# PEBBLE: Feedback-Efficient Interactive Reinforcement Learning via Relabeling Experience and Unsupervised Pre-training

Kimin Lee [* 1]    Laura Smith [* 1]    Pieter Abbeel [1]

## Abstract

Conveying complex objectives to reinforcement learning (RL) agents can often be difficult, involving meticulous design of reward functions that are sufficiently informative yet easy enough to provide. Human-in-the-loop RL methods allow practitioners to instead interactively teach agents through tailored feedback; however, such approaches have been challenging to scale since human feedback is very expensive. In this work, we aim to make this process more sample- and feedback-efficient. We present an off-policy, interactive RL algorithm that capitalizes on the strengths of both feedback and off-policy learning. Specifically, we learn a reward model by actively querying a teacher's preferences between two clips of behavior and use it to train an agent. To enable off-policy learning, we relabel all the agent's past experience when its reward model changes. We additionally show that pre-training our agents with unsupervised exploration substantially increases the mileage of its queries. We demonstrate that our approach is capable of learning tasks of higher complexity than previously considered by human-in-the-loop methods, including a variety of locomotion and robotic manipulation skills. We also show that our method is able to utilize real-time human feedback to effectively prevent reward exploitation and learn new behaviors that are difficult to specify with standard reward functions.

## 1. Introduction

Deep reinforcement learning (RL) has emerged as a powerful method whereby agents learn complex behaviors on their own through trial and error (Kohl & Stone, 2004;

*Equal contribution [1]University of California, Berkeley. Correspondence to: Kimin Lee <kiminlee@berkeley.edu>, Laura Smith <smithlaura@berkeley.edu>.

Kober & Peters, 2011; Kober et al., 2013; Silver et al., 2017; Andrychowicz et al., 2020; Kalashnikov et al., 2018; Vinyals et al., 2019). Scaling RL to many applications, however, is yet precluded by a number of challenges. One such challenge lies in providing a suitable reward function. For example, while it may be desirable to provide sparse rewards out of ease, they are often insufficient to train successful RL agents. Thus, to provide adequately dense signal, real-world problems may require extensive instrumentation, such as accelerometers to detect door opening (Yahya et al., 2017), thermal cameras to detect pouring (Schenck & Fox, 2017) or motion capture for object tracking (Kormushev et al., 2010; Akkaya et al., 2019; Peng et al., 2020).

Despite these costly measures, it may still be difficult to construct a suitable reward function due to reward exploitation. That is, RL algorithms often discover ways to achieve high returns by unexpected, unintended means. In general, there is nuance in how we might want agents to behave, such as obeying social norms, that are difficult to account for and communicate effectively through an engineered reward function (Amodei et al., 2016; Shah et al., 2019; Turner et al., 2020). A popular way to avoid reward engineering is through imitation learning, during which a learner distills information about its objectives or tries to directly follow an expert (Schaal, 1997; Ng et al., 2000; Abbeel & Ng, 2004; Argall et al., 2009). While imitation learning is a powerful tool, suitable demonstrations are often prohibitively expensive to obtain in practice (Calinon et al., 2009; Pastor et al., 2011; Akgun et al., 2012; Zhang et al., 2018).

In contrast, humans often learn fairly autonomously, relying on occasional external feedback from a teacher. Part of what makes a teacher effective is their ability to interactively guide students according to their progress, providing corrective or increasingly advanced instructions as needed. Such an interactive learning process is also alluring for artificial agents since the agent's behavior can naturally be tailored to one's preference (avoiding reward exploitation) without requiring extensive engineering. This approach is only feasible if the feedback is both practical for a human to provide and sufficiently high-bandwidth. As such, human-in-the-loop (HiL) RL (Knox & Stone, 2009; Christiano et al., 2017; MacGlashan et al., 2017) has not yet been widely adopted.
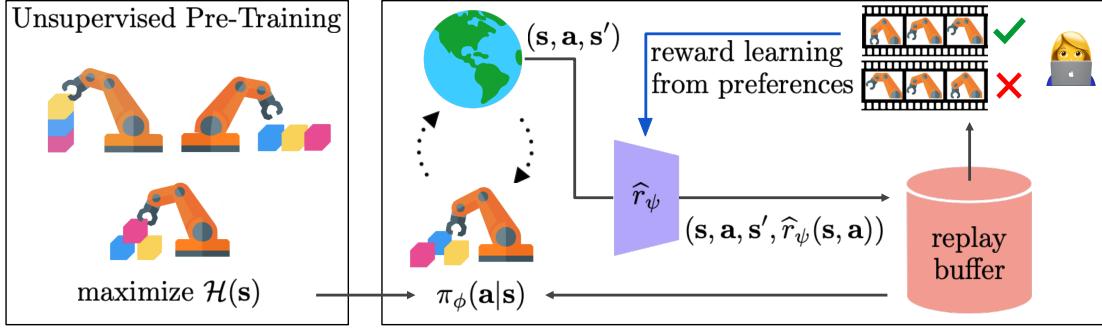
*Figure 1.* Illustration of our method. First, the agent engages in unsupervised pre-training during which it is encouraged to visit a diverse set of states so its queries can provide more meaningful signal than on randomly collected experience (left). Then, a teacher provides preferences between two clips of behavior, and we learn a reward model based on them. The agent is updated to maximize the expected return under the model. We also relabel all its past experiences with this model to maximize their utilization to update the policy (right).

In this work, we aim to substantially reduce the amount of human effort required for HiL learning. To this end, we present PEBBLE: unsupervised **PrE**-training and preference-**B**ased learning via rela**BeL**ing **E**xperience, a feedback-efficient RL algorithm by which learning is largely autonomous and supplemented by a practical number of binary labels (i.e. preferences) provided by a supervisor. Our method relies on two main, synergistic ingredients: *unsupervised pre-training* and *off-policy learning* (see Figure 1). For generality, we do not assume the agent is privy to rewards from its environment. Instead, we first allow the agent to explore using only intrinsic motivation ([Oudeyer et al., 2007](#); [Schmidhuber, 2010](#)) to diversify its experience and produce coherent behaviors. Collecting a breadth of experiences enables the teacher to provide more meaningful feedback, as compared to feedback on data collected in an indeliberate manner. The supervisor then steps in to teach the agent by expressing their preferences between pairs of clips of the agent's behavior ([Christiano et al., 2017](#)). The agent distills this information into a reward model and uses RL to optimize this inferred reward function.

Leveraging unsupervised pre-training increases the efficiency of the teacher's initial feedback; however, RL requires a large enough number of samples such that supervising the learning process is still quite expensive for humans. It is thus especially critical to enable off-policy algorithms that can reuse data to maximize the agent's, and thereby human's, efficiency. However, on-policy methods have typically been used thus far for HiL RL because of their ability to mitigate the effects of non-stationarity in reward caused by online learning. We show that by simply relabeling all of the agent's past experience every time the reward model is updated, we can make use and reuse of *all* the agent's collected experience to improve sample and feedback efficiency by a large margin. Source code and videos are available at `https://sites.google.com/view/icml21pebble`.

We summarize the main contributions of PEBBLE:

- For the first time, we show that *unsupervised pre-training* and *off-policy learning* can significantly improve the sample- and feedback-efficiency of HiL RL.

- PEBBLE outperforms prior preference-based RL baselines on complex locomotion and robotic manipulation tasks from DeepMind Control Suite (DMControl; [Tassa et al. 2018](#); [2020](#)) and Meta-world ([Yu et al., 2020](#)).

- We demonstrate that PEBBLE can learn behaviors for which a typical reward function is difficult to engineer very efficiently.

- We also show that PEBBLE can avoid reward exploitation, leading to more desirable behaviors compared to an agent trained with respect to an engineered reward function.

## 2. Related Work

**Learning from human feedback**. Several works have successfully utilized feedback from real humans to train agents where it is assumed that the feedback is available at *all* times ([Pilarski et al., 2011](#); [MacGlashan et al., 2017](#); [Arumugam et al., 2019](#)). Due to this high feedback frequency, these approaches are difficult to scale to more complex learning problems that require substantial agent experience.

Better suited to learning in complex domains is to learn a reward model so the agent can learn without a supervisor's perpetual presence. One simple yet effective direction in reward learning is to train a classifier that recognizes task success and use it as basis for a reward function ([Pinto & Gupta, 2016](#); [Levine et al., 2018](#); [Fu et al., 2018](#); [Xie et al., 2018](#)). Positive examples may be designated or reinforced through human feedback ([Zhang et al., 2019](#); [Singh et al., 2019](#); [Smith et al., 2020](#)). Another promising direction has focused on simply training a reward model via regres-

sion using unbounded real-valued feedback (Knox & Stone, 2009; Warnell et al., 2018), but this has been challenging to scale because it is difficult for humans to reliably provide a particular utility value for certain behaviors of the RL agent.

