

Causally Robust Reward Learning from Reason-Augmented Preference Feedback

Minjune Hwang, Yigit Korkmaz, Daniel Seita[†], Erdem Bıyık[†] ([†]: equal advising)
University of Southern California, Viterbi School of Engineering

Motivation

- Preference-based reinforcement learning shapes agent behaviors from user's binary preferences, which can leave it vulnerable to **causal confusion**.
- The learned reward can latch onto **spurious features** that cooccur with preferred trajectories, collapsing when those correlations disappear at test time.

Insight: A natural language reason can (1) clarify true causal signals behind preference and (2) improve generalization beyond spurious features.

Method: ReCouPLE

- We learn a trajectory encoder ϕ that maps a trajectory τ into the pretrained language encoder's embedding space.
- We represent each task's reward given its task description ℓ_{task} as **an inner product** of $\phi(\tau)$ and task embeddings $\theta = \text{LM}(\ell_{\text{task}})$:

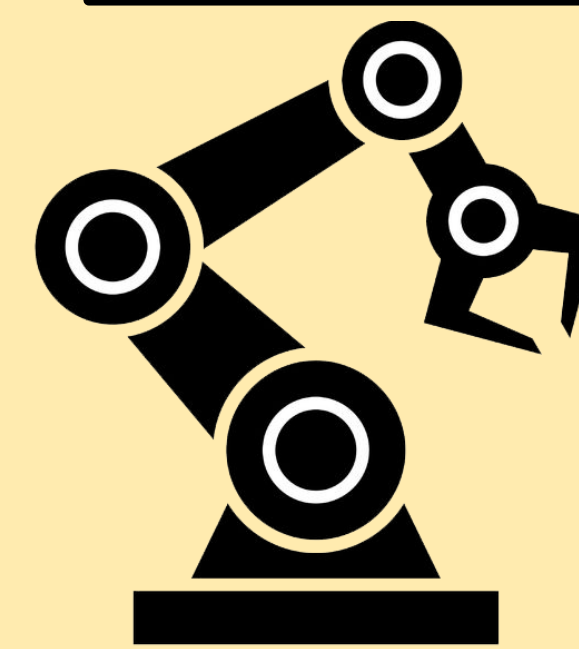
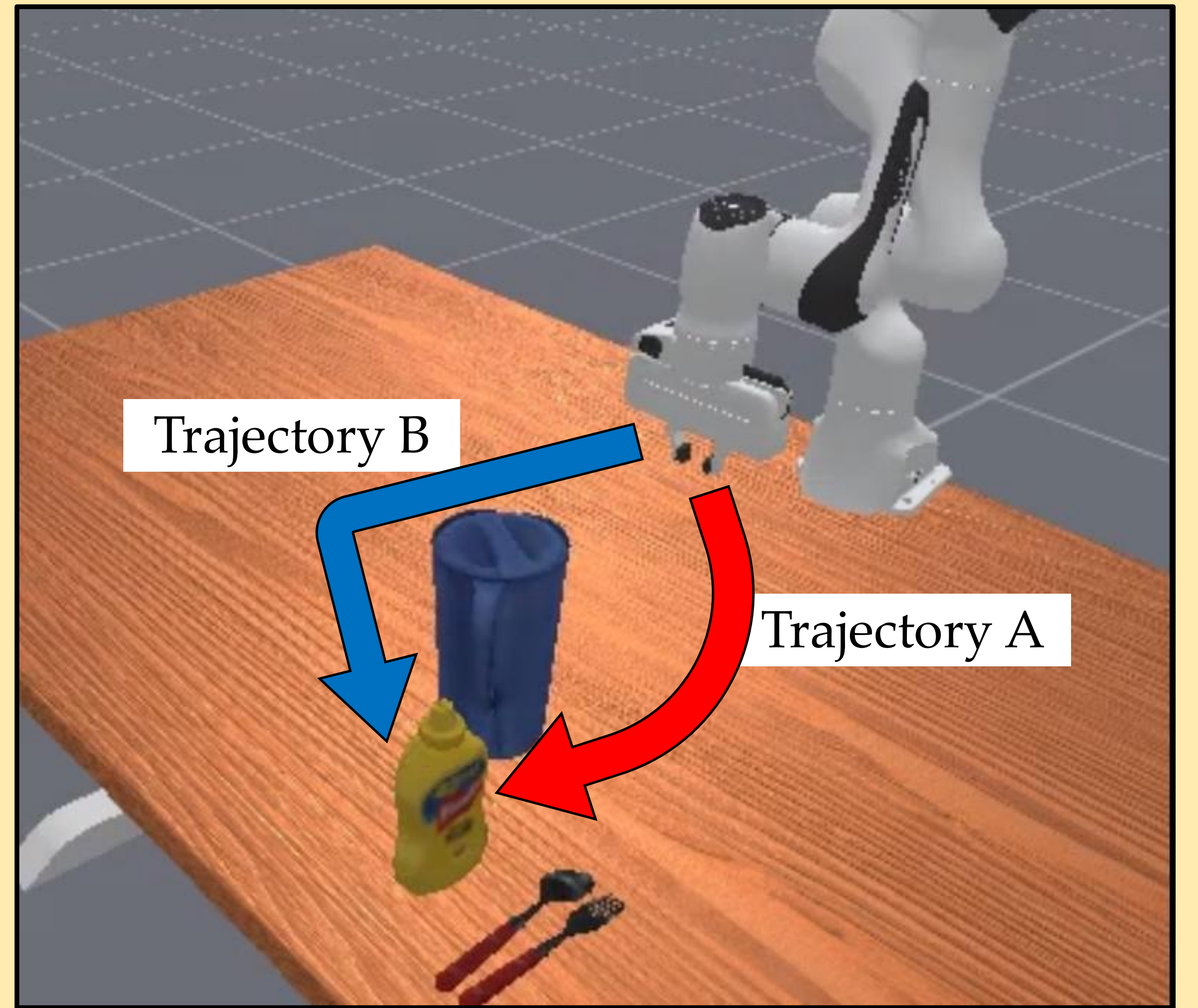
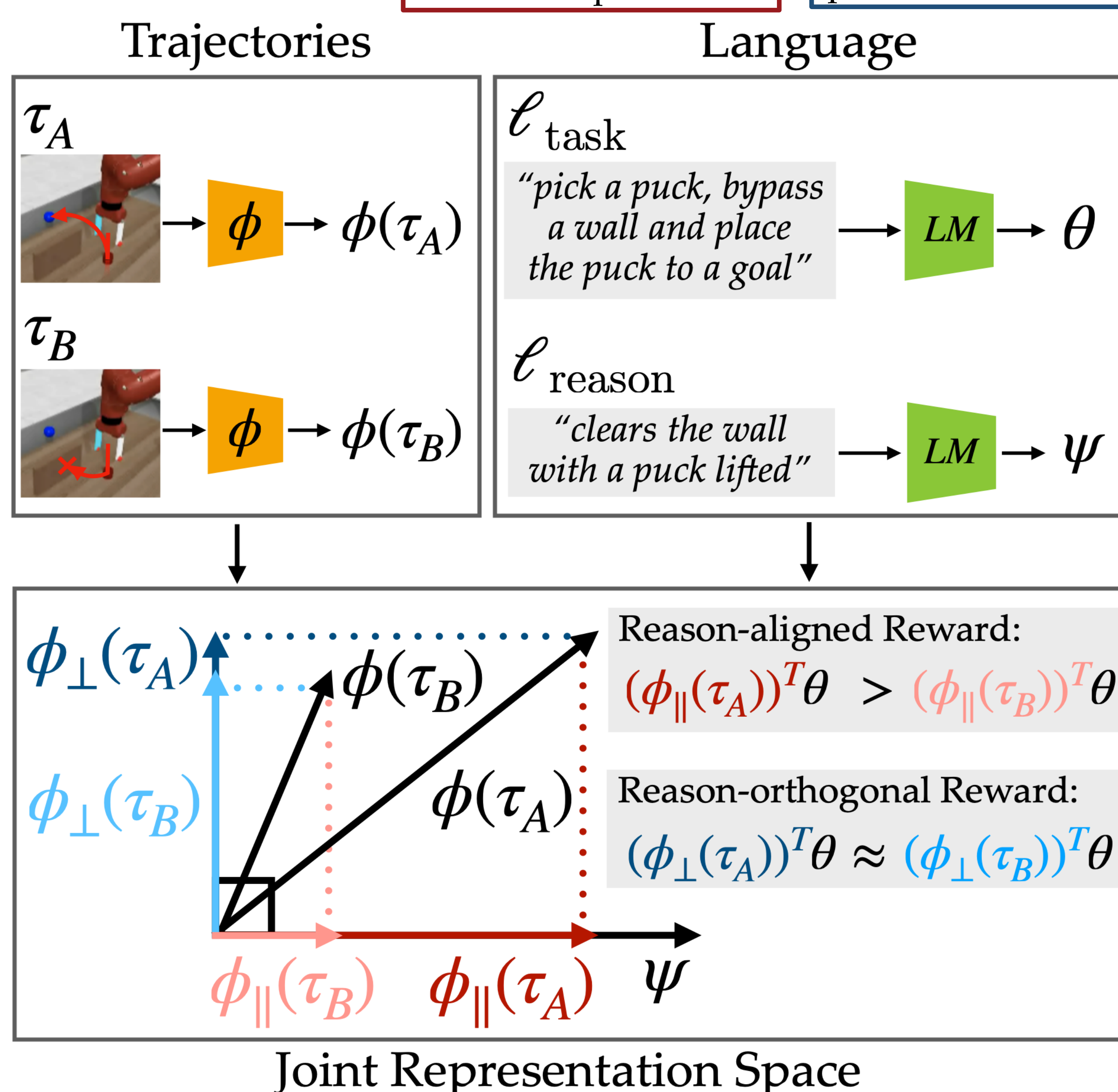
$$r(\tau, \ell_{\text{task}}) = \phi(\tau)^T \theta$$

