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

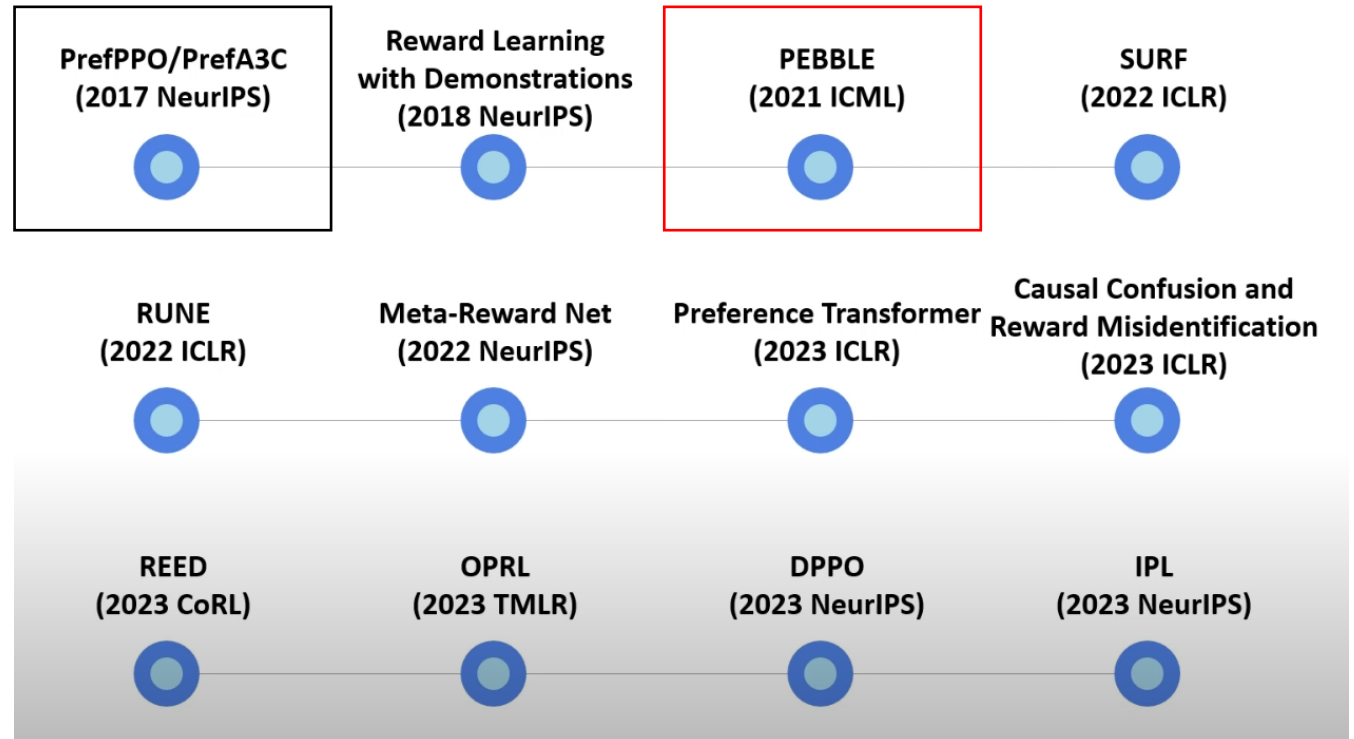Kimin Lee, Laura Smith, Pieter Abbeel
ICML 2021

2025. 01. 09
Learning Agents 강화학습 논문 리뷰 스터디
Minkyoung Kim

# Agenda

- Introduction

- Related work

- Preliminaries

- Method

- Experiment

- Conclusions

- introduction of PbRL
- Reward Ensemble and Sampling
- on-policy Algorithm (PPO)

# Introduction

- Scaling RL to many applications is yet preclude by a number of challenges; <u>providing a suitable reward function</u>
    - Sparse reward, dense reward
    - <span style="color:red">Reward exploitation: achieve high returns by unexpected, unintended means</span>
        - Imitation learning is one of the popular way to avoid reward engineering → expensive


- Human-in-the-loop(HiL) RL
  : Agent's behavior can be tailored to one's preference (avoid reward exploitation) without requiring extensive engineering.
    - given: feedback is both practical for a human to provide and sufficiently high-bandwidth
    - <u>How to reduce the amount of human effort required for HiL learning?</u>

