# Unsupervised-to-Online Reinforcement Learning

**Junsu Kim**[1,2*†]    **Seohong Park**[2*]    **Sergey Levine**[2]
[1]KAIST    [2]University of California, Berkeley
junsu.kim@kaist.ac.kr

## 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** (**U2O RL**), which replaces domain-specific *supervised* offline RL with *unsupervised* offline RL, as a 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 nine state-based and pixel-based environments, we empirically demonstrate that U2O RL achieves strong performance that matches or even outperforms previous offline-to-online RL approaches, 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? Indeed, offline-to-online RL 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 is in contrast to predominant pre-training recipes in large language models or visual representation learning, where they pre-train large models with self-supervised or *unsupervised* objectives to learn useful representations, which can facilitate learning a wide array of downstream tasks. Second, naïve offline-to-online RL is often brittle in practice [37, 51]. This is mainly because pre-trained offline RL agents suffer the distributional shift between the offline and online interaction data [37, 51] or experience feature collapse (Section 5.5), which necessitates specialized, potentially complicated techniques.
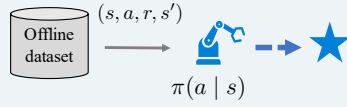
In this work, our central hypothesis is that *unsupervised pre-training of diverse policies* from offline data can serve as an effective data-driven recipe for *online* RL, and can be more effective than even
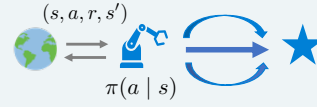
---

**Offline-to-Online RL**

(Supervised) Offline RL Pre-training

$(s, a, r, s')$

Online RL Fine-tuning

$(s, a, r, s')$

$\pi(a \mid s)$

**Unsupervsied-to-Online RL**

Unsupervised Offline RL Pre-training

$(s, a, s', r^{\text{int}}, z)$

$z \in \mathcal{Z}$

$\pi(a \mid s, z)$

Bridging

$(s, a, r, s')$

$\pi(a \mid s, z^*)$

Online RL Fine-tuning

$(s, a, r, s')$

$\pi(a \mid s, z^*)$

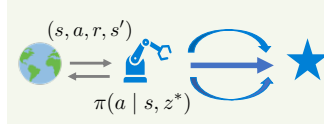Figure 1: **Illustration of U2O RL.** In this work, we propose to replace supervised offline RL with *unsupervised offline RL* in the offline-to-online RL framework. We call this scheme **unsupervised-to-online RL** (**U2O RL**). U2O RL consists of three stages: (1) unsupervised offline RL pre-training, (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.

domain-specific ("supervised") offline pre-training. We call this framework **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 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 [4, 8, 63, 24, 23, 25], improves over the performance of domain-specific specialist pre-training.

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

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

## 2 Related work

**Online RL from prior data**. Prior works have proposed several ways to leverage a previously collected offline dataset to accelerate online RL training. They can be categorized into two main groups: offline-to-online RL and off-policy online RL with offline data. Offline-to-online RL first pre-trains a policy and a value function with offline RL [35, 39, 14, 16, 32, 73, 81, 30, 33, 22, 31, 49, 60, 80], and then continues to fine-tune them with additional online interactions [37, 49,

51, 86, 38, 89, 46, 88]. Since naïve offline-to-online RL is often unstable in practice due to the distributional shift between the dataset and online interactions, prior works have proposed several techniques, such as balanced sampling [37], actor-critic alignment [86], adaptive conservatism [78], and return lower-bounding [51]. In this work, unlike offline-to-online RL, which trains a policy with the target task reward, we offline pre-train a multi-task policy with unsupervised (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 RL agent from scratch on top of a replay buffer filled with offline data, without any pre-training [2, 40, 42, 71]. While this simple approach often leads to improved stability and performance [2], it does not leverage the benefits of pre-training; in contrast, we do leverage pre-training by learning useful features via offline unsupervised RL, which we show leads to better fine-tuning performance in our experiments.

**Unsupervised RL.** The goal of unsupervised RL is to leverage unsupervised pre-training to facilitate downstream reinforcement learning. Prior works have mainly focused on unsupervised representation learning and unsupervised behavior learning. Unsupervised representation learning methods [67, 68, 54, 83, 50, 45, 17, 65, 64, 66] aim to extract useful (visual) representations from data. These representations are then fed into the policy to accelerate task learning. In this work, we focus on unsupervised behavior learning, which aims to pre-train policies that can be directly adapted to downstream tasks. Among unsupervised behavior learning methods, online unsupervised RL pre-trains useful policies by either maximizing state coverage [58, 59, 47, 41] or learning distinct skills [19, 9, 69, 57] via reward-free interactions with the environment. In this work, we consider *offline* unsupervised RL, which does not allow any environment interactions during the pre-training stage.

**Offline unsupervised RL.** Offline unsupervised RL methods focus on learning diverse policies (*i.e.*, skills) from the dataset, rather than exploration, as online interactions are not permitted in this problem setting. There exist three main approaches to offline unsupervised RL. Behavioral cloning-based approaches extract skills from an offline dataset by training a generative model (*e.g.*, variational autoencoders [29], Transformers [77], etc.) [1, 61, 70]. Offline goal-conditioned RL methods learn diverse goal-reaching behaviors with goal-conditioned reward functions [6, 10, 44, 55, 79, 84, 11, 12]. Offline unsupervised skill learning approaches learn diverse skills based on intrinsically defined reward functions [56, 76, 26]. Among these approaches, we use methods in the second and third categories (*i.e.*, goal- or skill-based unsupervised offline RL) as part of our framework.

Our goal in this work is to study how unsupervised offline RL, as opposed to supervised task-specific offline RL, can be employed to facilitate online RL fine-tuning. While somewhat similar unsupervised pre-training schemes have been explored in prior works, they either consider hierarchical RL (or zero-shot RL) with frozen learned skills without fine-tuning [1, 61, 76, 56, 26], assume online-only RL [36], or are limited to the specific setting of goal-conditioned RL [11, 12]. To the best of our knowledge, this is the first work that considers the *fine-tuning* of skill policies pre-trained with unsupervised offline RL *in the context of offline-to-online RL*. Through our experiments, we show that our fine-tuning framework leads to significantly better performance than previous approaches based on hierarchical RL, zero-shot RL, and standard offline-to-online RL.

## 3 Preliminaries

We formulate a decision making problem as a Markov decision process (MDP) [72], which is defined by a tuple of $(\mathcal{S}, \mathcal{A}, P, r, \rho, \gamma)$, where $\mathcal{S}$ is the state space, $\mathcal{A}$ is the action space, $P : \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$ is the transition dynamics, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ is the task reward function, $\rho \in \Delta(\mathcal{S})$ is the initial state distribution, and $\gamma \in (0, 1)$ is the discount factor. Our aim is to learn a policy $\pi : \mathcal{S} \to \Delta(\mathcal{A})$ that maximizes the expectation of cumulative task rewards, $\mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1})]$.