ReCouPLE treats the rationale embedding ψ as a **projection axis**, splitting the trajectory representation into reason-aligned and reason-orthogonal parts:

$$\phi_{\parallel}(\tau) = \left(\frac{\phi(\tau)^T \psi}{\|\psi\|_2^2} \right) \psi, \quad \phi_{\perp}(\tau) = \phi(\tau) - \phi_{\parallel}(\tau)$$

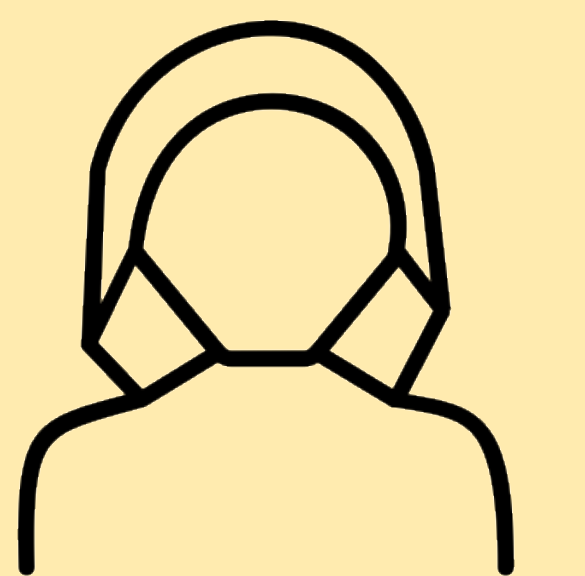
$$\begin{aligned} \text{Then, } r(\tau, \ell_{\text{task}}) &= \phi_{\parallel}(\tau)^T \theta + \phi_{\perp}(\tau)^T \theta \\ &= \underbrace{r_{\parallel}(\tau)}_{\text{reason-aligned}} + \underbrace{r_{\perp}(\tau)}_{\text{reason-orthogonal}} \end{aligned}$$

Should mainly determine preference Should not affect preference as much



Could you tell me which trajectory you prefer?

