

RePReL

Integrating Relational Planning and Reinforcement Learning for Effective Abstraction



Sriraam Natarajan



Harsha Kokel



Balaraman Ravindran



Arjun Manoharan



Prasad Tadepalli



RL and Planning

- Reinforcement Learning
 - + *Proven successful in complex games*
 - + *Adaptive and robust against uncertainties*
 - *Relies on huge amount of data*
 - *Not effective for long-horizon problems*
 - *Generalization to different task*
- Planning
 - + *Not data but prior-knowledge*
 - + *Better generalization*
 - *May not capture uncertainties*

Planner + RL

- Policy Sketches¹ for modular policies that address multi-task RL
- PLANQ-learning² uses planner to shape reward function that guides the Q-learner.
- PEORL³ (Planning–Execution–Observation–Reinforcement-Learning) framework uses symbolic planner to guide the exploration and learning in RL
- TMP-RL⁵ and PDDL+HER⁶ framework use integrated approach for robotic systems with uncertainty and continuous state space
- Taskable-RL⁷ formalizes the high-level planner and low-level RL executioner setup

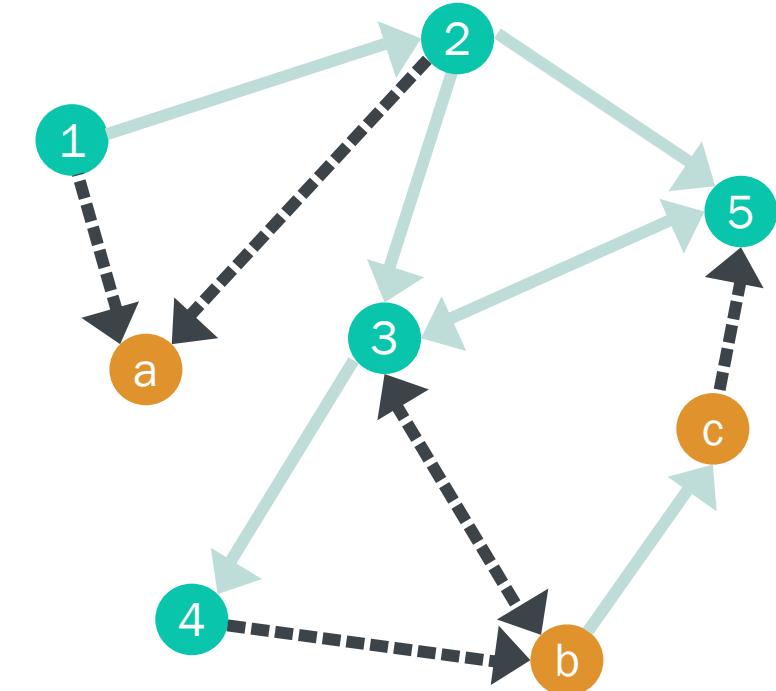
¹Andreas, Klein, and Levine ICML 2017

²Grounds and Kudenko, AAMAS 2008; ³Yang, Lyu, Liu, and Gustafson IJCAI 2018; ⁵Jiang, Yang, Zhang, and Stone, IROS 2019;

⁶Eppe, Nguyen, and Wermter, Front. in Rob. and AI 2019; ⁷Illanes, Yan, Icarte, and McIlraith ICAPS 2020

