

RePreL: Integrating Relational Planning and Reinforcement Learning for Effective Abstraction

Harsha Kokel, Arjun Manoharan, Sriraam Natarajan, Balaraman Ravindran, Prasad Tadepalli



Aim Learn to act in relational domains, with varying task, number of objects and relations.

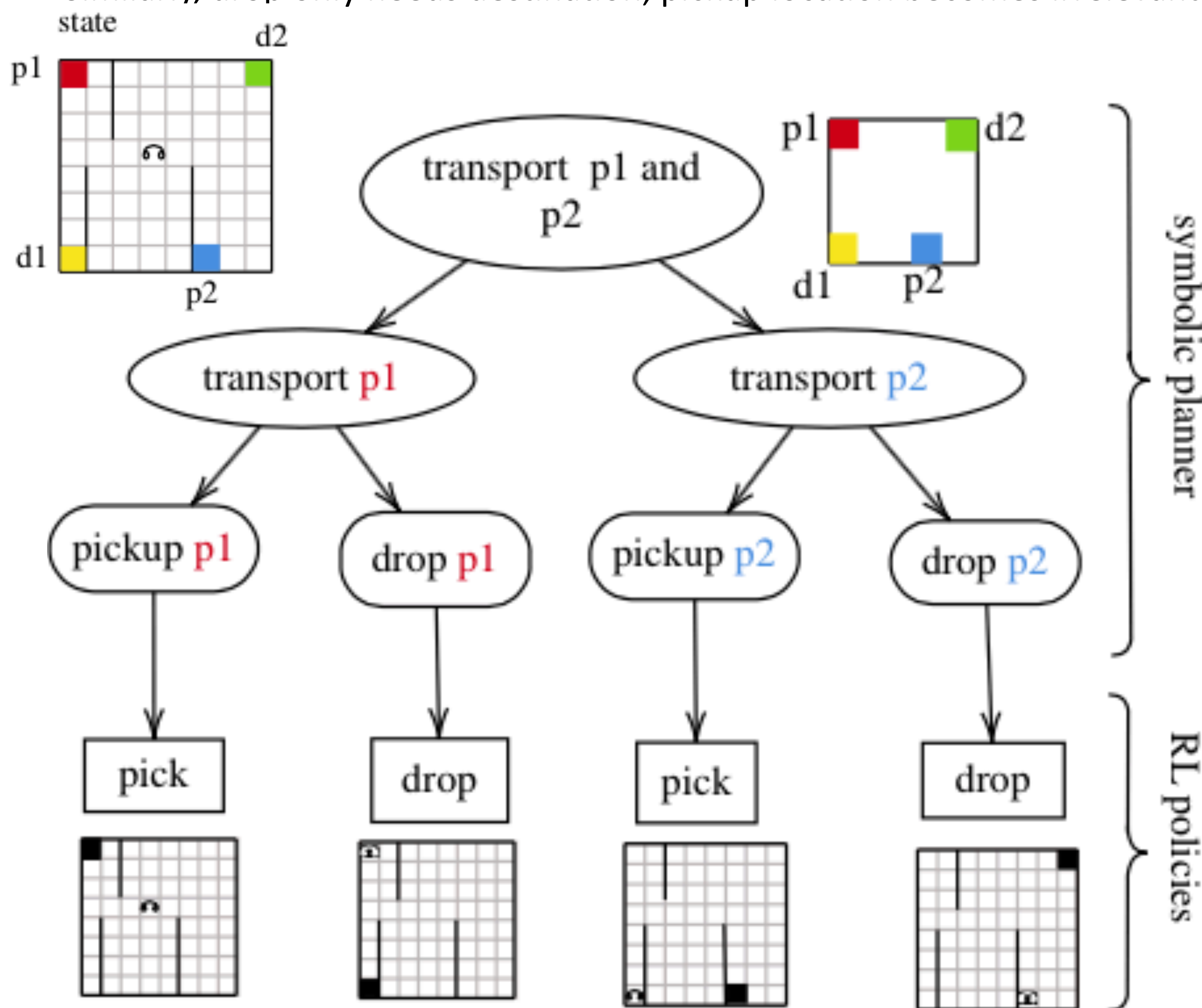
Motivation Humans are quite natural at creating an abstract representation for efficient planning. They use different abstractions for different tasks. Can we leverage their domain knowledge and provide such abstractions to the RL agent?

