 236

251

253

# 4.2 BRIDGING OFFLINE UNSUPERVISED RL AND ONLINE SUPERVISED RL

After finishing unsupervised offline policy pre-training, our next step is to convert the learned multitask skill policy into a task-specific policy that can be fine-tuned to maximize a given downstream reward function r (Figure 1 (bottom middle)). There exist two challenges in this step: (1) we need a mechanism to identify the skill vector z that best solves the given task and (2) we need to reconcile the gap between intrinsic rewards and downstream task rewards for seamless online fine-tuning.

242 Skill identification. Since we chose to use a successor feature- or goal-based unsupervised pre-243 training method in the previous phase, the first challenge is relatively straightforward. For goal-244 oriented tasks (e.g., AntMaze (Fu et al., 2020) and Kitchen (Gupta et al., 2020)), we assume the 245 task goal g to be available, and we either directly use g (for goal-conditioned methods) or infer the 246 skill  $z^*$  that corresponds to g based on a predefined conversion formula (for successor feature-based methods that support such a conversion (Touati et al., 2022; Park et al., 2024c)). For general reward-247 maximization tasks, we employ successor feature-based unsupervised pre-training methods, and use 248 the following linear regression to find the skill latent vector  $z^*$  that best approximates the downstream 249 task reward function  $r : S \times A \times S \rightarrow \mathbb{R}$  (Touati et al., 2022; Park et al., 2024c): 250

$$z^* = \operatorname*{arg\,min}_{z \in \mathcal{Z}} \mathbb{E}_{(s,a,s') \sim \mathcal{D}_{\text{reward}}} \left[ \left( r(s,a,s') - f(s,a,s')^\top z \right)^2 \right],\tag{6}$$

where f is the feature network in the successor feature framework (Section 3) and  $\mathcal{D}_{reward}$  is a reward-labeled dataset. This reward-labeled dataset  $\mathcal{D}_{reward}$  can be either the full offline dataset  $\mathcal{D}_{off}$ (if it is fully reward-labeled), a subset of the offline dataset (if it is partially reward-labeled), or a newly collected dataset with additional environment interactions. In our experiments, we mainly use a small number (*e.g.*, 0.2% for DMC tasks) of reward-labeled samples from the offline dataset for  $\mathcal{D}_{reward}$ , following previous works (Touati et al., 2022; Park et al., 2024c), but we do not require  $\mathcal{D}_{reward}$  to be a subset of  $\mathcal{D}_{off}$  (see Appendix A.4).

260 **Reward scale matching.** After identifying the best skill latent vector  $z^*$ , our next step is to bridge 261 the gap between intrinsic and extrinsic rewards. Since these two reward functions can have very 262 different scales, naïve online adaptation can lead to abrupt shifts in target Q-values, potentially 263 causing significant performance drops in the early stages of online fine-tuning. While one can employ 264 sophisticated reward-shaping techniques to deal with this issue (Ng et al., 1999; Gleave et al., 2021), 265 in this work, we propose a simple yet effective reward scale matching technique that we find to be 266 effective enough in practice. Specifically, we compute the running mean and standard deviation of intrinsic rewards during the pre-training phase, and normalize the intrinsic rewards with the calculated 267 statistics. Similarly, during the fine-tuning phase, we compute the statistics of task rewards and 268 normalize the task rewards so that they have the same scale and mean as normalized intrinsic rewards. 269 This way, we can prevent abrupt shifts in reward scales without altering the optimal policy for the



Figure 2: Environments. We evaluate U2O RL on eleven state-based or pixel-based environments.

downstream task. In our experiments, we find that this simple technique is crucial for achieving good performance, especially in environments with dense rewards (Q6 in Section A.8).

### 280 4.3 **ONLINE FINE-TUNING**

Our final step is to fine-tune the skill policy with online environment interactions (Figure 1 (bottom 282 right)). This step is straightforward: since we have found  $z^*$  in the previous stage, we can simply 283 fix the skill vector  $z^*$  in the policy  $\pi_{\theta}(a \mid s, z^*)$  and the Q-function  $Q_{\phi}(s, a, z^*)$ , and fine-tune them with the same (offline) RL algorithm used in the first phase (e.g., IQL (Kostrikov et al., 2022), 285 TD3 (Fujimoto et al., 2018)) with additional online interaction data. While one can employ existing 286 specialized techniques for offline-to-online RL for better online adaptation in this phase, we find in our experiments that, thanks to rich representations learned by unsupervised pre-training, simply 288 using the same (offline) RL algorithm is enough to achieve strong performance that matches or even 289 outperforms state-of-the-art offline-to-online RL techniques.

290 291

292

275 276 277

278

279

281

284

287

#### WHY IS U2O RL POTENTIALLY BETTER THAN OFFLINE-TO-ONLINE RL? 44

293 Our main claim is that *unsupervised* offline RL is better than supervised offline RL for online finetuning. However, this might sound counterintuitive. Especially, if we know the downstream task 294 ahead of time, how can unsupervised offline RL potentially lead to better performance than supervised 295 offline RL, despite the fact that the former does not use any task information during the offline phase? 296

297 We hypothesize that this is because unsupervised multi-task offline RL enables better *feature learning* 298 than supervised single-task offline RL. By training the agent on a number of diverse intrinsically 299 defined tasks, it gets to acquire rich knowledge about the environment, dynamics, and potential tasks in the form of representations, which helps improve and facilitate the ensuing task-specific 300 online fine-tuning. This resembles the recent observation in machine learning that large-scale 301 unsupervised pre-training improves downstream task performances over task-specific supervised 302 pre-training (Brown et al., 2020; Devlin et al., 2019; Radford et al., 2019; He et al., 2021; 2020; 303 Hénaff et al., 2020). In our experiments, we empirically show that U2O RL indeed learns better 304 features than its supervised counterpart (Q4 in Section 5). 305

Another benefit of U2O RL is that it does not use any task-specific information during pre-training. 306 This is appealing because we can reuse a single pre-trained policy for a number of different down-307 stream tasks. Moreover, it enables leveraging potentially large, task-agnostic offline data during 308 pre-training, which is often cheaper to collect than task-specific, curated datasets (Lynch et al., 2019). 309

310

#### 311 5 **EXPERIMENTS**

312

313 In our experiments, we evaluate the performance of U2O RL in the context of offline-to-online RL. We 314 aim to answer the following research questions: (1) Is U2O RL better than previous offline-to-online 315 RL strategies? (2) Can a single pre-trained model from U2O be fine-tuned to solve multiple tasks? (3) What makes unsupervised offline RL result in better fine-tuning performance than supervised 316 offline RL? (4) Which components of U2O RL are important? 317

318 Environments and offline datasets. In our experiments, we consider eleven tasks across five 319 benchmarks (Figure 2). ExORL (Yarats et al., 2022) is a benchmark suite that consists of offline 320 datasets collected by exploratory policies (e.g., RND (Burda et al., 2019)) on the DeepMind Control 321 Suite (Tassa et al., 2018). We consider four embodiments (Walker, Cheetah, Quadruped, and Jaco), each of which has four tasks. AntMaze (Fu et al., 2020; Jiang et al., 2023) is a navigation task, 322 whose goal is to control an 8-DoF quadruped agent to reach a target position. We consider the two 323 most challenging mazes with the largest sizes, large and ultra, and two types of offline datasets,



Figure 3: Online fine-tuning plots of U2O RL and previous offline-to-online RL frameworks (8 seeds). Across the benchmarks, our U2O RL mostly shows consistently better performance than standard offline-toonline RL and off-policy online RL with offline data.

