

Query-Policy Misalignment in Preference-Based Reinforcement Learning

Xiao Hu^{*1}, Jianxiong Li^{*1}, Xianyuan Zhan^{†1}, Qing-Shan Jia^{†1}, Ya-Qin Zhang^{†1}

¹ Tsinghua University, Beijing, China

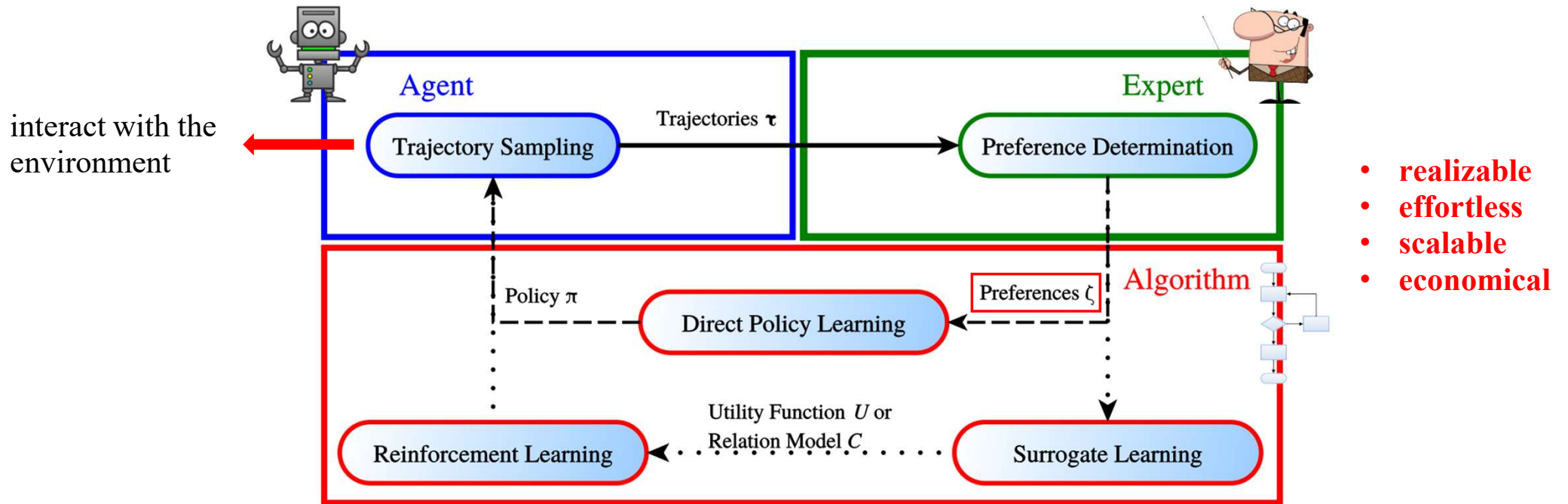
{hu-x21, li-jx21}@mails.tsinghua.edu.cn,