RePreL

- Hierarchical Framework with Planner and RL
- Plan sequence of operators at High level and learn to execute each operator at lower level
- Advantages
 - Planner provides compositionality
 - Different RL agent for execution allows separate state representations
- Abstraction is most useful when we have multiple relations and objects, so consider Relational MDP.
- Use D-FOCI statements to identify relevant objects and relations for a given task.

Toy Example

- Taxi domain with multiple passengers
- Acquire the sequence of passenger pickup and drop operators from planner
- Learn to execute the operators using RL
- Use different goal-conditioned RL for each operator
- Pickup operator only needs pickup location, drop location is irrelevant
- Similarly, drop only needs destination, pickup location becomes irrelevant



Model-agnostic Abstraction

Definition 4 (Li, Walsh, and Littman (2006)). A *model-agnostic abstraction* $\phi(s)$ is such that for any action a and abstract state \bar{s} , $\phi(s_1) = \phi(s_2)$ if and only if

$$\sum_{\{s'_1 | \phi(s'_1) = \bar{s}\}} R_o(s_1, a, s'_1) = \sum_{\{s'_2 | \phi(s'_2) = \bar{s}\}} R_o(s_2, a, s'_2)$$

$$\sum_{\{s'_1 | \phi(s'_1) = \bar{s}\}} P_o(s_1, a, s'_1) = \sum_{\{s'_2 | \phi(s'_2) = \bar{s}\}} P_o(s_2, a, s'_2)$$

- This is the bisimulation condition, where the abstraction function preserves the transition and the reward distribution of ground MDP in the abstract MDP

D-FOCI Statements

- First Order Conditional Influence (FOCI) statement is of the form "if *condition* then X_1 influence X_2 "
- FOCI statements encode the information that literal X_2 is influenced only by the literals in X_1 when the stated condition is satisfied
- We present Dynamic FOCI (D-FOCI) statements,

$$\text{operator} : \{p(X_1), q(X_1)\} \xrightarrow{+1} q(X_1)$$
- With +1, the D-FOCI encode the influence of literals in current time step on the literal in next time-step
- With operator, D-FOCI statement constraints the task being executed
- Given the current state and operator being executed we can obtain complete set of relevant state variables by grounding and rolling out the D-FOCI statements. This forms our abstract state
- Guarantees safe model-agnostic abstraction, if MDP satisfies the D-FOCI statement with fixed-depth unrolling

state : $\{at(p1, r), taxi-at(13), dest(p1, y), \neg at-dest(p1), \neg in-taxi(p1), at(p2, b), dest(p2, g), \neg at-dest(p2), \neg in-taxi(p1)\}$

$[at-p1, dest-p1, in-taxi-p1, at-dest-p2, at-p2, dest-p2, in-taxi-p2, at-dest-p2, taxi-at, move]$

operator $\langle pickup(P), \{P/p1, L/r\} \rangle$

D-FOCI : $\{taxi-at(L1), move(Dir)\} \xrightarrow{+1} taxi-at(L2)$

$\{taxi-at(L1), move(Dir)\} \rightarrow R$

pickup(P):