357 diverse and play. Kitchen (Gupta et al., 2020; Fu et al., 2020) is a robotic manipulation task, 358 where the goal is to control a 9-DoF Franka robot arm to achieve four subtasks sequentially. We 359 consider two types of offline datasets from the D4RL suite (Fu et al., 2020), partial and mixed. 360 Visual Kitchen (Gupta et al., 2020; Fu et al., 2020; Park et al., 2024c) is a pixel-based variant of 361 the Kitchen environment, where an agent must achieve four subtasks purely from  $64 \times 64 \times 3$  pixel 362 observations instead of low-dimensional state information. Adroit (Fu et al., 2020) is a dexterous 363 manipulation benchmark, where the goal is to control a 24-DoF robot hand to twirl a pen or open a door. OGBench-Cube (Park et al., 2024a) is an additional manipulation benchmark whose goal is to 364 control a 6-DoF UR5e robot arm to perform pick-and-place manipulation of multiple cubes from an 365 unlabeled, diverse dataset. We use a single-task version of OGBench-Cube to make it compatible 366 with our offline-to-online RL setting. We provide further details in Appendix C.1. 367

368 Implementation. In our experiments, we mainly employ HILP (Park et al., 2024c) as the unsupervised offline policy pre-training algorithm in U2O RL. For the offline RL backbone, we use 369 TD3 (Fujimoto et al., 2018) for ExORL and IQL (Kostrikov et al., 2021) for others following previous 370 works (Touati et al., 2022; Park et al., 2024c). Since both IQL and TD3+BC (Fujimoto & Gu, 2021; 371 Tarasov et al., 2023a) have been known to achieve strong performance in the offline-to-online RL 372 setting (Tarasov et al., 2023b), we use them for the online fine-tuning phase in U2O RL as well. For 373 sparse-reward tasks (AntMaze, Kitchen, and Adroit), we do not apply reward scale matching. For 374 AntMaze, Kitchen, and Adroit, we report normalized scores, following Fu et al. (2020). In our experi-375 ments, we use 8 random seeds and report standard deviations with shaded areas, unless otherwise 376 stated. We refer the reader to Appendix C for the full implementation details and hyperparameters.

377

354

355

356

Q1. Is U2O RL better than previous offline-to-online RL frameworks?

Table 1: Comparison between U2O RL and previous offline-to-online RL methods. We denote how performances change before and after online fine-tuning with arrows. Baseline scores except RLPD (Ball et al., 2023) are taken from Nakamoto et al. (2023); Wang et al. (2023a). Scores within the 5% of the best score are highlighted in bold, as in Kostrikov et al. (2022). We use 8 random seeds for each task for U2O RL.

202							
302	Task	antmaze-ultra-diverse	antmaze-ultra-play	antmaze-large-diverse	antmaze-large-play	kitchen-partial	kitchen-mixed
383	CQL	-	-	$25 \rightarrow 87$	$34 \rightarrow 76$	71 → <b>75</b>	$56 \rightarrow 50$
384	IQL	$13 \rightarrow 29$	$17 \rightarrow 29$	$40 \rightarrow 59$	$41 \rightarrow 51$	$40 \rightarrow 60$	$48 \rightarrow 48$
50-	AWAC	-	-	$00 \rightarrow 00$	$00 \rightarrow 00$	$01 \rightarrow 13$	$02 \rightarrow 12$
385	O3F	-	-	$59 \rightarrow 28$	$68 \rightarrow 01$	$11 \rightarrow 22$	$06 \rightarrow 33$
	ODT	-	-	$00 \rightarrow 01$	$00 \rightarrow 00$	-	-
386	CQL+SAC	-	-	$36 \rightarrow 00$	$21 \rightarrow 00$	$71 \rightarrow 00$	$59 \rightarrow 01$
007	Hybrid RL	-	-	$\rightarrow 00$	$\rightarrow 00$	$\rightarrow 00$	$\rightarrow 01$
387	SAC+od	-	-	$\rightarrow 00$	$\rightarrow 00$	$\rightarrow 07$	$\rightarrow 00$
200	SAC	-	-	$\rightarrow 00$	$\rightarrow 00$	$\rightarrow 03$	$\rightarrow 02$
300	IQL+od	$\rightarrow 04$	$\rightarrow 05$	$\rightarrow 71$	$\rightarrow 56$	$\rightarrow 74$	$\rightarrow 61$
389	FamO2O	-	-	$\rightarrow 64$	$\rightarrow 61$	-	-
000	RLPD	$00 \rightarrow 00$	$00 \rightarrow 00$	$00 \rightarrow 94$	$00 \rightarrow 81$	-	-
390	Cal-QL	$05 \rightarrow 05$	$15 \rightarrow 13$	$33 \rightarrow 95$	$26 \rightarrow 90$	$67 \rightarrow 79$	$38 \rightarrow 80$
201	U2O (Ours)	$22 \rightarrow 54$	$17 \rightarrow 58$	$11 \rightarrow 95$	$14 \rightarrow 88$	48  ightarrow 75	$48 \rightarrow 74$

<sup>391</sup> 392

We begin our experiments by comparing our approach, unsupervised-to-online RL, with two previous 394 offline-to-online RL frameworks (Section 2): offline-to-online RL (O2O RL) and off-policy online 395 RL with offline data (Online w/ Off Data). To recall, offline-to-online RL (Lee et al., 2022; Nair 396 et al., 2020; Nakamoto et al., 2023; Yu & Zhang, 2023; Lei et al., 2023) first pre-trains a policy 397 with supervised offline RL using the task reward, and then continues training it with online rollouts. 398 Off-policy online RL (Ball et al., 2023; Luo et al., 2024; Song et al., 2023) trains a policy from scratch 399 on top of a replay buffer filled with offline data. Here, we use the same offline RL backbone (i.e., 400 TD3 for ExORL and IOL for AntMaze, Kitchen, and Adroit) to ensure apples-to-apples comparisons 401 between the three frameworks. We will compare U2O RL with previous specialized offline-to-online 402 RL techniques in Q2 of Section 5.