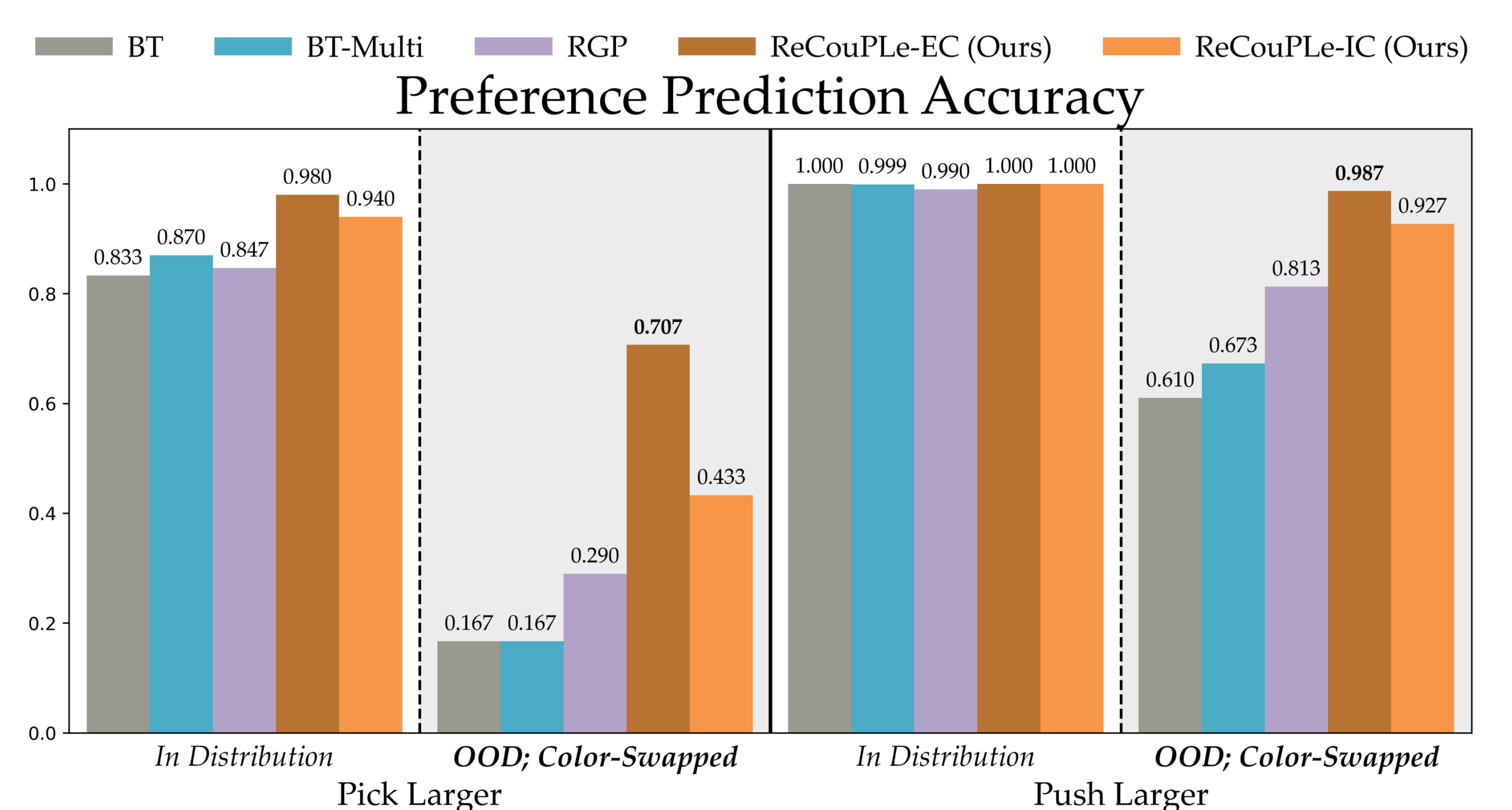
I prefer Trajectory A, **because it smoothly avoids the obstacle.**