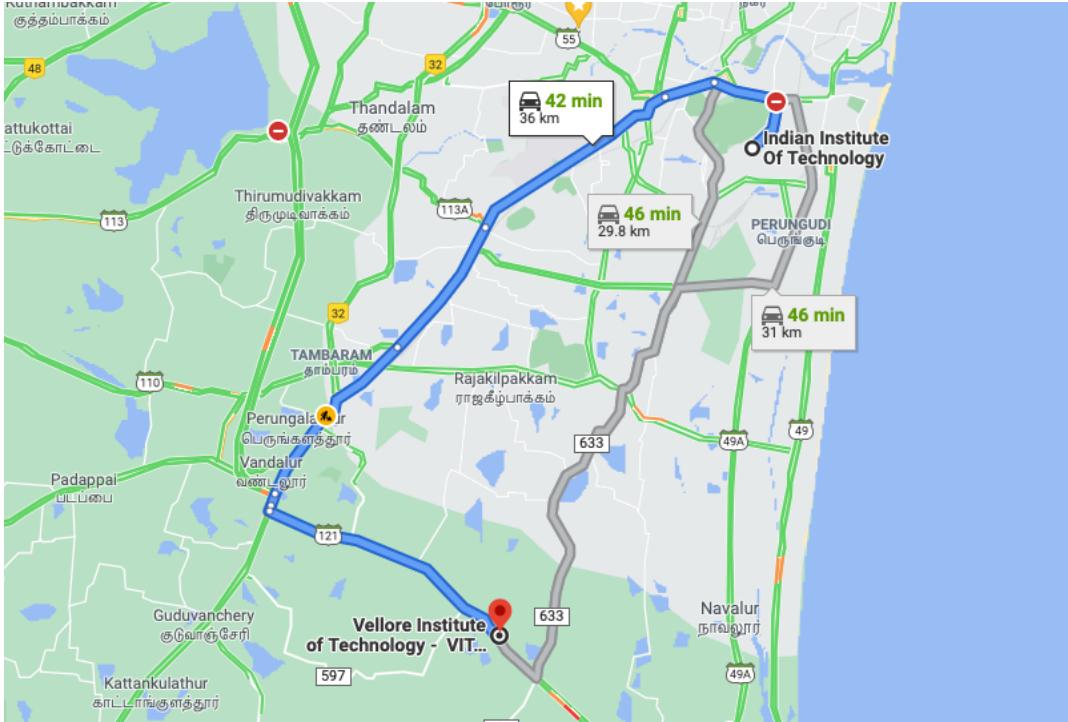
Motivation

Relational
domains



Motivation

Abstract Representations



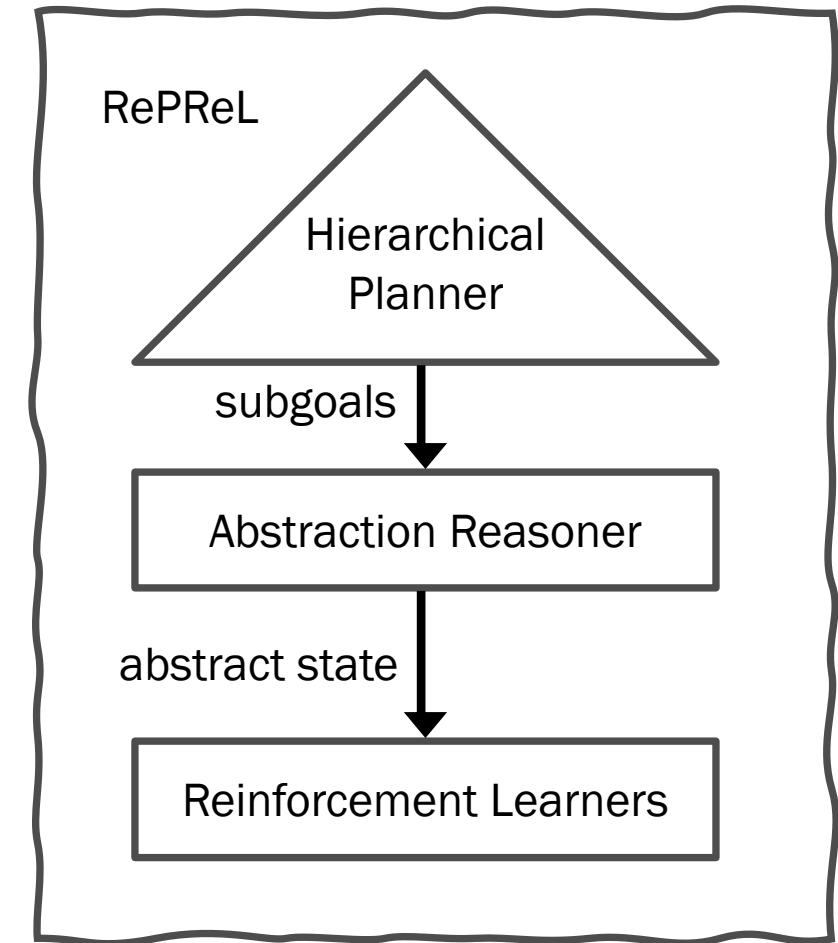
Planning



Execution

RePReL

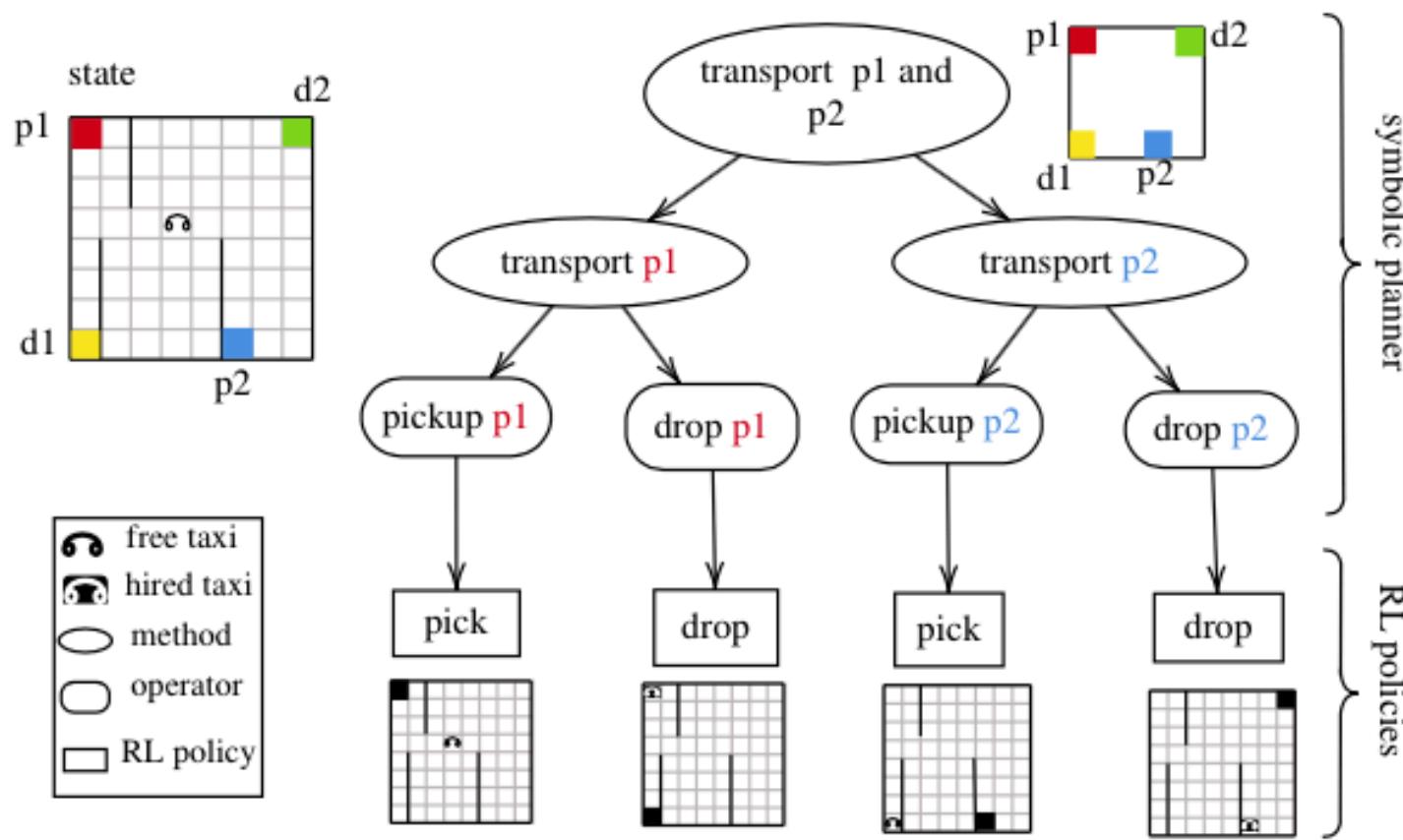
- Plan the sequence of high level subgoals and learn to execute each subgoal at lower level



RePREL

Goal directed relational MDP:

