

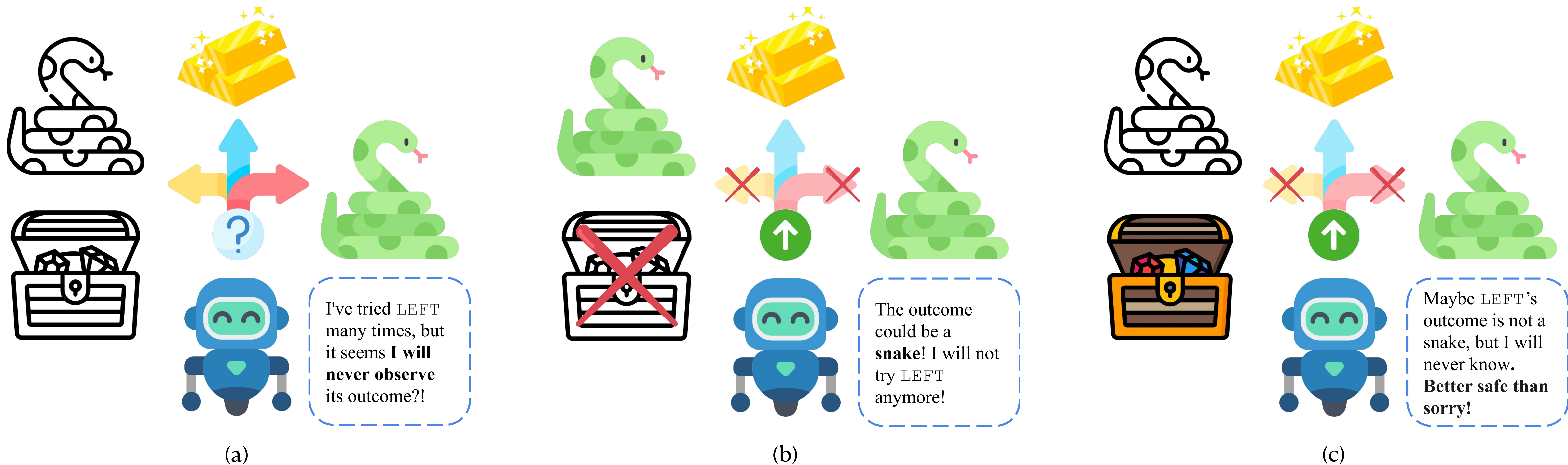
Model-Based Exploration in Monitored Markov Decision Processes

Alireza Kazemipour, Matthew E. Taylor, Michael Bowling

The 42nd Conference on Machine Learning (ICML 2025)

With feedback, **Optimism** Leads to Efficient Exploration

Core Problem: Can we use optimism in the absence of feedback?