Much easier for humans is to make *relative* judgments, i.e., comparing behaviors as better or worse. Preference-based learning is thus an attractive alternative because the supervision is easy to provide yet information-rich (Akrour et al., 2011; Pilarski et al., 2011; Akrour et al., 2012; Wilson et al., 2012; Sugiyama et al., 2012; Wirth & Fürnkranz, 2013; Wirth et al., 2016; Sadigh et al., 2017; Biyik & Sadigh, 2018; Leike et al., 2018; Biyik et al., 2020). Christiano et al. (2017) scaled preference-based learning to utilize modern deep learning techniques—they learn a reward function, modeled with deep neural networks, that is consistent with the observed preferences and use it to optimize an agent using RL. They choose on-policy RL methods (Schulman et al., 2015; Mnih et al., 2016) since they are more robust to the non-stationarity in rewards caused by online learning. Although they demonstrate that preference-based learning provides a fairly efficient (requiring feedback on less than 1% of the agent's experience) means of distilling information from feedback, they rely on notoriously sample-inefficient on-policy RL, so a large burden can yet be placed on the human. Subsequent works have aimed to improve the efficiency of this method by introducing additional forms of feedback such as demonstrations (Ibarz et al., 2018) or non-binary rankings (Cao et al., 2020). Our proposed approach similarly focuses on developing a more sample- and feedback-efficient preference-based RL algorithm *without* adding any additional forms of supervision. Instead, we enable off-policy learning as well as utilize *unsupervised* pre-training to substantially improve efficiency.

**Unsupervised pre-training for RL**. Unsupervised pre-training has been studied for extracting strong behavioral priors that can be utilized to solve downstream tasks efficiently in the context of RL (Daniel et al., 2016; Florensa et al., 2018; Achiam et al., 2018; Eysenbach et al., 2019; Sharma et al., 2020). Specifically, agents are encouraged to expand the boundary of seen states by maximizing various intrinsic rewards, such as prediction errors (Houthooft et al., 2016; Pathak et al., 2017; Burda et al., 2019), count-based state novelty (Bellemare et al., 2016; Tang et al., 2017; Ostrovski et al., 2017), mutual information (Eysenbach et al., 2019) and state entropy (Hazan et al., 2019; Lee et al., 2019; Hao & Pieter, 2021). Such unsupervised pre-training methods allow learning diverse behaviors without extrinsic rewards, effectively facilitating accelerated learning of downstream tasks. In this work, we show that unsupervised pre-training enables a teacher to provide more meaningful signal by showing them a diverse set of behaviors.

# 3. Preliminaries

**Reinforcement learning**. We consider a standard RL framework where an agent interacts with an environment in discrete time. Formally, at each timestep $t$, the agent receives a state $\mathbf{s}_t$ from the environment and chooses an action $\mathbf{a}_t$ based on its policy $\pi$. The environment returns a reward $r_t$ and the agent transitions to the next state $\mathbf{s}_{t+1}$. The return $\mathcal{R}_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ is the discounted sum of rewards from timestep $t$ with discount factor $\gamma \in [0, 1)$. RL then maximizes the expected return from each state $\mathbf{s}_t$.

**Soft Actor-Critic**. SAC (Haarnoja et al., 2018) is an off-policy actor-critic method based on the maximum entropy RL framework (Ziebart, 2010), which encourages exploration and greater robustness to noise by maximizing a weighted objective of the reward and the policy entropy. To update the parameters, SAC alternates between a soft policy evaluation and a soft policy improvement. At the soft policy evaluation step, a soft Q-function, which is modeled as a neural network with parameters $\theta$, is updated by minimizing the following soft Bellman residual:

$$\mathcal{L}_{\text{critic}}^{\text{SAC}} = \mathbb{E}_{\tau_t \sim \mathcal{B}}\Big[ \big(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - r_t - \gamma \bar{V}(\mathbf{s}_{t+1})\big)^2 \Big], \quad (1)$$

$$\text{with} \quad \bar{V}(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi}\big[Q_{\bar{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - \alpha \log \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)\big],$$

where $\tau_t = (\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, r_t)$ is a transition, $\mathcal{B}$ is a replay buffer, $\bar{\theta}$ are the delayed parameters, and $\alpha$ is a temperature parameter. At the soft policy improvement step, the policy $\pi_\phi$ is updated by minimizing the following objective:

$$\mathcal{L}_{\text{act}}^{\text{SAC}} = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{B}, \mathbf{a}_t \sim \pi_\phi}\Big[\alpha \log \pi_\phi(\mathbf{a}_t|\mathbf{s}_t) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t)\Big]. \quad (2)$$

SAC enjoys good sample-efficiency relative to its on-policy counterparts by reusing its past experiences. However, for the same reason, SAC is not robust to a non-stationary reward function.

**Reward learning from preferences**. We follow the basic framework for learning a reward function $\hat{r}_\psi$ from preferences in which the function is trained to be consistent with observed human feedback (Wilson et al., 2012; Christiano et al., 2017). In this framework, a segment $\sigma$ is a sequence of observations and actions $\{\mathbf{s}_k, \mathbf{a}_k, ..., \mathbf{s}_{k+H}, \mathbf{a}_{k+H}\}$. We elicit preferences $y$ for segments $\sigma^0$ and $\sigma^1$, where $y$ is a distribution indicating which segment a human prefers, i.e., $y \in \{(0, 1), (1, 0), (0.5, 0.5)\}$. The judgment is recorded in a dataset $\mathcal{D}$ as a triple $(\sigma^0, \sigma^1, y)$. By following the Bradley-Terry model (Bradley & Terry, 1952), we model a preference predictor using the reward function $\hat{r}_\psi$ as follows:

$$P_\psi[\sigma^1 \succ \sigma^0] = \frac{\exp \sum_t \hat{r}_\psi(\mathbf{s}_t^1, \mathbf{a}_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}_\psi(\mathbf{s}_t^i, \mathbf{a}_t^i)}, \quad (3)$$

where $\sigma^i \succ \sigma^j$ denotes the event that segment $i$ is preferable to segment $j$. Intuitively, this can be interpreted as

assuming the probability of preferring a segment depends exponentially on the sum over the segment of an underlying reward function. While $\widehat{r}_\psi$ is not a binary classifier, learning $\widehat{r}_\psi$ amounts to binary classification with labels $y$ provided by a supervisor. Concretely, the reward function, modeled as a neural network with parameters $\psi$, is updated by minimizing the following loss:

$$\mathcal{L}^{\texttt{Reward}} = - \mathop{\mathbb{E}}_{(\sigma^0, \sigma^1, y) \sim \mathcal{D}} \Big[ y(0) \log P_\psi[\sigma^0 \succ \sigma^1] \\ + y(1) \log P_\psi[\sigma^1 \succ \sigma^0] \Big]. \quad (4)$$

## 4. PEBBLE

In this section, we present PEBBLE: unsupervised **PrE**-training and preference-**B**ased learning via rela**BeL**ing **E**xperience, an off-policy actor-critic algorithm for HiL RL. Formally, we consider a policy $\pi_\phi$, Q-function $Q_\theta$ and reward function $\widehat{r}_\psi$, which are updated by the following processes (see Algorithm 2 for the full procedure):

- *Step 0 (unsupervised pre-training)*: We pre-train the policy $\pi_\phi$ only using intrinsic motivation to explore and collect diverse experiences (see Section 4.1).

- *Step 1 (reward learning)*: We learn a reward function $\widehat{r}_\psi$ that can lead to the desired behavior by getting feedback from a teacher (see Section 4.2).

- *Step 2 (agent learning)*: We update the policy $\pi_\phi$ and Q-function $Q_\theta$ using an off-policy RL algorithm with relabeling to mitigate the effects of a non-stationary reward function $\widehat{r}_\psi$ (see Section 4.3).

- Repeat *Step 1* and *Step 2*.

### 4.1. Accelerating Learning via Unsupervised Pre-training

In our setting, we assume the agent is given feedback in the form of preferences between segments. In the beginning of training, though, a naive agent executes a random policy, which does not provide good state coverage nor coherent behaviors. The agent's queries are thus quite limited and likely difficult for human teachers to judge. As a result, it requires many samples (and thus queries) for these methods to show initial progress. Recent work has addressed this issue by means of providing demonstrations; however, this is not ideal since these are notoriously hard to procure (Ibarz et al., 2018). Instead, our insight is to produce informative queries at the start of training by utilizing unsupervised pre-training to collect diverse samples solely through intrinsic motivation (Oudeyer et al., 2007; Schmidhuber, 2010).