zhanxianyuan@air.tsinghua.edu.cn, jiaqs@tsinghua.edu.cn

ICLR 2024 Spotlight

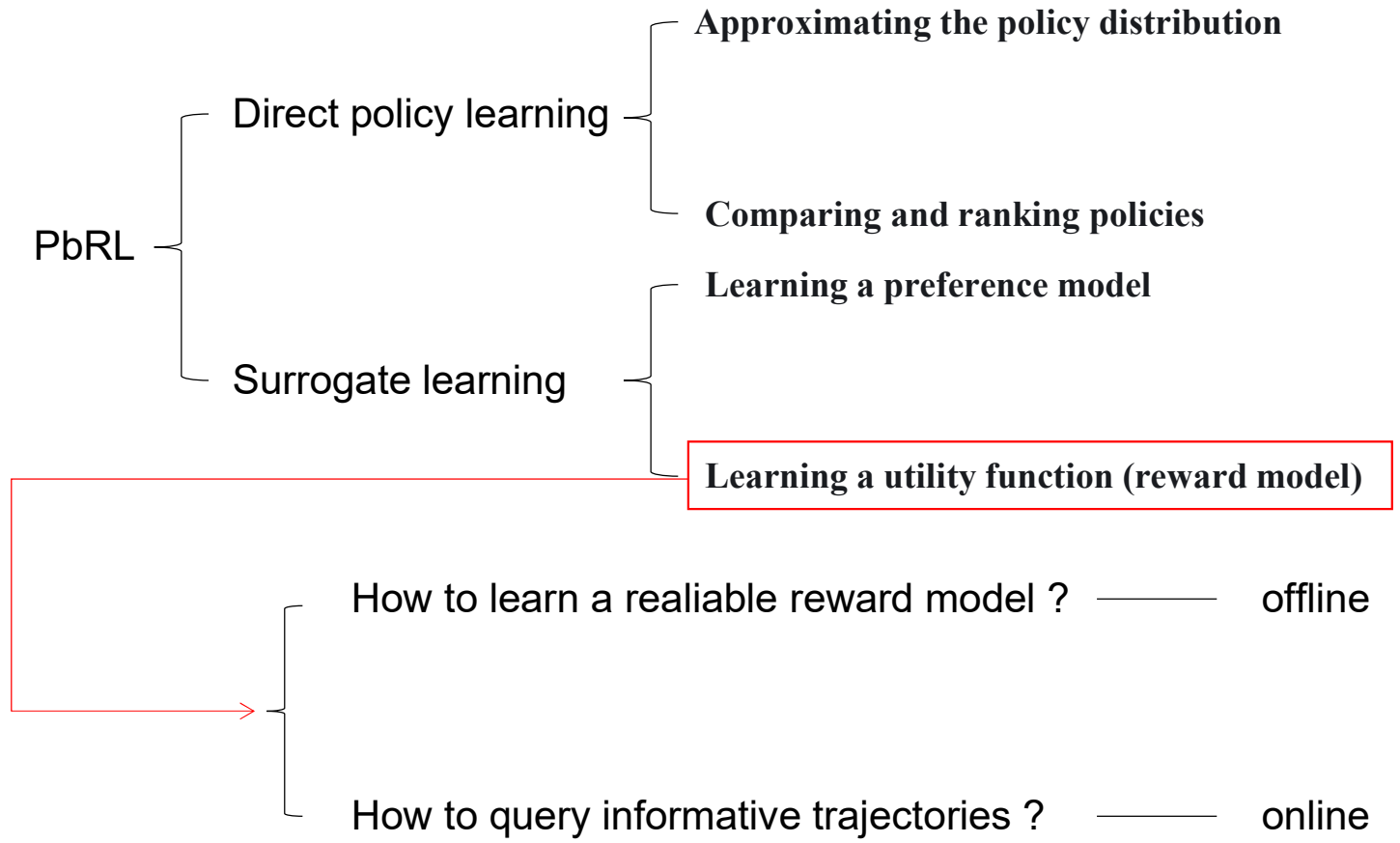
Background

Preference-Based Reinforcement Learning (PbRL)



The agent does not receive a reward after an interacting step, but selects some historical trajectories to query the preferences of the expert.

Background



Background

How to learn a reliable reward model ?

trajectory $\sigma : \{s_k, a_k, \dots, s_{k+L-1}, a_{k+L-1}\}$

preference predictor $P_\psi [\sigma^1 \succ \sigma^0] = \frac{\exp \sum_t \hat{r}_\psi (s_t^1, a_t^1)}{\sum_{i \in \{0,1\}} \exp \sum_t \hat{r}_\psi (s_t^i, a_t^i)}$

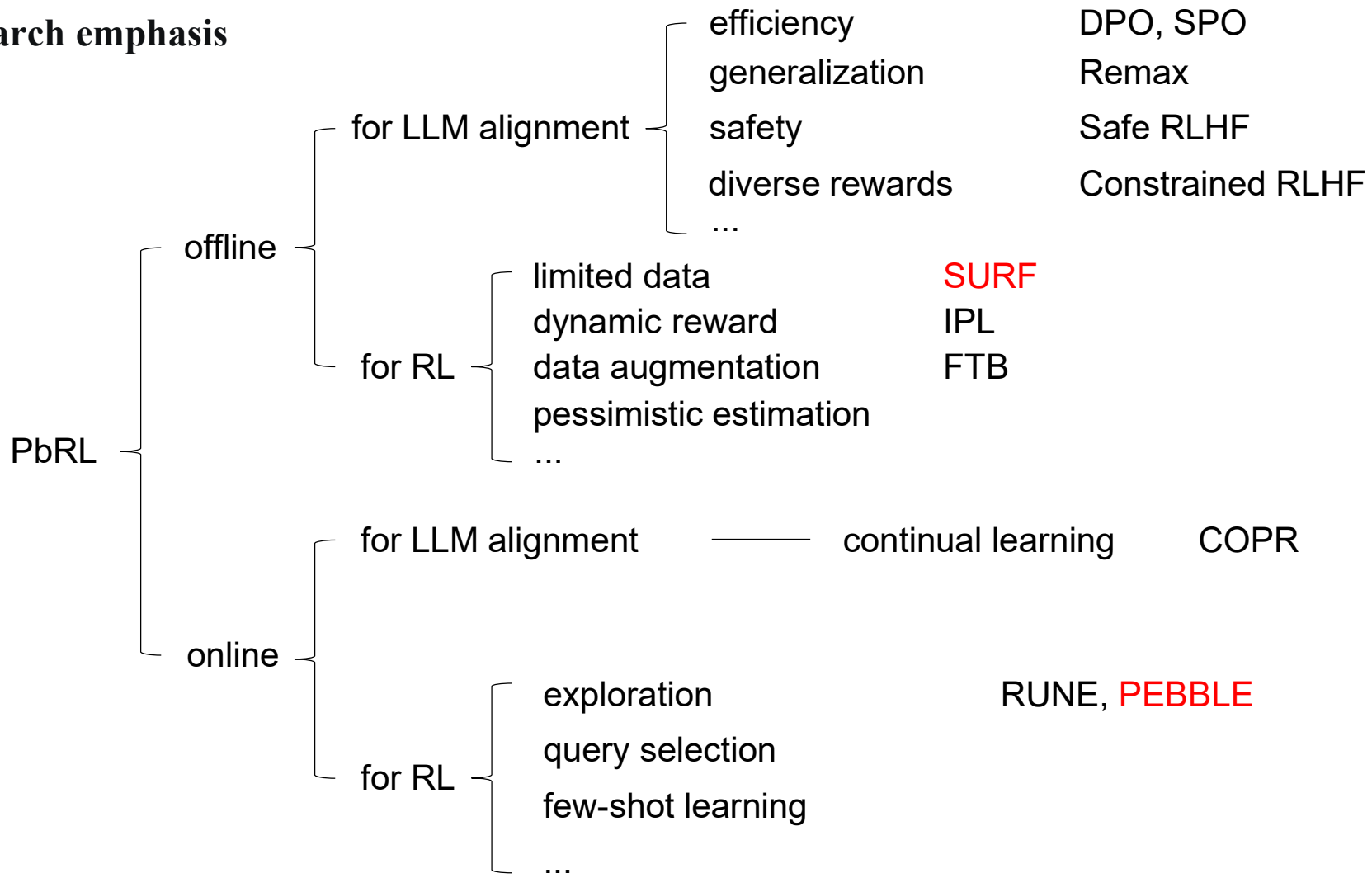
loss $\mathcal{L}^{\text{reward}} = - \mathbb{E}_{(\sigma^0, \sigma^1, y) \sim \mathcal{D}^\sigma} [(1 - y) \log P_\psi [\sigma^0 \succ \sigma^1] + y \log P_\psi [\sigma^1 \succ \sigma^0]]$