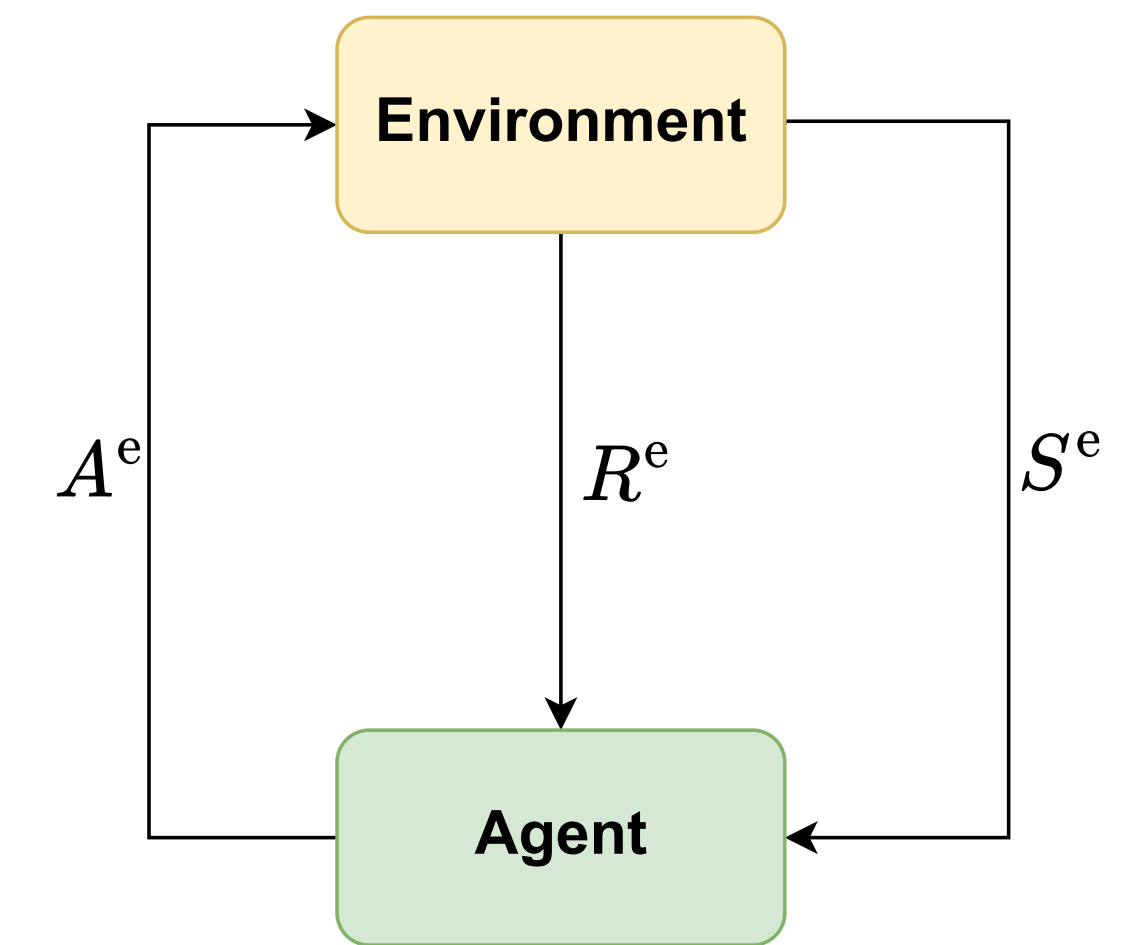


Without feedback, **pessimism** leads to worst-case Guarantee