Specifically, we encourage our agent to visit a wider range of states by using the state entropy $\mathcal{H}(\mathbf{s}) =$

---

**Algorithm 1** EXPLORE: Unsupervised exploration

1: Initialize parameters of $Q_\theta$ and $\pi_\phi$ and a replay buffer $\mathcal{B} \leftarrow \emptyset$
2: **for** each iteration **do**
3:     **for** each timestep $t$ **do**
4:         Collect $\mathbf{s}_{t+1}$ by taking $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t|\mathbf{s}_t)$
5:         Compute intrinsic reward $r_t^{\texttt{int}} \leftarrow r^{\texttt{int}}(\mathbf{s}_t)$ as in (5)
6:         Store transitions $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, r_t^{\texttt{int}})\}$
7:     **end for**
8:     **for** each gradient step **do**
9:         Sample minibatch $\{(\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}_{j+1}, r_j^{\texttt{int}})\}_{j=1}^B \sim \mathcal{B}$
10:        Optimize $\mathcal{L}_{\texttt{critic}}^{\texttt{SAC}}$ in (1) and $\mathcal{L}_{\texttt{act}}^{\texttt{SAC}}$ in (2) with respect to $\theta$ and $\phi$
11:     **end for**
12: **end for**
13: **return** $\mathcal{B}, \pi_\phi$

---

$-\mathbb{E}_{\mathbf{s} \sim p(\mathbf{s})}[\log p(\mathbf{s})]$ as an intrinsic reward (Hazan et al., 2019; Lee et al., 2019; Hao & Pieter, 2021; Seo et al., 2021). By updating the agent to maximize the sum of expected intrinsic rewards, it can efficiently explore an environment and learn how to generate diverse behaviors. However, this intrinsic reward is intractable to compute in most settings. To handle this issue, we employ the simplified version of particle-based entropy estimator (Beirlant et al., 1997; Singh et al., 2003) (see the supplementary material for more details):

$$\widehat{\mathcal{H}}(\mathbf{s}) \propto \sum_i \log(||\mathbf{s}_i - \mathbf{s}_i^k||),$$

where $\widehat{\mathcal{H}}$ denotes the particle-based entropy estimator and $\mathbf{s}_i^k$ is the $k$-th nearest neighbor ($k$-NN) of $\mathbf{s}_i$. This implies that maximizing the distance between a state and its nearest neighbor increases the overall state entropy. Inspired by this, we define the intrinsic reward of the current state $\mathbf{s}_t$ as the distance between $\mathbf{s}_t$ and its $k$-th nearest neighbor by following the idea of Hao & Pieter (2021) that treats each transition as a particle:

$$r^{\texttt{int}}(\mathbf{s}_t) = \log(||\mathbf{s}_t - \mathbf{s}_t^k||). \quad (5)$$

In our experiments, we compute $k$-NN distances between a sample and all samples in the replay buffer and normalize the intrinsic reward by dividing it by a running estimate of the standard deviation. The full procedure of unsupervised pre-training is summarized in Algorithm 1.

### 4.2. Selecting Informative Queries

As previously mentioned, we learn our reward function by modeling the probability that a teacher prefers one sampled segment over another as proportional to the exponentiated sum of rewards over the segment (see Section 3). Ideally, one should solicit preferences so as to maximize *expected value of information* (EVOI; Savage 1972): the improvement of an agent caused by optimizing with respect to the resulting reward model (Viappiani, 2012; Akrour et al., 2012).

**Algorithm 2** PEBBLE

**Require:** frequency of teacher feedback $K$
**Require:** number of queries $M$ per feedback session
1: Initialize parameters of $Q_\theta$ and $\widehat{r}_\psi$
2: Initialize a dataset of preferences $\mathcal{D} \leftarrow \emptyset$
3: // EXPLORATION PHASE
4: $\mathcal{B}, \pi_\phi \leftarrow$ EXPLORE() in Algorithm 1
5: // POLICY LEARNING
6: **for** each iteration **do**
7:      // REWARD LEARNING
8:      **if** iteration % $K == 0$ **then**
9:          **for** $m$ in $1 \ldots M$ **do**
10:            $(\sigma^0, \sigma^1) \sim$ SAMPLE() (see Section 4.2)
11:            Query instructor for $y$
12:            Store preference $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\sigma^0, \sigma^1, y)\}$
13:          **end for**
14:          **for** each gradient step **do**
15:            Sample minibatch $\{(\sigma^0, \sigma^1, y)_j\}_{j=1}^{D} \sim \mathcal{D}$
16:            Optimize $\mathcal{L}^{\text{Reward}}$ in (4) with respect to $\psi$
17:          **end for**
18:          Relabel entire replay buffer $\mathcal{B}$ using $\widehat{r}_\psi$
19:      **end if**
20:      **for** each timestep $t$ **do**
21:          Collect $\mathbf{s}_{t+1}$ by taking $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$
22:          Store transitions $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, \widehat{r}_\psi(\mathbf{s}_t))\}$
23:      **end for**
24:      **for** each gradient step **do**
25:          Sample random minibatch $\{(\tau_j)\}_{j=1}^{B} \sim \mathcal{B}$
26:          Optimize $\mathcal{L}^{\text{SAC}}_{\text{critic}}$ in (1) and $\mathcal{L}^{\text{SAC}}_{\text{act}}$ in (2) with respect to $\theta$ and $\phi$, respectively
27:      **end for**
28: **end for**



*Figure 2.* Examples from the environments we test on. We consider learning a variety of complex locomotion and manipulation skills through interacting with a scripted or human trainer.

Computing the EVOI is intractable since it involves taking an expectation over all possible trajectories induced by the updated policy. To handle this issue, several approximations have been explored by prior works to sample queries that are likely to change the reward model (Daniel et al., 2014; Christiano et al., 2017; Ibarz et al., 2018). In this work, we consider the sampling schemes employed by Christiano et al. (2017): (1) uniform sampling and (2) ensemble-based sampling, which selects pairs of segments with high variance across ensemble reward models. We explore an additional third method, entropy-based sampling, which seeks to disambiguate pairs of segments nearest the decision boundary. That is, we sample a large batch of segment pairs and select pairs that maximize $\mathcal{H}(P_\psi)$. We evaluate the effects of these sampling methods in Section 5.

**4.3. Using Off-policy RL with Non-Stationary Reward**

Once we learn a reward function $\widehat{r}_\psi$, we can update the policy $\pi_\phi$ and Q-function $Q_\theta$ using any RL algorithm. A caveat is that the reward function $\widehat{r}_\psi$ may be non-stationary because we update it during training. Christiano et al. (2017) used on-policy RL algorithms, TRPO (Schulman et al., 2015) and A2C (Mnih et al., 2016), to address this issue. However, their poor sample-efficiency leads to poor feedback-efficiency of the overall HiL method.
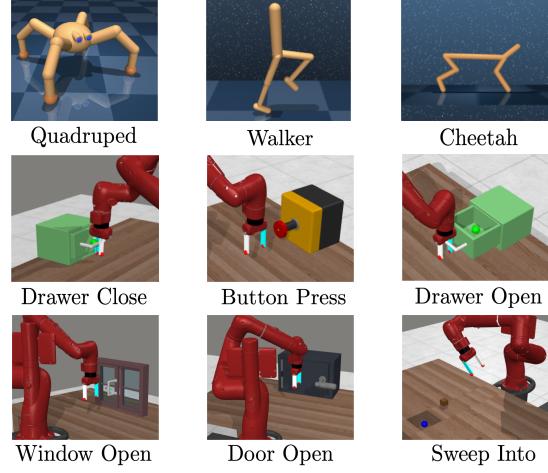
In this work, we use an off-policy RL algorithm, which provides for sample-efficient learning by reusing past experiences that are stored in the replay buffer. However, the learning process of off-policy RL algorithms can be unstable because previous experiences in the replay buffer are labeled with previous learned rewards. To handle this issue, we relabel all of the agent's past experience every time we update the reward model. We find that this simple technique stabilizes the learning process and provides large gains in performance (see Figure 5(a) for supporting results). The full procedure of PEBBLE is summarized in Algorithm 2.

## 5. Experiments

We design our experiments to investigate the following:

1. How does PEBBLE compare to existing methods in terms of sample and feedback efficiency?

2. What is the contribution of each of the proposed techniques in PEBBLE?

3. Can PEBBLE learn novel behaviors for which a typical reward function is difficult to engineer?

4. Can PEBBLE mitigate the effects of reward exploitation?

**5.1. Setups**

We evaluate PEBBLE on several continuous control tasks involving locomotion and robotic manipulation from DeepMind Control Suite (DMControl; Tassa et al. 2018; 2020) and Meta-world (Yu et al., 2020). In order to verify the efficacy of our method, we first focus on having an agent solve a range of tasks without being able to directly observe the
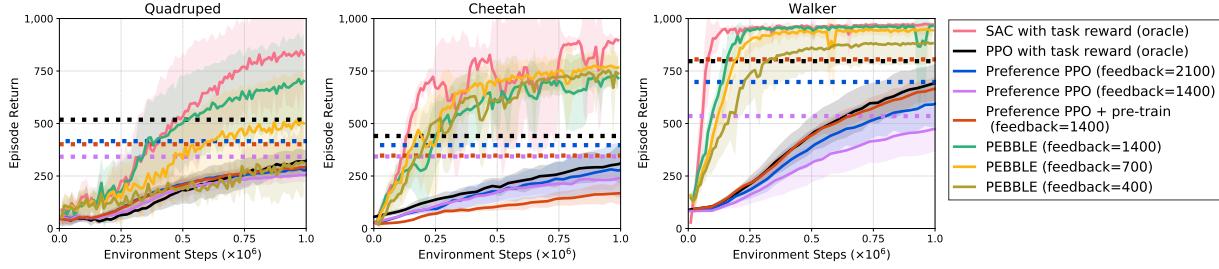
*Figure 3.* Learning curves on locomotion tasks as measured on the ground truth reward. The solid line and shaded regions represent the mean and standard deviation, respectively, across ten runs. Asymptotic performance of PPO and Preference PPO is indicated by dotted lines of the corresponding color.

ground truth reward function. Instead, similar to Christiano et al. (2017) and Ibarz et al. (2018), the agent learns to perform a task only by getting feedback from a scripted teacher that provides preferences between trajectory segments according to the true, underlying task reward. Because this scripted teacher's preferences are immediately generated by a ground truth reward, we are able to evaluate the agent quantitatively by measuring the true average return and do more rapid experiments. For all experiments, we report the mean and standard deviation across ten runs.