**Offline RL and implicit Q-learning (IQL).** The goal of offline RL is to learn a policy solely from an offline dataset $\mathcal{D}_{\texttt{off}}$, which consists of transition tuples $(s, a, s', r)$. One straightforward approach to offline RL is to simply employ an off-policy RL algorithm (*e.g.* TD3 [15]). For instance, we can minimize the following temporal difference (TD) loss to learn an action-value function (Q-function) from data:

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

3

where $Q_\phi$ denotes the parameterized action-value function, and $Q_{\bar\phi}$ represents the target action-value function [48], whose parameter $\bar\phi$ is updated via Polyak averaging [62] using $\phi$. We can then train a policy $\pi$ to maximize $\mathbb{E}_{a\sim\pi}[Q_\phi(s,a)]$.

While this simple off-policy TD learning can be enough when the dataset has sufficiently large state-action coverage, offline datasets in practice often have limited coverage, which makes the agent susceptible to value overestimation and exploitation, as the agent cannot get corrective feedback from the environment [39]. To address this issue, Kostrikov et al. [31] proposed implicit Q-learning (IQL), which fits an optimal action-value function without querying out-of-distribution actions: IQL replaces the $\arg\max$ operator, which potentially allows the agent to exploit Q-values from out-of-distribution actions, with an expectile loss that implicitly approximates the maximum value. Specifically, IQL minimizes the following losses:

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

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

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 [52] 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 framework 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 [56] in our experiments. HILP consists of two phases. In the first phase, HILP trains a feature network $\xi : \mathcal{S} \to \mathcal{Z}$ 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 $s, g \in \mathcal{S}$, 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 [55] (see Park et al. [56] 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 [31, 15]:

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

where $z$ is sampled from the unit ball, $\{z \in \mathcal{Z} : \|z\| = 1\}$. Intuitively, Equation 5 encourages the agent to learn behaviors that move in every possible latent direction, resulting in diverse state-spanning skills [56]. 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 [7, 3].

## 4 Unsupervised-to-online RL (U2O RL)

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 framework **unsupervised-to-online RL** (**U2O RL**). In this section, we first describe the three stages of our U2O RL framework (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 $\mathcal{D}_{\text{off}}$. 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 $r^{\text{int}}(s, a, s', z)$, where $z$ is a skill latent vector sampled from a latent space $\mathcal{Z} = \mathbb{R}^d$. Even if $\mathcal{D}_{\text{off}}$ contains reward labels, we do not use any reward information in this phase.

---

**Algorithm 1** U2O RL: Unsupervised-to-Online Reinforcement Learning

---
**Require**: offline dataset $\mathcal{D}_{\mathtt{off}}$, reward-labeled dataset $\mathcal{D}_{\mathtt{reward}}$, empty replay buffer $\mathcal{D}_{\mathtt{on}}$, offline pre-training steps $N_{\mathtt{PT}}$, online fine-tuning steps $N_{\mathtt{FT}}$, skill latent space $\mathcal{Z}$
Initialize the parameters of policy $\pi_\theta$ and action-value function $Q_\phi$
**for** $t = 0, 1, 2, \ldots N_{\mathtt{PT}} - 1$ **do**
    Sample transitions $(s, a, s')$ from $\mathcal{D}_{\mathtt{off}}$
    Sample latent vector $z \in \mathcal{Z}$ and compute intrinsic rewards $r^{\mathtt{int}}$
    Update policy $\pi_\theta(a \mid s, z)$ and $Q_\phi(s, a, z)$ using normalized intrinsic rewards $\tilde{r}^{\mathtt{int}}$
**end for**
Compute the best latent vector $z^*$ with Equation 6 using samples $(s, a, s', r)$ from $\mathcal{D}_{\mathtt{reward}}$
**for** $t = 0, 1, 2, \ldots N_{\mathtt{FT}} - 1$ **do**
    Collect transition $(s, a, s', r)$ via environment interaction with $\pi_\theta$ and add to replay buffer $\mathcal{D}_{\mathtt{on}}$
    Sample transitions $(s, a, s', r)$ from $\mathcal{D}_{\mathtt{off}} \cup \mathcal{D}_{\mathtt{on}}$
    Update policy $\pi_\theta(a \mid s, z^*)$ and $Q_\phi(s, a, z^*)$ using normalized task rewards $\tilde{r}$
**end for**

---

Among existing unsupervised offline policy pre-training methods (Section 2), we opt to employ successor feature-based methods [7, 3, 82, 76, 56] or offline goal-conditioned RL methods [6, 55] for our unsupervised pre-training, since they provide a convenient mechanism to identify the best skill latent vector given a downstream task, which we will utilize in the next phase. More concretely, we mainly choose to employ HILP [56] (Section 3) as an unsupervised offline policy pre-training method in our experiments for its state-of-the-art performance in previous benchmarks [56]. We note, however, that any other unsupervised offline successor feature-based skill learning methods [76] or offline goal-conditioned RL methods [55] can also be used in place of HILP (see Appendix A.1).

### 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 multi-task 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.

**Skill identification.** Since we chose to use a successor feature- or goal-based unsupervised pre-training method in the previous phase, the first challenge is relatively straightforward. For goal-oriented tasks (*e.g.*, AntMaze [13] and Kitchen [20]), we assume the task goal $g$ to be available, and we either directly use $g$ (for goal-conditioned methods) or infer the skill $z^*$ that corresponds to $g$ based on a predefined conversion formula (for successor feature-based methods that support such a conversion [76, 56]). For general reward-maximization tasks, we employ successor feature-based unsupervised pre-training methods, and use the following linear regression to find the skill latent vector $z^*$ that best approximates the downstream task reward function $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$:

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

where $f$ is the feature network in the successor feature framework (Section 3) and $\mathcal{D}_{\mathtt{reward}}$ is a reward-labeled dataset. This reward-labeled dataset $\mathcal{D}_{\mathtt{reward}}$ can be either the full offline dataset $\mathcal{D}_{\mathtt{off}}$ (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 $\mathcal{D}_{\mathtt{reward}}$, following previous works [76, 56], but we do not require $\mathcal{D}_{\mathtt{reward}}$ to be a subset of $\mathcal{D}_{\mathtt{off}}$ (see Appendix A.4).

**Reward scale matching.** After identifying the best skill latent vector $z^*$, our next step is to bridge the gap between intrinsic and extrinsic rewards. Since these two reward functions can have very different scales, naïve online adaptation can lead to abrupt shifts in target Q-values, potentially causing significant performance drops in the early stages of online fine-tuning. While one can employ sophisticated reward-shaping techniques to deal with this issue [53, 18], in this work, we propose a simple yet effective reward scale matching technique that we find to be effective enough in practice. Specifically, we compute the running mean and standard deviation of intrinsic rewards during the

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

pre-training phase, and normalize the intrinsic rewards with the calculated statistics. Similarly, during the fine-tuning phase, we compute the statistics of task rewards and normalize the task rewards so that they have the same scale and mean as normalized intrinsic rewards. This way, we can prevent abrupt shifts in reward scales without altering the optimal policy for the downstream task. In our experiments, we find that this simple technique is crucial for achieving good performance, especially in environments with dense rewards (Section 5.7).

### 4.3 Online fine-tuning

Our final step is to fine-tune the skill policy with online environment interactions (Figure 1 (bottom right)). This step is straightforward: since we have found $z^*$ in the previous stage, we can simply *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 [31], TD3 [15]) with additional online interaction data. While one can employ existing 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 using the same (offline) RL algorithm is enough to achieve strong performance that matches or even outperforms state-of-the-art offline-to-online RL techniques.