# Introduction

- PEBBLE
  - unsupervised **PrE**-training and preference-**B**ased learning via rela**B**e**L**ing **E**xperience
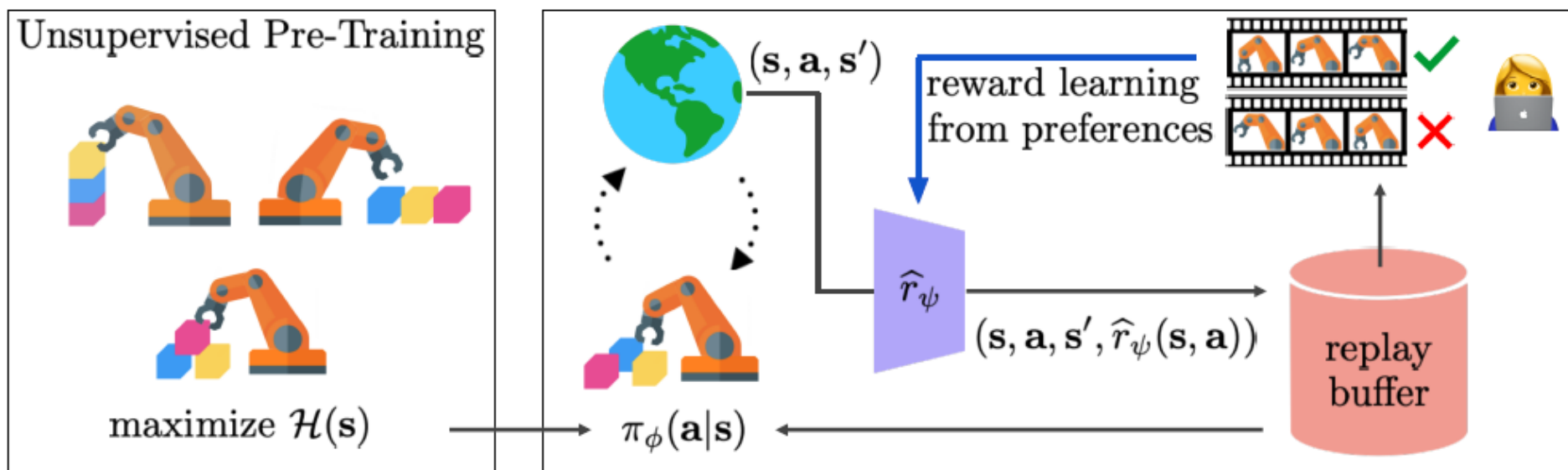  - unsupervised pre-training + off-policy



*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).
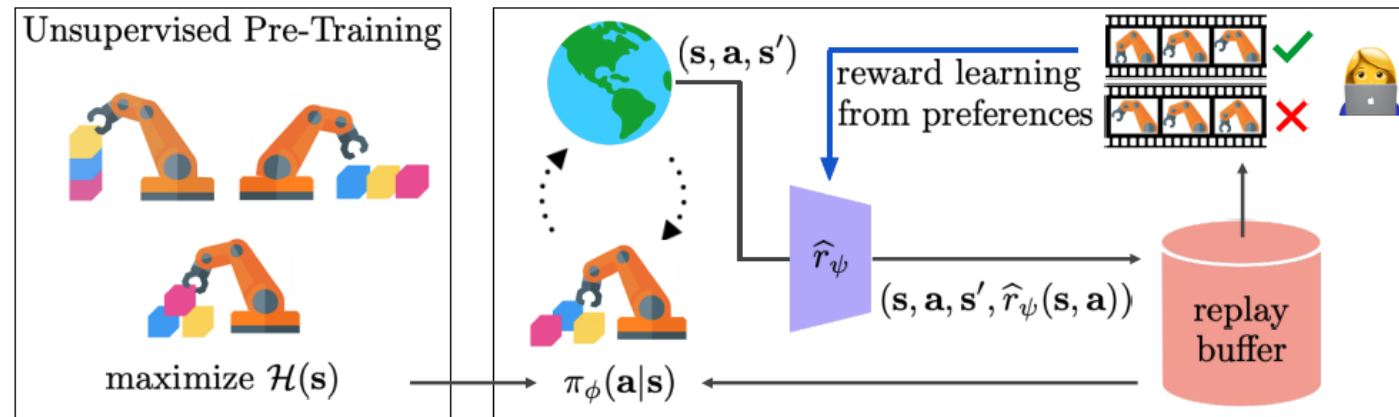
# Introduction

- **Unsupervised Pre-Training**
    - collective a breadth of experience enables the teacher to provide more meaningful feedback (compared to collected data in a indeliberate manner)
    - increase the efficiency of the teacher's initial feedback

- **Off-policy**
    - reuse data to maximize the agent's and human's efficiency
    - to mitigate the effects of non-stationarity
        - relabeling all of the agent's past experiences, every time the reward model is updated

# Introduction

- **Contribution of PEBBLE**

- unsupervised pre-training and off-policy learning can significantly improve the sample- and feedback-efficiency of HiL RL.