$\langle S, A, P, R, \gamma, G \rangle$



RePReLU

Definition 3. The subgoal RMDP M_o for each operator o is defined by the tuple $\langle S, A, P_o, R_o, \gamma \rangle$ consisting of states S , actions A , transition function P_o , reward function R_o , and discount factor γ . State and Actions remain same as the original RMDP. The reward function R_o and transition probability distribution function P_o are defined as follows:

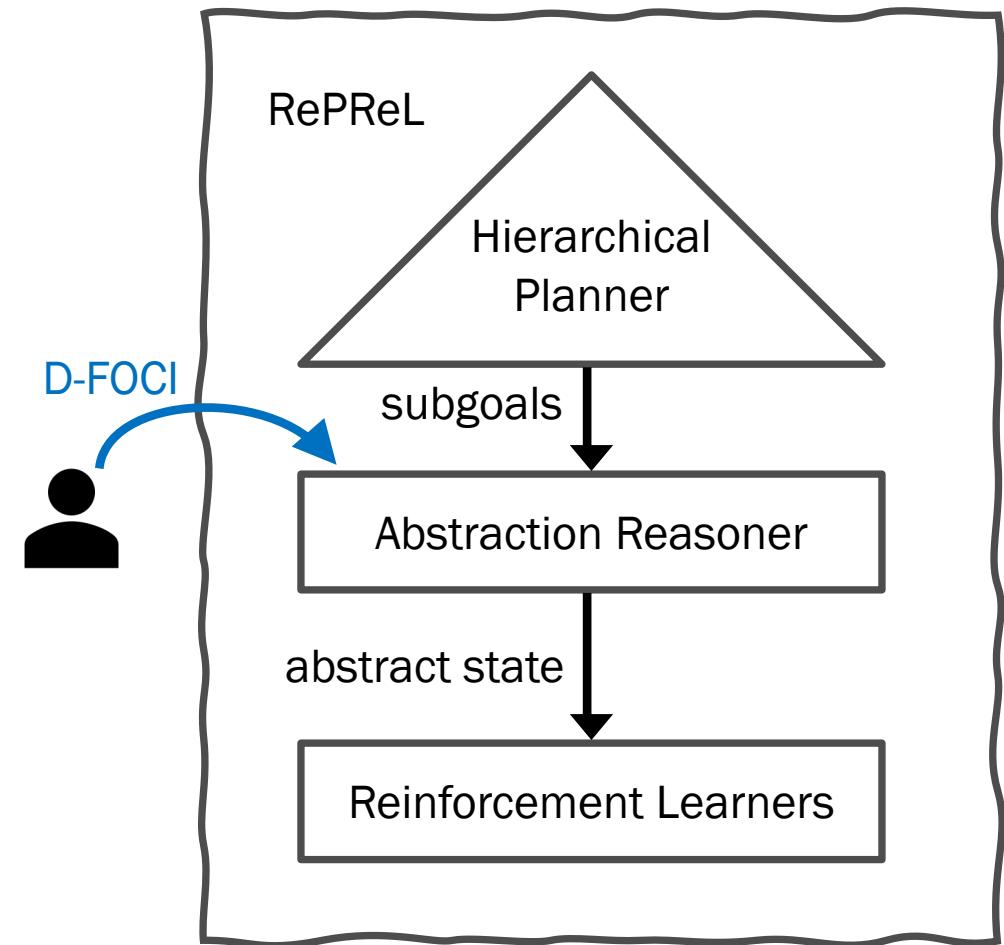
$$R_o(s, a, s') = \begin{cases} t_R + R(s, a, s') & \text{if } s' \in \beta(o) \text{ and } s \notin \beta(o) \\ 0 & \text{if } s' \in \beta(o) \text{ and } s \in \beta(o) \\ R(s, a, s') & \text{otherwise} \end{cases}$$

$$P_o(s, a, s') = \begin{cases} 0 & \text{if } s \in \beta(o) \text{ and } s' \notin \beta(o) \\ 1 & \text{if } s \in \beta(o) \text{ and } s' \in \beta(o) \\ P(s, a, s') & \text{otherwise} \end{cases}$$

with $R(s, a, s')$ indicating the reward function from the original GRMDP definition. t_R is a fixed terminal reward.

RePReL

- Plan the sequence of high level subgoals and learn to execute each subgoal at lower level
- Advantage:
 - *Compositionality*
 - *Task specific state representations*
- Dynamic First Order Conditional Influence (D-FOCI) statements to obtain task-specific abstract representations