### 4.4 Why is U2O RL potentially better than offline-to-online RL?

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

We hypothesize that this is because unsupervised multi-task offline RL enables better *feature learning* than supervised single-task offline RL. By training the agent on a number of diverse intrinsically 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 online fine-tuning. This resembles the recent observation in machine learning that large-scale unsupervised pre-training improves downstream task performances over task-specific supervised pre-training [4, 8, 63, 24, 23, 25]. In our experiments, we empirically show that U2O RL indeed learns better features than its supervised counterpart (Section 5.5).

Another benefit of U2O RL is that it does not use any task-specific information during pre-training. This is appealing because we can reuse a single pre-trained policy for a number of different downstream tasks. Moreover, it enables leveraging potentially very large, task-agnostic offline data during pre-training, which is often much cheaper to collect than task-specific, curated datasets [43].

## 5 Experiments

In our experiments, we evaluate the performance of U2O RL in the context of offline-to-online RL. We aim to answer the following research questions: (1) Is U2O RL better than previous offline-to-online RL strategies? (2) What makes unsupervised offline RL result in better fine-tuning performance than supervised offline RL? (3) Which components of U2O RL are important?
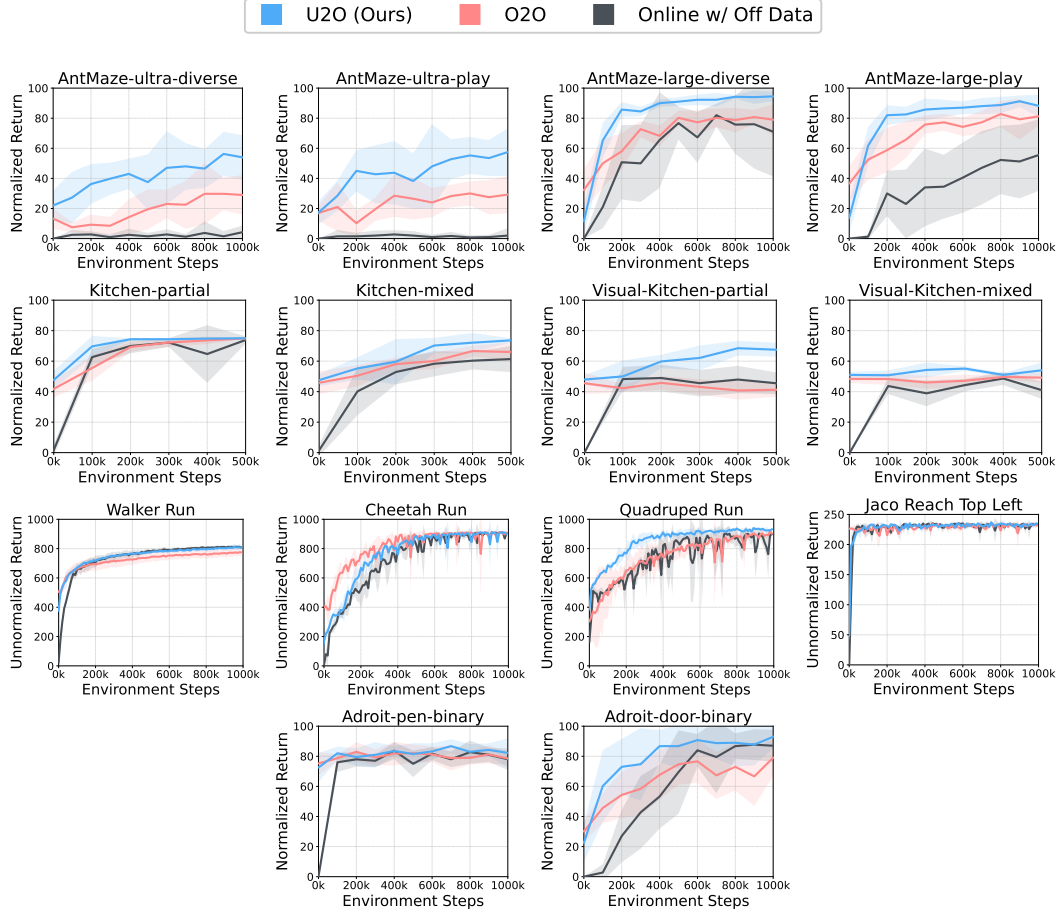
Figure 3: **Online fine-tuning plots of U2O RL and previous offline-to-online RL frameworks (8 seeds).** Across the benchmarks, our U2O RL mostly shows consistently better performance than standard offline-to-online RL and off-policy online RL with offline data.

## 5.1 Experimental setup

**Environments and offline datasets.** In our experiments, we consider nine tasks across five benchmarks (Figure 2). **ExORL** [85] is a benchmark suite that consists of offline datasets collected by exploratory policies (*e.g.*, RND [5]) on the DeepMind Control Suite [75]. We consider four embodiments (Walker, Cheetah, Quadruped, and Jaco), each of which has four tasks. **AntMaze** [13, 27] is a navigation task, whose goal is to control an 8-DoF quadruped agent to reach a target position. We consider the two most challenging mazes with the largest sizes, large and ultra, and two types of offline datasets, diverse and play. **Kitchen** [20, 13] is a robotic manipulation task, where the goal is to control a 9-DoF Franka robot arm to achieve four subtasks sequentially. The environment has four subtasks in total. We consider two types of offline datasets from the D4RL suite [13], partial and mixed. **Visual Kitchen** [20, 13, 56] is a pixel-based variant of the Kitchen environment, where an agent must achieve four subtasks purely from $64 \times 64 \times 3$ pixel observations instead of low-dimensional state information. **Adroit** [13] is a dexterous manipulation benchmark, where the goal is to control a 24-DoF robot hand to twirl a pen or open a door. We provide further details in Appendix B.1

**Implementation.** In our experiments, we mainly employ HILP [56] as the unsupervised offline policy pre-training algorithm in U2O RL. For the offline RL backbone, we use TD3 [15] for ExORL and IQL [30] for others following previous works [76, 56]. Since both IQL and TD3+BC [14, 73] have been known to achieve strong performance in the offline-to-online RL setting [74], we use them for the online fine-tuning phase in U2O RL as well. For sparse-reward tasks (AntMaze, Kitchen, and Adroit), we do not apply reward scale matching. For AntMaze, Kitchen, and Adroit, we report

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 [2] are taken from Nakamoto et al. [51], Wang et al. [78]. Scores within the 5% of the best score are highlighted in bold, as in Kostrikov et al. [31]. We use 8 random seeds for each task for U2O RL.