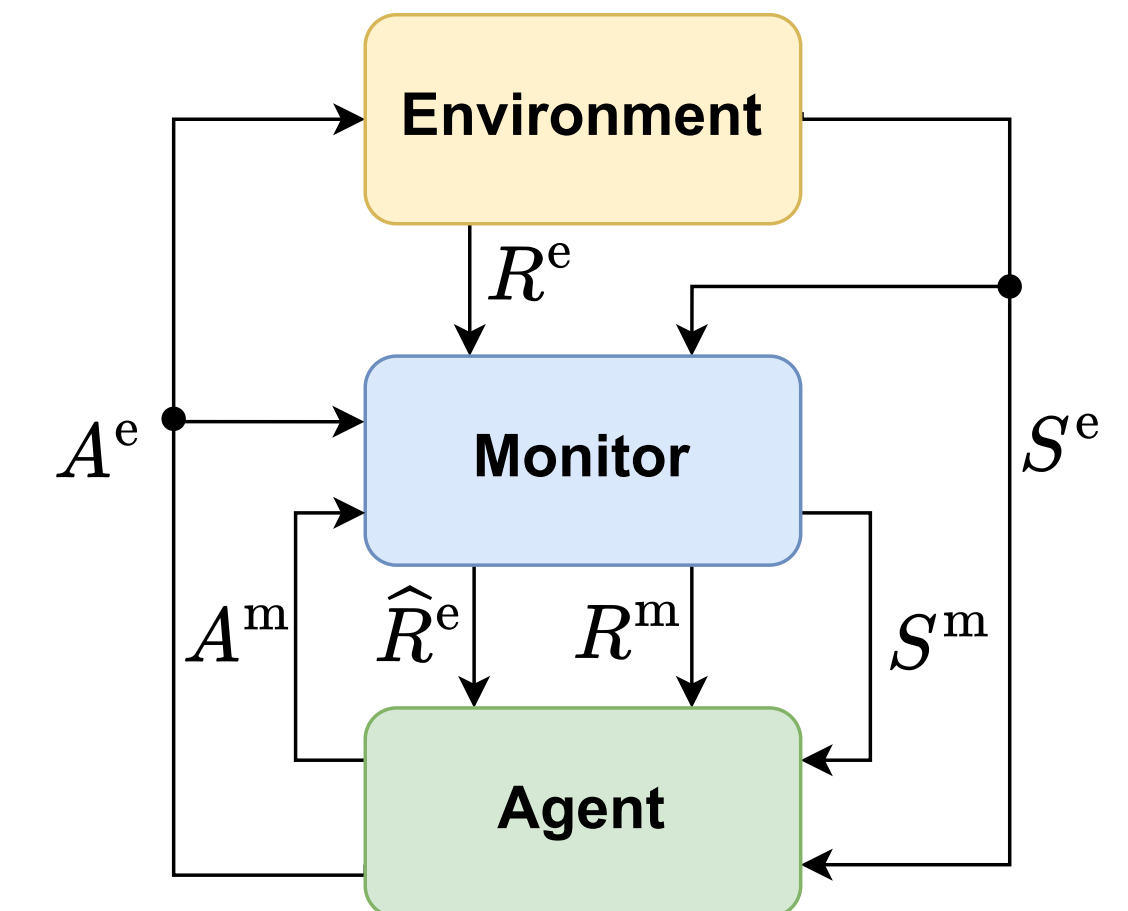
Monitored Markov Decision Processes (Mon-MDPs)