How to query informative trajectories ?

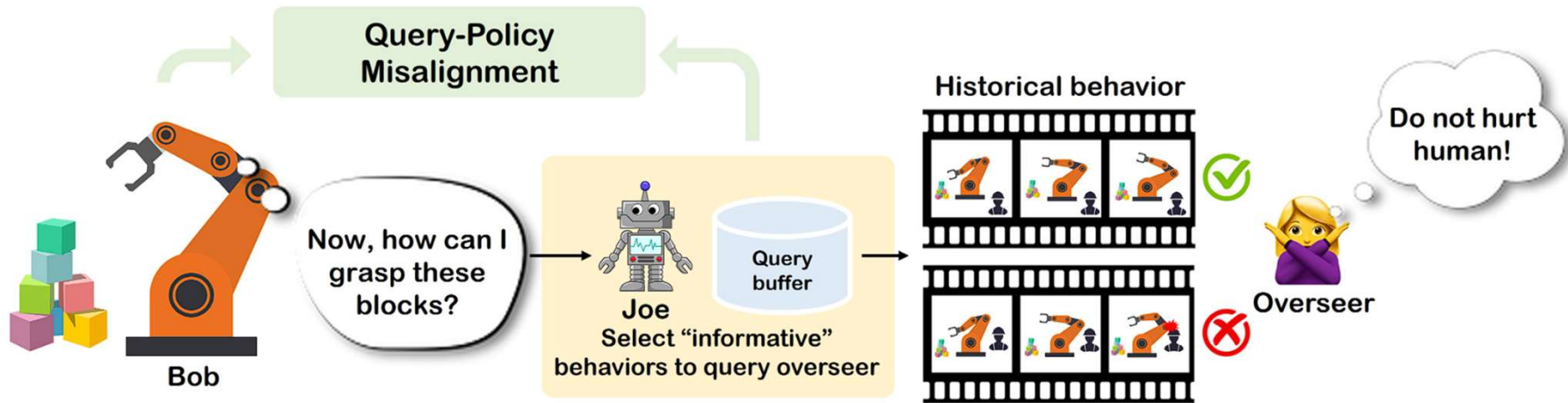
- random query
- entropy-based query
- disagreement-based query