| Task | antmaze-ultra-diverse | antmaze-ultra-play | antmaze-large-diverse | antmaze-large-play | kitchen-partial | kitchen-mixed |
|---|---|---|---|---|---|---|
| CQL | - | - | 25 → 87 | 34 → 76 | 71 → **75** | 56 → 50 |
| IQL | 13 → 29 | 17 → 29 | 40 → 59 | 41 → 51 | 40 → 60 | 48 → 48 |
| AWAC | - | - | 00 → 00 | 00 → 00 | 01 → 13 | 02 → 12 |
| O3F | - | - | 59 → 28 | 68 → 01 | 11 → 22 | 06 → 33 |
| ODT | - | - | 00 → 00 | 00 → 00 | - | - |
| CQL+SAC | - | - | 36 → 00 | 21 → 00 | 71 → 00 | 59 → 01 |
| Hybrid RL | - | - | → 00 | → 00 | → 00 | → 01 |
| SAC+od | - | - | → 00 | → 00 | → 07 | → 00 |
| SAC | - | - | → 00 | → 00 | → 03 | → 02 |
| IQL+od | → 04 | → 05 | → 71 | → 56 | → 74 | → 61 |
| FamO2O | - | - | → 64 | → 61 | - | - |
| RLPD | 00 → 00 | 00 → 00 | 00 → 94 | 00 → 81 | - | - |
| Cal-QL | 05 → 05 | 15 → 13 | 33 → **95** | 26 → **90** | 67 → 79 | 38 → **80** |
| **U2O (Ours)** | 22 → **54** | 17 → **58** | 11 → **95** | 14 → **88** | 48 → **75** | 48 → 74 |

normalized scores, following Fu et al. [13]. In our experiments, we use **8 random seeds** and report standard deviations with shaded areas, unless otherwise stated. We refer the reader to Appendix B for the full implementation details and hyperparameters.

## 5.2 Is U2O RL better than previous offline-to-online RL frameworks?

We begin our experiments by comparing our approach, unsupervised-to-online RL, with two previous offline-to-online RL *frameworks* (Section 2): **offline-to-online RL** (**O2O RL**) and **off-policy online RL with offline data** (**Online w/ Off Data**). To recall, offline-to-online RL [37, 49, 51, 86, 38] first pre-trains a policy with *supervised* offline RL using the task reward, and then continues training it with online rollouts. Off-policy online RL [2, 42, 71] trains a policy from scratch on top of a replay buffer filled with offline data. Here, we use the *same* offline RL backbone (*i.e.*, TD3 for ExORL and IQL for AntMaze, Kitchen, and Adroit) to ensure apples-to-apples comparisons between the three frameworks. We will compare U2O RL with previous specialized offline-to-online RL techniques in Section 5.3.

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}`).

## 5.3 How does U2O RL compare to previous specialized offline-to-online RL techniques?

Next, we compare U2O RL with 13 previous specialized offline-to-online RL methods, including **CQL** [33], **IQL** [31], **AWAC** [49], **O3F** [46], **ODT** [89], **CQL+SAC** [33, 21], **Hybrid RL** [71], **SAC+od (offline data)** [21, 2], **SAC** [21], **IQL+od (offline data)** [31, 2], **RLPD** [2], **FamO2O** [78] and **Cal-QL** [51]. We show the comparison results in Table 1, where we take the scores from Nakamoto et al. [51], Wang et al. [78] for the tasks that are common to ours. Since Cal-QL achieves the best performance in the table, we additionally make a comparison with Cal-QL on `antmaze-ultra-{diverse, play}` as well, by running their official implementation with tuned hyperparameters.

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.
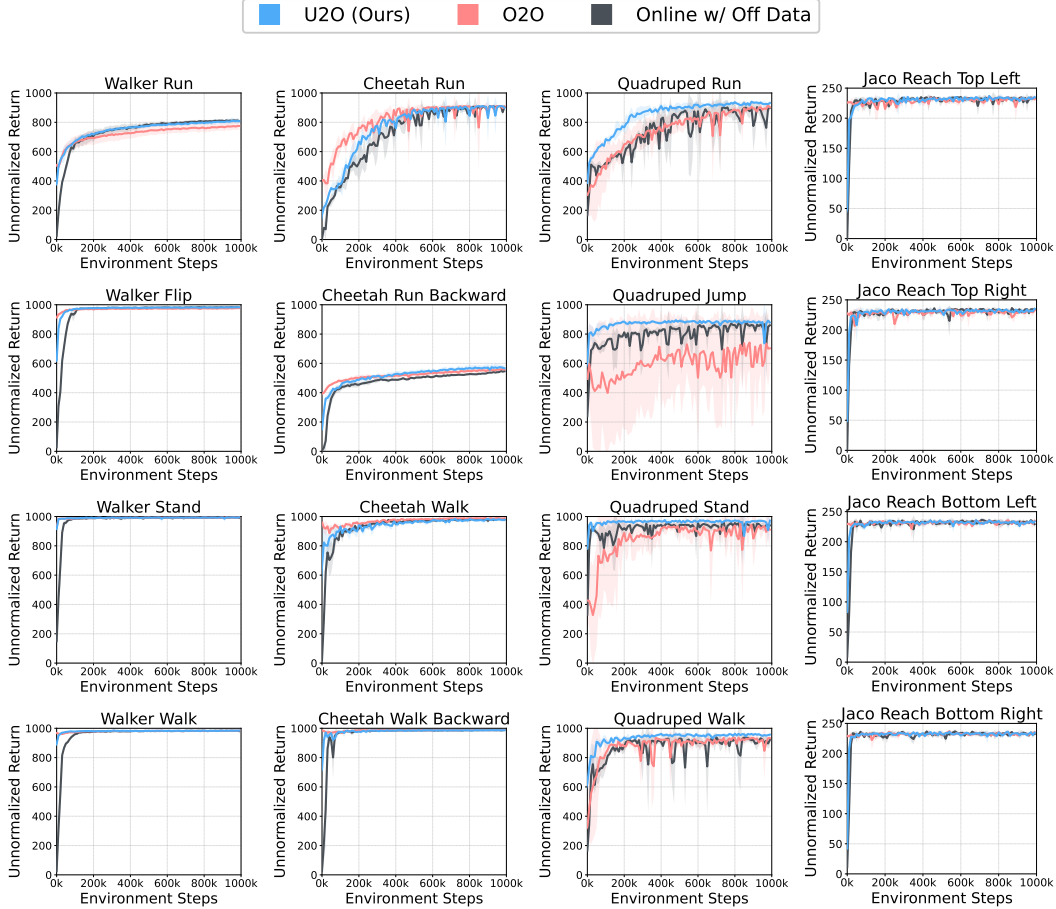
Figure 4: **Learning curves during online RL fine-tuning (8 seeds).** A single pre-trained model from U2O can be fine-tuned to solve multiple downstream tasks. Across the embodiments and tasks, our U2O RL matches or outperforms standard offline-to-online RL and off-policy online RL with offline data even though U2O RL uses a single task-agnostic pre-trained model.

## 5.4 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 4. The results show that, for example, 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 task-agnostic pre-trained model, the performance of U2O RL matches or even outperforms O2O RL, which pre-trains a model with task-specific rewards.