$\{taxi-at(L1), at(P, L), in-taxi(P)\} \xrightarrow{+1} in-taxi(P)$

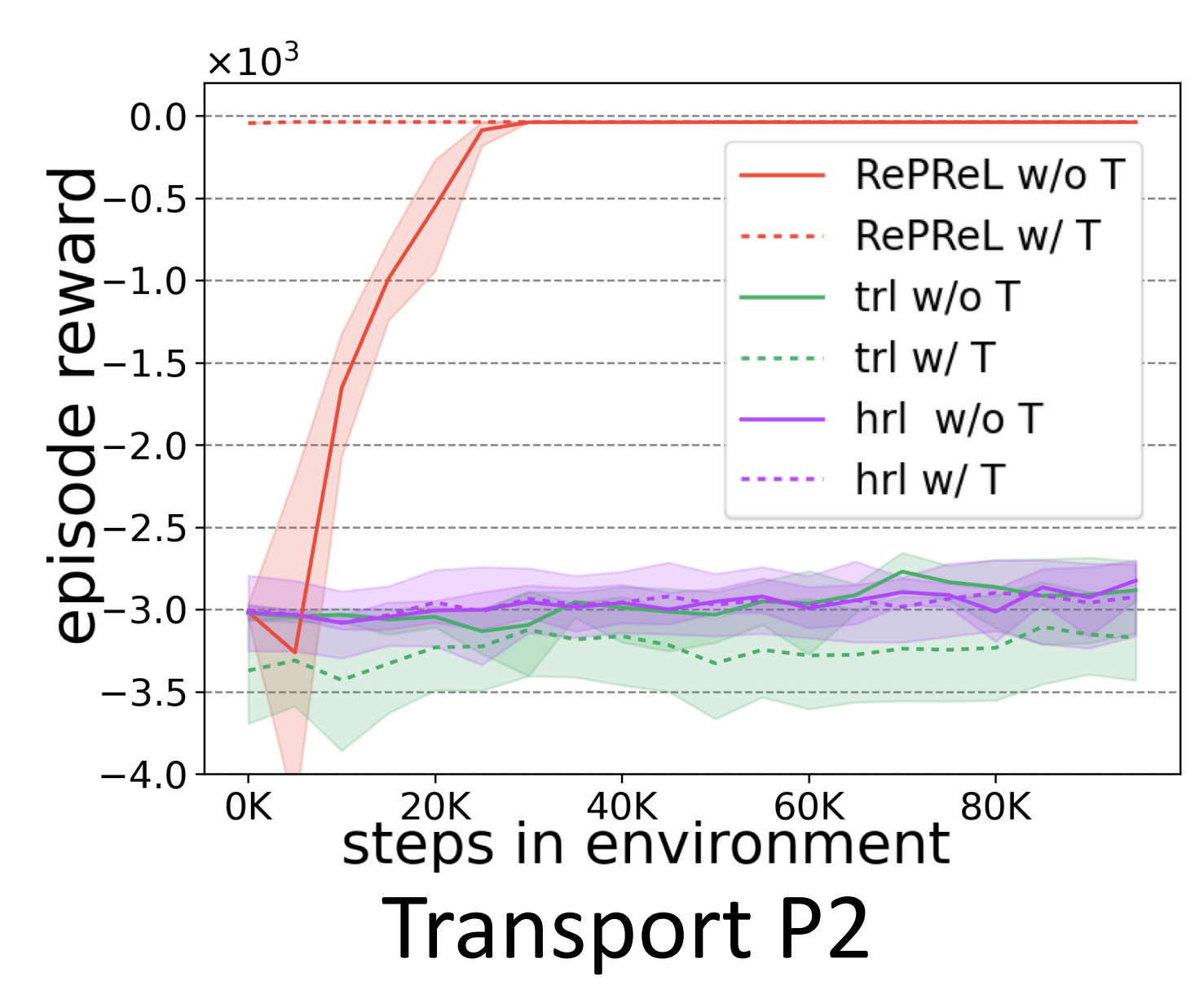