A framework where the reward can be unobservable

- The **monitor is another MDP** in addition to the environment
- The goal is to maximize the joint reward $\sum_{t=0}^{\infty} \gamma^t (R_t^e + R_t^m)$
- But the agent observes \widehat{R}_t^e instead of the true environment reward R_t^e
- For some state-action $\widehat{R}_t^e = \perp$ (NaN) at all times t !



MDP



Mon-MDP

Monitored Model-Based Interval Estimation with Exploration Bonus

Monitored MBIE-EB

- **Innovation 1:** Extending MBIE-EB to Mon-MDPs

$$\tilde{R}_{opt}(s, a) = \bar{R}^e(s^e, a^e) + \underbrace{\frac{\beta^e}{N(s^e, a^e)}}_{\text{bonus for } r^e} + \bar{R}^m(s^m, a^m) + \underbrace{\frac{\beta^m}{N(s^m, a^m)}}_{\text{bonus for } r^m} + \underbrace{\frac{\beta}{N(s, a)}}_{\text{bonus for } p}$$

- **Innovation 2:** Pessimism instead of the optimism

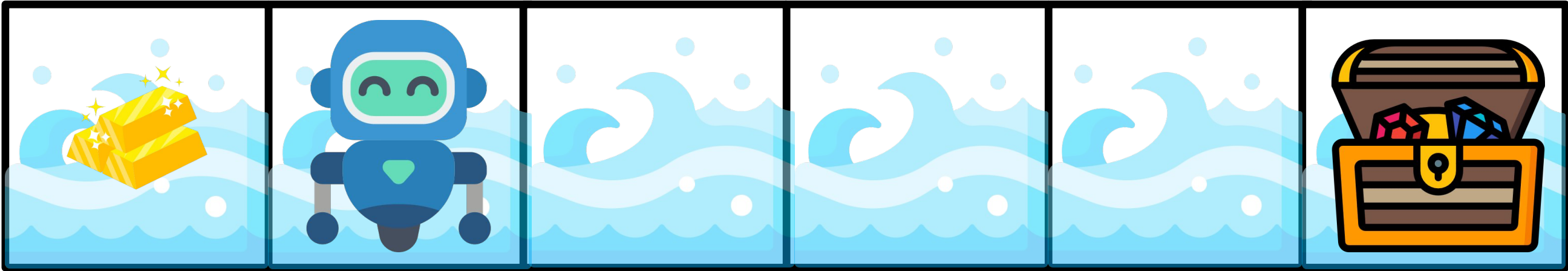
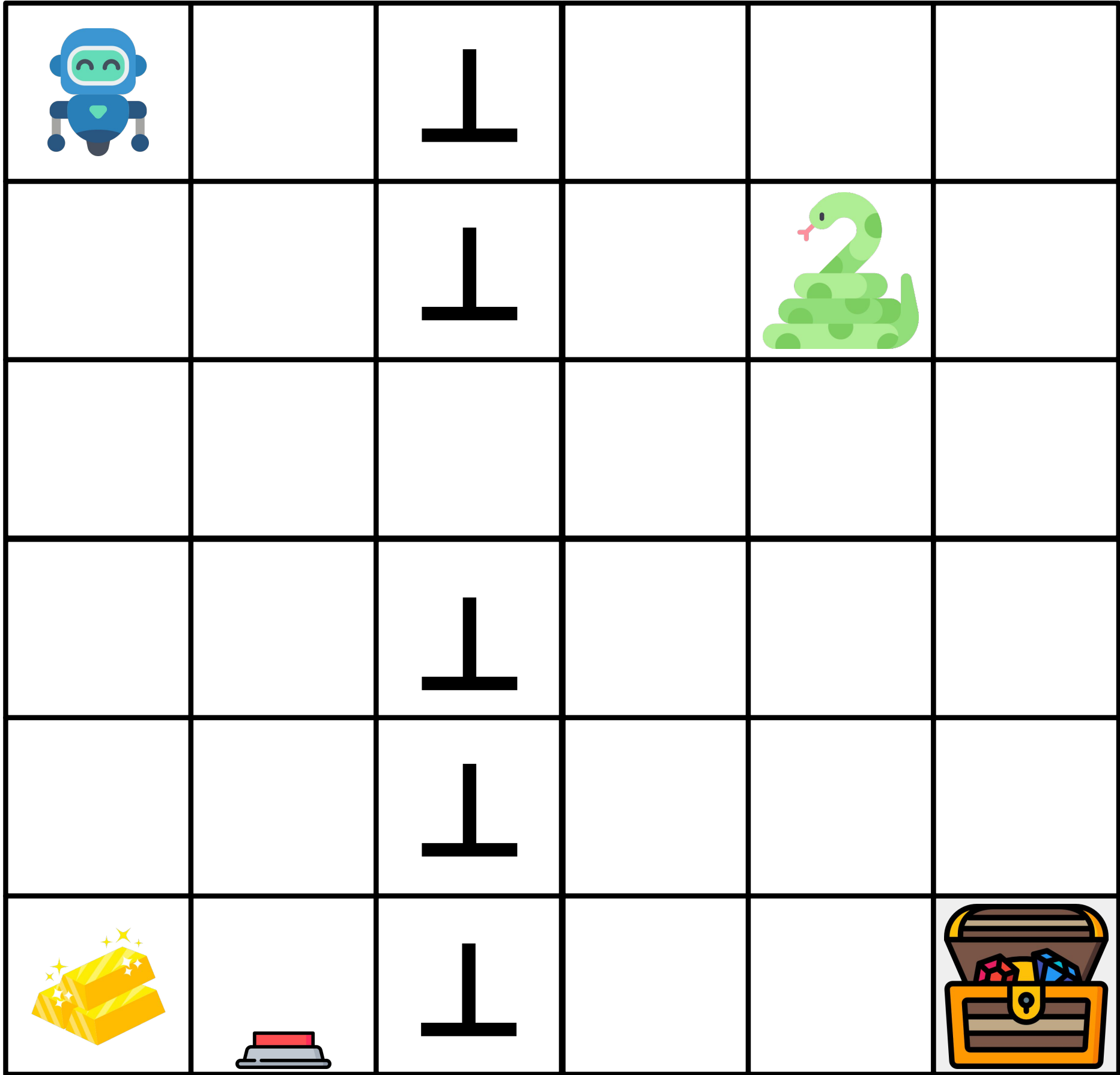
$$\tilde{R}_{opt}(s, a) = r_{\min}^e + \bar{R}^m(s^m, a^m) + \underbrace{\frac{\beta^m}{N(s^m, a^m)}}_{\text{bonus for } r^m} + \underbrace{\frac{\beta}{N(s, a)}}_{\text{bonus for } p}$$