## 5.5 Why does U2O RL often outperform supervised offline-to-online RL?

In the above experiments, we showed that our unsupervised-to-online RL framework often even outperforms previous supervised offline-to-online RL methods. We hypothesized in Section 4.4 that this is because unsupervised offline pre-training yields better *representations* that facilitate online task adaptation. To empirically verify this hypothesis, we measure the quality of the value function representations using the method proposed by Kumar et al. [34]. Specifically, we define the value features $\zeta_\phi(s, a)$ as the penultimate layer of the value function $Q_\phi$, *i.e.*, $Q_\phi(s, a) = w_\phi^\top \zeta_\phi(s, a)$, and measure the dot product between consecutive state-action pairs, $\zeta_\phi(s, a)^\top \zeta_\phi(s', a')$ [34]. Intuitively,
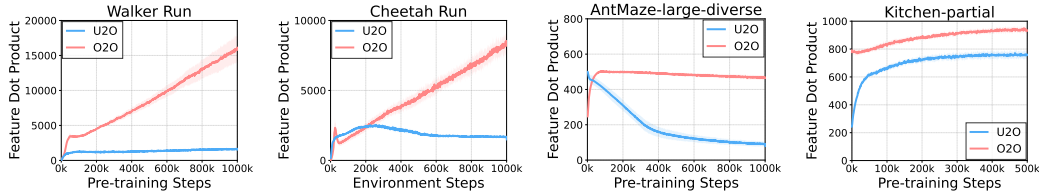
Figure 5: **Feature dot products during offline RL pre-training (lower is better, 8 seeds).** The plots show that *unsupervised* offline pre-training effectively prevents feature collapse (co-adaptation), yielding better representations than supervised offline pre-training.

this dot product represents the degree to which these two representations are "collapsed" (or "co-adapted"), which is known to be correlated with poor performance [34] (*i.e.*, the lower the better).

Figure 5 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 [76, 82]) in Appendix A.1.

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

In this work, we focus on the *fine-tuning* of offline pre-trained skill policies, but this is not the only way to leverage pre-trained skills for downstream tasks. To see how our fine-tuning scheme compares to other alternative strategies, we compare U2O RL with three previously considered approaches: **hierarchical RL** (**HRL**, *e.g.*, **OPAL** [1], **SPiRL** [61]) [1, 61, 76, 56, 26], **zero-shot RL** [76, 56], and **PEX** [87]. HRL additionally trains a high-level policy that 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 pre-trained multi-task policies and a newly initialized policy with a multiplexer that chooses the best policy.

Figure 6 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.
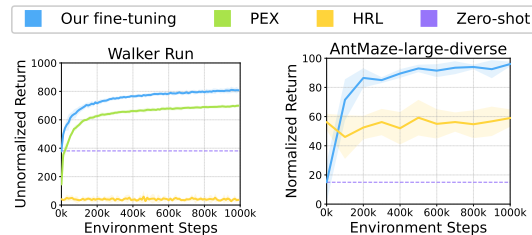


Figure 6: **Fine-tuning is better than previous strategies, such as hierarchical RL, zero-shot RL, and PEX (8 seeds).**

### 5.7 Ablation study

**Reward scale matching.** In Section 4.2, we propose a simple reward scale matching technique that bridges the gap between intrinsic rewards and downstream task rewards. We ablate this component, and report the results in Figure 7. The results suggest that our reward scale matching technique effectively prevents a performance drop at the beginning of the online fine-tuning stage, leading to substantially better final performance on dense-reward tasks (*e.g.*, Walker Run and Cheetah Run).

**Value transfer vs. policy transfer.** In our U2O RL framework, we transfer *both* the value function and policy from unsupervised pre-training to supervised fine-tuning. To dissect the importance of each component, we conduct an ablation study, in which we compare four settings: (1) without any transfer, (2) value-only transfer, (3) policy-only transfer, and (4) full transfer. Figure 8 demonstrates the ablation results on Walker and AntMaze. The results suggest that both value transfer and policy
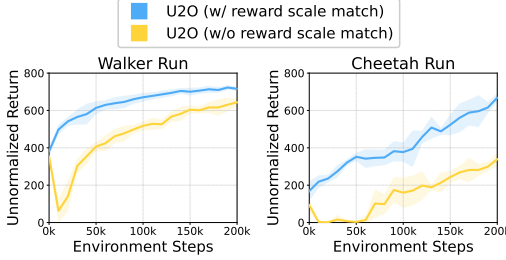
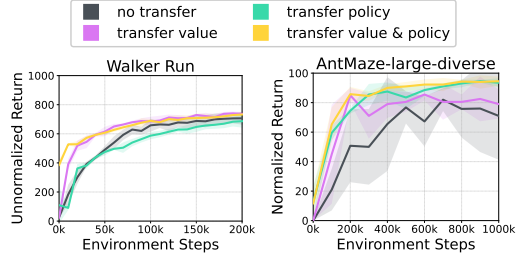Figure 7: **Ablation study of reward scale matching (4 seeds).**

Figure 8: **Ablation study of value transfer and policy transfer (4 seeds).**

transfer matter in general, but value transfer is more important than policy transfer. This aligns with our findings in Section 5.5 as well as Kumar et al. [34], which says that the quality of value features often correlates with the performance of TD-based RL algorithms.

We refer to the reader to Appendix for further analysis including (1) combining U2O RL with other offline unsupervised skill learning methods (Appendix A.1), (2) U2O RL on a different type of datasets (Appendix A.2), (3) comparisons between U2O RL and pure representation learning schemes (Appendix A.3), (4) U2O RL without reward samples in the bridging phase (Appendix A.4), (5) an ablation study with different skill identification strategies (Appendix A.5), and (6) additional results with different online RL strategies (Appendix A.6).

# 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 U2O RL framework 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.

**Do we expect U2O RL to be *always* better than O2O RL?** While we showed strong results of U2O RL throughout the paper, our framework does have limitations. One limitation is that U2O RL may not be as effective when the offline dataset is monolithic and heavily tailored toward the downstream task (see Appendix A.2), which leaves little room for improvement for unsupervised offline RL compared to task-specific supervised offline RL. We believe U2O RL can be most effective (compared to standard offline-to-online RL) when the dataset is highly diverse so that unsupervised offline RL methods can learn a variety of behaviors and thus learn better features and representations. Nonetheless, given the recent successes in large-scale self-supervised and unsupervised pre-training from unlabeled data, we believe our U2O RL framework serves as a step toward a general recipe for scalable data-driven decision-making.

## Acknowledgments and Disclosure of Funding

## References

[1] A. Ajay, A. Kumar, P. Agrawal, S. Levine, and O. Nachum. Opal: Offline primitive discovery for accelerating offline reinforcement learning. In *International Conference on Learning Representations*, 2021.

[2] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. Efficient online reinforcement learning with offline data. In *International Conference on Machine Learning*, 2023.

[3] A. Barreto, W. Dabney, R. Munos, J. J. Hunt, T. Schaul, H. P. van Hasselt, and D. Silver. Successor features for transfer in reinforcement learning. In *Advances in Neural Information Processing Systems*, 2017.

[4] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020.

[5] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. Exploration by random network distillation. In *International Conference on Learning Representations*, 2019.

[6] Y. Chebotar, K. Hausman, Y. Lu, T. Xiao, D. Kalashnikov, J. Varley, A. Irpan, B. Eysenbach, R. C. Julian, C. Finn, et al. Actionable models: Unsupervised offline reinforcement learning of robotic skills. In *International Conference on Machine Learning*, 2021.

[7] P. Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5:613–624, 1993.

[8] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.

[9] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.

[10] B. Eysenbach, T. Zhang, S. Levine, and R. R. Salakhutdinov. Contrastive learning as goal-conditioned reinforcement learning. In *Advances in Neural Information Processing Systems*, 2022.

[11] K. Fang, P. Yin, A. Nair, and S. Levine. Planning to practice: Efficient online fine-tuning by composing goals in latent space. In *International Conference on Intelligent Robots and Systems*, 2022.