- PEBBLE outperform prior preference-based RL baselines on complex locomotion and robotics manipulation tasks from DeepMind Control Suite and Meta-world

- PEBBLE can learn behavior for which a typical reward function is difficult to engineer very efficiency.

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

# Related Work

- **Learning from human feedback**
  - assumption: feedback is available at all time (high feedback frequency)

    → difficult to scale to more complex learning problems

  - learn a reward model, agent can learn without supervisor's perpetual presence.

    → e.g. reward model; classifier, regression

    → difficult for human to reliably provide a particular utility value

  - relative judgement(e.g. comparing behavior as better or worse) is mush easier for human

    → preference RL
    - on-policy RL: more robust to the non-stationarity in rewards caused by online learning
      **but, sample-inefficient on-policy RL**

# Related Work

- **Unsupervised pre-training for RL**
  - agents are encouraged to expand the boundary of seen state by maximizing various intrinsic rewards
    - e.g. prediction errors, count-based state novelty, mutual information
  - allows learning diverse behaviors without extrinsic reward
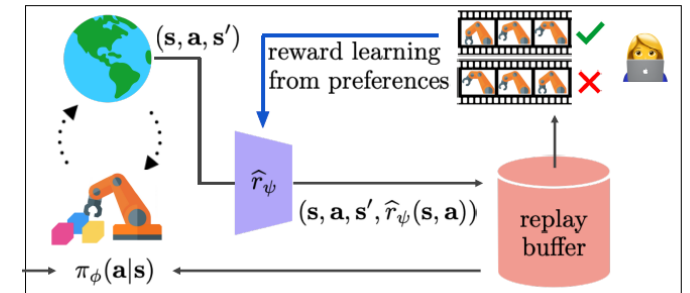
# Preliminaries

- **Soft Actor-Critic(SAC)**
    - <u>off-policy</u> actor-critic method based on maximum entropy RL Framework
    - <u>good sample-efficiency</u> relative to its on-policy counterparts by reusing its past experiences.
    - However, SAC <span style="color:red">is not robust to a non-stationary reward function.</span>

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

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

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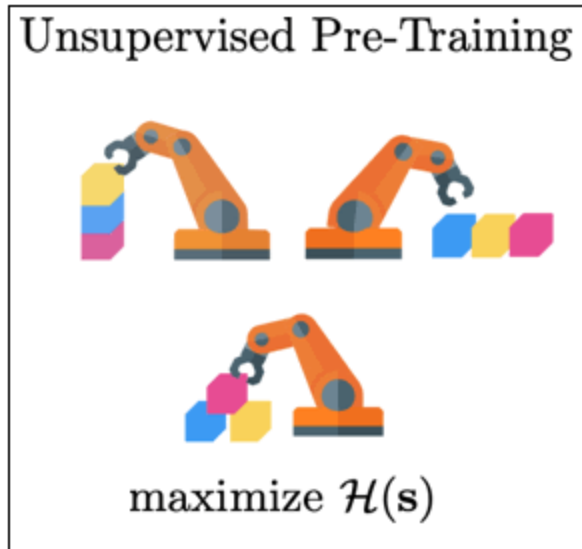
# Preliminaries

- **Reward learning from preference**
  - learning reward function $\widehat{r}_\psi$ from preferences in which the function is trained to be consistent with observed human feedback.
  - following the Bradley-Terry model,

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

where $\sigma^i \succ \sigma^j$ denotes the event that segment $i$ is preferable to segment $j$.

# PEBBLE



Unsupervised Pre-Training

maximize $\mathcal{H}(\mathbf{s})$

# PEBBLE

- policy $\pi_\theta$, Q-function $Q_\theta$, reward function $\widehat{r}_\psi$

**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: $\mathcal{B}, \pi_\phi \leftarrow \text{EXPLORE}()$ 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:       $(\sigma^0, \sigma^1) \sim \text{SAMPLE}()$ (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 $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$ and $\phi$, respectively
27:    **end for**
28: **end for**

- *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*.