Figure 3 shows the online fine-tuning curves on 14 different tasks. The results suggest that U2O
RL generally leads to better performance than both offline-to-online RL and off-policy online RL
across the environments, despite not using any task information during pre-training. Notably, U2O
RL significantly outperforms these two previous frameworks in the most challenging AntMaze tasks
(antmaze-ultra-{diverse, play}).

### 408 409 Q2. How does U2O RL compare to previous specialized offline-to-online RL techniques?

410 Next, we compare U2O RL with 13 previous specialized offline-to-online RL methods, including CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2022), AWAC (Nair et al., 2020), O3F (Mark 411 et al., 2022), ODT (Zheng et al., 2022), CQL+SAC (Kumar et al., 2020; Haarnoja et al., 2018), 412 Hybrid RL (Song et al., 2023), SAC+od (offline data) (Haarnoja et al., 2018; Ball et al., 2023), 413 SAC (Haarnoja et al., 2018), IQL+od (offline data) (Kostrikov et al., 2022; Ball et al., 2023), 414 FamO2O (Wang et al., 2023a), RLPD (Ball et al., 2023), and Cal-QL (Nakamoto et al., 2023). We 415 show the comparison results in Table 1, where we take the scores from Nakamoto et al. (2023); Wang 416 et al. (2023a) for the tasks that are common to ours. Since Cal-QL achieves the best performance in 417 the table, we additionally make a comparison with Cal-QL on antmaze-ultra-{diverse, play} 418 as well, by running their official implementation with tuned hyperparameters. 419

Table 1 shows that U2O RL achieves strong performance that matches or sometimes even outperforms previous offline-to-online RL methods, even though U2O RL does *not* use any task information during offline pre-training nor any specialized offline-to-online techniques. In particular, in the most challenging antmaze-ultra tasks, U2O RL outperforms the previous best method (Cal-QL) by a significant margin. This is very promising because, even if U2O RL does not necessarily outperform the state-of-the-art methods on every single task (though it is at least on par with the previous best methods), U2O RL enables reusing a single unsupervised pre-trained policy for multiple downstream tasks, unlike previous offline-to-online RL methods that perform *task-specific* pretraining.

427 428

# Q3. Can a single pre-trained model from U2O be fine-tuned to solve multiple tasks?

One important advantage of U2O RL is that it can reuse a single task-agnostic dataset for multiple
 different downstream tasks, unlike standard offline-to-online RL. To demonstrate this, we train U2O
 RL with four different tasks from the same task-agnostic ExORL dataset on each DMC environment, and report the full training curves in Figure 7 of Appendix A.1. The results show that, for example,

437

438

439

440

441

442

443

444

445

446

447

a single pre-trained model on the Walker domain can be fine-tuned for all four tasks (Walker Run,
Walker Flip, Walker Stand, and Walker Walk). Note that even though U2O RL uses a single taskagnostic pre-trained model, the performance of U2O RL matches or even outperforms O2O RL,
which pre-trains a model with task-specific rewards.

Q4. Why does U2O RL often outperform supervised offline-to-online RL?



Figure 4: Feature dot products during offline RL pre-training (lower is better, 8 seeds). The plots show that *unsupervised* offline pre-training effectively prevents feature collapse (co-adaptation), yielding better representations than supervised offline pre-training.

448 In the above experiments, we showed that U2O RL often even outperforms previous supervised 449 offline-to-online RL methods. We hypothesized in Section 4.4 that this is because unsupervised 450 offline pre-training yields better *representations* that facilitate online task adaptation. To empirically 451 verify this hypothesis, we measure the quality of the value function representations using the method 452 proposed by Kumar et al. (2022). Specifically, we define the value features  $\zeta_{\phi}(s, a)$  as the penultimate 453 layer of the value function  $Q_{\phi}$ , *i.e.*,  $Q_{\phi}(s, a) = w_{\phi}^{+} \zeta_{\phi}(s, a)$ , and measure the dot product between 454 consecutive state-action pairs,  $\zeta_{\phi}(s, a)^{\top}\zeta_{\phi}(s', a')$  (Kumar et al., 2022). Intuitively, this dot product 455 represents the degree to which these two representations are "collapsed" (or "co-adapted"), which is 456 known to be correlated with poor performance (Kumar et al., 2022) (*i.e.*, the lower the better). 457

Figure 4 compares the dot product metrics of unsupervised offline RL (in U2O RL) and supervised offline RL (in O2O RL) on four benchmark tasks. The results suggest that our unsupervised multi-task pre-training effectively prevents feature co-adaptation and thus indeed yields better representations across the environments. This highlights the benefits of unsupervised offline pre-training, and (partially) explains the strong online fine-tuning performance of U2O RL. We additionally provide further analyses with different offline unsupervised RL algorithms (*e.g.*, graph Laplacian-based successor feature learning (Touati et al., 2022; Wu et al., 2019b)) in Appendix A.2.

### 464 465 Q5. Is fine-tuning better than other alternative strategies (*e.g.*, hierarchical RL)?

In this work, we focus on the *fine-tuning* of of-466 fline pre-trained skill policies, but this is not 467 the only way to leverage pre-trained skills for 468 downstream tasks. To see how our fine-tuning 469 scheme compares to other alternative strategies, 470 we compare U2O RL with three previously con-471 sidered approaches: hierarchical RL (HRL, 472 e.g., OPAL (Ajay et al., 2021), SPiRL (Pertsch 473 et al., 2021)) (Ajay et al., 2021; Pertsch et al., 474 2021; Touati et al., 2022; Park et al., 2024c; Hu 475 et al., 2023), zero-shot RL (Touati et al., 2022; 476 Park et al., 2024c), and **PEX** (Zhang et al., 2023). 477 HRL additionally trains a high-level policy that



Figure 5: Fine-tuning is better than previous strategies, such as hierarchical RL, zero-shot RL, and PEX (8 seeds).

combines fixed pre-trained skills in a sequential manner. Zero-shot RL simply finds the skill policy
that best solves the downstream task, with no fine-tuning or hierarchies. PEX combines fixed pretrained multi-task policies and a newly initialized policy with a multiplexer that chooses the best
policy.

Figure 5 shows the comparison results on top of the same pre-trained unsupervised skill policy. Since
PEX is not directly compatible with IQL, we evaluate PEX only on the tasks with TD3 (*e.g.*, ExORL
tasks). The plots suggest that our fine-tuning strategy leads to significantly better performance than
previous approaches. This is because pre-trained offline skill policies are often not perfect (due to
the limited coverage or suboptimality of the dataset), and thus using a fixed offline policy is often



# not sufficient to achieve strong performance in downstream tasks. We provide additional results in Appendix A.6.

Figure 6: Online RL learning curves with expert-only datasets (4 seeds).