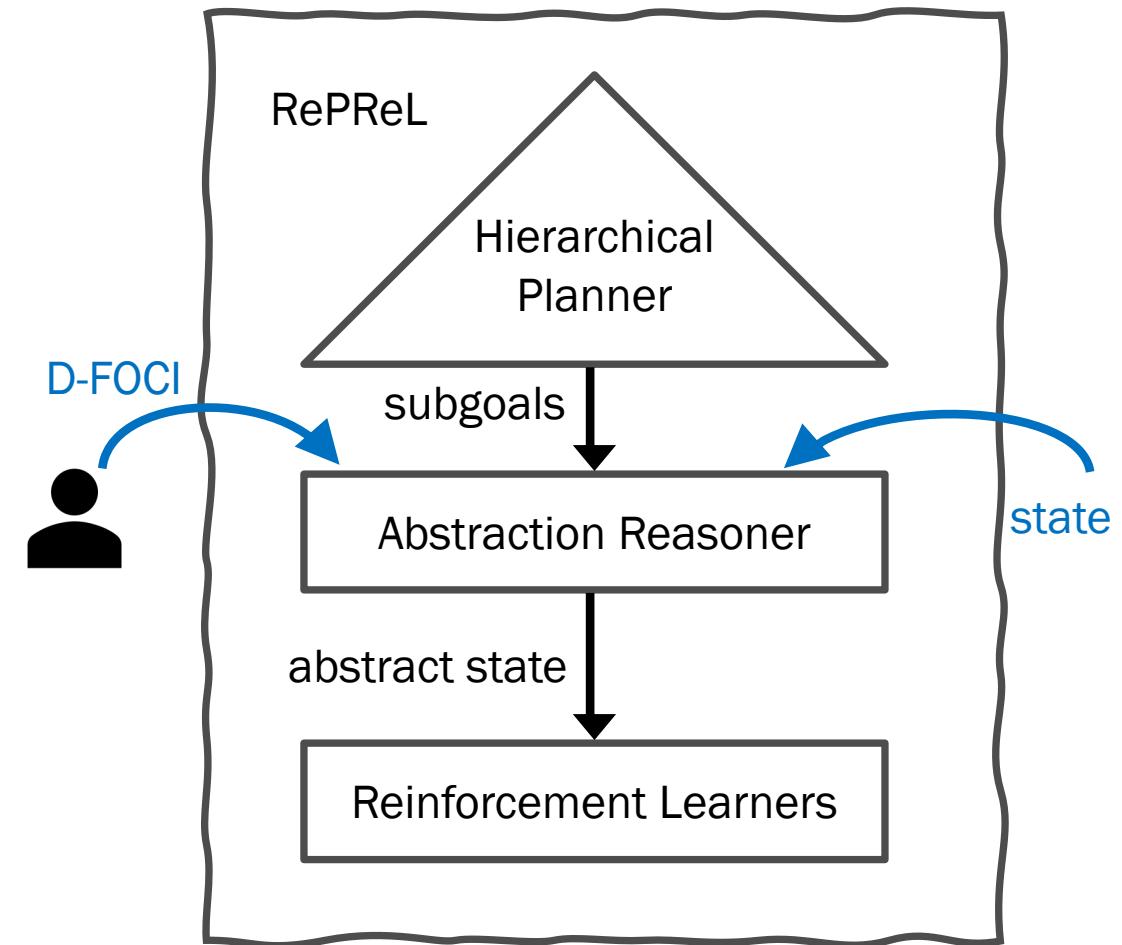
D-FOCI

First Order Conditional Influence (FOCI) statements

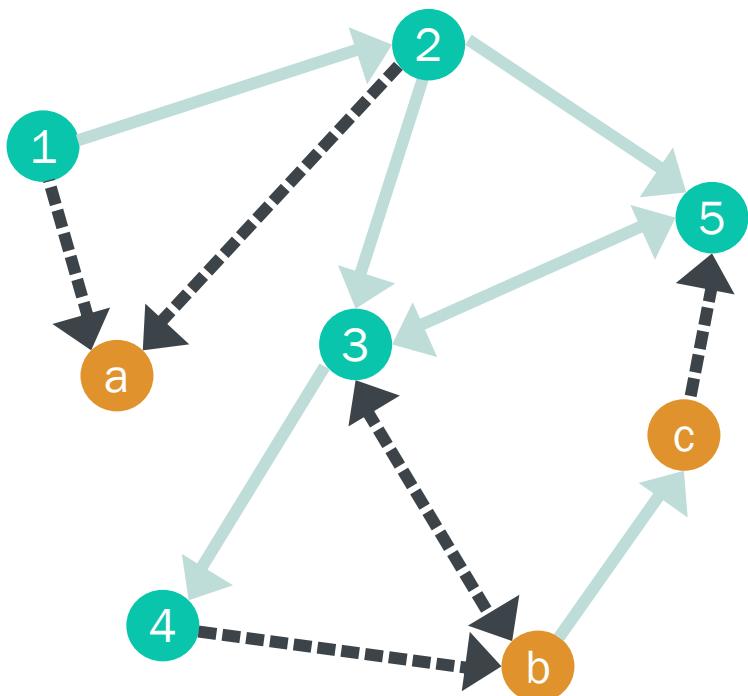
if *condition* then X_1 influence X_2

Dynamic FOCI statements

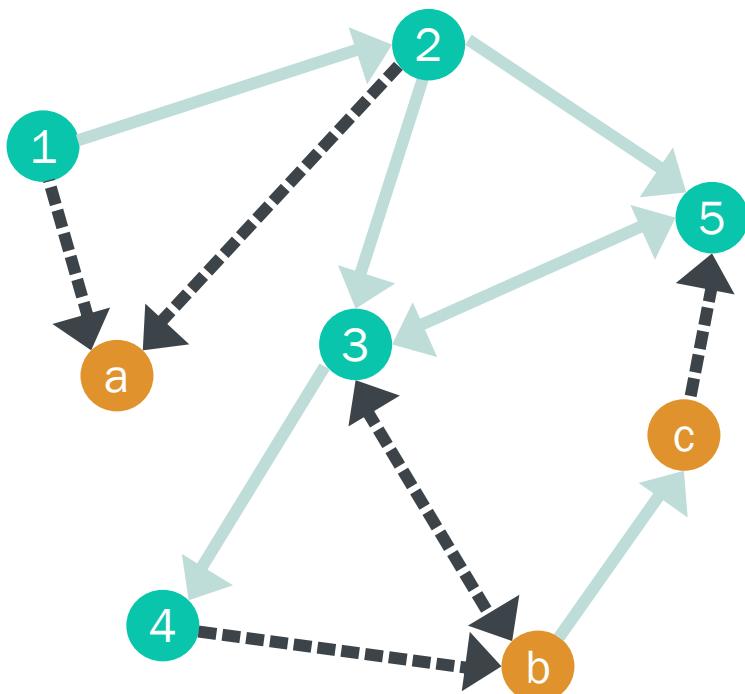