Background

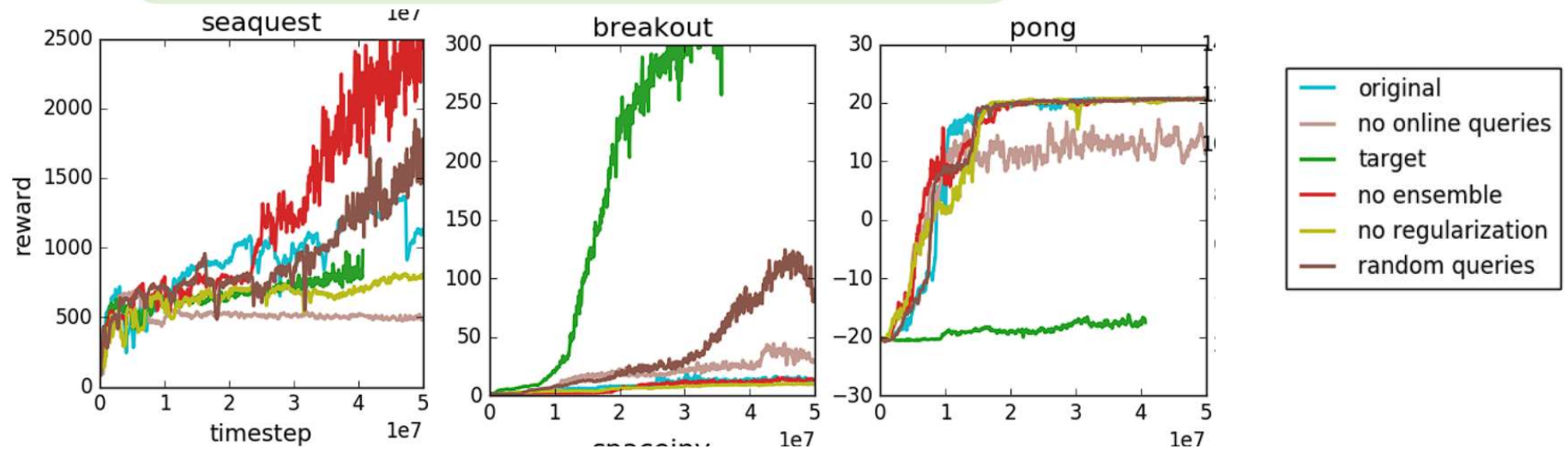
Research emphasis



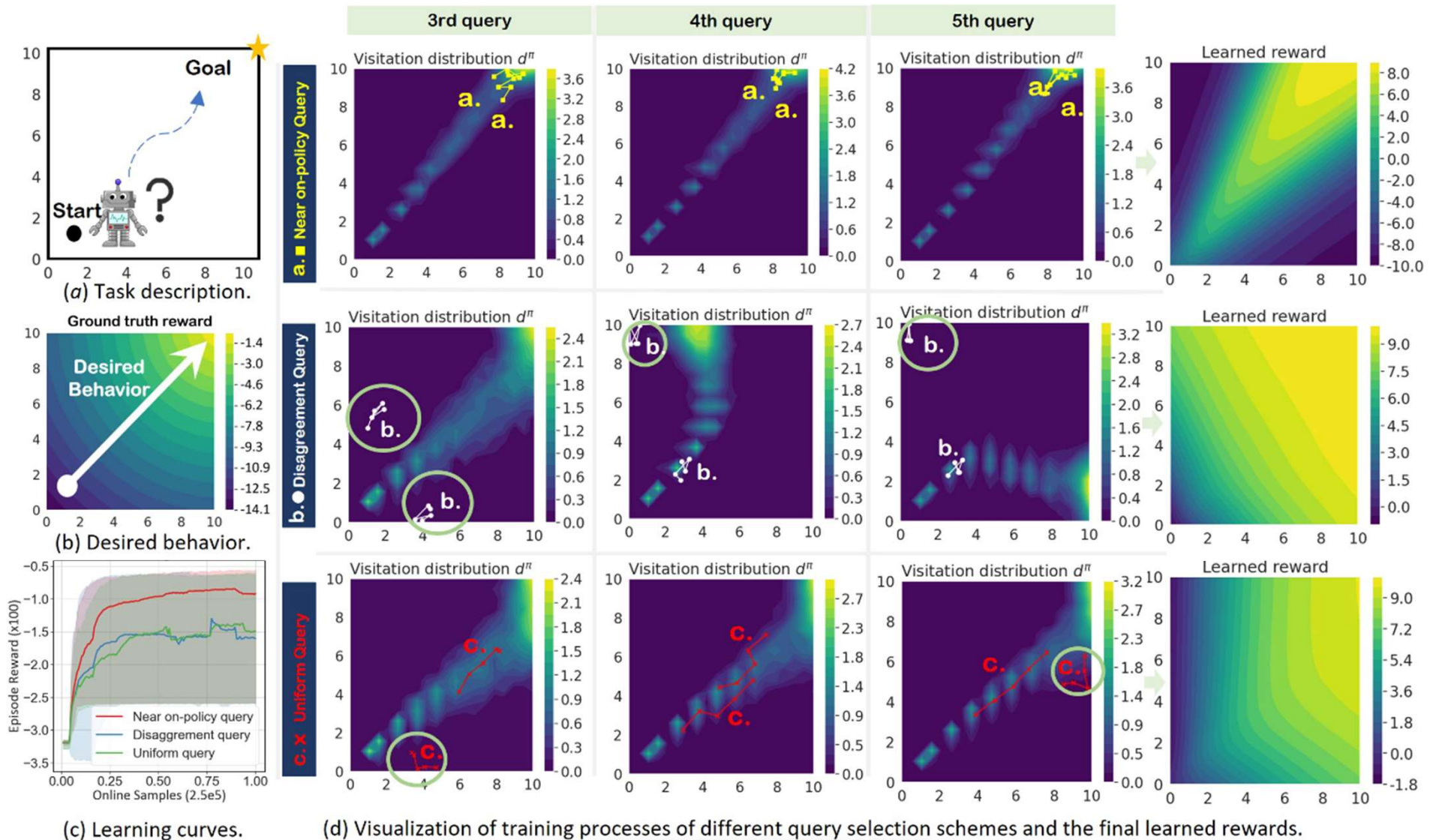
Motivation



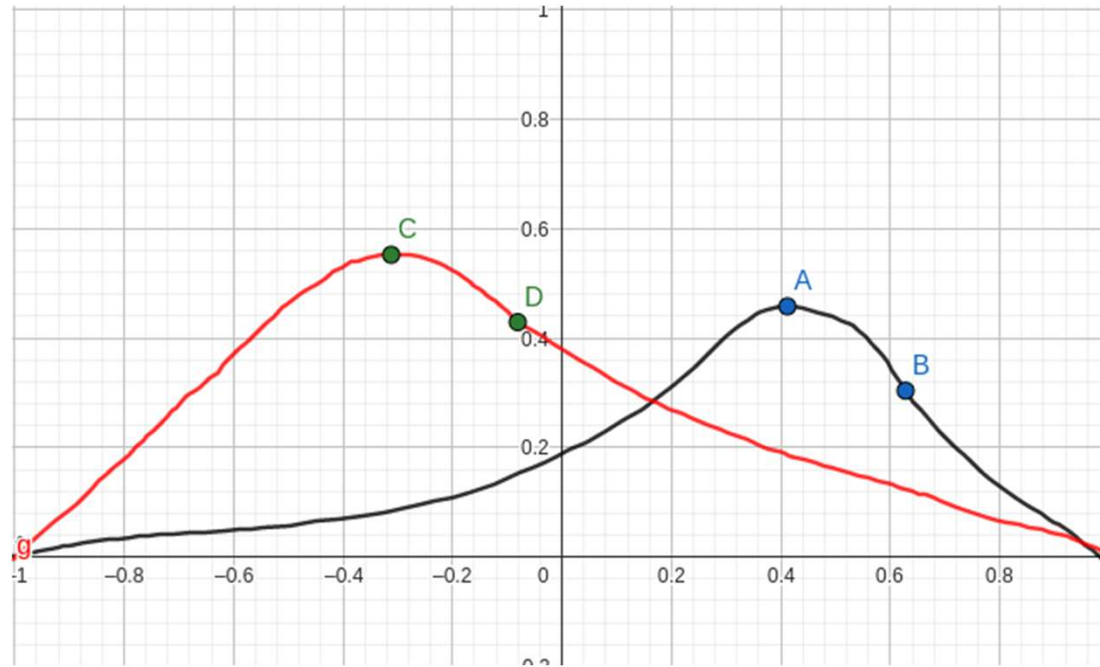
Bob still don't know how to pick up these blocks!



Motivation



Motivation



Theorem 1. Given the two conditions $\|\hat{r}_\psi - r\|_{d^\pi} \leq \epsilon$ and $\|Q_{\hat{r}_\psi}^\pi - \hat{Q}_{\hat{r}_\psi}^\pi\|_{d^\pi} \leq \alpha$, the value approximation error $\|Q_r^\pi - \hat{Q}_{\hat{r}_\psi}^\pi\|_{d^\pi}$ is upper bounded as:

$$\|Q_r^\pi - \hat{Q}_{\hat{r}_\psi}^\pi\|_{d^\pi} \leq \frac{\epsilon}{1 - \gamma} + \alpha \quad (5)$$