# Q6. Negative results: When is U2O RL better than O2O RL?

501 While we showed strong results of U2O RL throughout the paper, U2O RL is not *always* better 502 than O2O RL. Specifically, U2O RL may not be as effective when the dataset is monolithic (e.g., consists of only expert trajectories, has less diversity, etc.). To empirically show this, we conduct an 504 additional experiment with a different dataset in antmaze-large and antmaze-ultra that consists of monolithic, expert trajectories (we collect a 1M-sized dataset by rolling out an offline pre-trained 505 policy) as well as kitchen-complete dataset, which also consists of expert trajectories. Figure 6 506 shows the (negative) results, which suggest that U2O RL is not particularly better than O2O RL 507 on these monolithic, optimal datasets. However, we emphasize that U2O RL achieves similar final 508 performance to O2O RL even on these datasets, and has the unique strength that a single unsupervised 509 pre-trained model can be fine-tuned to many different reward functions, unlike standard O2O RL. 510

We refer to the reader to Appendix for further analysis including (1) combining U2O RL with other offline unsupervised skill learning methods (Appendix A.2), (2) comparisons between U2O RL and pure representation learning schemes (Appendix A.3), (3) U2O RL without reward samples in the bridging phase (Appendix A.4), (4) an ablation study with different skill identification strategies (Appendix A.5), (5) additional results with different online RL strategies (Appendix A.6), (6) comparison to O2O RL combined with methods for mitigating the feature collapse issue (Appendix A.7), and (7) ablation studies on each component in U2O RL such as reward scale matching, value transfer and policy transfer (Appendix A.8).

518 519 520

521

486

487

488 489

490

491

492

493

494

495 496

497

498 499

500

# 6 CONCLUSION

In this work, we investigated how unsupervised pre-training of diverse policies enables better online
 fine-tuning than standard supervised offline-to-online RL. We showed that our unsupervised-to online recipe often achieves even better performance and stability than previous offline-to-online
 RL approaches, thanks to the rich representations learned by pre-training on diverse tasks. We also
 demonstrated that U2O RL enables reusing a single offline pre-trained policy for multiple downstream
 tasks.

Limitation. As shown in Q7 of Section 5, U2O RL is not necessarily better than O2O RL when the offline dataset is monolithic and heavily tailored toward the downstream task. We believe U2O RL is most effective (compared to standard offline-to-online RL) when the dataset is highly diverse so that the unsupervised offline RL method can learn a variety of behaviors and thus learn better features and representations. Given the recent successes in large-scale self-supervised and unsupervised pre-training from unlabeled data, we believe U2O RL serves as a step toward a general recipe for scalable data-driven decision-making.

534 535

# REPRODUCIBILITY STATEMENT

536 537

538 We provide implementation details in Section 5 and Appendix C including hyperparameters. We 539 also provide pseudo-code in Appendix B and have attached our source code to the OpenReview submission page.

### 540 REFERENCES 541

547

554

567

577

- Anurag Ajay, Aviral Kumar, Pulkit Agrawal, Sergey Levine, and Ofir Nachum. Opal: Offline 542 primitive discovery for accelerating offline reinforcement learning. In International Conference on 543 Learning Representations, 2021. 544
- Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient online reinforcement learning 546 with offline data. In International Conference on Machine Learning, 2023.
- 548 André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. In Advances in Neural 549 Information Processing Systems, 2017. 550
- 551 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 552 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are 553 few-shot learners. In Advances in Neural Information Processing Systems, 2020.
- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network 555 distillation. In International Conference on Learning Representations, 2019. 556
- Yevgen Chebotar, Karol Hausman, Yao Lu, Ted Xiao, Dmitry Kalashnikov, Jacob Varley, Alex Irpan, 558 Benjamin Eysenbach, Ryan C Julian, Chelsea Finn, et al. Actionable models: Unsupervised offline 559 reinforcement learning of robotic skills. In International Conference on Machine Learning, 2021. 560
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. 561 Neural computation, 5:613-624, 1993. 562
- 563 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of* 565 the North American Chapter of the Association for Computational Linguistics: Human Language 566 Technologies, 2019.
- Ben Eysenbach, Russ R Salakhutdinov, and Sergey Levine. Search on the replay buffer: Bridging 568 planning and reinforcement learning. In Advances in Neural Information Processing Systems, 569 2019a. 570
- 571 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you 572 need: Learning skills without a reward function. In International Conference on Learning 573 Representations, 2019b.
- 574 Benjamin Eysenbach, Tianjun Zhang, Sergey Levine, and Russ R Salakhutdinov. Contrastive learning 575 as goal-conditioned reinforcement learning. In Advances in Neural Information Processing Systems, 576 2022.
- 578 Kuan Fang, Patrick Yin, Ashvin Nair, and Sergey Levine. Planning to practice: Efficient online 579 fine-tuning by composing goals in latent space. In International Conference on Intelligent Robots 580 and Systems, 2022.
- 581 Kuan Fang, Patrick Yin, Ashvin Nair, Homer Rich Walke, Gengchen Yan, and Sergey Levine. 582 Generalization with lossy affordances: Leveraging broad offline data for learning visuomotor tasks. 583 In Conference on Robot Learning, 2023. 584
- 585 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep 586 data-driven reinforcement learning. arXiv preprint arXiv:2004.07219, 2020.
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. In 588 Advances in Neural Information Processing Systems, 2021. 589
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-591 critic methods. In International Conference on Machine Learning, 2018. 592
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In International Conference on Machine Learning, 2019.