# PEBBLE

- Accelerating Learning via Unsupervised Pre-training
  - collect diverse samples through intrinsic motivation
  - visit a sider range of state by using state entropy(particle-based entropy estimator)
    - compute $k$-NN distance between a sample and all samples in replay buffer $B$

1: Initialize parameters of $Q_\theta$ and $\widehat{r}_\psi$
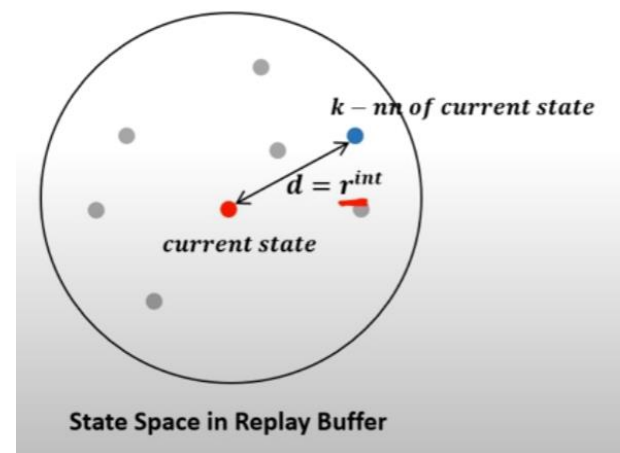2: Initialize a dataset of preferences $\mathcal{D} \leftarrow \emptyset$

**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 $r_t^{\text{int}} \leftarrow r^{\text{int}}(\mathbf{s}_t)$ as in (5)
6:         Store transitions $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}, r_t^{\text{int}})\}$
7:     **end for**
8:     **for** each gradient step **do**
9:         Sample minibatch $\{(\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}_{j+1}, r_j^{\text{int}})\}_{j=1}^B \sim \mathcal{B}$
10:        Optimize $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$ and $\phi$
11:     **end for**
12: **end for**
13: **return** $\mathcal{B}, \pi_\phi$

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

$s_i^k$ = the $k$-th nearest neighbor(k-NN) of $s_i$

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



$k - nn$ of current state

$d = r^{int}$

current state

**State Space in Replay Buffer**

policy $\pi_\theta$, Q-function $Q_\theta$, reward function $\widehat{r}_\psi$

# PEBBLE

- Selecting Informative Queries

**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 $s_{t+1}$ by taking $a_t \sim \pi_\phi(a_t|s_t)$
22:         Store transitions $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s_t, a_t, s_{t+1}, \widehat{r}_\psi(s_t))\}$
23:     **end for**
24:     **for** each gradient step **do**
25:         Sample random minibatch $\{(\tau_j)\}_{j=1}^{B} \sim \mathcal{B}$
26:         Optimize $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$ and $\phi$, respectively
27:     **end for**
28: **end for**

- *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*.

policy $\pi_\theta$, Q-function $Q_\theta$, reward function $r_\psi$

# PEBBLE

- Selecting Informative Queries
  - uniform sampling
  - ensemble-based sampling
    - selects pairs of segment with high variance across ensemble models.
  - ✓ entropy-based sampling (selected)
    - seeks to disambiguate pairs of segments nearest the decision boundary

**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 $1 \dots M$ **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 $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$ and $\phi$, respectively
27:     **end for**
28: **end for**

policy $\pi_\theta$, Q-function $Q_\theta$, reward function $\widehat{r}_\psi$

# PEBBLE

- Using Off-policy RL with Non-stationary Reward
  - relabel all of the agent's past experience ever time we update the reward model

**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: $\mathcal{B}, \pi_\phi \leftarrow \text{EXPLORE}()$ 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:           $(\sigma^0, \sigma^1) \sim \text{SAMPLE}()$ (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 $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$ and $\phi$, respectively
27:     **end for**
28: **end for**

- *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*.

policy $\pi_\theta$, Q-fun

# PEBBLE

**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 $\mathcal{L}_{\text{critic}}^{\text{SAC}}$ in (1) and $\mathcal{L}_{\text{act}}^{\text{SAC}}$ in (2) with respect to $\theta$
      and $\phi$, respectively
27:    **end for**
28: **end for**