sub-task : $X_1 \xrightarrow{+1} X_2$



Example

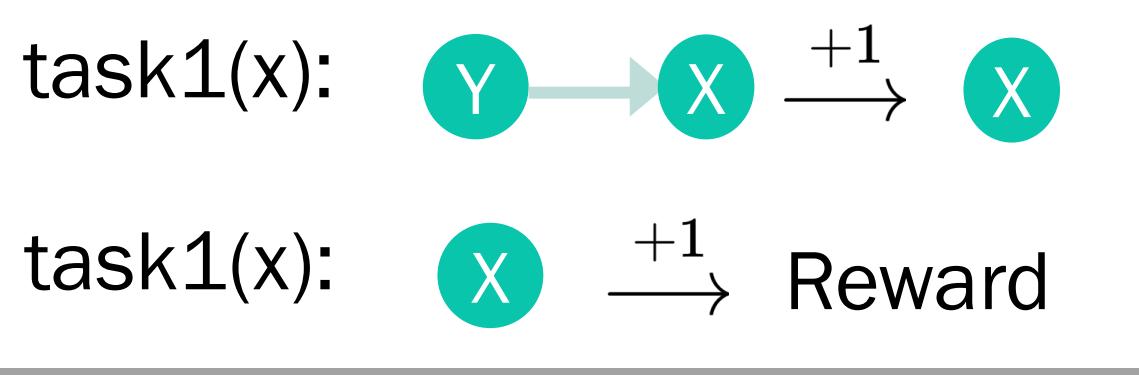
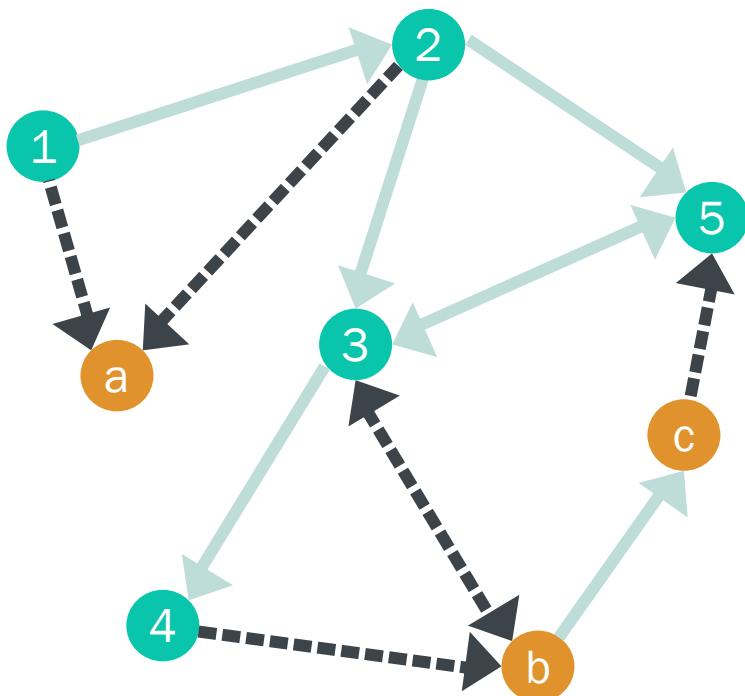


Example

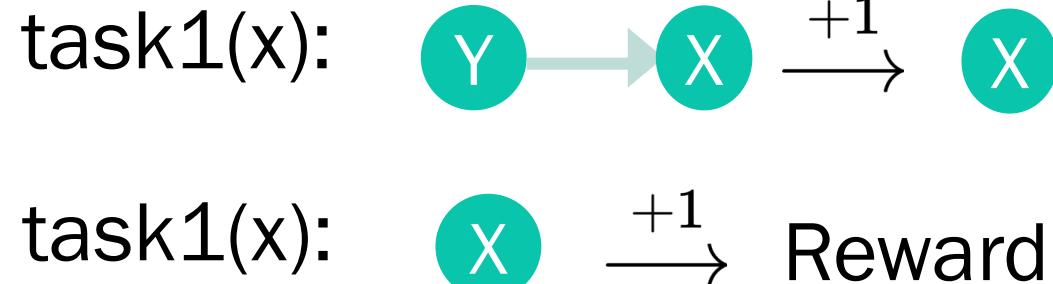
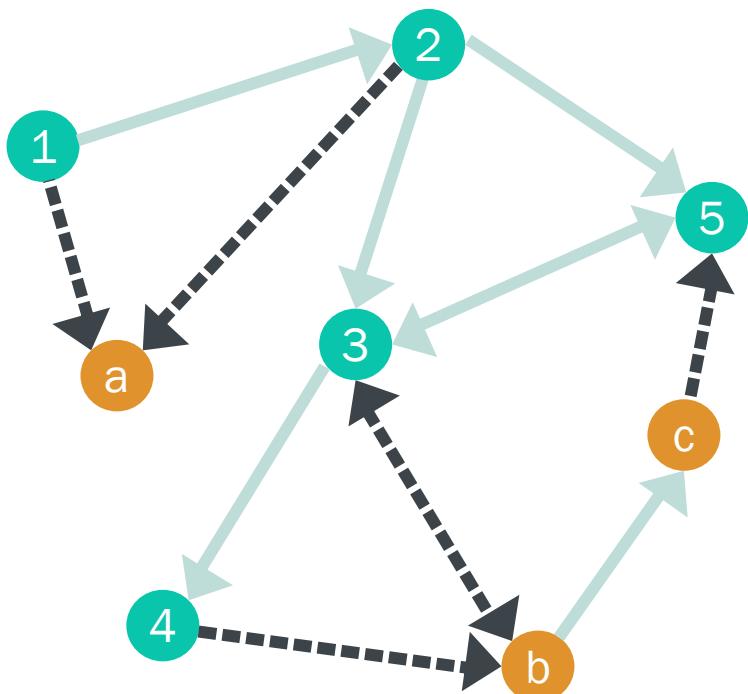