- **Innovation 3:** Explore to observe rewards

$$\tilde{R}_{obs}(s, a) = \text{KL-UCB}(0, N(s, a)) \mathbb{I}\{N(s^e, a^e) = 0\} + \frac{\beta^{obs}}{\sqrt{N(s, a)}},$$

Results

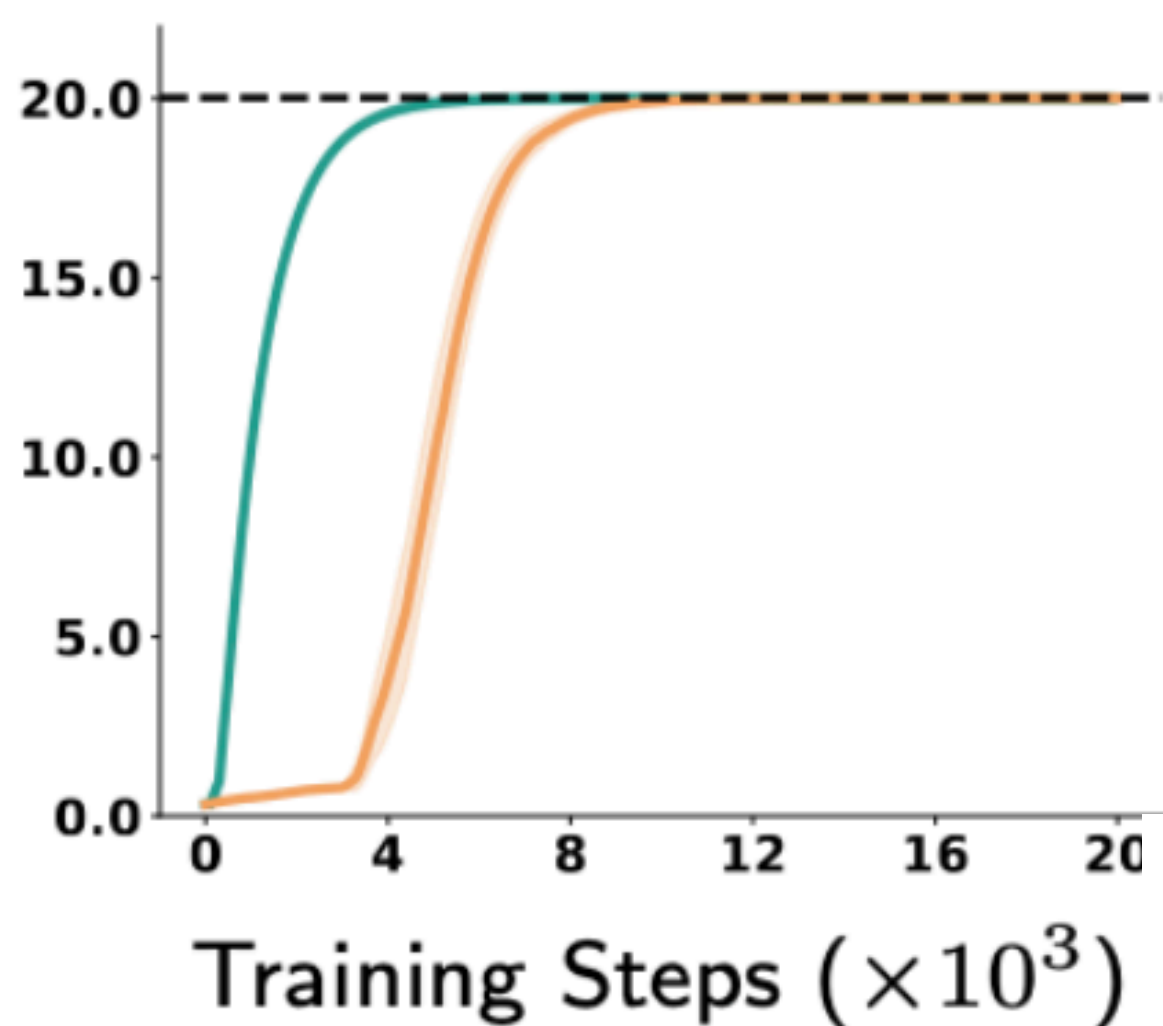
Environments



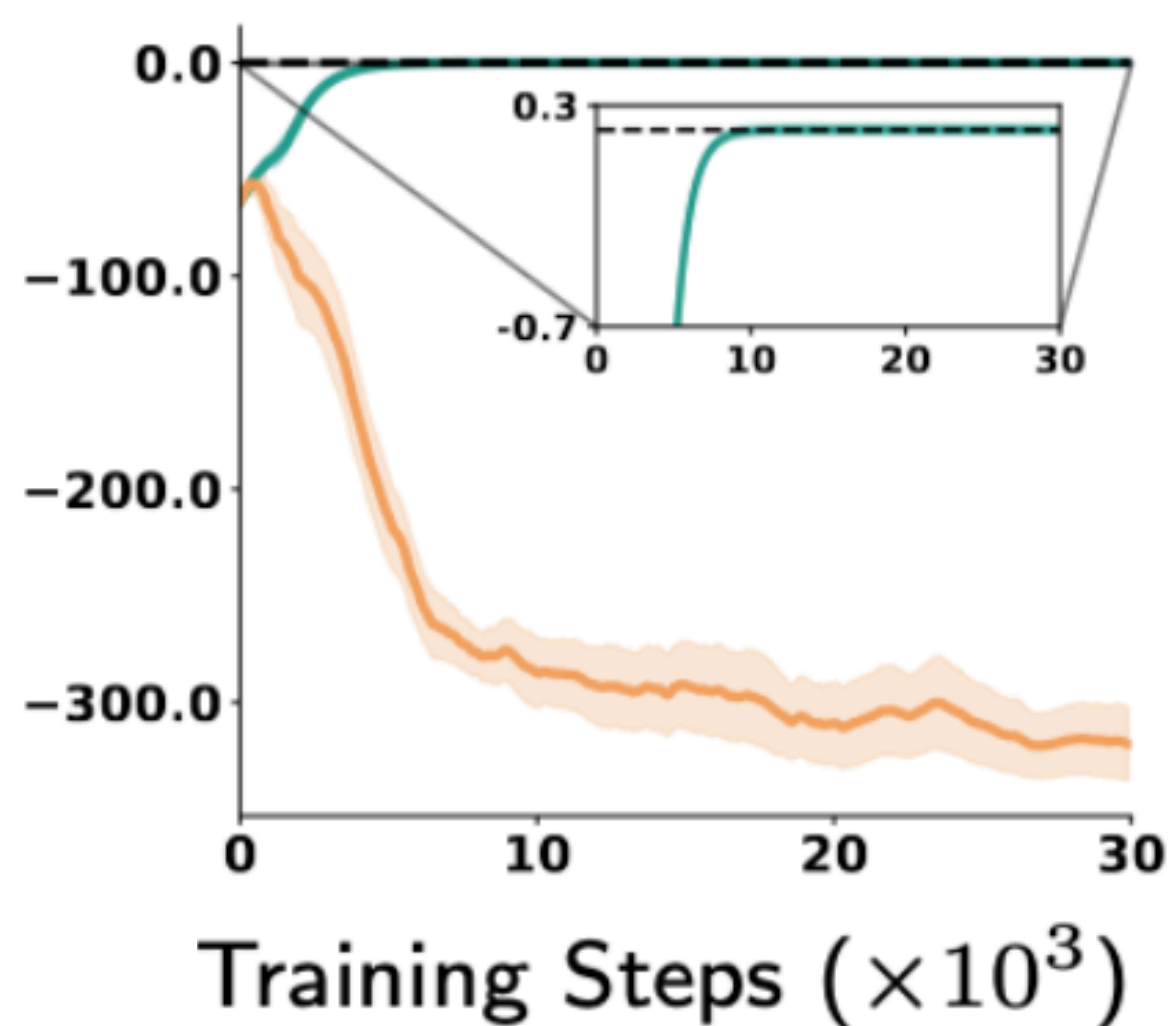