595 data via latent intentions. In International Conference on Machine Learning, 2023. 596 Adam Gleave, Michael Dennis, Shane Legg, Stuart Russell, and Jan Leike. Quantifying differences 597 in reward functions. In International Conference on Learning Representations, 2021. 598 Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. arXiv 600 preprint arXiv:1611.07507, 2016. 601 Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy 602 learning: Solving long-horizon tasks via imitation and reinforcement learning. In Conference on 603 Robot Learning, 2020. 604 605 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy 606 maximum entropy deep reinforcement learning with a stochastic actor. In International Conference on Machine Learning, 2018. 607 608 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine. 609 Idql: Implicit q-learning as an actor-critic method with diffusion policies. arXiv preprint 610 arXiv:2304.10573, 2023. 611 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for 612 unsupervised visual representation learning. In Proceedings of the IEEE/CVF Conference on 613 Computer Vision and Pattern Recognition, 2020. 614 615 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked 616 autoencoders are scalable vision learners. arXiv preprint arXiv:2111.06377, 2021. 617 Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Ali Eslami, and 618 Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. In 619 International Conference on Machine Learning, 2020. 620 621 Hao Hu, Yiqin Yang, Jianing Ye, Ziqing Mai, and Chongjie Zhang. Unsupervised behavior extraction 622 via random intent priors. In Advances in Neural Information Processing Systems, 2023. 623 Zhengyao Jiang, Tianjun Zhang, Michael Janner, Yueying Li, Tim Rocktäschel, Edward Grefenstette, 624 and Yuandong Tian. Efficient planning in a compact latent action space. In International Conference 625 on Learning Representations, 2023. 626 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In International 627 Conference on Learning Representations, 2015. 628 629 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In International Conference 630 on Learning Representations, 2014. 631 Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement learning 632 with fisher divergence critic regularization. In International Conference on Machine Learning, 633 2021. 634 635 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit 636 q-learning. In International Conference on Learning Representations, 2022. 637 Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy 638 q-learning via bootstrapping error reduction. In Advances in Neural Information Processing 639 Systems, 2019. 640 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline 641 reinforcement learning. Advances in Neural Information Processing Systems, 2020. 642 643 Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron Courville, George Tucker, and Sergey Levine. 644 DR3: Value-based deep reinforcement learning requires explicit regularization. In International 645 Conference on Learning Representations, 2022. 646

Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement learning from passive

647 Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In *Reinforcement learning: State-of-the-art*, pp. 45–73. Springer, 2012.

648 649 650	Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel Pinto, and Pieter Abbeel. Urlb: Unsupervised reinforcement learning benchmark. In Advances in Neural Information Processing Systems Datasets and Benchmarks Track, 2021.
652 653 654	Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In <i>Conference on Robot Learning</i> , 2022.
655 656 657	Kun Lei, Zhengmao He, Chenhao Lu, Kaizhe Hu, Yang Gao, and Huazhe Xu. Uni-o4: Unifying online and offline deep reinforcement learning with multi-step on-policy optimization. <i>arXiv</i> preprint arXiv:2311.03351, 2023.
658 659 660	Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. <i>arXiv preprint arXiv:2005.01643</i> , 2020.
661 662	Qiyang Li, Jason Zhang, Dibya Ghosh, Amy Zhang, and Sergey Levine. Accelerating exploration with unlabeled prior data. In <i>Advances in Neural Information Processing Systems</i> , 2023.
663 664 665	Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. <i>arXiv preprint arXiv:2103.04551</i> , 2021.
666 667 668	Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal, Chelsea Finn, Abhishek Gupta, and Sergey Levine. Serl: A software suite for sample-efficient robotic reinforcement learning. <i>arXiv preprint arXiv:2401.16013</i> , 2024.
669 670 671	Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In <i>Conference on Robot Learning</i> , 2019.
672 673 674 675	Jason Yecheng Ma, Jason Yan, Dinesh Jayaraman, and Osbert Bastani. Offline goal-conditioned reinforcement learning via <i>f</i> -advantage regression. In <i>Advances in Neural Information Processing Systems</i> , 2022.
676 677 678	Yecheng Jason Ma, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. In <i>International Conference on Machine Learning</i> , 2023a.
679 680 681	Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training. In <i>International Conference on Learning Representations</i> , 2023b.
682 683 684	Max Sobol Mark, Ali Ghadirzadeh, Xi Chen, and Chelsea Finn. Fine-tuning offline policies with optimistic action selection. In <i>Deep Reinforcement Learning Workshop NeurIPS</i> 2022, 2022.
685 686 687 688	Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Discovering and achieving goals via world models. In <i>Advances in Neural Information Processing Systems</i> , 2021.
689 690 691	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. <i>arXiv preprint arXiv:1312.5602</i> , 2013.
692 693 694	Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. <i>arXiv preprint arXiv:2006.09359</i> , 2020.
695 696	Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. <i>arXiv preprint arXiv:2203.12601</i> , 2022.
697 698 699 700	Mitsuhiko Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. In <i>Advances in Neural Information Processing Systems</i> , 2023.
700	Soroush Nasiriany, Vitchyr Pong, Steven Lin, and Sergey Levine. Planning with goal-conditioned policies. In Advances in Neural Information Processing Systems, 2019.

702 703 704	Whitney K Newey and James L Powell. Asymmetric least squares estimation and testing. <i>Econometrica: Journal of the Econometric Society</i> , pp. 819–847, 1987.
705 706 707	Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In <i>International Conference on Machine Learning</i> , 1999.
708 709 710	Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsurprising effectiveness of pre-trained vision models for control. In <i>International Conference on Machine Learning</i> , 2022.
711 712 713	Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking offline goal-conditioned rl. <i>arXiv preprint arXiv:2410.20092</i> , 2024a.
714 715 716	Seohong Park, Dibya Ghosh, Benjamin Eysenbach, and Sergey Levine. Hiql: Offline goal-conditioned rl with latent states as actions. In <i>Advances in Neural Information Processing Systems</i> , 2024b.
717 718	Seohong Park, Tobias Kreiman, and Sergey Levine. Foundation policies with hilbert representations. <i>arXiv preprint arXiv:2402.15567</i> , 2024c.
719 720 721	Seohong Park, Oleh Rybkin, and Sergey Levine. METRA: Scalable unsupervised RL with metric- aware abstraction. In <i>International Conference on Learning Representations</i> , 2024d.
721 722 723	Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In <i>International Conference on Machine Learning</i> , 2017.
724 725	Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement. In <i>International Conference on Machine Learning</i> , 2019.
720 727 728	Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. <i>arXiv preprint arXiv:1910.00177</i> , 2019.
729 730 731	Karl Pertsch, Youngwoon Lee, and Joseph Lim. Accelerating reinforcement learning with learned skill priors. In <i>Conference on Robot Learning</i> , 2021.
732 733	Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. <i>SIAM journal on control and optimization</i> , 1992.
734 735 736	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 2019.
737 738	Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter Abbeel. Masked world models for visual control. In <i>Conference on Robot Learning</i> , 2022a.
739 740 741	Younggyo Seo, Kimin Lee, Stephen James, and Pieter Abbeel. Reinforcement learning with action- free pre-training from videos. In <i>International Conference on Machine Learning</i> , 2022b.
742 743 744	Younggyo Seo, Junsu Kim, Stephen James, Kimin Lee, Jinwoo Shin, and Pieter Abbeel. Multi-view masked world models for visual robotic manipulation. In <i>International Conference on Machine Learning</i> , 2023.
745 746 747 748	Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, and Sergey Levine. Time-contrastive networks: Self-supervised learning from video. In <i>International Conference on Robotics and Automation</i> , 2018.
749 750	Rutav Shah and Vikash Kumar. Rrl: Resnet as representation for reinforcement learning. In <i>International Conference on Machine Learning</i> , 2021.
751 752 753	Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In <i>International Conference on Learning Representations</i> , 2020.
754 755	Avi Singh, Huihan Liu, Gaoyue Zhou, Albert Yu, Nicholas Rhinehart, and Sergey Levine. Parrot: Data-driven behavioral priors for reinforcement learning. In <i>International Conference on Learning Representations</i> , 2021.