task1(x):
Y → X →⁺¹ X

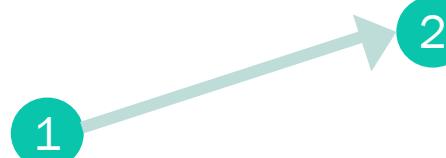
Example



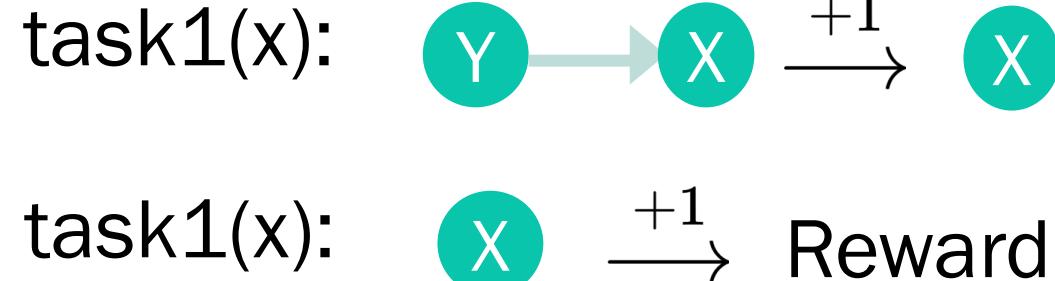
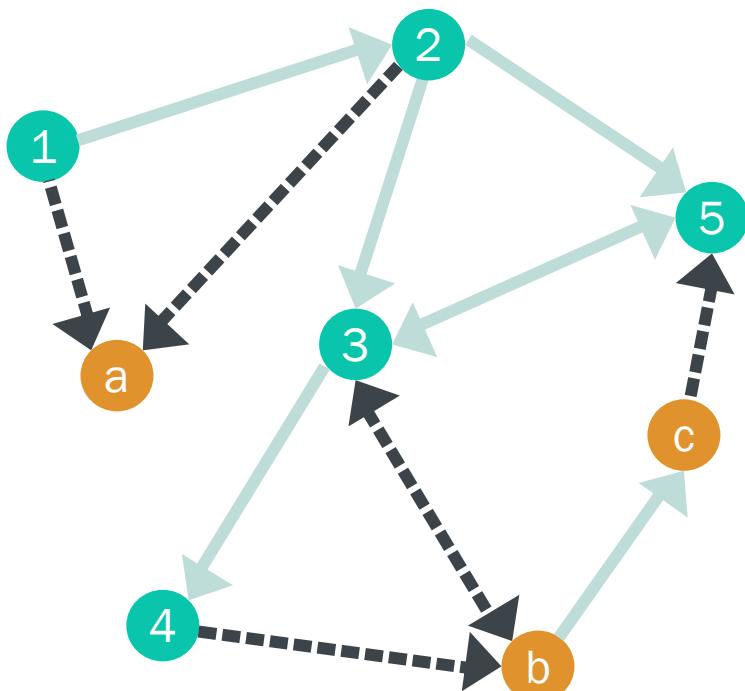
Example



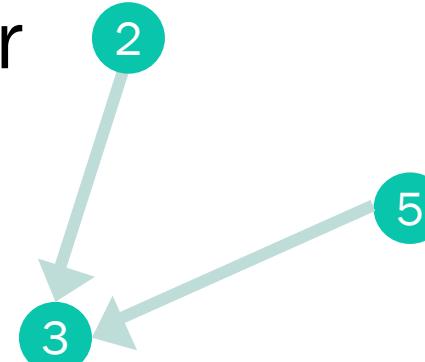
Abstract state for
task1(2):



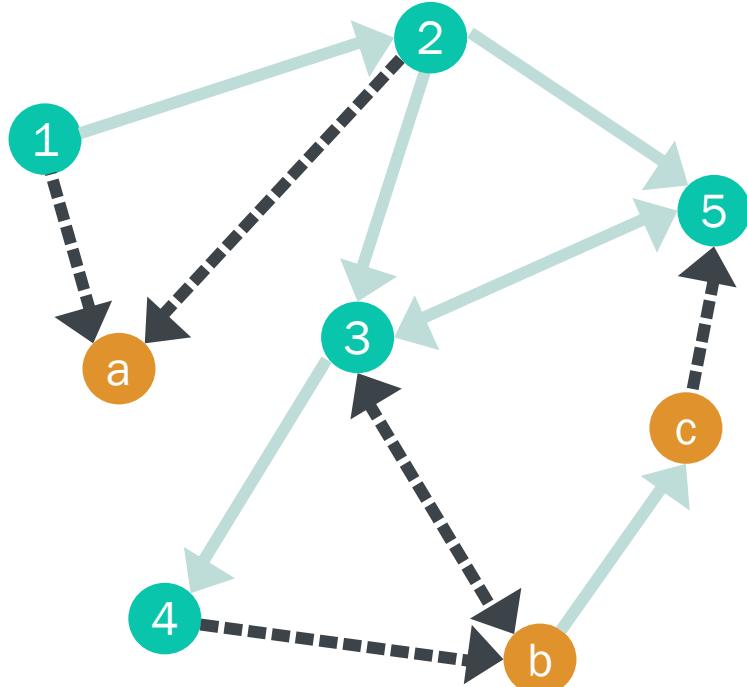
Example



Abstract state for
task1(3):



Example



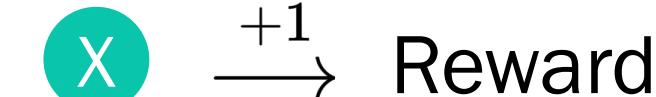
task1(x):



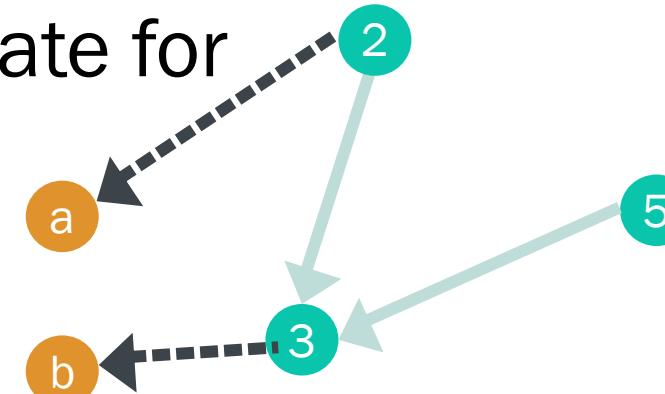
task1(x):



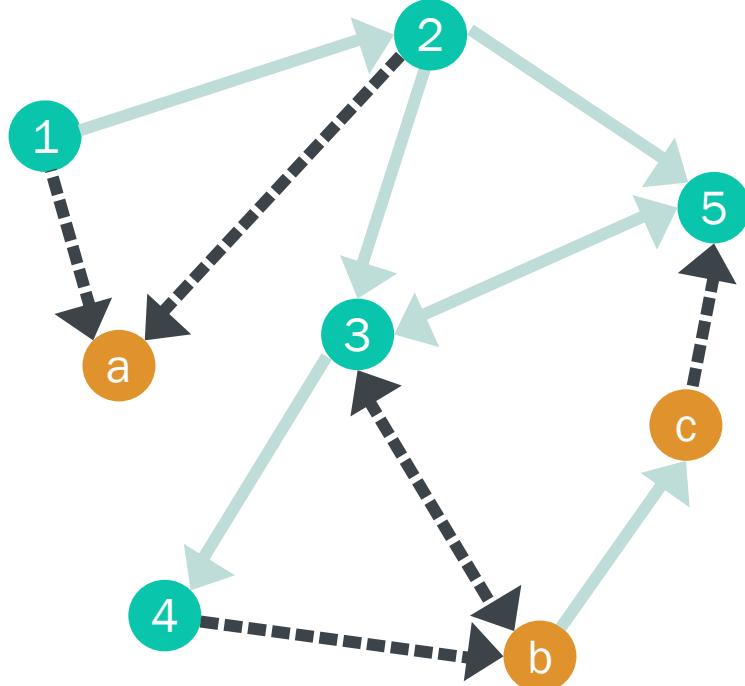
task1(x):