- *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*.

policy $\pi_\theta$, Q-function $Q_\theta$, reward function $\widehat{r}_\psi$

# Experiment

- used scripted teacher
  - provides preferences between trajectory segments according to the true task reward.
- each trajectory segment is presented to the human as a 1 second video clips
  - maximum of one hour of human time is required.
- pre-train an agent for 10K timesteps in all learning curves.

locomotion task from DeepMind Control Suite

robotics manipulation tasks from Meta-world



Quadruped    Walker    Cheetah

Drawer Close    Button Press    Drawer Open

Window Open    Door Open    Sweep Into

# Experiment

- How does PEBBLE compare to existing method in terms of sample and feedback efficiency?
  - prePPO가 더 많은 feedback을 사용
  - PEBBLE(green, 1400)은 SAC(pink)과 거의 성능이 유사
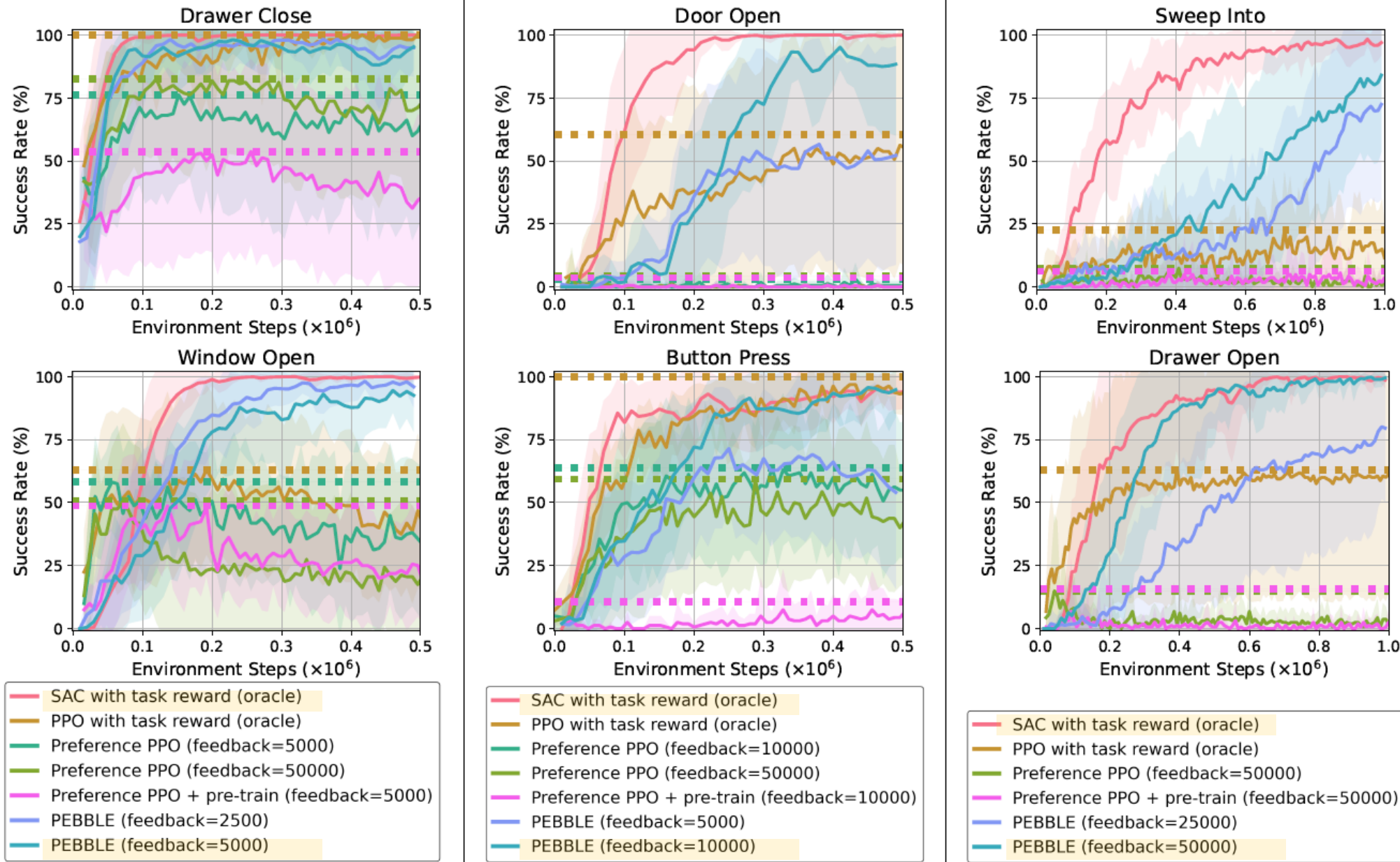  - prePPO(purple, 1400)은 PPO(black)에 도달하지 않음



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.