We also run experiments with actual human trainers (the authors) to show the benefits of human-in-the-loop RL. First, we show that human trainers can teach novel behaviors (e.g., waving a leg), which are not defined in original benchmarks. Second, we show that agents trained with the hand-engineered rewards from benchmarks can perform the task in an undesirable way (i.e., the agent exploits a misspecified reward function), while agents trained using human feedback can perform the same task in the desired way. For all experiments, each trajectory segment is presented to the human as a 1 second video clip, and a maximum of one hour of human time is required.

For evaluation, we compare to Christiano et al. (2017), which is the current state-of-the-art approach using the same type of feedback. The primary differences in our method are (1) the introduction of unsupervised pre-training, (2) the accommodation of *off-policy* RL, and (3) entropy-based sampling. We re-implemented Christiano et al. (2017) using the state-of-the-art on-policy RL algorithm: PPO (Schulman et al., 2017). We use the same reward learning framework and ensemble disagreement-based sampling as they proposed. We refer to this baseline as Preference PPO.

As an upper bound, since we evaluate against the task reward function, we also compare to SAC (Haarnoja et al., 2018) and PPO using the same ground truth reward. For our method, we pre-train an agent for 10K timesteps and include these pre-training steps in all learning curves. We do not alter any hyperparameters of the original SAC algorithm and

use an ensemble of three reward models. Unless stated otherwise, we use entropy-based sampling. More experimental details including model architectures, sampling schemes, and reward learning are in the supplementary material.

### 5.2. Benchmark Tasks with Unobserved Rewards

**Locomotion tasks from DMControl**. Figure 3 shows the learning curves of PEBBLE with 1400, 700 or 400 pieces of feedback[1] and that of Preference PPO with 2100 or 1400 pieces of feedback on three complex environments: Cheetah-run, Walker-walk and Quadruped-walk. Note that we explicitly give preference PPO an advantage by providing it with more feedback. We find that given a budget of 1400 queries, PEBBLE (green) reaches the same performance as SAC (pink) while Preference PPO (purple) is unable to match PPO (black). That PEBBLE requires less feedback than Preference PPO to match its respective oracle performance corroborates that PEBBLE is indeed more feedback-efficient. These results demonstrate that PEBBLE can enable the agent to solve the tasks without directly observing the ground truth reward function.

For further analysis, we incorporated our pre-training with Preference PPO (red) and find that it improves performance for Quadruped and Walker. We emphasize that our insight of using pre-training is able to improve both methods in terms of feedback-efficiency and asymptotic performance, but PEBBLE is uniquely positioned to benefit as it is able to utilize unsupervised experience for policy learning.

**Robotic manipulation tasks from Meta-world**. One application area in which HiL methods could have significant real-world impact is robotic manipulation, since learning often requires extensive engineering in the real world (Yahya et al., 2017; Schenck & Fox, 2017; Kormushev et al., 2010; Rusu et al., 2017; Akkaya et al., 2019; Peng et al., 2020). However, the common approach is to perform goal-conditioned learning with classifiers (Singh et al., 2019),

---

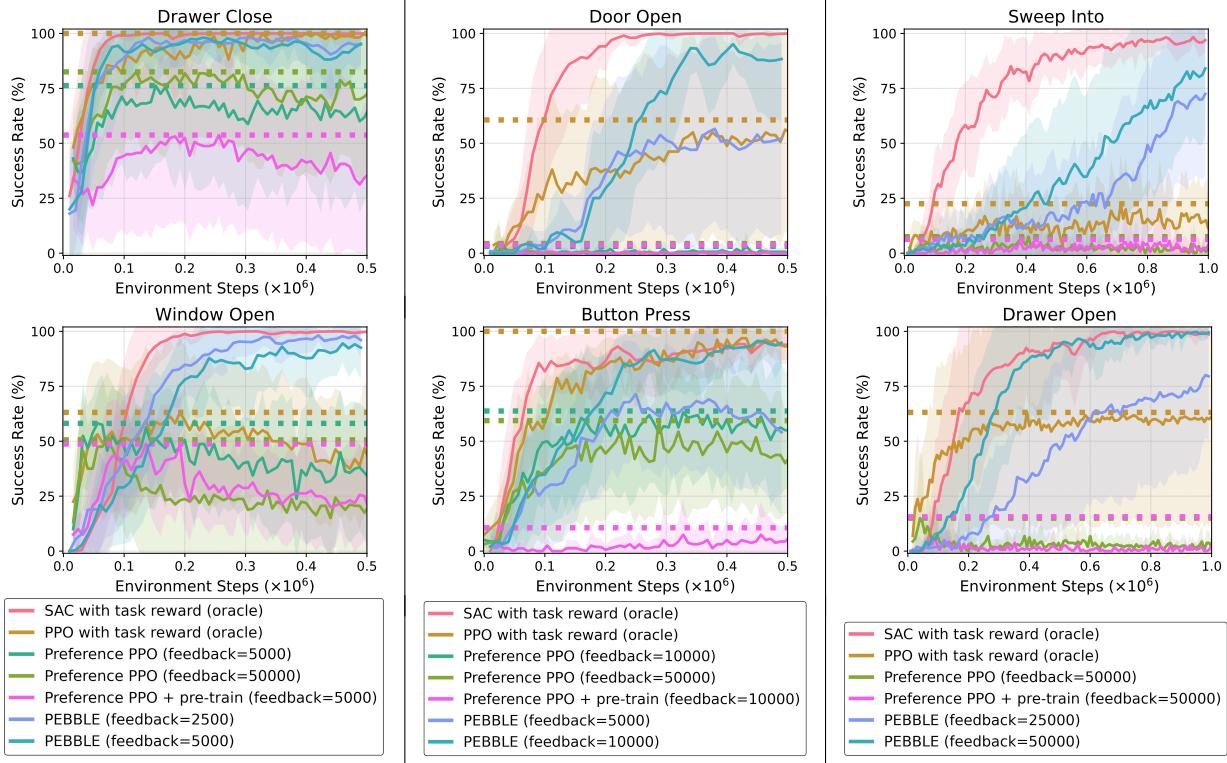[1]One piece of feedback corresponds to one preference query.

*Figure 4.* Learning curves on robotic manipulation tasks as measured on the success rate. The solid line and shaded regions represent the mean and standard deviation, respectively, across ten runs. Asymptotic performance of PPO and Preference PPO is indicated by dotted lines of the corresponding color.
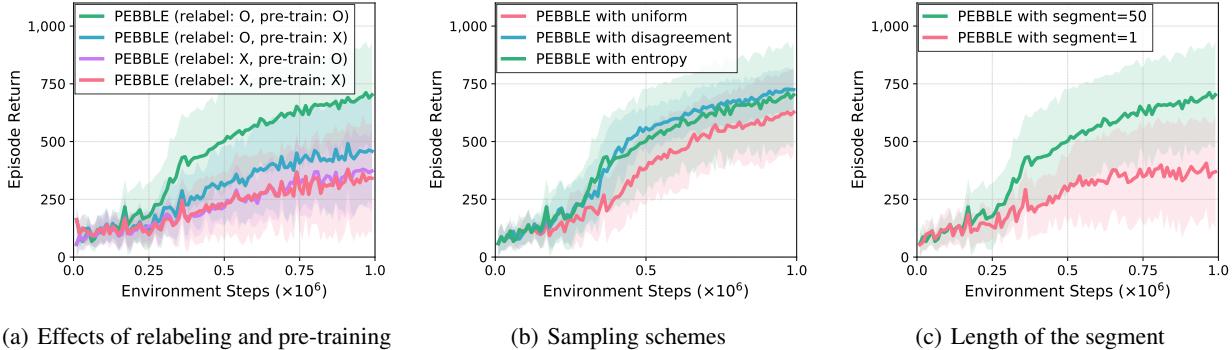


(a) Effects of relabeling and pre-training    (b) Sampling schemes    (c) Length of the segment

*Figure 5.* Ablation study on Quadruped-walk. (a) Contribution of each technique in PEBBLE, i.e., relabeling the replay buffer (relabel) and unsupervised pre-training (pre-train). (b) Effects of sampling schemes to select queries. (c) PEBBLE with varying the length of the segment. The results show the mean and standard deviation averaged over ten runs.

which can only capture limited information about what goal states *are*, and not about *how* they can be achieved. To study how we can utilize preference-based learning to perform more complex skills, we also consider six tasks covering a range of fundamental robotic manipulation skills from Meta-world (see Figure 2). As shown in Figure 4, PEB-BLE matches the performance of SAC using the ground truth reward and outperforms Preference PPO, given comparable (and more) feedback, on every task. By demonstrating

the applicability of PEBBLE to learning a variety of robotic manipulation tasks, we believe that we are taking an important step towards anyone (non-experts included) being able to teach robots in real-world settings.