Abstract state for
task1(3):



Example



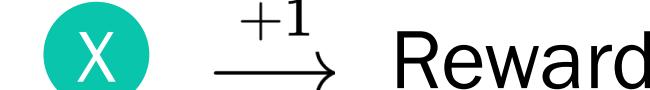
task1(x):



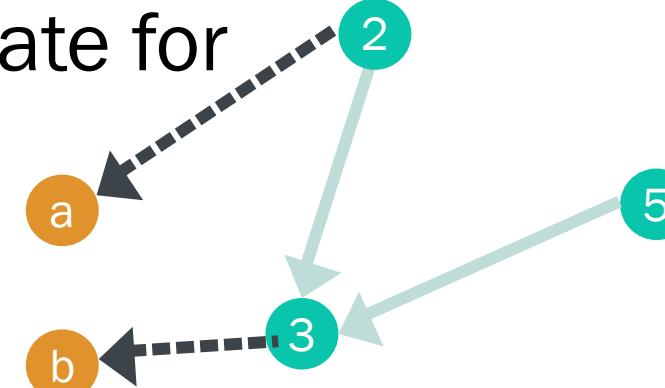
task1(x):



task1(x):



Abstract state for
task1(3):



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*

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

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

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*

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

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



Abstraction

$$\{\text{taxi-at(L1)}, \text{move(Dir)}\} \xrightarrow{+1} \text{taxi-at(L2)}$$
$$\{\text{taxi-at(L1)}, \text{move(Dir)}\} \longrightarrow R$$

pickup(P):

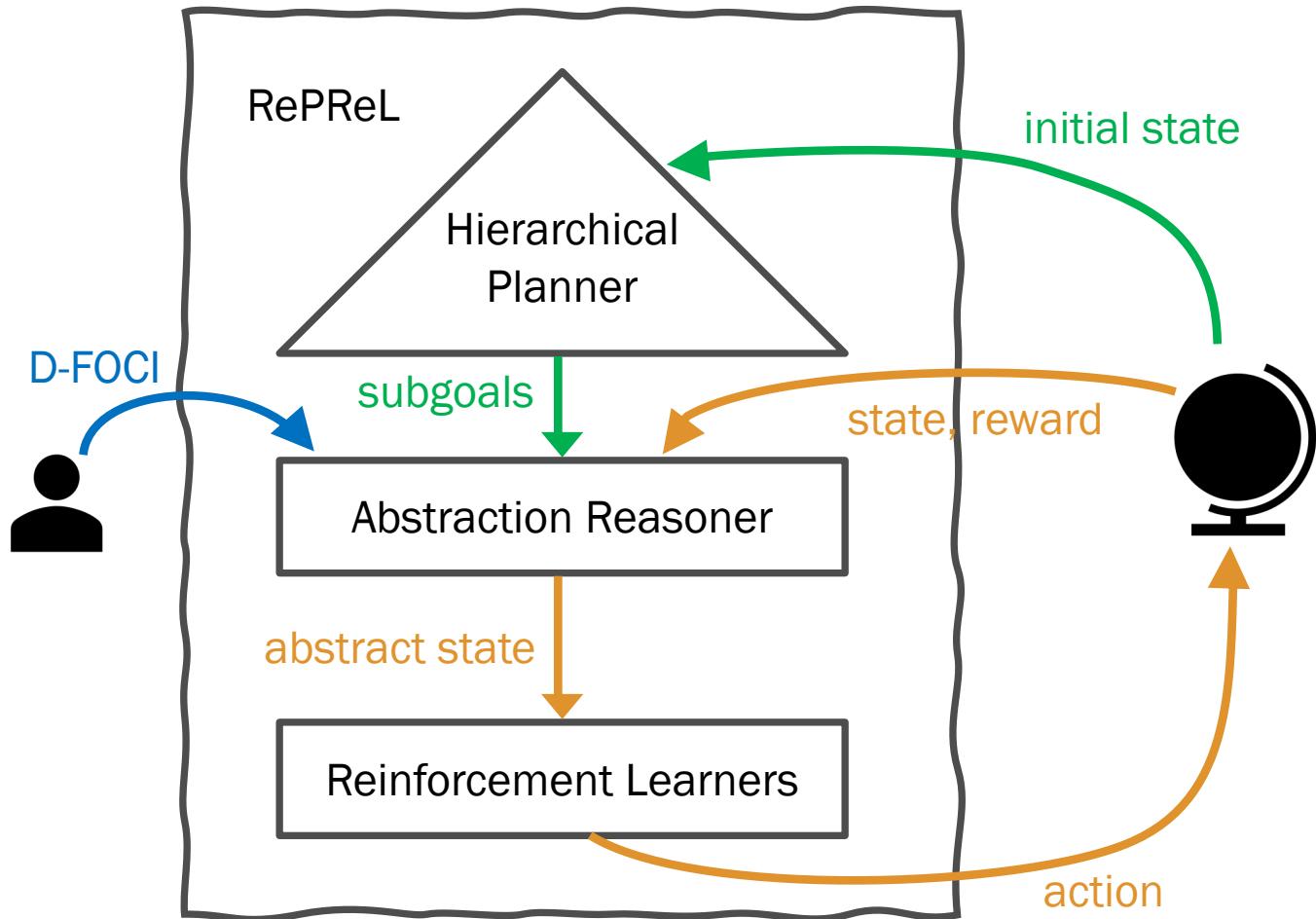
$$\{\text{taxi-at(L1)}, \text{at}(P, L), \text{in-taxi}(P)\} \xrightarrow{+1} \text{in-taxi}(P)$$

pickup(P): $\text{in-taxi}(P) \longrightarrow R_o$

D-FOCI for Taxi pickup task

RePReL Learning

- Get high level plan
- For each subgoal
 - Loop till the state is reached
 - Get the action
 - Get the path
 - Take a step
 - Next state
 - Update the state