near on-policy selection can boost the performance

Method

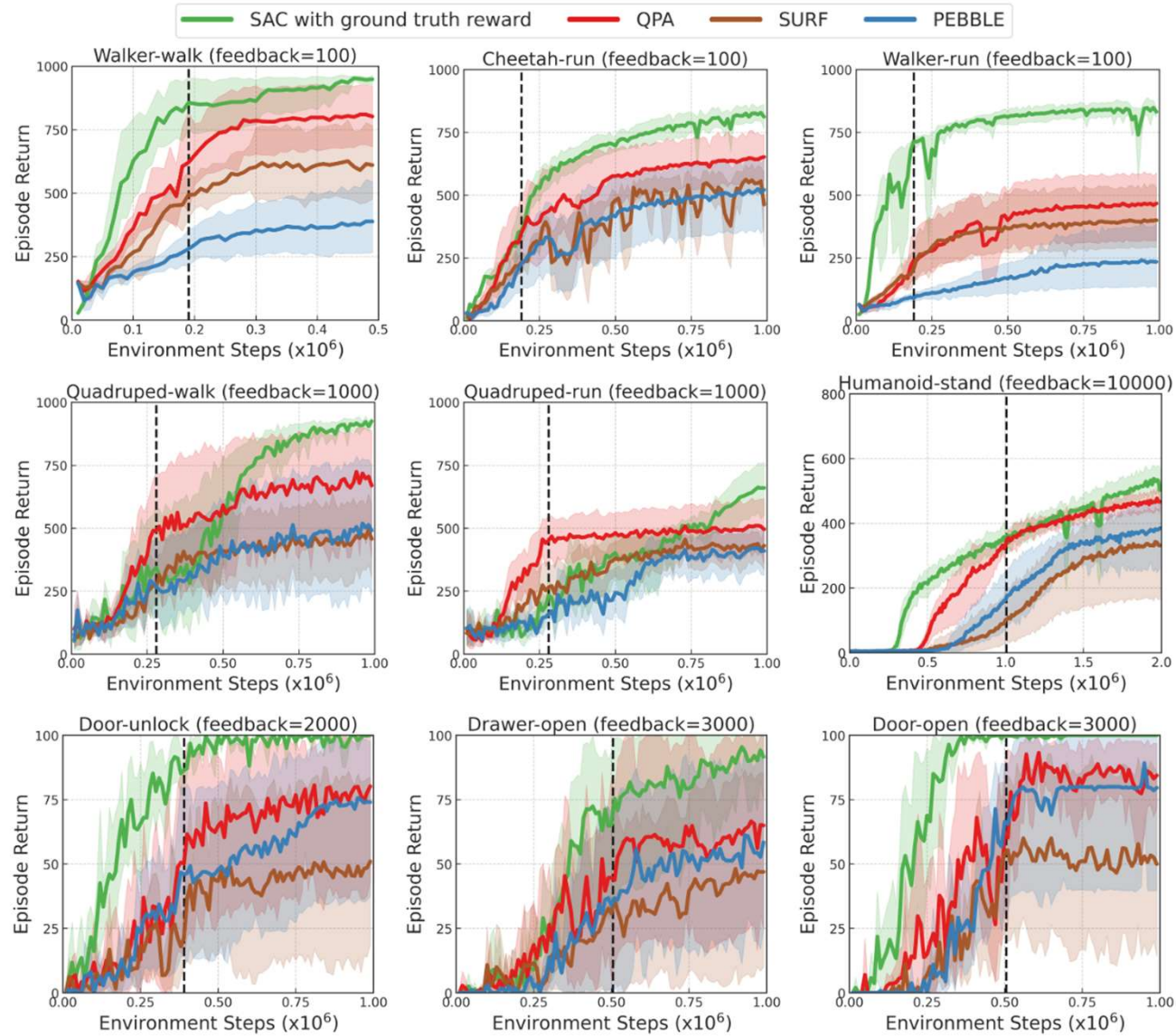
Algorithm 1: QPA

Input : Frequency of overseer feedback K , number of queries per feedback session M
Near on-policy buffer size N , data augmentation ratio τ

Initialize : Initialize replay buffer \mathcal{D} , query buffer \mathcal{D}^σ , near on-policy buffer \mathcal{D}^{on} with size N