### 5.3. Ablation Study

**Contribution of each technique**. In order to evaluate the individual effects of each technique in PEBBLE, we incre-

Clock-wise windmill      Counter clock-wise windmill

Quadruped waving its left front leg      Quadruped waving its right front leg      Hopper backflip
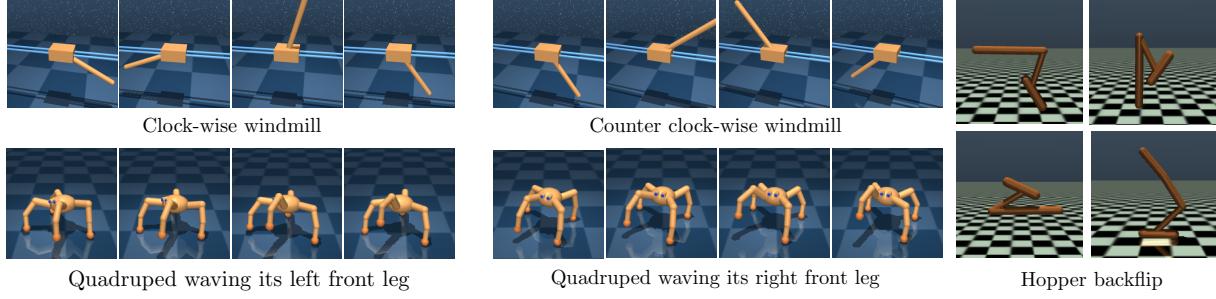
*Figure 6.* Novel behaviors trained using feedback from human trainers. The corresponding videos and examples of selected queries are available at the supplementary material.
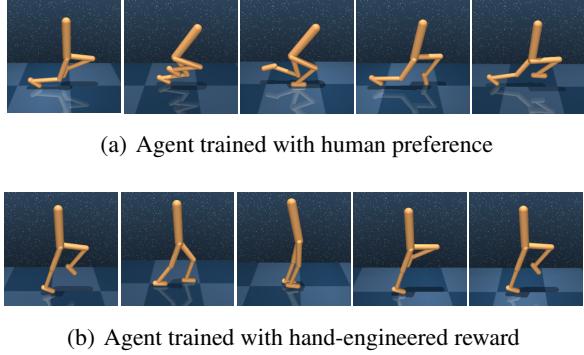


(a) Agent trained with human preference



(b) Agent trained with hand-engineered reward

*Figure 7.* Five frames from agents trained with (a) human preference and (b) hand-engineered reward from DMControl benchmark.

mentally apply unsupervised pre-training and relabeling. Figure 5(a) shows the learning curves of PEBBLE with 1400 queries on Quadruped-walk. First, we remark that relabeling significantly improves performance because it enables the agent to be robust to changes in its reward model. By additionally utilizing unsupervised pre-training, both sample-efficiency and asymptotic performance of PEB-BLE are further improved because showing diverse behaviors to a teacher can induce a better-shaped reward. This shows that PEBBLE's key ingredients are fruitfully wed, and their unique combination is crucial to our method's success.

**Effects of sampling schemes**. We also analyze the effects of different sampling schemes to select queries. Figure 5(b) shows the learning curves of PEBBLE with three different sampling schemes: uniform sampling, disagreement sampling and entropy sampling on Quadruped-walk. For this complex domain, we find that the uncertainty-based sampling schemes (using ensemble disagreement or entropy) are superior to the naive uniform sampling scheme. However, we note that they did not lead to extra gains on relatively simple environments, like Walker and Cheetah, similar to observations from Ibarz et al. (2018) (see the supplementary material for more results).

**Comparison with step-wise feedback**. We also measure the performance of PEBBLE by varying the length of segments. Figure 5(c) shows that feedback from longer segments (green curve) provide more meaningful signal than step-wise feedback (red curve). We believe that this is because longer segments can provide more context in reward learning.

### 5.4. Human Experiments

**Novel behaviors**. We show that agents can perform various novel behaviors based on human feedback using PEBBLE in Figure 6. Specifically, we demonstrate (a) the Cart agent swinging a pole (using 50 queries), (b) the Quadruped agent waving a front leg (using 200 queries), and (c) the Hopper performing a backflip (using 50 queries). We note that the human is indeed able to guide the agent in a controlled way, as evidenced by training the same agent to perform several variations of the same task (e.g., waving different legs or spinning in opposite directions). The videos of all behaviors and examples of selected queries are available in the supplementary material.

**Reward exploitation**. One concern in utilizing hand-engineered rewards is that an agent can exploit unexpected sources of reward, leading to unintended behaviors. Indeed, we find that the Walker agent learns to walk using only one leg even though it achieves the maximum scores as shown in Figure 7(b). However, using 200 human queries, we were able to train the Walker to walk in a more natural, human-like manner (using both legs) as shown in Figure 7(a). This result clearly shows the advantage of HiL RL to avoid reward exploitation.

## 6. Discussion

In this work, we present PEBBLE, a feedback-efficient algorithm for HiL RL. By leveraging unsupervised pre-training and off-policy learning, we show that sample- and feedback-efficiency of HiL RL can be significantly improved and this framework can be applied to tasks of higher complexity

than previously considered by previous methods, including a variety of locomotion and robotic manipulation skills. Additionally, we demonstrate that PEBBLE can learn novel behaviors and avoid reward exploitation, leading to more desirable behaviors compared to an agent trained with respect to an engineered reward function. We believe by making preference-based learning more tractable, PEBBLE may facilitate broadening the impact of RL beyond settings in which experts can carefully craft reward functions to those in which laypeople can likewise utilize the advances of robot learning in the real world.

## Acknowledgements

## References

Abbeel, P. and Ng, A. Y. Apprenticeship learning via inverse reinforcement learning. In *International Conference on Machine Learning*, 2004.

Achiam, J., Edwards, H., Amodei, D., and Abbeel, P. Variational option discovery algorithms. *arXiv preprint arXiv:1807.10299*, 2018.

Akgun, B., Cakmak, M., Yoo, J., and Thomaz, A. Trajectories and keyframes for kinesthetic teaching: A human-robot interaction perspective. In *International Conference on Human-Robot Interaction*, 2012.

Akkaya, I., Andrychowicz, M., Chociej, M., Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., et al. Solving rubik's cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.

Akrour, R., Schoenauer, M., and Sebag, M. Preference-based policy learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2011.

Akrour, R., Schoenauer, M., and Sebag, M. April: Active preference learning-based reinforcement learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2012.

Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.

Andrychowicz, O. M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.

Argall, B. D., Chernova, S., Veloso, M., and Browning, B. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57(5):469–483, 2009.

Arumugam, D., Lee, J. K., Saskin, S., and Littman, M. L. Deep reinforcement learning from policy-dependent human feedback. *arXiv preprint arXiv:1902.04257*, 2019.

Beirlant, J., Dudewicz, E. J., Györfi, L., and Van der Meulen, E. C. Nonparametric entropy estimation: An overview. *International Journal of Mathematical and Statistical Sciences*, 6(1):17–39, 1997.

Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems*, 2016.

Biyik, E. and Sadigh, D. Batch active preference-based learning of reward functions. In *Conference on Robot Learning*, 2018.

Biyik, E., Huynh, N., Kochenderfer, M. J., and Sadigh, D. Active preference-based gaussian process regression for reward learning. In *Robotics: Science and Systems*, 2020.

Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.

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

Calinon, S., Evrard, P., Gribovskaya, E., Billard, A., and Kheddar, A. Learning collaborative manipulation tasks by demonstration using a haptic interface. In *International Conference on Advanced Robotics*, 2009.

Cao, Z., Wong, K., and Lin, C.-T. Human preference scaling with demonstrations for deep reinforcement learning. *arXiv preprint arXiv:2007.12904*, 2020.

Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, 2017.

Daniel, C., Viering, M., Metz, J., Kroemer, O., and Peters, J. Active reward learning. In *Robotics: Science and Systems*, 2014.

Daniel, C., Neumann, G., Kroemer, O., and Peters, J. Hierarchical relative entropy policy search. *The Journal of Machine Learning Research*, 17(1):3190–3239, 2016.

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

Florensa, C., Duan, Y., and Abbeel, P. Stochastic neural networks for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2018.

Fu, J., Singh, A., Ghosh, D., Yang, L., and Levine, S. Variational inverse control with events: A general framework for data-driven reward definition. In *Advances in Neural Information Processing Systems*, 2018.

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

Hao, L. and Pieter, A. Behavior from the void: Unsupervised active pre-training. In *International Conference on Machine Learning*, 2021.