# Experiment

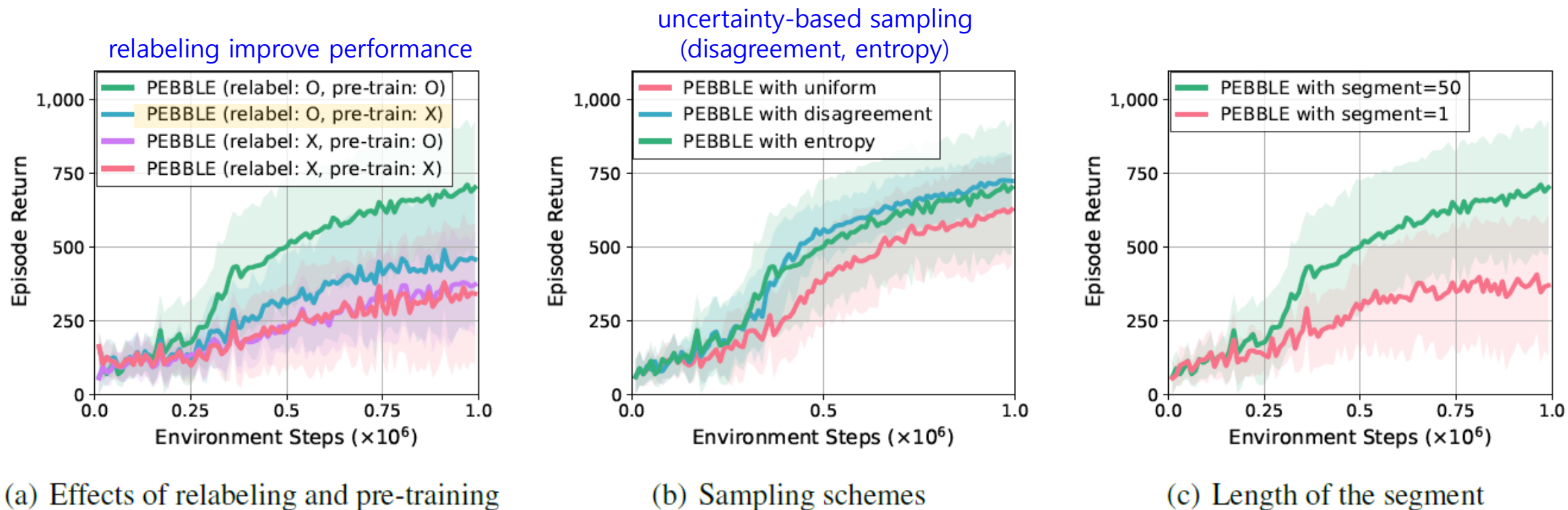- What is the contribution of each of proposed techniques in PEBBLE?
  - pretrain
    - prePPO(red, 1400) + pretrain이 prePPO(purple, 1400)보다 성능이 좋음 @Quadruped, Walker



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.

*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.

# Experiment

- What is the contribution of each of proposed techniques in PEBBLE?



(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.