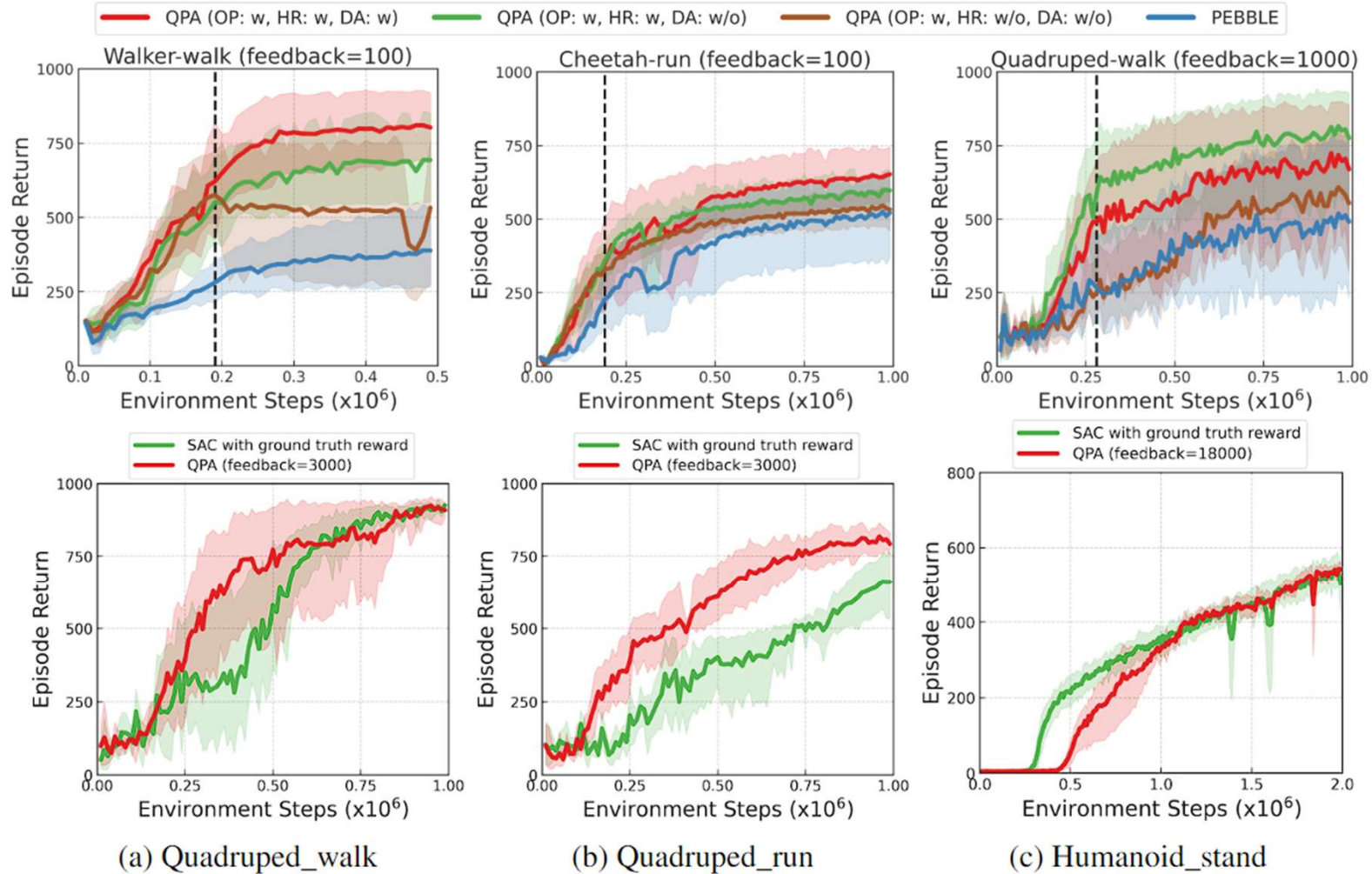
```
1 (Option) Unsupervised pretraining [17]
2 for each iteration do
3   Collect and store new experience  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s, a, r, s')\}$ ,  $\mathcal{D}^{\text{on}} \leftarrow \mathcal{D}^{\text{on}} \cup \{(s, a, r, s')\}$ 
4   if iteration %  $K == 0$  then
5     /* Near on-policy query selection (see Section 5.1) */
6      $\{(\sigma^0, \sigma^1)\}_{i=1}^M \sim \mathcal{D}^{\text{on}}$ 
7     Query for preferences  $\{y\}_{i=1}^M$ , and store preference  $\mathcal{D}^\sigma \leftarrow \mathcal{D}^\sigma \cup \{(\sigma^0, \sigma^1, y)\}_{i=1}^M$ 
8     for each gradient step do
9       Sample a minibatch preferences  $\mathcal{B} \leftarrow \{(\sigma^0, \sigma^1, y)\}_{i=1}^h \sim \mathcal{D}^\sigma$ 
10      /* Data augmentation for reward learning (see Section 5.3) */
11      Generate augmented preferences  $\hat{\mathcal{B}} \leftarrow \{(\hat{\sigma}^0, \hat{\sigma}^1, y)\}_{i=1}^{h \times \tau}$  based on  $\mathcal{B}$ 
12      Optimize  $\mathcal{L}^{\text{reward}}$  in Eq. (4) w.r.t.  $\hat{r}_\psi$  using  $\hat{\mathcal{B}}$ 
13   for each gradient step do
14     /* Hybrid experience replay (see Section 5.2) */
15     Sample minibatch  $\mathcal{D}_{\text{mini}} \leftarrow \{(s, a, r, s')\}_{i=1}^{\frac{n}{2}} \sim \mathcal{D}$ ,  $\mathcal{D}_{\text{mini}}^{\text{on}} \leftarrow \{(s, a, r, s')\}_{i=1}^{\frac{n}{2}} \sim \mathcal{D}^{\text{on}}$ 
16     Optimize SAC agent using  $\mathcal{D}_{\text{mini}} \cup \mathcal{D}_{\text{mini}}^{\text{on}}$ 
```