[12] K. Fang, P. Yin, A. Nair, H. R. Walke, G. Yan, and S. Levine. Generalization with lossy affordances: Leveraging broad offline data for learning visuomotor tasks. In *Conference on Robot Learning*, 2023.

[13] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.

[14] S. Fujimoto and S. S. Gu. A minimalist approach to offline reinforcement learning. In *Advances in Neural Information Processing Systems*, 2021.

[15] S. Fujimoto, H. Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, 2018.

[16] S. Fujimoto, D. Meger, and D. Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, 2019.

[17] D. Ghosh, C. A. Bhateja, and S. Levine. Reinforcement learning from passive data via latent intentions. In *International Conference on Machine Learning*, 2023.

[18] A. Gleave, M. Dennis, S. Legg, S. Russell, and J. Leike. Quantifying differences in reward functions. In *International Conference on Learning Representations*, 2021.

[19] K. Gregor, D. J. Rezende, and D. Wierstra. Variational intrinsic control. *arXiv preprint arXiv:1611.07507*, 2016.

[20] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning. In *Conference on Robot Learning*, 2020.

[21] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, 2018.

[22] P. Hansen-Estruch, I. Kostrikov, M. Janner, J. G. Kuba, and S. Levine. Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint arXiv:2304.10573*, 2023.

[23] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.

[24] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.

[25] O. J. Hénaff, A. Srinivas, J. De Fauw, A. Razavi, C. Doersch, S. A. Eslami, and A. van den Oord. Data-efficient image recognition with contrastive predictive coding. In *International Conference on Machine Learning*, 2020.

[26] H. Hu, Y. Yang, J. Ye, Z. Mai, and C. Zhang. Unsupervised behavior extraction via random intent priors. In *Advances in Neural Information Processing Systems*, 2023.

[27] Z. Jiang, T. Zhang, M. Janner, Y. Li, T. Rocktäschel, E. Grefenstette, and Y. Tian. Efficient planning in a compact latent action space. In *International Conference on Learning Representations*, 2023.

[28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.

[29] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.

[30] I. Kostrikov, R. Fergus, J. Tompson, and O. Nachum. Offline reinforcement learning with fisher divergence critic regularization. In *International Conference on Machine Learning*, 2021.

[31] I. Kostrikov, A. Nair, and S. Levine. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations*, 2022.

[32] A. Kumar, J. Fu, M. Soh, G. Tucker, and S. Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In *Advances in Neural Information Processing Systems*, 2019.

[33] A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 2020.

[34] A. Kumar, R. Agarwal, T. Ma, A. Courville, G. Tucker, and S. Levine. DR3: Value-based deep reinforcement learning requires explicit regularization. In *International Conference on Learning Representations*, 2022.

[35] S. Lange, T. Gabel, and M. Riedmiller. Batch reinforcement learning. In *Reinforcement learning: State-of-the-art*, pages 45–73. Springer, 2012.

[36] M. Laskin, D. Yarats, H. Liu, K. Lee, A. Zhan, K. Lu, C. Cang, L. Pinto, and P. Abbeel. Urlb: Unsupervised reinforcement learning benchmark. In *Advances in Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.

[37] S. Lee, Y. Seo, K. Lee, P. Abbeel, and J. Shin. Offline-to-online reinforcement learning via balanced replay and pessimistic q-ensemble. In *Conference on Robot Learning*, 2022.

[38] K. Lei, Z. He, C. Lu, K. Hu, Y. Gao, and H. Xu. Uni-o4: Unifying online and offline deep reinforcement learning with multi-step on-policy optimization. *arXiv preprint arXiv:2311.03351*, 2023.

[39] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.

[40] Q. Li, J. Zhang, D. Ghosh, A. Zhang, and S. Levine. Accelerating exploration with unlabeled prior data. In *Advances in Neural Information Processing Systems*, 2023.

[41] H. Liu and P. Abbeel. Behavior from the void: Unsupervised active pre-training. *arXiv preprint arXiv:2103.04551*, 2021.

[42] J. Luo, Z. Hu, C. Xu, Y. L. Tan, J. Berg, A. Sharma, S. Schaal, C. Finn, A. Gupta, and S. Levine. Serl: A software suite for sample-efficient robotic reinforcement learning. *arXiv preprint arXiv:2401.16013*, 2024.

[43] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning latent plans from play. In *Conference on Robot Learning*, 2019.

[44] J. Y. Ma, J. Yan, D. Jayaraman, and O. Bastani. Offline goal-conditioned reinforcement learning via $f$-advantage regression. In *Advances in Neural Information Processing Systems*, 2022.

[45] Y. J. Ma, S. Sodhani, D. Jayaraman, O. Bastani, V. Kumar, and A. Zhang. VIP: Towards universal visual reward and representation via value-implicit pre-training. In *International Conference on Learning Representations*, 2023.

[46] M. S. Mark, A. Ghadirzadeh, X. Chen, and C. Finn. Fine-tuning offline policies with optimistic action selection. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.

[47] R. Mendonca, O. Rybkin, K. Daniilidis, D. Hafner, and D. Pathak. Discovering and achieving goals via world models. In *Advances in Neural Information Processing Systems*, 2021.

[48] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.

[49] A. Nair, A. Gupta, M. Dalal, and S. Levine. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.

[50] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.

[51] M. Nakamoto, S. Zhai, A. Singh, M. Sobol Mark, Y. Ma, C. Finn, A. Kumar, and S. Levine. Cal-ql: Calibrated offline rl pre-training for efficient online fine-tuning. In *Advances in Neural Information Processing Systems*, 2023.

[52] W. K. Newey and J. L. Powell. Asymmetric least squares estimation and testing. *Econometrica: Journal of the Econometric Society*, pages 819–847, 1987.

[53] A. Y. Ng, D. Harada, and S. Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *International Conference on Machine Learning*, 1999.

[54] S. Parisi, A. Rajeswaran, S. Purushwalkam, and A. Gupta. The unsurprising effectiveness of pre-trained vision models for control. In *International Conference on Machine Learning*, 2022.

[55] S. Park, D. Ghosh, B. Eysenbach, and S. Levine. Hiql: Offline goal-conditioned rl with latent states as actions. In *Advances in Neural Information Processing Systems*, 2024.

[56] S. Park, T. Kreiman, and S. Levine. Foundation policies with hilbert representations. *arXiv preprint arXiv:2402.15567*, 2024.

[57] S. Park, O. Rybkin, and S. Levine. METRA: Scalable unsupervised RL with metric-aware abstraction. In *International Conference on Learning Representations*, 2024.

[58] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *International Conference on Machine Learning*, 2017.

[59] D. Pathak, D. Gandhi, and A. Gupta. Self-supervised exploration via disagreement. In *International Conference on Machine Learning*, 2019.

[60] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.

[61] K. Pertsch, Y. Lee, and J. Lim. Accelerating reinforcement learning with learned skill priors. In *Conference on Robot Learning*, 2021.

[62] B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 1992.

[63] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.

[64] Y. Seo, D. Hafner, H. Liu, F. Liu, S. James, K. Lee, and P. Abbeel. Masked world models for visual control. In *Conference on Robot Learning*, 2022.

[65] Y. Seo, K. Lee, S. James, and P. Abbeel. Reinforcement learning with action-free pre-training from videos. In *International Conference on Machine Learning*, 2022.