# Experiment

- What is the contribution of each of proposed techniques in PEBBLE?
  - uncertainty-based sampling schemes not lead to extra gain on relatively simple environment.
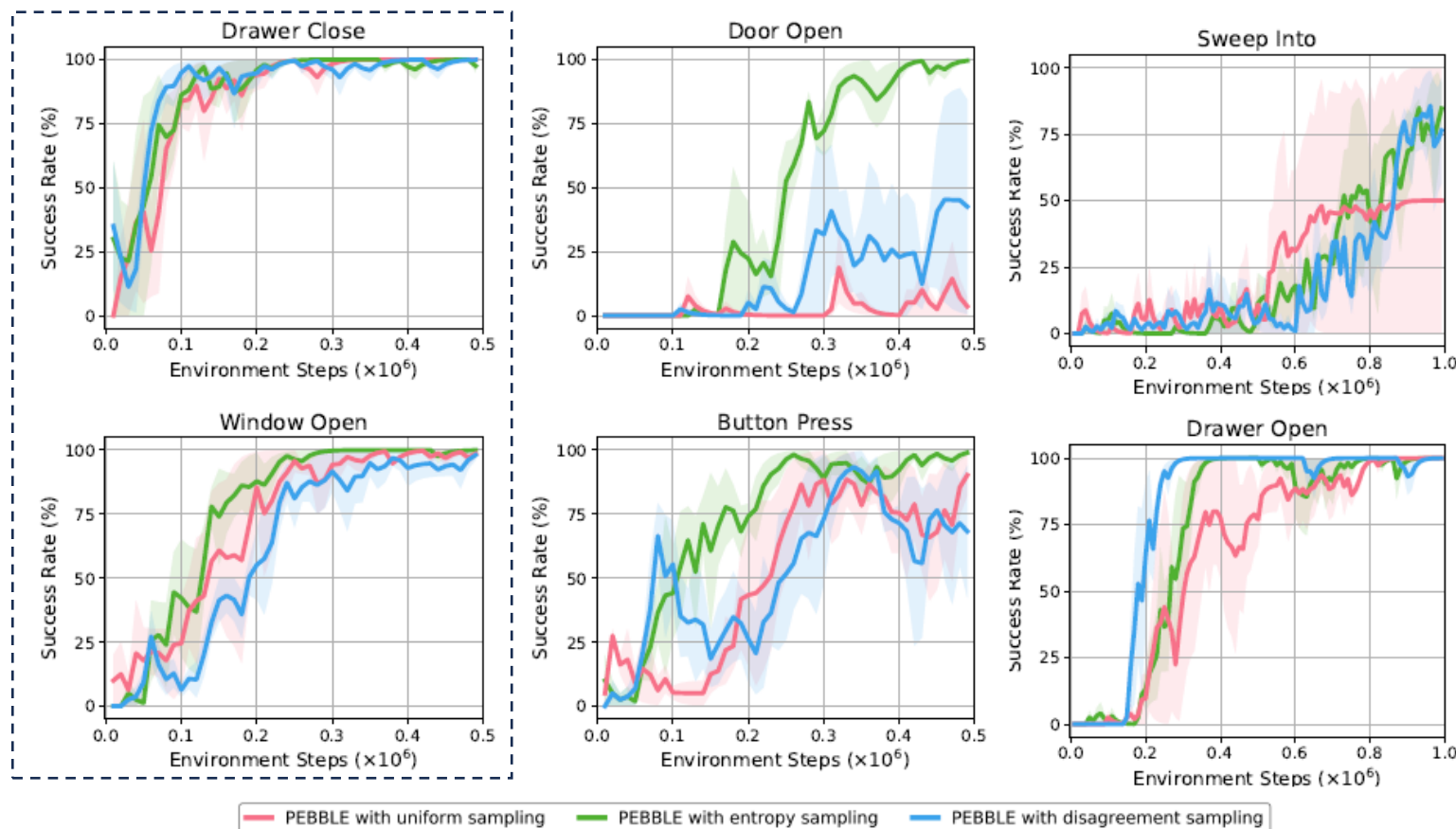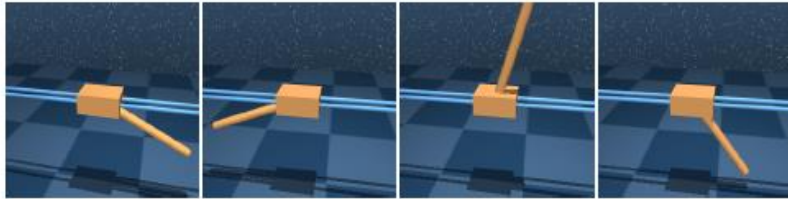    - e.g. Walker, Cheetah, Drawer Close, Window Open



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.