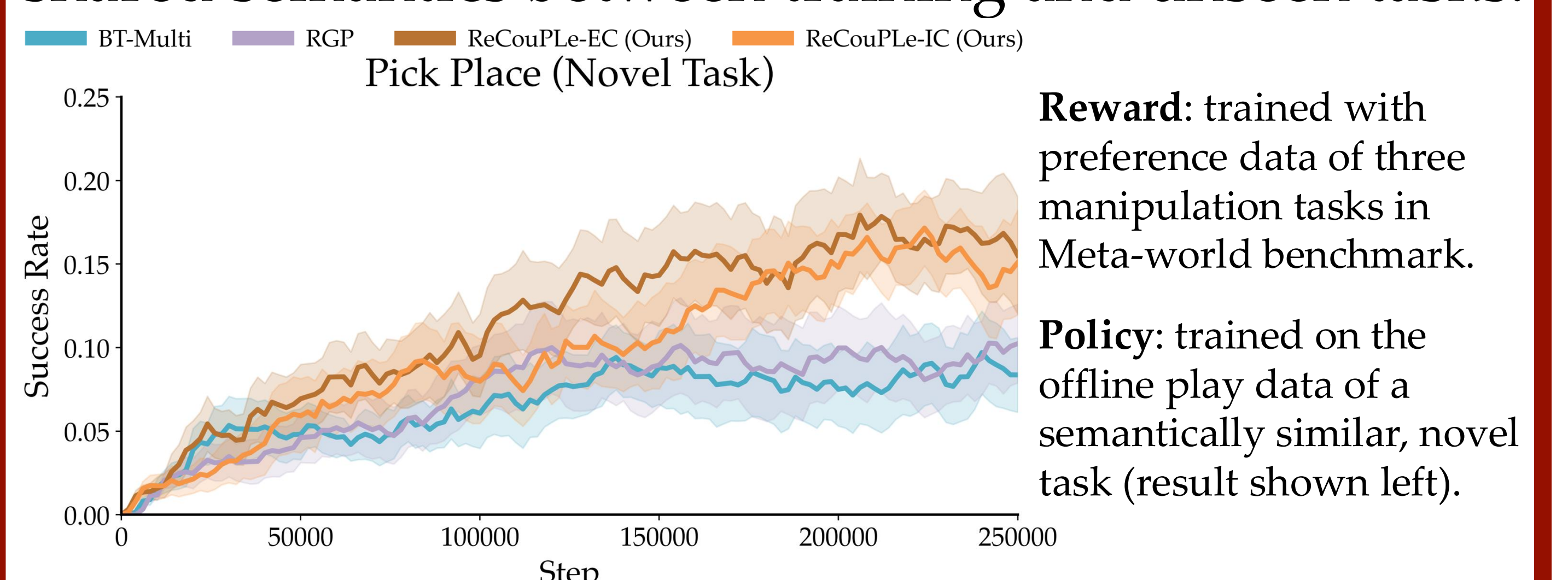
How can we use **this information**, i.e., **reasoning**?

Results

1. ReCouPLE robustly predicts user's preference under the distribution shifts.



2. ReCouPLE transfers to a novel task, thanks to shared semantics between training and unseen tasks.



3. Ablation shows necessities of each component.

Model	2-task		4-task	
	ID	OOD	ID	OOD
ReCouPLE	0.995	0.872	1.000	0.878
ReCouPLE-no-consistency	0.980	0.726	0.977	0.745
ReCouPLE-no-consistency-no-ratio	0.987	0.727	0.990	0.730