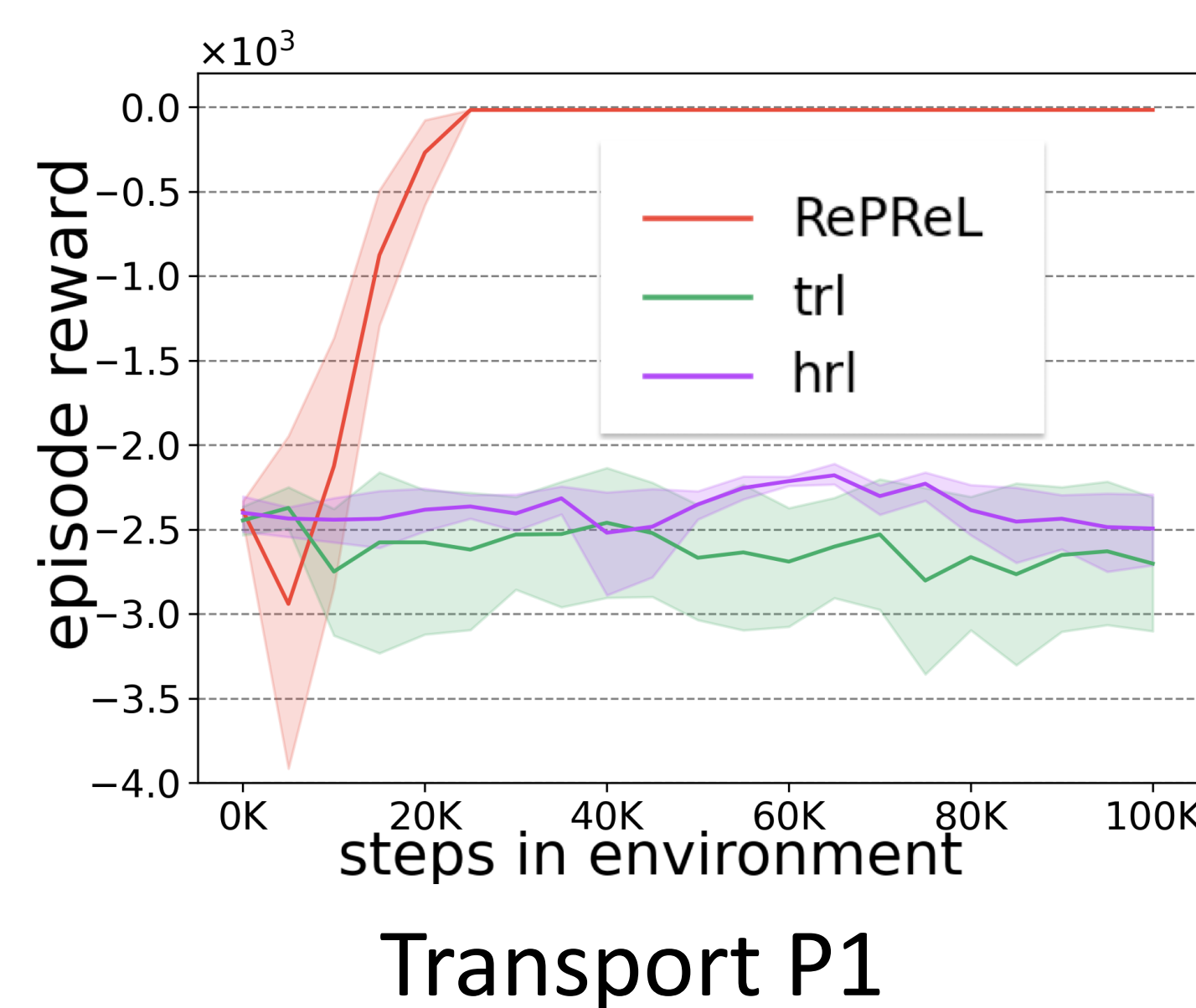
pickup(P): $in-taxi(P) \rightarrow R_o$