# Experiment

- What is the contribution of each of proposed techniques in PEBBLE?
  - uncertainty-based sampling schemes not lead to extra gain on relatively simple environment.
    - e.g. Walker, Cheetah, Drawer Close, Window Open



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.
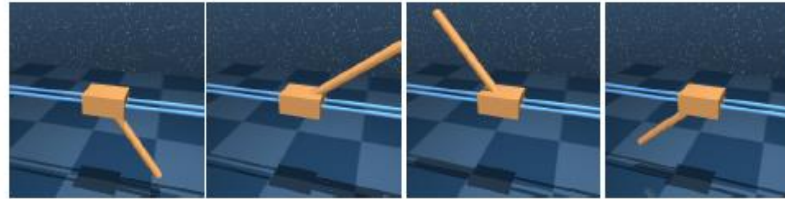
23

# Experiment

- Can PEBBLE learn novel behaviors for which a typical reward function is difficult to engineer?
  - https://youtu.be/w88IYSeV7PQ?si=eVuJAyDgLN1MXl8O



50 queries

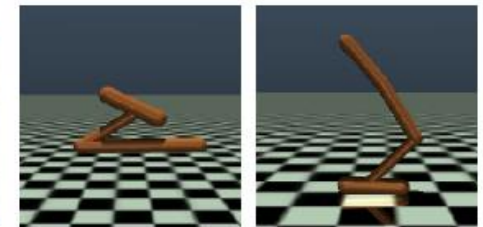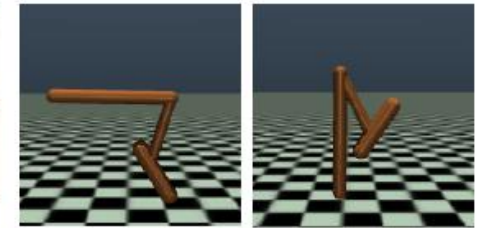Clock-wise windmill

50 queries

Counter clock-wise windmill

Quadruped waving its left front leg

200 queries

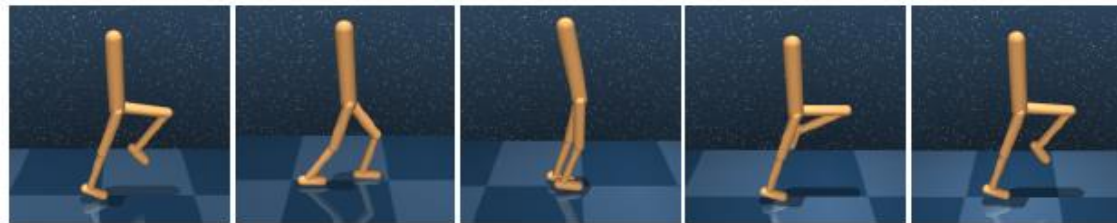Quadruped waving its right front leg

200 queries

Hopper backflip

500 queries

# Experiment

- Can PEBBLE mitigate the effects of reward exploitation?



(a) Agent trained with human preference   200 queries

(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.

https://sites.google.com/view/icml21pebble

# Conclusion

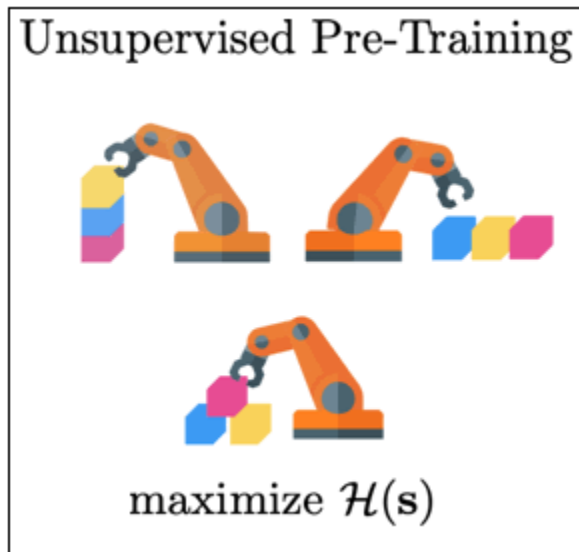### Unsupervised pre-training via state entropy maximization

We pre-train the policy only using an intrinsic motivation to explore and collect diverse experiences.

### Reward learning from preferences

We learn a reward function that can lead to the desired behavior by getting feedback from a teacher.

### Off-policy RL with non-stationary reward

We update the agent using an off-policy RL algorithm with relabeling to mitigate the effects of a non-stationary reward function.



Unsupervised Pre-Training

maximize $\mathcal{H}(\mathbf{s})$

# References

- [Open DMQA Seminar] RLHF-Preference-based Reinforcement Learning, https://www.youtube.com/watch?v=Vzno0oBbm6w

- PEBBLE: Feedback-Efficient Interactive Reinforcement Learning via Relabeling Experience and Unsupervised Pre-training, https://sites.google.com/view/icml21pebble