Hazan, E., Kakade, S., Singh, K., and Van Soest, A. Provably efficient maximum entropy exploration. In *International Conference on Machine Learning*, 2019.

Houthooft, R., Chen, X., Duan, Y., Schulman, J., De Turck, F., and Abbeel, P. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems*, 2016.

Ibarz, B., Leike, J., Pohlen, T., Irving, G., Legg, S., and Amodei, D. Reward learning from human preferences and demonstrations in atari. In *Advances in Neural Information Processing Systems*, 2018.

Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, 2018.

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

Knox, W. B. and Stone, P. Interactively shaping agents via human reinforcement: The tamer framework. In *International Conference on Knowledge Capture*, 2009.

Kober, J. and Peters, J. Policy search for motor primitives in robotics. *Machine learning*, 84(1-2):171–203, 2011.

Kober, J., Bagnell, J. A., and Peters, J. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.

Kohl, N. and Stone, P. Policy gradient reinforcement learning for fast quadrupedal locomotion. In *International Conference on Robotics and Automation*, 2004.

Kormushev, P., Calinon, S., and Caldwell, D. Robot motor skill coordination with EM-based reinforcement learning. In *International Conference on Intelligent Robots and Systems*, 2010.

Lee, L., Eysenbach, B., Parisotto, E., Xing, E., Levine, S., and Salakhutdinov, R. Efficient exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.

Leike, J., Krueger, D., Everitt, T., Martic, M., Maini, V., and Legg, S. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.

Levine, S., Pastor, P., Krizhevsky, A., Ibarz, J., and Quillen, D. Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection. *The International Journal of Robotics Research*, 37(4-5):421–436, 2018.

MacGlashan, J., Ho, M. K., Loftin, R., Peng, B., Roberts, D., Taylor, M. E., and Littman, M. L. Interactive learning from policy-dependent human feedback. In *International Conference on Machine Learning*, 2017.

Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, 2016.

Ng, A. Y., Russell, S. J., et al. Algorithms for inverse reinforcement learning. In *International Conference on Machine Learning*, 2000.

Ostrovski, G., Bellemare, M. G., Oord, A. v. d., and Munos, R. Count-based exploration with neural density models. In *International Conference on Machine Learning*, 2017.

Oudeyer, P.-Y., Kaplan, F., and Hafner, V. V. Intrinsic motivation systems for autonomous mental development. *IEEE transactions on evolutionary computation*, 11(2): 265–286, 2007.

Pastor, P., Righetti, L., Kalakrishnan, M., and Schaal, S. Online movement adaptation based on previous sensor experiences. In *International Conference on Intelligent Robots and Systems*, 2011.

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

Peng, X. B., Coumans, E., Zhang, T., Lee, T.-W., Tan, J., and Levine, S. Learning agile robotic locomotion skills by imitating animals. In *Robotics: Science and Systems*, 2020.

Pilarski, P. M., Dawson, M. R., Degris, T., Fahimi, F., Carey, J. P., and Sutton, R. S. Online human training of a myoelectric prosthesis controller via actor-critic reinforcement learning. In *International Conference on Rehabilitation Robotics*, 2011.

Pinto, L. and Gupta, A. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *International Conference on Robotics and Automation*, 2016.

Rusu, A., Večerík, M., Rothörl, T., Heess, N., Pascanu, R., and Hadsell, R. Sim-to-real robot learning from pixels with progressive nets. In *Conference on Robot Learning*, 2017.

Sadigh, D., Dragan, A. D., Sastry, S., and Seshia, S. A. Active preference-based learning of reward functions. In *Robotics: Science and Systems*, 2017.

Savage, L. J. *The foundations of statistics*. Courier Corporation, 1972.

Schaal, S. Learning from demonstration. In *Advances in Neural Information Processing Systems*, 1997.

Schenck, C. and Fox, D. Visual closed-loop control for pouring liquids. In *International Conference on Robotics and Automation*, 2017.

Schmidhuber, J. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010.

Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust region policy optimization. In *International Conference on Machine Learning*, 2015.

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Seo, Y., Chen, L., Shin, J., Lee, H., Abbeel, P., and Lee, K. State entropy maximization with random encoders for efficient exploration. In *International Conference on Machine Learning*, 2021.

Shah, R., Krasheninnikov, D., Alexander, J., Abbeel, P., and Dragan, A. Preferences implicit in the state of the world. In *International Conference on Learning Representations*, 2019.

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

Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.

Singh, A., Yang, L., Hartikainen, K., Finn, C., and Levine, S. End-to-end robotic reinforcement learning without reward engineering. In *Robotics: Science and Systems*, 2019.

Singh, H., Misra, N., Hnizdo, V., Fedorowicz, A., and Demchuk, E. Nearest neighbor estimates of entropy. *American journal of mathematical and management sciences*, 23 (3-4):301–321, 2003.

Smith, L., Dhawan, N., Zhang, M., Abbeel, P., and Levine, S. Avid: Learning multi-stage tasks via pixel-level translation of human videos. In *Robotics: Science and Systems*, 2020.

Sugiyama, H., Meguro, T., and Minami, Y. Preference-learning based inverse reinforcement learning for dialog control. In *Conference of the International Speech Communication Association*, 2012.

Tang, H., Houthooft, R., Foote, D., Stooke, A., Chen, X., Duan, Y., Schulman, J., De Turck, F., and Abbeel, P. # exploration: A study of count-based exploration for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2017.

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

Tassa, Y., Tunyasuvunakool, S., Muldal, A., Doron, Y., Liu, S., Bohez, S., Merel, J., Erez, T., Lillicrap, T., and Heess, N. dm_control: Software and tasks for continuous control. *arXiv preprint arXiv:2006.12983*, 2020.

Turner, A. M., Ratzlaff, N., and Tadepalli, P. Avoiding side effects in complex environments. *arXiv preprint arXiv:2006.06547*, 2020.

Viappiani, P. Monte carlo methods for preference learning. In *International Conference on Learning and Intelligent Optimization*, 2012.

Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575 (7782):350–354, 2019.

Warnell, G., Waytowich, N., Lawhern, V., and Stone, P. Deep tamer: Interactive agent shaping in high-dimensional state spaces. In *Conference on Artificial Intelligence*, 2018.

Wilson, A., Fern, A., and Tadepalli, P. A bayesian approach for policy learning from trajectory preference queries. In *Advances in Neural Information Processing Systems*, 2012.

Wirth, C. and Fürnkranz, J. Preference-based reinforcement learning: A preliminary survey. In *ECML/PKDD Workshop on Reinforcement Learning from Generalized Feedback: Beyond Numeric Rewards*, 2013.

Wirth, C., Fürnkranz, J., and Neumann, G. Model-free preference-based reinforcement learning. In *Conference on Artificial Intelligence*, 2016.

Xie, A., Singh, A., Levine, S., and Finn, C. Few-shot goal inference for visuomotor learning and planning. In *Conference on Robot Learning*, 2018.

Yahya, A., Li, A., Kalakrishnan, M., Chebotar, Y., and Levine, S. Collective robot reinforcement learning with distributed asynchronous guided policy search. In *International Conference on Intelligent Robots and Systems*, 2017.

Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on Robot Learning*, 2020.

Zhang, M., Vikram, S., Smith, L., Abbeel, P., Johnson, M., and Levine, S. Solar: Deep structured representations for model-based reinforcement learning. In *International Conference on Machine Learning*, 2019.

Zhang, T., McCarthy, Z., Jow, O., Lee, D., Goldberg, K., and Abbeel, P. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *International Conference on Robotics and Automation*, 2018.

Ziebart, B. D. Modeling purposeful adaptive behavior with the principle of maximum causal entropy. 2010.

# Appendix

## A. State Entropy Estimator

To approximate state entropy, we employ the simplified version of particle-based entropy estimator (Beirlant et al., 1997; Singh et al., 2003). Specifically, let $\mathbf{s}$ be a random variable with a probability density function $p$ whose support is a set $\mathcal{S} \subset \mathbb{R}^q$. Then its differential entropy is given as $\mathcal{H}(\mathbf{s}) = -\mathbb{E}_{\mathbf{s} \sim p(\mathbf{s})}[\log p(\mathbf{s})]$. When the distribution $p$ is not available, this quantity can be estimated given $N$ i.i.d realizations of $\{\mathbf{s}_i\}_{i=1}^N$ (Beirlant et al., 1997). However, since it is difficult to estimate $p$ with high-dimensional data, particle-based $k$-nearest neighbors ($k$-NN) entropy estimator (Singh et al., 2003) can be employed:

$$\widehat{\mathcal{H}}(\mathbf{s}) = \frac{1}{N} \sum_{i=1}^N \log \frac{N \cdot ||\mathbf{s}_i - \mathbf{s}_i^k||_2^q \cdot \widehat{\pi}^{\frac{q}{2}}}{k \cdot \Gamma(\frac{q}{2}+1)} + C_k \tag{6}$$

$$\propto \frac{1}{N} \sum_{i=1}^N \log ||\mathbf{s}_i - \mathbf{s}_i^k||_2, \tag{7}$$

where $\widehat{\pi}$ is the ratio of a circle's circumference to its diameter, $\mathbf{s}_i^k$ is the $k$-NN of $\mathbf{s}_i$ within a set $\{\mathbf{s}_i\}_{i=1}^N$, $C_k = \log k - \Psi(k)$ a bias correction term, $\Psi$ the digamma function, $\Gamma$ the gamma function, $q$ the dimension of $\mathbf{s}$, and the transition from (6) to (7) always holds for $q > 0$. Then, from Equation 7, we define the intrinsic reward of the current state $\mathbf{s}_t$ as follows:

$$r^{\text{int}}(\mathbf{s}_t) = \log(||\mathbf{s}_t - \mathbf{s}_t^k||).$$

| Hyperparameter | Value | Hyperparameter | Value |
|---|---|---|---|
| Initial temperature | 0.1 | Hidden units per each layer | 1024 (DMControl), 256 (Meta-world) |
| Learning rate | 0.0003 (Meta-world), 0.001 (cheetah) | Batch Size | 1024 (DMControl), 512 (Meta-world) |
| | 0.0001 (qauadruped), 0.0005 (walker) | Optimizer | Adam (Kingma & Ba, 2015) |
| Critic target update freq | 2 | Critic EMA $\tau$ | 0.005 |
| $(\beta_1, \beta_2)$ | (.9, .999) | Discount $\gamma$ | .99 |