Results

Discounted Test Return

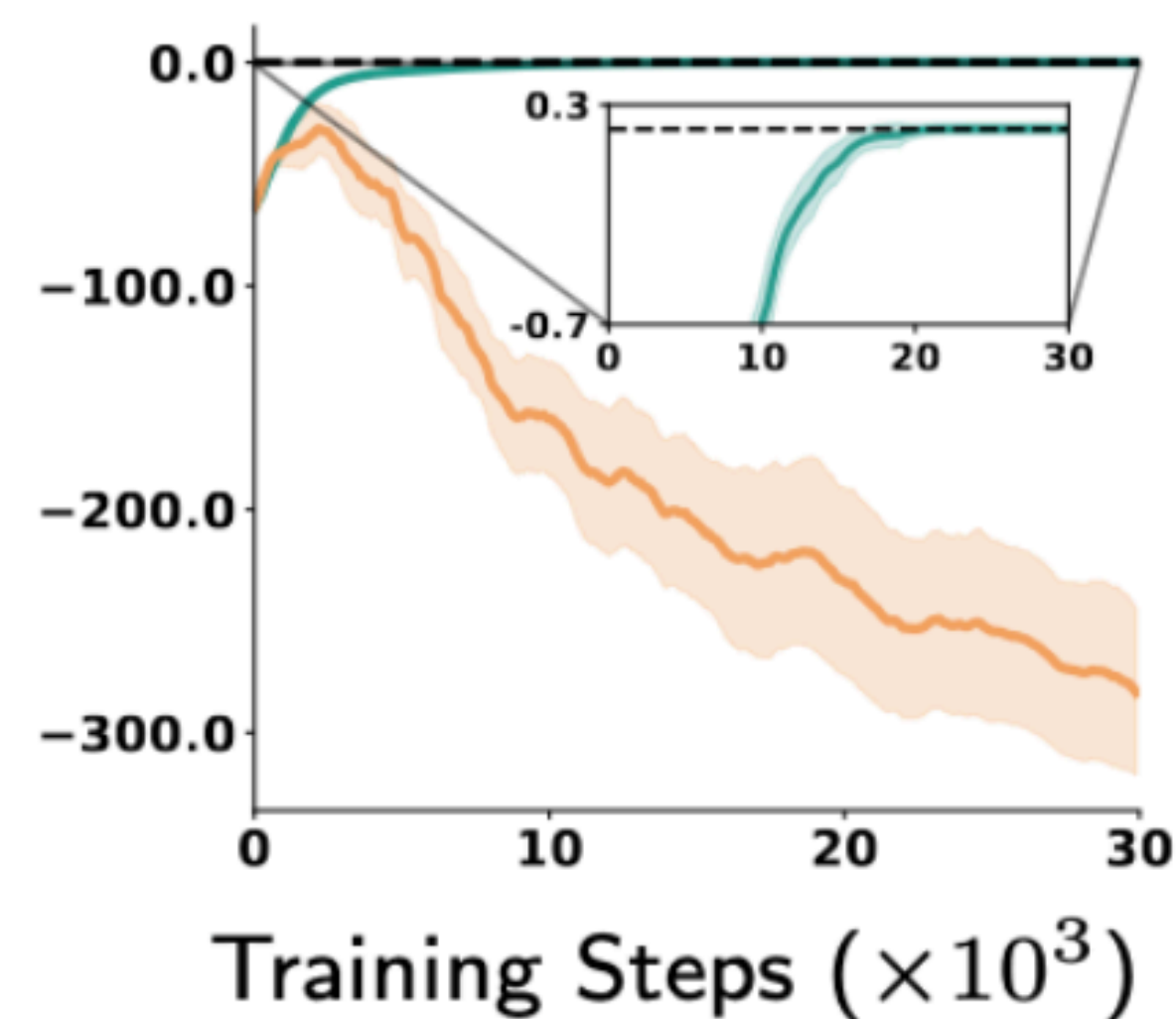
Monitored MBIE-EB Directed- E^2 Known Monitor Minimax-Optimal