756 757 758	Yuda Song, Yifei Zhou, Ayush Sekhari, J Andrew Bagnell, Akshay Krishnamurthy, and Wen Sun. Hybrid RL: Using both offline and online data can make RL efficient. In <i>International Conference</i> <i>on Learning Representations</i> , 2023.
759 760	Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT Press, 2018.
761 762 763	Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the mini- malist approach to offline reinforcement learning. In <i>Advances in Neural Information Processing</i> <i>Systems</i> , 2023a.
764 765 766 767	Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov. Corl: Research-oriented deep offline reinforcement learning library. In <i>Advances in Neural</i> <i>Information Processing Systems</i> , 2023b.
768 769 770	Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. <i>arXiv preprint arXiv:1801.00690</i> , 2018.
771 772	Ahmed Touati, Jérémy Rapin, and Yann Ollivier. Does zero-shot reinforcement learning exist? In <i>International Conference on Learning Representations</i> , 2022.
773 774 775 776	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , 2017.
777 778 779	Shenzhi Wang, Qisen Yang, Jiawei Gao, Matthieu Lin, Hao Chen, Liwei Wu, Ning Jia, Shiji Song, and Gao Huang. Train once, get a family: State-adaptive balances for offline-to-online reinforcement learning. In <i>Advances in Neural Information Processing Systems</i> , 2023a.
780 781 782	Tongzhou Wang, Antonio Torralba, Phillip Isola, and Amy Zhang. Optimal goal-reaching rein- forcement learning via quasimetric learning. In <i>International Conference on Machine Learning</i> , 2023b.
783 784 785	Ziyu Wang, Alexander Novikov, Konrad Zolna, Josh S Merel, Jost Tobias Springenberg, Scott E Reed, Bobak Shahriari, Noah Siegel, Caglar Gulcehre, Nicolas Heess, et al. Critic regularized regression. In <i>Advances in Neural Information Processing Systems</i> , 2020.
786 787 788	Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. arXiv preprint arXiv:1911.11361, 2019a.
789 790	Yifan Wu, George Tucker, and Ofir Nachum. The laplacian in RL: Learning representations with efficient approximations. In <i>International Conference on Learning Representations</i> , 2019b.
791 792 793	Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for motor control. <i>arXiv preprint arXiv:2203.06173</i> , 2022.
794 795 796	Rui Yang, Lin Yong, Xiaoteng Ma, Hao Hu, Chongjie Zhang, and Tong Zhang. What is essential for unseen goal generalization of offline goal-conditioned rl? In <i>International Conference on Machine Learning</i> , 2023.
797 798 799	Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro Lazaric, and Lerrel Pinto. Don't change the algorithm, change the data: Exploratory data for offline reinforcement learning. <i>arXiv preprint arXiv:2201.13425</i> , 2022.
800 801 802	Zishun Yu and Xinhua Zhang. Actor-critic alignment for offline-to-online reinforcement learning. In <i>International Conference on Machine Learning</i> , 2023.
803 804	Haichao Zhang, We Xu, and Haonan Yu. Policy expansion for bridging offline-to-online reinforcement learning. In <i>International Conference on Learning Representations</i> , 2023.
805 806 807 808	Kai Zhao, Yi Ma, Jinyi Liu, HAO Jianye, Yan Zheng, and Zhaopeng Meng. Improving offline-to- online reinforcement learning with q-ensembles. In <i>ICML Workshop on New Frontiers in Learning,</i> <i>Control, and Dynamical Systems</i> , 2023.
809	Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In <i>International Conference on Machine Learning</i> , 2022.



Figure 7: Learning curves during online RL fine-tuning (8 seeds). A single pre-trained model from U2O can be fine-tuned to solve multiple downstream tasks. Across the embodiments and tasks, our U2O RL matches or outperforms standard offline-to-online RL and off-policy online RL with offline data even though U2O RL uses a single task-agnostic pre-trained model.





Figure 8: U2O RL with Laplacian-based successor Figure feature learning (8 seeds).

Figure 9: U2O RL with goal-conditioned IQL (8 seeds).

While we employ HILP (Park et al., 2024c) as an offline unsupervised skill learning method in U2O RL in our main experiments, our recipe can be combined with other offline unsupervised skill

learning methods as well. To show this, we replace HILP with a graph Laplacian-based successor
feature method (Touati et al., 2022; Wu et al., 2019b) or goal-conditioned IQL (GC-IQL) (Kostrikov
et al., 2022; Park et al., 2024b), and report the results in Figures 8 and 9, respectively. The results
demonstrate that U2O RL with different unsupervised RL methods also improves performance over
standard offline-to-online RL.

Additionally, we show that other unsupervised 870 skill learning methods also lead to better value 871 representations. We measure the same fea-872 ture dot product metric in Section 5 with the 873 graph Laplacian-based successor feature learn-874 ing method and report the results in Figure 10. The results suggest that this unsupervised RL 875 method also prevents feature co-adaptation, lead-876 ing to better features. 877



Figure 10: Feature dot product analysis with Laplacian-based successor feature learning (8 seeds).

# 878 879

880

088

A.3 Do we need to use unsupervised RL for pre-training representations?

In Sections 4.4 and 5, we hypothesized and empirically showed that U2O RL is often better than O2O RL because it learns better representations. This leads to the following natural question: do we *need* to use offline unsupervised *reinforcement learning*, as opposed to general representation learning? To answer this question, we consider two pure representation learn-