*Table 1.* Hyperparameters of the SAC algorithm. Most hyperparameters values are unchanged across environments with the exception for learning rate.

| Hyperparameter | Value | Hyperparameter | Value |
|---|---|---|---|
| GAE parameter $\lambda$ | 0.92 | Hidden units per each layer | 1024 (DMControl), 256 (Meta-world) |
| Learning rate | 0.0003 (Meta-world), 0.0001 (quadruped) | Batch Size | 512 (cheetah), 128 (Otherwise) |
| | $5e^{-5}$ (quadruped, Walker) | # of timesteprs per rollout | 100 (cheetah, Walker), 500 (quadruped) |
| # of environments per worker | 16 (quadruped, cheetah), 32 (Walker) | PPO clip range | 0.2 |
| Entropy bonus | 0.0 | Discount $\gamma$ | .99 |

*Table 2.* Hyperparameters of the PPO algorithm. Most hyperparameters values are unchanged across environments with the exception for learning rate.

## B. Experimental Details

**Training details**. For our method, we use the publicly released implementation repository of the SAC algorithm (https://github.com/denisyarats/pytorch_sac) with a full list of hyperparameters in Table 1. On the DMControl environments, we use segments of length 50 and a frequency of teacher feedback ($K$ in Algorithm 2) of 20K timesteps, which corresponds to roughly 20 episodes. We choose the number of queries per feedback session $M = 140, 70, 40$ for the maximum budget of 1400, 700, 400 on Walker and Cheetah, and choose $M = 70, 35, 20$ for the maximum budget of 1400, 700, 400 on Quadruped. For Meta-world, we use segments of length 10 and set $M = 64, K = 2400$ for the maximum budget of 2500, 5000, and 10000 on Drawer Close, Window Open, Door Open, and Button Press and $M = 128, K = 4800$ for maximum budget of 25000, 50000 on Sweep Into and Drawer Open.

For preference PPO, we use the publicly released implementation repository of the PPO algorithm (`https://github.com/DLR-RM/stable-baselines3`) with a full list of hyperparameters in Table 2. We choose the number of queries per feedback session $M = 70, 45$ for the maximum budget of 2100, 1400 on the DMControl environments. For the reward model, we use same setups for our method. For Meta-world, we use segments of length 10 and set $M = 256, K = 2400$ for all environments and budgets of feedback.

**Reward model**. For the reward model, we use a three-layer neural network with 256 hidden units each, using leaky ReLUs. To improve the stability in reward learning, we use an ensemble of three reward models, and bound the output using tanh function. Each model is trained by optimizing the cross-entropy loss defined in (4) using ADAM learning rule (Kingma & Ba, 2015) with the initial learning rate of 0.0003.

**Environments**. We follow the standard evaluation protocol for the benchmark locomotion tasks from DMControl. The Meta-world single-task benchmark involves training and testing on a single instantiation (fixed reset and goal) of the task. To constitute a more realistic single-task manipulation setting, we randomize the reset and goal positions in all our experiments. We also use new reward function, which are nicely normalized and make the tasks stable.

## C. Effects of Sampling Schemes

Figures 8 and 9 show the learning curves of PEBBLE with various sampling schemes. For Quadruped, we find that the uncertainty-based sampling schemes (using ensemble disagreement or entropy) are superior to the naive uniform sampling scheme. However, they did not lead to extra gains on relatively simple environments, like Walker and Cheetah, similar to observations from Ibarz et al. (2018). Similarly, on the robotic manipulation tasks, we find little difference in performance for simpler tasks (Drawer Close, Window Open). However, we find that the uncertainty-based sampling schemes generally fare better on the other environments.
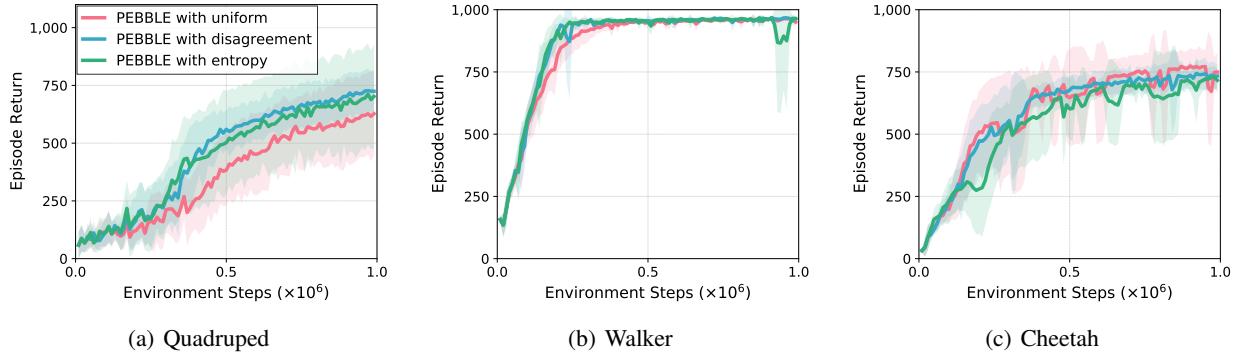


(a) Quadruped

(b) Walker

(c) Cheetah

*Figure 8.* Learning curves of PEBBLE with 1400 pieces of feedback by varying sampling schemes. The solid line and shaded regions represent the mean and standard deviation, respectively, across ten runs.