Experiment



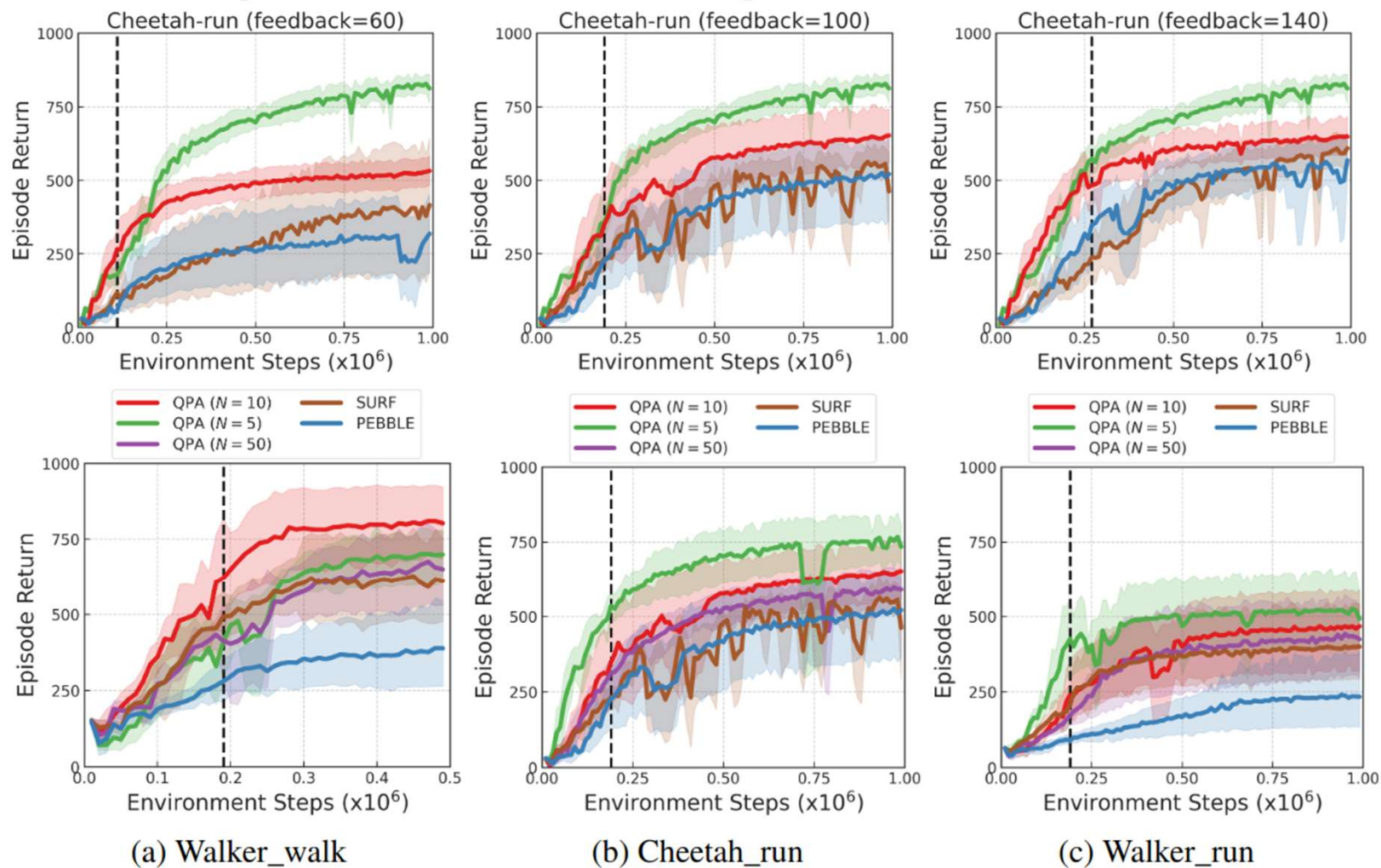
Experiment

Ablation Study



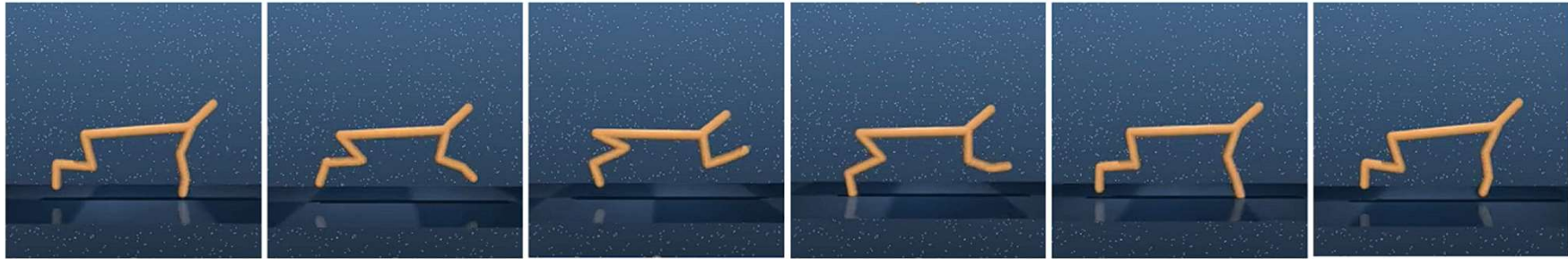
Experiment

Ablation Study

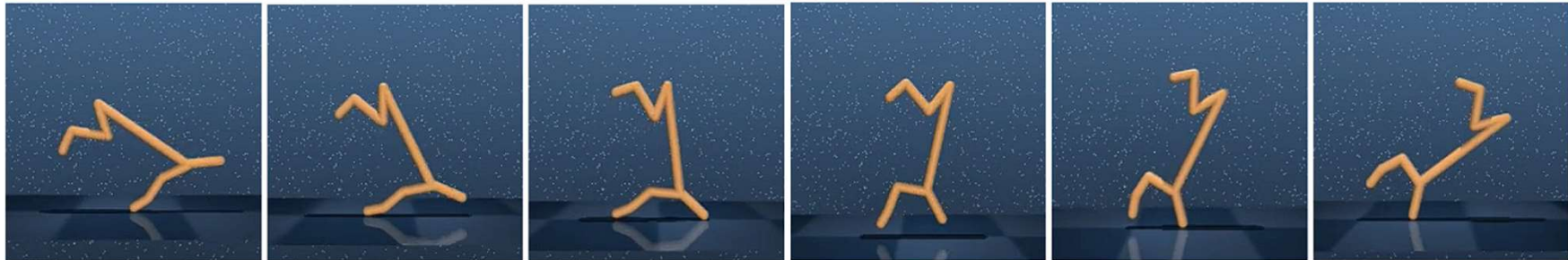


Experiment

Human Experiments



(a) Agent trained with human preferences



(b) Agent trained with hand-engineered preferences

Discussion

PbRL

limited preferences in offline datasets

active query schemes

noisy preferences

Weakly Supervised Learning

semi-supervised learning

active learning

learning with noisy labels

Difficulties in connecting weakly supervised learning and PbRL

- pair-wise query / relative reward
- distribution shift
- sequence-level preference

Thanks