# Table 2: Comparison between U2O RL and pure representation learning algorithms (4 seeds).

Task	antmaze-large-diverse
U2O (HILP, Q Ours)	$\textbf{94.50} \pm \textbf{3.16}$
U2O (HILP, $\xi$ )	$5.50 \pm 1.91$
Temporal contrastive learning	$37.50 \pm 15.00$

ing algorithms as alternatives to unsupervised RL: temporal contrastive learning (Eysenbach et al., 890 2022) and Hilbert (metric) representation learning (Park et al., 2024c), where the latter is equivalent 891 to directly taking  $\xi$  in the HILP framework (Equation 4) (note that the original U2O RL takes the Q 892 function of HILP, not the Hilbert representation  $\xi$  itself, which is used to train the Q function). To 893 evaluate their fine-tuning performances, for the temporal contrastive representation, we fine-tune both 894 the Q function and policy with contrastive RL (Eysenbach et al., 2022); for the Hilbert representation, 895 we take the pre-trained representation, add one new layer, and use it as the initialization of the 896 Q function. Table 2 shows the results on antmaze-large-diverse. Somewhat intriguingly, the 897 results suggest that it is important to use the full unsupervised RL procedure, and pure representation learning methods result in much worse performance in this case. This becomes more evident if we 899 compare U2O RL (HILP Q, ours) and U2O RL (HILP  $\xi$ ), given that they are rooted in the same Hilbert representation. We believe this is because, if we simply use an off-the-shelf representation 900 learning, there exists a discrepancy in training objectives between pre-training (e.g., metric learning) 901 and fine-tuning (Q-learning). On the other hand, in U2O RL, we pre-train a representation with 902 unsupervised Q-learning (though with a different reward function), and thus the discrepancy between 903 pre-training and fine-tuning becomes less severe. 904

905 906

## A.4 CAN WE DO "BRIDGING" WITHOUT ANY REWARD-LABELED DATA?

907 In the bridging phase of U2O RL (Section 4.2), 908 we assume a (small) reward-labeled dataset 909  $\mathcal{D}_{reward}$ . In our experiments, we sample a small 910 number of transitions (e.g., 0.2% in the case of 911 DMC) from the offline dataset and label them 912 with the ground-truth reward function, as in prior 913 works (Touati et al., 2022; Park et al., 2024c). 914 However, these samples do not necessarily have 915 to come from the offline dataset. To show this, we conduct an additional experiment where we 916 do not assume access to any of the existing re-917



Figure 11: U2O RL without using reward-labeling in the offline dataset (8 seeds).

ward samples or the ground-truth reward function in the bridging phase. Specifically, we collect 10K

online samples with random skills and perform the linear regression in Equation 6 only using the
 collected online transitions. We report the performances of U2O (without offline samples) and O2O
 in Figure 11. The results show that U2O still works and outperforms the supervised offline-to-online
 RL baseline.

922 923

924

935

936

937

938

939

940

941

942 943

944

945

946 947

948

949

950

951

952 953

954

965

966 967 A.5 How do different strategies of skill identification affect performance?

925 To understand how skill identification strategies 926 affect online RL performance, we compare our 927 strategy in Section 4.2 with an alternative strategy that simply selects a random latent vector 928 z from the skill space. Figure 12 shows that 929 the skill identification with a randomly selected 930 latent vector performs worse than our strategy. 931 This is likely because modulating the policy with 932 the best latent vector helps boost task-relevant 933 exploration and information. 934



Figure 12:Ablation Figure 13: Comparisonstudy of skill identifica- with PEX and zero-shottion (4 seeds).RL (4 seeds).

A.6 ADDITIONAL EXPERIMENTS ON FINE-TUNING STRATEGIES

We additionally provide experimental results of fine-tuning strategies on a different task (*i.e.*, Cheetah Run). Figure 13 shows that our fine-tuning strategy outperforms previous strategies, such as zero-shot RL and PEX. This result further supports the effectiveness of fine-tuning.



Figure 14: Comparison with O2O RL + DR3 (4 seeds).

# A.7 How does U2O perform compared to O2O combined with methods for mitigating feature collapse?

To further understand the effectiveness of U2O RL, we compare the performance of U2O RL with that of O2O RL combined with DR3 (Kumar et al., 2022), a regularizer that regularizes feature dot products to prevent the feature collapse issue. The result in Figure 14 shows that simply adding the DR3 regularizer is not as effective as U2O RL. We believe this is likely because the full unsupervised RL procedure can lead to much richer representations than simply adding a regularizer.

# A.8 HOW DOES EACH COMPONENT IN U2O RL AFFECT PERFORMANCE?



Figure 15: Ablation study of reward scale matching (4 seeds).

Figure 16: Ablation study of value transfer and policy transfer (4 seeds).

968 Reward scale matching. In Section 4.2, we propose a simple reward scale matching technique that
969 bridges the gap between intrinsic rewards and downstream task rewards. We ablate this component,
970 and report the results in Figure 15. The results suggest that our reward scale matching technique
971 effectively prevents a performance drop at the beginning of the online fine-tuning stage, leading to
972 substantially better final performance on dense-reward tasks (*e.g.*, Walker Run and Cheetah Run).

972 Value transfer vs. policy transfer. In U2O RL, we transfer *both* the value function and policy from 973 unsupervised pre-training to supervised fine-tuning. To dissect the importance of each component, we 974 conduct an ablation study, in which we compare four settings: (1) without any transfer, (2) value-only 975 transfer, (3) policy-only transfer, and (4) full transfer. Figure 16 demonstrates the ablation results 976 on Walker and AntMaze. The results suggest that both value transfer and policy transfer matter in general, but value transfer is more important than policy transfer. This aligns with our findings in 977 Q4 of Section 5 as well as Kumar et al. (2022), which says that the quality of value features often 978 correlates with the performance of TD-based RL algorithms. 979

# 980 981

1001

В ALGORITHM TABLE

984	Algorithm 1 U2O RL: Unsupervised-to-Online Reinforcement Learning
985	<b>Require</b> : offline dataset $\mathcal{D}_{off}$ , reward-labeled dataset $\mathcal{D}_{reward}$ , empty replay buffer $\mathcal{D}_{on}$ , offline
986	pre-training steps $N_{\text{PT}}$ , online fine-tuning steps $N_{\text{FT}}$ , skill latent space $Z$
987	Initialize the parameters of policy $\pi_{\theta}$ and action-value function $Q_{\phi}$
988	for $t = 0, 1, 2, \dots N_{PT} - 1$ do
080	Sample transitions $(s, a, s')$ from $\mathcal{D}_{off}$
000	Sample latent vector $z \in \mathcal{Z}$ and compute intrinsic rewards $r^{\texttt{int}}$
990	Update policy $\pi_{\theta}(a \mid s, z)$ and $Q_{\phi}(s, a, z)$ using normalized intrinsic rewards $\tilde{r}^{int}$
991	end for
992	Compute the best latent vector $z^*$ with Equation 6 using samples $(s, a, s', r)$ from $\mathcal{D}_{reward}$
993	for $t = 0, 1, 2, \dots N_{\text{FT}} - 1$ do
994	Collect transition $(s, a, s', r)$ via environment interaction with $\pi_{\theta}$ and add to replay buffer $\mathcal{D}_{on}$
995	Sample transitions $(s, a, s', r)$ from $\mathcal{D}_{off} \cup \mathcal{D}_{on}$
996	Update policy $\pi_{\theta}(a \mid s, z^*)$ and $Q_{\phi}(s, a, z^*)$ using normalized task rewards $\tilde{r}$
997	end for