## D. Examples of Selected Queries

Figure 10, 11 and 12 show some examples from the selected queries to teach the agents.
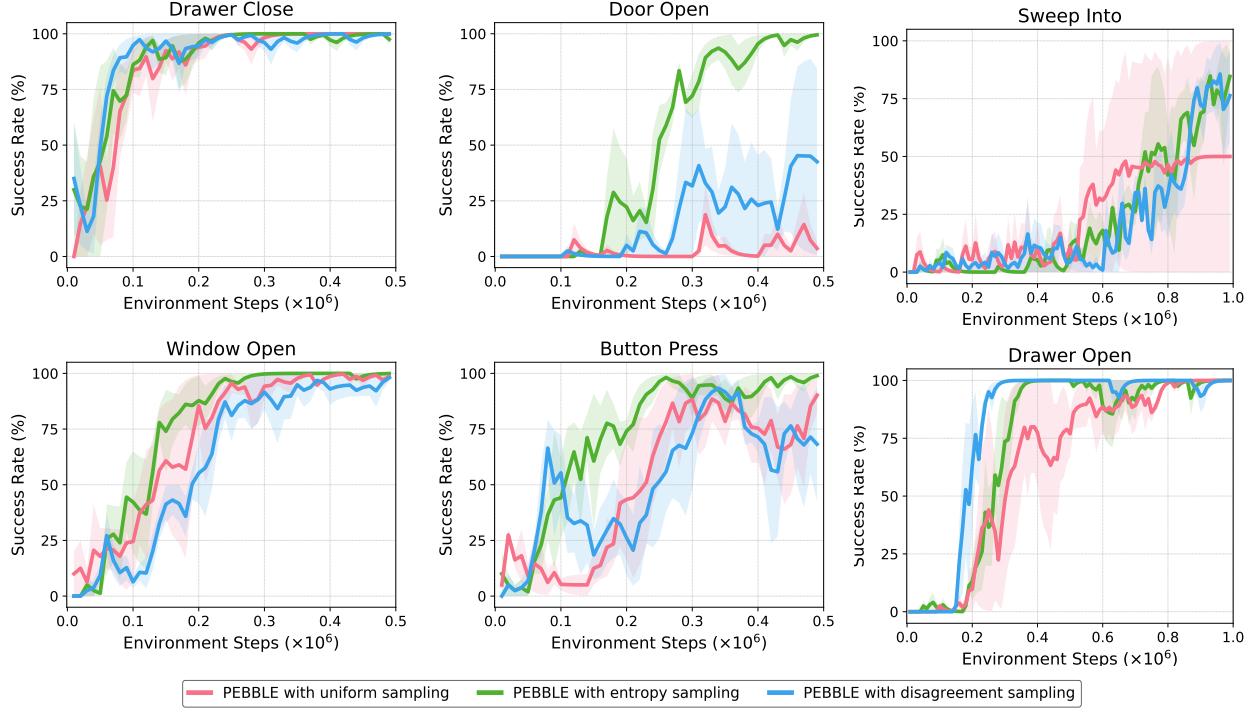
*Figure 9.* Learning curves of PEBBLE with various sampling schemes on the Meta-world tasks. The solid line and shaded regions represent the mean and standard deviation, respectively, across ten runs.
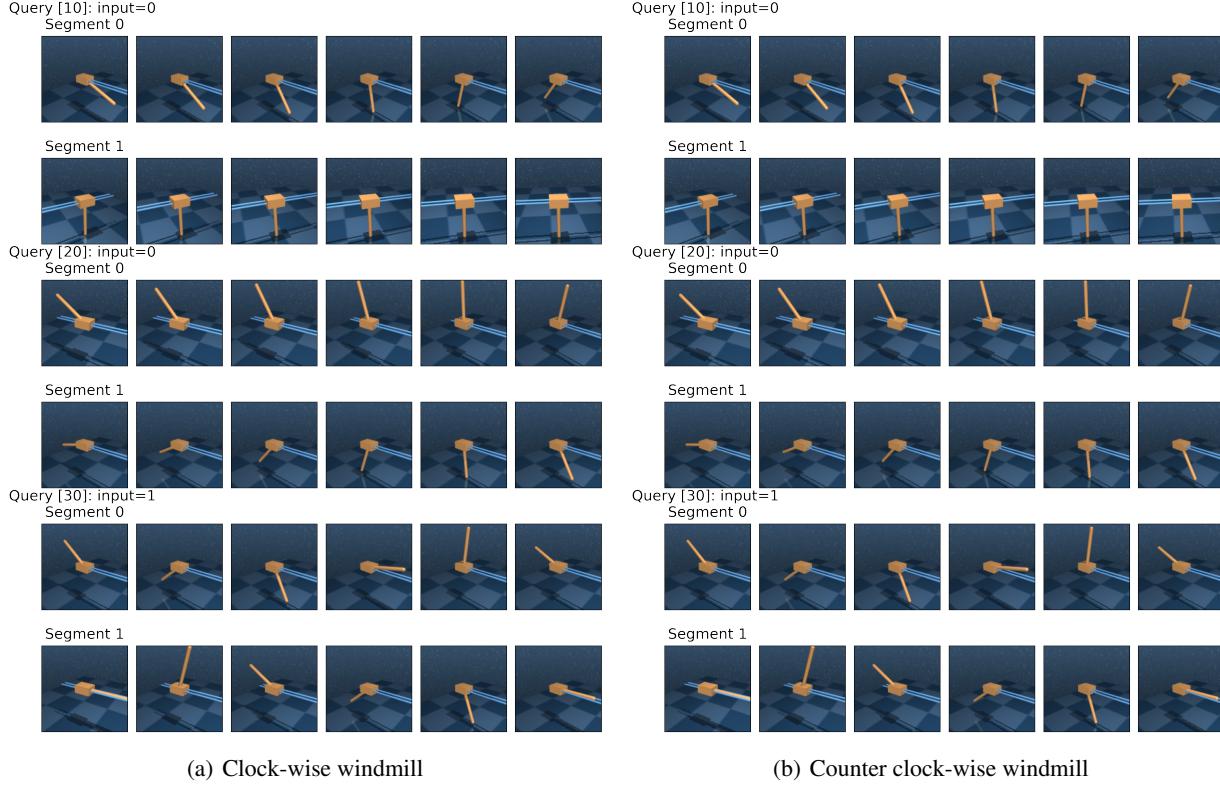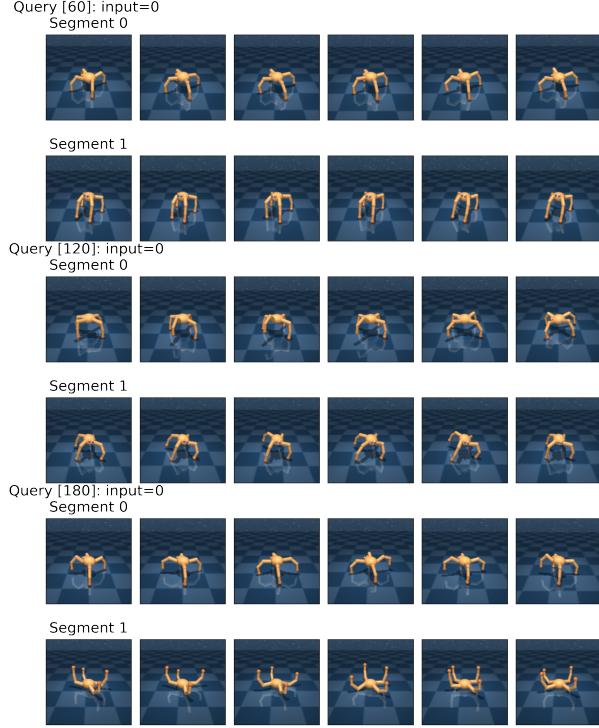


(a) Clock-wise windmill

(b) Counter clock-wise windmill

*Figure 10.* Examples from the selected queries to teach the Cart agent.

Query [60]: input=0
Segment 0

Segment 1

Query [120]: input=0
Segment 0

Segment 1

Query [180]: input=0
Segment 0

Segment 1

(a) Waving left front leg

Query [60]: input=0
Segment 0

Segment 1

Query [120]: input=1
Segment 0

Segment 1

Query [180]: nothing
Segment 0

Segment 1

(b) Waving right front leg

*Figure 11.* Examples from the selected queries to teach the Quadruped agent.

Query [2]: input=1
Segment 0

Segment 1

Query [9]: input=0
Segment 0

Segment 1

Query [22]: input=0
Segment 0

Segment 1

Query [36]: input=0
Segment 0

Segment 1

Query [38]: input=0
Segment 0

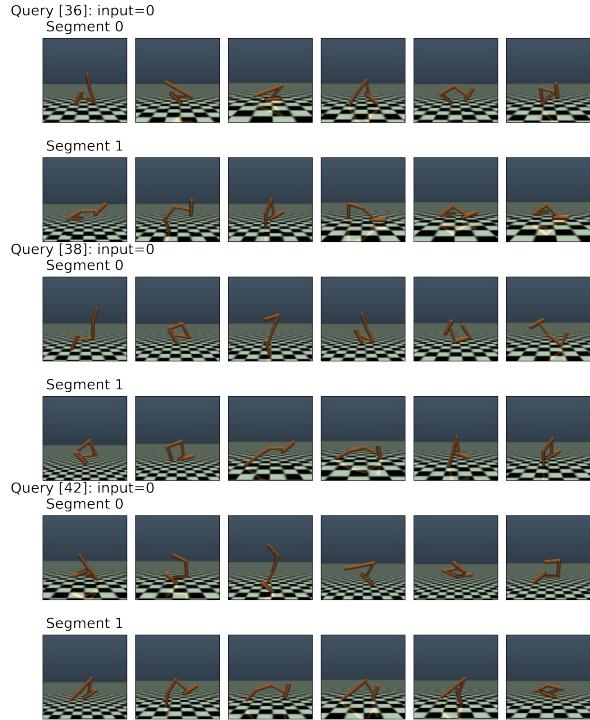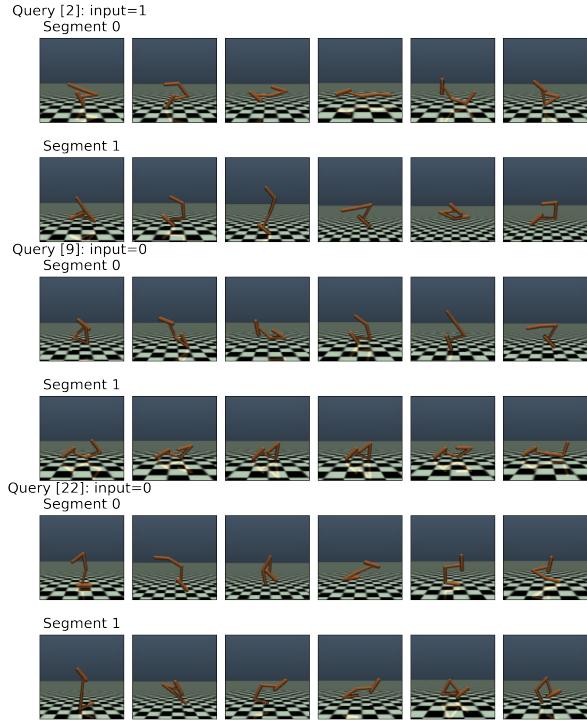Segment 1

Query [42]: input=0
Segment 0

Segment 1

*Figure 12.* Examples from the selected queries to teach the Hopper agent.