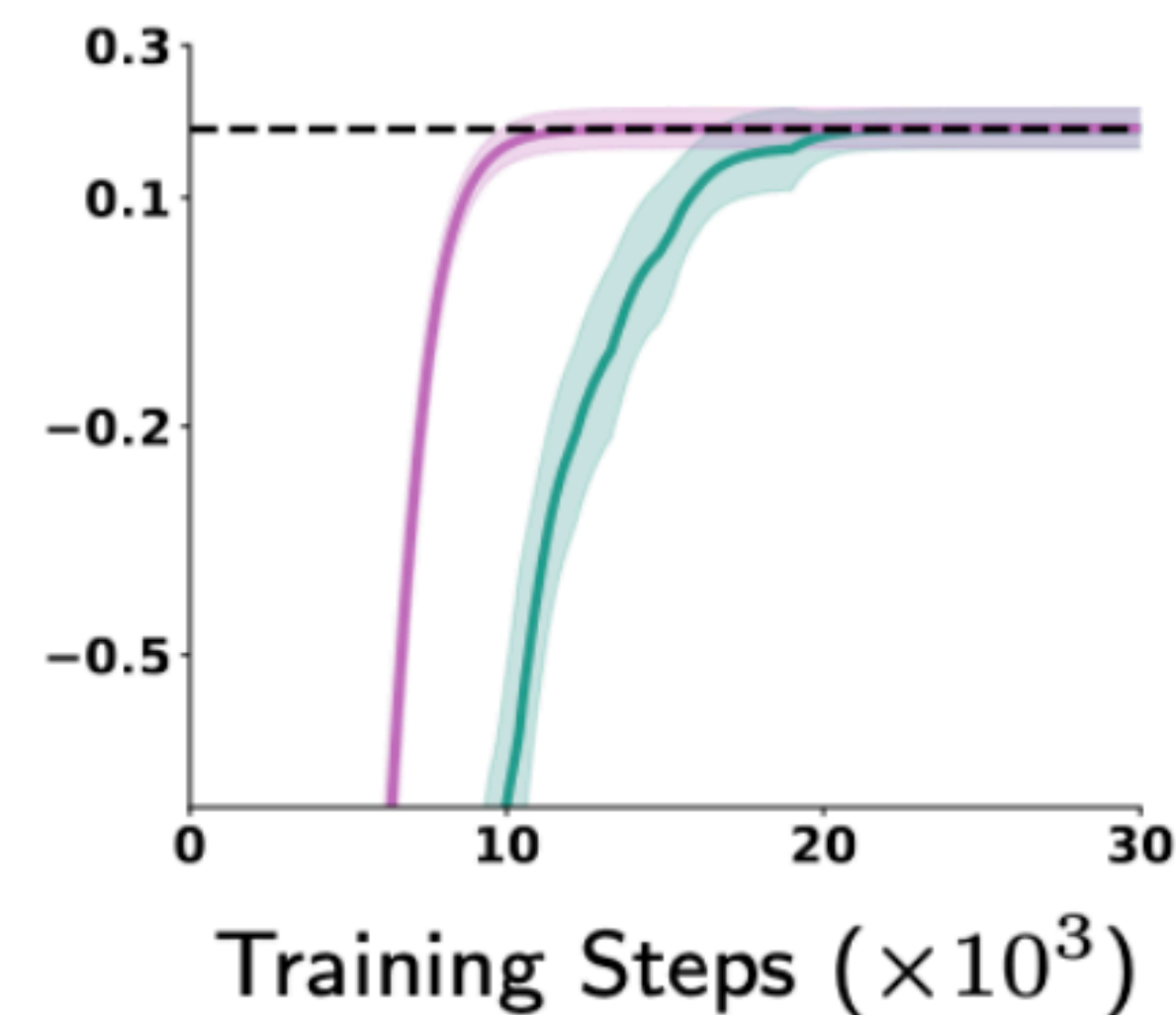
(a) River Swim



(b) Bottleneck (100%)



(c) Bottleneck (5%)



(d) Bottleneck (5%)