1000 **EXPERIMENTAL DETAILS** С

1002 For offline RL pre-training, we use 1M training steps for ExORL, AntMaze, and Adroit and 500K 1003 steps for Kitchen, following Park et al. (2024c). For online fine-tuning, we use 1M additional environment steps for ExORL, AntMaze, and Adroit and 500K steps for Kitchen with an update-to-1004 data ratio of 1. We implement U2O RL based on the official implementation of HILP (Park et al., 1005 2024c). We evaluate the normalized return with 50 episodes every 10k online steps for ExORL tasks, and every 100k online steps for AntMaze, Kitchen, and Adroit tasks. We run our experiments on 1007 A5000 or RTX 3090 GPUs. Each run takes at most 40 hours (e.g. Visual Kitchen). We provide our 1008 implementation in the supplementary material. 1009

1010

C.1 ENVIRONMENTS AND DATASETS 1011

1012 ExORL (Yarats et al., 2022). In the ExORL benchmark, we consider four embodiments, Walker, 1013 Cheetah, Quadruped, and Jaco. Each embodiment has four tasks: Walker has {Run, Flip, Stand, 1014 Walk}, Cheetah has {Run, Run Backward, Walk, Walk Backward}, Quadruped has {Run, Jump, 1015 Stand, Walk}, and Jaco has {Reach Top Left, Reach Top Right, Reach Bottom Left, Reach Bottom 1016 Right}. For all the tasks in Walker, Cheetah, and Quadruped, the maximum return is 1000, and Jaco 1017 has 250. Each embodiment has an offline dataset, which is collected by running exploratory agents such as RND (Burda et al., 2019), and then annotated with task reward function. We use the first 5M 1018 transitions of the offline dataset following the prior work (Touati et al., 2022; Park et al., 2024c). The 1019 maximum episode length is 250 (Jaco) or 1000 (others). 1020

1021 AntMaze (Fu et al., 2020; Jiang et al., 2023). In AntMaze, a quadruped agent aims at reaching the (pre-defined) target position in a maze and gets a positive reward when the agent arrives at a pre-1023 defined neighborhood of the target position. We consider two types of Maze: antmaze-large (Fu et al., 2020) and antmaze-ultra (Jiang et al., 2023), where the latter has twice the size of the 1024 former. Each maze has two types of offline datasets: play and diverse. The dataset consists 1025 of 999 trajectories with an episode length of 1000. In each trajectory, an agent is initialized at a

random location in the maze and is directed to an arbitrary location, which may not be the same as the target goal. At the evaluation, antmaze-large has a maximum episode length of 1000, and antmaze-ultra has 2000. We report normalized scores by multiplying the returns by 100.

1029 Kitchen (Gupta et al., 2020; Fu et al., 2020). In the Kitchen environment, a Franka robot should 1030 achieve four sub-tasks, microwave, slide cabinet, light switch, and kettle. Each task has 1031 a success criterion determined by an object configuration. Whenever the agent achieves a sub-task, 1032 a task reward of 1 is given, where the maximum return is 4. We consider two types of offline 1033 datasets: mixed and partial. We report normalized scores by multiplying the returns by 100. For 1034 Visual-Kitchen, we follow the same camera configuration as Mendonca et al. (2021), Park et al. 1035 (2024d), and Park et al. (2024c), to render  $64 \times 64$  RGB observations, which are used instead of 1036 low-dimensional states. We report normalized scores by multiplying the returns by 25.

Adroit (Fu et al., 2020). In Adroit, a 24-DoF Shadow Hand robot should be controlled to achieve a desired task. We consider two tasks: pen-binary and door-binary, following prior works (Ball et al., 2023; Li et al., 2023). The maximum episode lengths of pen-binary and door-binary are 100 and 200. respectively. We report normalized scores by multiplying the returns by 100.

OGBench-Cube (Park et al., 2024a). In the OGBench-Cube, a 6-DoF UR5e robot arm should be controlled to arrange multiple cubes into the desired configuration. The maximum episode lengths for cube-single and cube-double are 200 and 500, respectively. To make the originally goal-conditioned tasks compatible with regular (single-task) offline-to-online RL, we fix the task\_id to 2, and define the task reward as the negative of the number of the unmatched cubes. We report binary success rates multiplied by 100.

# C.2 HYPERPARAMETERS

1048

1049 1050

1051

## Table 3: Hyperparameters of unsupervised RL pre-training in ExORL.

1052		
1053	Hyperparameter	Value
1054	Learning rate	0.0005 (feature), 0.0001 (others)
1055	Optimizer	Adam (Kingma & Ba, 2015)
1056	Minibatch size	1024
1057	Feature MLP dimensions	(512, 512)
1058	Value MLP dimensions	(1024, 1024, 1024)
1059	Policy MLP dimensions	(1024, 1024, 1024)
1060	TD3 target smoothing coefficient	0.01
1061	TD3 discount factor $\gamma$	0.98
1062	Latent dimension State complex for latent vector informed	5U 10000
1063	State samples for fatent vector interence Successor feature loss	$\Omega$ loss
1064	Hilbert representation discount factor	0.96 (Walker) 0.98 (others)
1065	Hilbert representation expectile	0.50 (Warker), 0.50 (others)
1066	Hilbert representation target smoothing coefficient	0.005
1067		
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		

1080		
1081		
1082		
1083		
1084		
1005		
1005		
1000		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094		
1095		
1096		
1097		
1098		
1099	Table 4: Hyperparameters of unsupervised RL pre	-training in AntMaze, Kitchen, and Adroit.
1100		
1101		
1102	Hyperparameter	Value
1103	Learning rate	0.0003
1104	Optimizer	Adam (Kingma & Ba, 2015)
1105	Minibatch size	256 (Adroit), 512 (others)
1106	Value MLP dimensions	(256, 256, 256) (Adroit), (512, 512, 512) (others)
1107	Policy MLP dimensions	(256, 256, 256) (Adroit), (512, 512, 512) (others)
1108	Target smoothing coefficient	0.005
1109	Discount factor $\gamma$	0.99
1110	Latent dimension	32
1110	Hilbert representation discount factor	0.99
1111	Hilbert representation expectile	0.95
1112	Hilbert representation target smoothing coefficient	0.005
1113	HILP IQL expectile	0.9 (AntMaze), 0.7 (others)
1114	HILP AWR temperature	0.5 (Kitchen) 3 (Adroit-door), 10 (others)
1115		
1116		
1117		
1118		
1119		
1120		
1121		
1122		
1123		
1124		
1125		
1126		
1127		
1128		
1129		
1130		
1131		
1132		
1100		