[66] Y. Seo, J. Kim, S. James, K. Lee, J. Shin, and P. Abbeel. Multi-view masked world models for visual robotic manipulation. In *International Conference on Machine Learning*, 2023.

[67] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, and S. Levine. Time-contrastive networks: Self-supervised learning from video. In *International Conference on Robotics and Automation*, 2018.

[68] R. Shah and V. Kumar. Rrl: Resnet as representation for reinforcement learning. In *International Conference on Machine Learning*, 2021.

[69] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020.

[70] A. Singh, H. Liu, G. Zhou, A. Yu, N. Rhinehart, and S. Levine. Parrot: Data-driven behavioral priors for reinforcement learning. In *International Conference on Learning Representations*, 2021.

[71] Y. Song, Y. Zhou, A. Sekhari, J. A. Bagnell, A. Krishnamurthy, and W. Sun. Hybrid RL: Using both offline and online data can make RL efficient. In *International Conference on Learning Representations*, 2023.

[72] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT Press, 2018.

[73] D. Tarasov, V. Kurenkov, A. Nikulin, and S. Kolesnikov. Revisiting the minimalist approach to offline reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.

[74] D. Tarasov, A. Nikulin, D. Akimov, V. Kurenkov, and S. Kolesnikov. Corl: Research-oriented deep offline reinforcement learning library. In *Advances in Neural Information Processing Systems*, 2023.

[75] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

[76] A. Touati, J. Rapin, and Y. Ollivier. Does zero-shot reinforcement learning exist? In *International Conference on Learning Representations*, 2022.

[77] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.

[78] S. Wang, Q. Yang, J. Gao, M. Lin, H. Chen, L. Wu, N. Jia, S. Song, and G. Huang. Train once, get a family: State-adaptive balances for offline-to-online reinforcement learning. In *Advances in Neural Information Processing Systems*, 2023.

[79] T. Wang, A. Torralba, P. Isola, and A. Zhang. Optimal goal-reaching reinforcement learning via quasimetric learning. In *International Conference on Machine Learning*, 2023.

[80] Z. Wang, A. Novikov, K. Zolna, J. S. Merel, J. T. Springenberg, S. E. Reed, B. Shahriari, N. Siegel, C. Gulcehre, N. Heess, et al. Critic regularized regression. In *Advances in Neural Information Processing Systems*, 2020.

[81] Y. Wu, G. Tucker, and O. Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.

[82] Y. Wu, G. Tucker, and O. Nachum. The laplacian in RL: Learning representations with efficient approximations. In *International Conference on Learning Representations*, 2019.

[83] T. Xiao, I. Radosavovic, T. Darrell, and J. Malik. Masked visual pre-training for motor control. *arXiv preprint arXiv:2203.06173*, 2022.

[84] R. Yang, L. Yong, X. Ma, H. Hu, C. Zhang, and T. Zhang. What is essential for unseen goal generalization of offline goal-conditioned rl? In *International Conference on Machine Learning*, 2023.

[85] D. Yarats, D. Brandfonbrener, H. Liu, M. Laskin, P. Abbeel, A. Lazaric, and L. Pinto. Don't change the algorithm, change the data: Exploratory data for offline reinforcement learning. *arXiv preprint arXiv:2201.13425*, 2022.

[86] Z. Yu and X. Zhang. Actor-critic alignment for offline-to-online reinforcement learning. In *International Conference on Machine Learning*, 2023.

[87] H. Zhang, W. Xu, and H. Yu. Policy expansion for bridging offline-to-online reinforcement learning. In *International Conference on Learning Representations*, 2023.

[88] K. Zhao, Y. Ma, J. Liu, H. Jianye, Y. Zheng, and Z. Meng. Improving offline-to-online reinforcement learning with q-ensembles. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023.

[89] Q. Zheng, A. Zhang, and A. Grover. Online decision transformer. In *International Conference on Machine Learning*, 2022.

# Appendices

## A   Additional Experiments

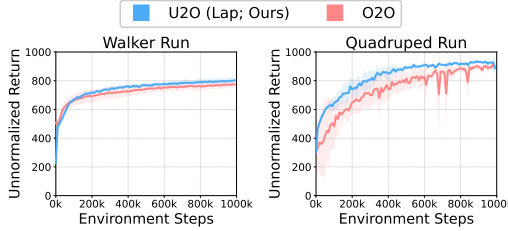### A.1   Can U2O RL be combined with other offline unsupervised RL methods?



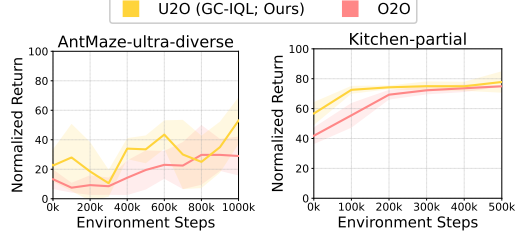Figure 9: **U2O RL with Laplacian-based successor feature learning (8 seeds).**



Figure 10: **U2O RL with goal-conditioned IQL (8 seeds).**

While we employ HILP [56] as an offline unsupervised skill learning method in the U2O RL framework in our main experiments, our framework 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 [76, 82] or goal-conditioned IQL (GC-IQL) [31, 55], and report the results in Figures 9 and 10, 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 skill learning methods also lead to better value representations. We measure the same feature dot product metric in Section 5.5 with the graph Laplacian-based successor feature learning method and report the results in Figure 11. The results suggest that this unsupervised RL method also prevents feature co-adaptation, leading to better features.
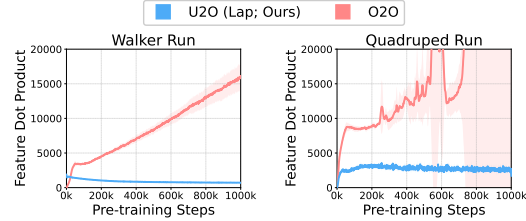


Figure 11: **Feature dot product analysis with Laplacian-based successor feature learning (8 seeds).**

### A.2   When is U2O RL better than O2O RL?

In Section 6, we discussed that U2O RL is advantageous over O2O RL especially when the offline dataset allows the underlying unsupervised RL method to learn diverse behaviors and thus capture diverse features. Conversely, we do not particularly expect U2O RL to be better than O2O RL when the dataset is monolithic (*e.g.*, consists of only expert trajectories, has less diversity, etc.). To empirically show this, we conduct an additional experiment with a different dataset in `AntMaze-large` and `AntMaze-ultra` that consists of monolithic, ex-



Figure 12: **Online RL learning curves with expert-only datasets (4 seeds).**

pert trajectories (we collect a 1M-sized dataset by rolling out an offline pre-trained policy). Figure 12 shows the results. As expected, we find that the performance difference between U2O and O2O is marginal in this monolithic setting.
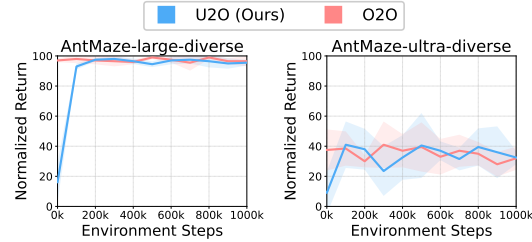
### A.3 Do we need to use unsupervised RL for pre-training representations?