abstract state: $\{at(p1, r), taxi-at(13), \neg in-taxi(p1), move(Dir)\}$

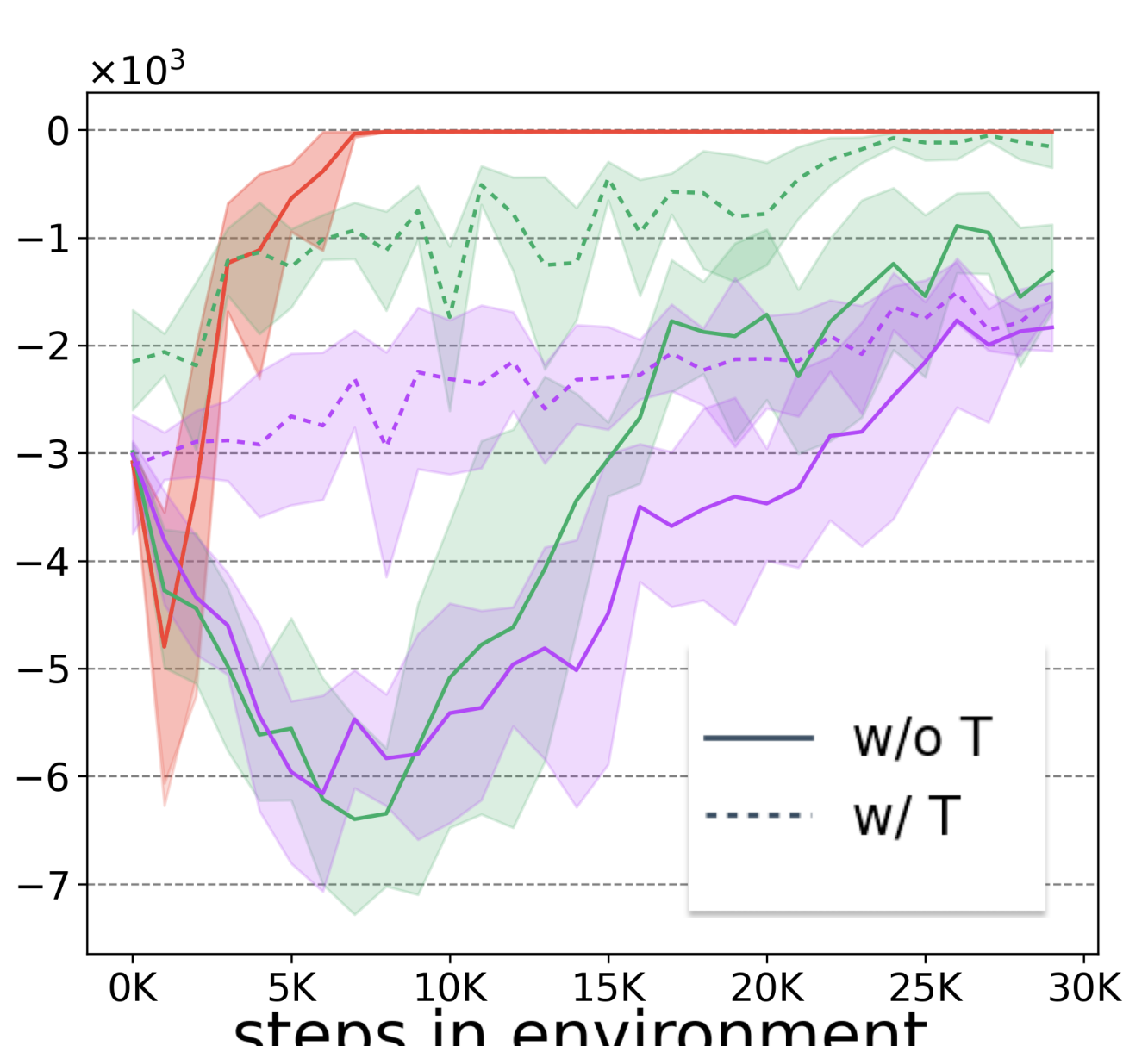
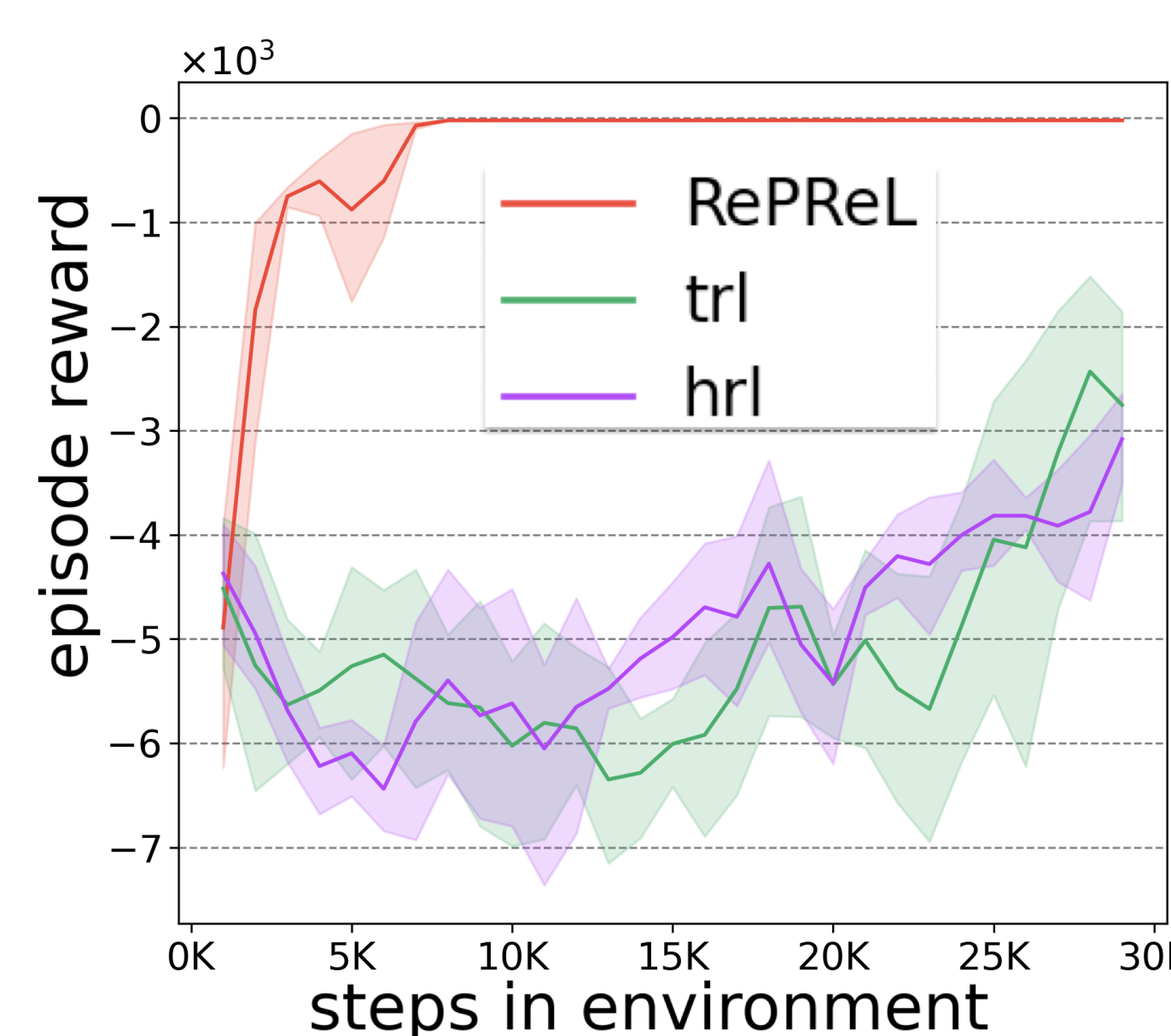
$[at, taxi-at, in-taxi, move]$

Experiments

Taxi domain



Office domain



Human knowledge can help provide effective abstractions

that enable efficient learning and effective generalization across tasks and objects.