In Sections 4.4 and 5.5, 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) | **94.50 $\pm$ 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 [10] and Hilbert (metric) representation learning [56], where the latter is equivalent to directly taking $\xi$ in the HILP framework (Equation 4) (note that the original U2O RL takes the Q function of HILP, not the Hilbert representation $\xi$ itself, which is used to train the Q function). To evaluate their fine-tuning performances, for the temporal contrastive representation, we fine-tune both the Q function and policy with contrastive RL [10]; for the Hilbert representation, we take the pre-trained representation, add one new layer, and use it as the initialization of the Q function. Table 2 shows the results on `antmaze-large-diverse`. Somewhat intriguingly, the 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 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 learning, there exists a discrepancy in training objectives between pre-training (*e.g.*, metric learning) and fine-tuning (Q-learning). On the other hand, in U2O RL, we pre-train a representation with unsupervised Q-learning (though with a different reward function), and thus the discrepancy between pre-training and fine-tuning becomes less severe.

### A.4 Can we do "bridging" without any reward-labeled data?

In the bridging phase of U2O RL (Section 4.2), we assume a (small) reward-labeled dataset $\mathcal{D}_{\texttt{reward}}$. In our experiments, we sample a small number of transitions (*e.g.*, 0.2% in the case of DMC) from the offline dataset and label them with the ground-truth reward function, as in prior works [76, 56]. However, these samples do not necessarily have to come from the offline dataset. To show this, we conduct an additional experiment where we do not assume access to any of the existing reward samples or the ground-truth



Figure 13: **U2O RL without using reward-labeling in the offline dataset (8 seeds).**

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 13. The results show that U2O still works and outperforms the supervised offline-to-online RL baseline.
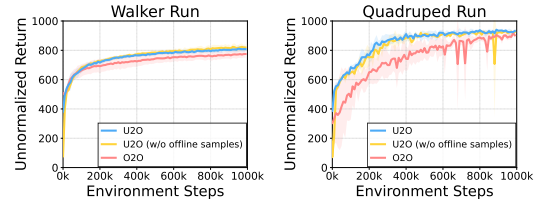
### A.5 How do different strategies of skill identification affect performance?

To understand how skill identification strategies affect online RL performance, we compare our strategy in Section 4.2 with an alternative strategy that simply selects a random latent vector $z$ from the skill space. Figure 14 shows that the skill identification with a randomly selected latent vector performs worse than our strategy. This is likely because modulating the policy with the best latent vector helps boost task-relevant exploration and information.
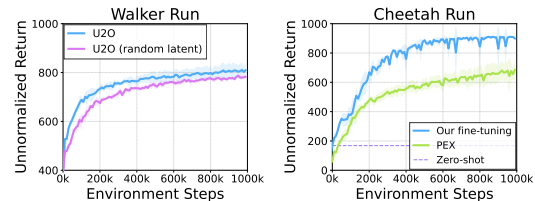


Figure 14: **Ablation study of skill identification (4 seeds).**



Figure 15: **Comparison with PEX and zero-shot 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 15 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.

## B   Experimental Details

For offline RL pre-training, we use 1M training steps for ExORL, AntMaze, and Adroit and 500K steps for Kitchen, following Park et al. [56]. For online fine-tuning, we use 1M additional environment steps for ExORL, AntMaze, and Adroit and 500K steps for Kitchen with an update-to-data ratio of 1. We implement U2O RL based on the official implementation of HILP [56]. 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 A5000 or RTX 3090 GPUs. Each run takes at most 40 hours (*e.g.* Visual Kitchen). We provide our implementation in the supplementary material.

### B.1   Environments and Datasets

**ExORL [85].**   In the ExORL benchmark, we consider four embodiments, Walker, Cheetah, Quadruped, and Jaco. Each embodiment has four tasks: Walker has {Run, Flip, Stand, Walk}, Cheetah has {Run, Run Backward, Walk, Walk Backward}, Quadruped has {Run, Jump, Stand, Walk}, and Jaco has {Reach Top Left, Reach Top Right, Reach Bottom Left, Reach Bottom Right}. For all the tasks in Walker, Cheetah, and Quadruped, the maximum return is 1000, and Jaco has 250. Each embodiment has an offline dataset, which is collected by running exploratory agents such as RND [5], and then annotated with task reward function. We use the first 5M transitions of the offline dataset following the prior work [76, 56]. The maximum episode length is 250 (Jaco) or 1000 (others).

**AntMaze [13, 27].**   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-defined neighborhood of the target position. We consider two types of Maze: `antmaze-large` [13] and `antmaze-ultra` [27], where the latter has twice the size of the former. Each maze has two types of offline datasets: `play` and `diverse`. The dataset consists 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.

**Kitchen [20, 13].**   In the Kitchen environment, a Franka robot should achieve four sub-tasks, `microwave`, `slide cabinet`, `light switch`, and `kettle`. Each task has a success criterion determined by an object configuration. Whenever the agent achieves a sub-task, a task reward of 1 is given, where the maximum return is 4. We consider two types of offline datasets: `mixed` and `partial`. We report normalized scores by multiplying the returns by 100. For Visual-Kitchen, we follow the same camera configuration as Mendonca et al. [47], Park et al. [57], and Park et al. [56], to render $64 \times 64$ RGB observations, which are used instead of low-dimensional states. We report normalized scores by multiplying the returns by 25.

**Adroit [13].**   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 [2, 40]. 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.

## B.2 Hyperparameters

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

| Hyperparameter | Value |
|---|---|
| Learning rate | 0.0005 (feature), 0.0001 (others) |
| Optimizer | Adam [28] |
| Minibatch size | 1024 |
| Feature MLP dimensions | $(512, 512)$ |
| Value MLP dimensions | $(1024, 1024, 1024)$ |
| Policy MLP dimensions | $(1024, 1024, 1024)$ |
| TD3 target smoothing coefficient | 0.01 |
| TD3 discount factor $\gamma$ | 0.98 |
| Latent dimension | 50 |
| State samples for latent vector inference | 10000 |
| Successor feature loss | Q loss |
| Hilbert representation discount factor | 0.96 (Walker), 0.98 (others) |
| Hilbert representation expectile | 0.5 |
| Hilbert representation target smoothing coefficient | 0.005 |

Table 4: Hyperparameters of unsupervised RL pre-training in AntMaze, Kitchen, and Adroit.

| Hyperparameter | Value |
|---|---|
| Learning rate | 0.0003 |
| Optimizer | Adam [28] |
| Minibatch size | 256 (Adroit), 512 (others) |
| Value MLP dimensions | $(256, 256, 256)$ (Adroit), $(512, 512, 512)$ (others) |
| Policy MLP dimensions | $(256, 256, 256)$ (Adroit), $(512, 512, 512)$ (others) |
| Target smoothing coefficient | 0.005 |
| Discount factor $\gamma$ | 0.99 |
| Latent dimension | 32 |
| Hilbert representation discount factor | 0.99 |
| Hilbert representation expectile | 0.95 |
| Hilbert representation target smoothing coefficient | 0.005 |
| HILP IQL expectile | 0.9 (AntMaze), 0.7 (others) |
| HILP AWR temperature | 0.5 (Kitchen) 3 (Adroit-door), 10 (others) |