RePReL Learning

- Get high level plan
- For each subgoal
 - *Loop till the subgoal is achieved*
 - Get the abstract state
 - Get the policy for that subgoal
 - Take a step and observe reward, next state
 - Update the policy using abstract state

Algorithm 1 RePReL Learning Algorithm

INPUT: $\mathfrak{P}(O, M)$, goal set g , env, t_R, F
OUTPUT: RL policies $\pi_o, \forall o \in O$

```
1:  $\pi_o \leftarrow 0, \forall o \in O$            ▷ initialize RL policy for each operator
2: for each episode do
3:    $s \leftarrow$  get state from env
4:    $\Pi \leftarrow \mathfrak{P}(s, g)$                   ▷ get high-level plan
5:   for  $o_g$  in  $\Pi$  do
6:      $\pi \leftarrow \pi_o$                       ▷ get resp. RL policy
7:      $\hat{s} \leftarrow \text{GetAbstractState}(s, o_g, F)$ 
8:     done  $\leftarrow \hat{s} \in \beta(o_g)$           ▷ check terminal state
9:     while not done do
10:       $a \leftarrow \pi(\hat{s})$                 ▷ get action
11:       $s' \leftarrow \text{env.step}(a)$           ▷ take step in env
12:       $r \leftarrow R(s, a, s')$               ▷ get step reward
13:       $\hat{s}' \leftarrow \text{GetAbstractState}(s, o_g, F)$ 
14:      done  $\leftarrow \hat{s}' \in \beta(o_g)$       ▷ check terminal next state
15:      if done then
16:         $r = r + t_R$                   ▷ add terminal reward
17:      end if
18:       $\pi.\text{update}(\hat{s}, a, \hat{s}', r)$     ▷ update policy
19:       $s, \hat{s} \leftarrow s', \hat{s}'$ 
20:    end while
21:  end for
22: end for
23: return  $\pi_o, \forall o \in O$ 
```

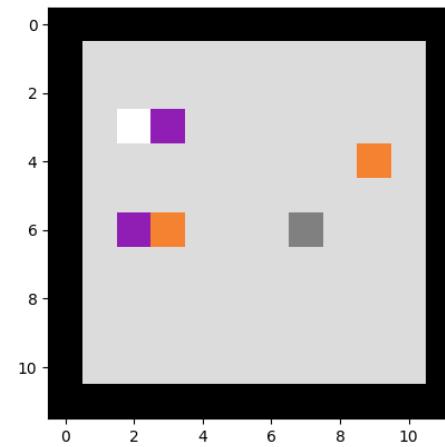
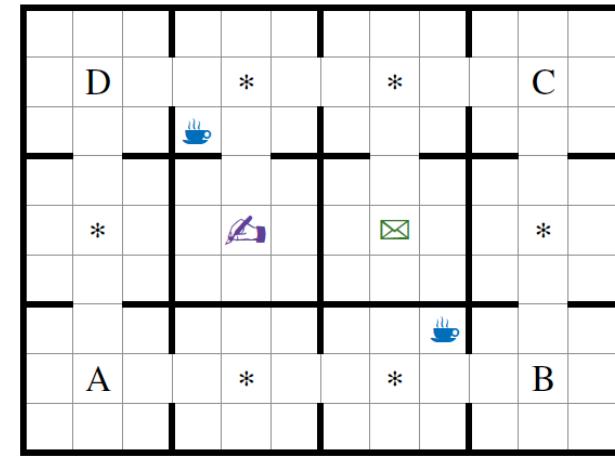
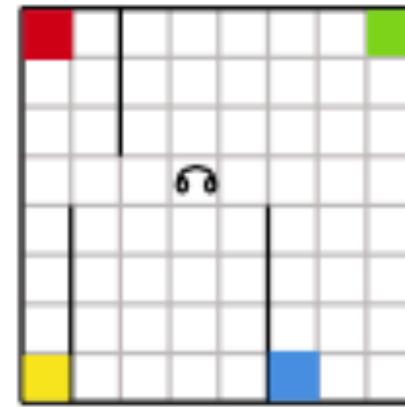
Experiments

- Domains

- *Office World*
- *Craft World*
- *Relational Taxi*
- *Relational Box World*

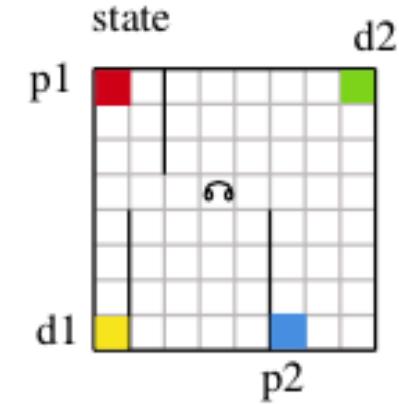
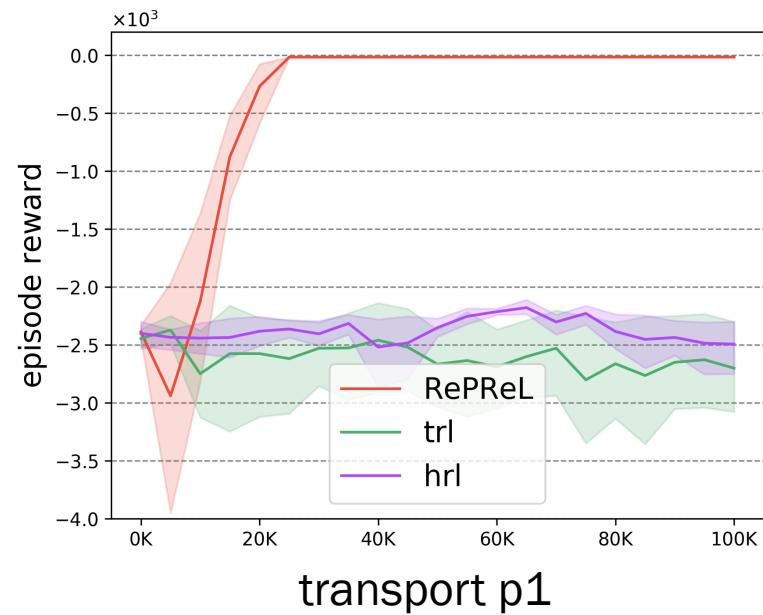
- Baselines

- *HRL (options framework)*
- *TRL (Taskable RL, Illanes et al. 2020)*



Experiments

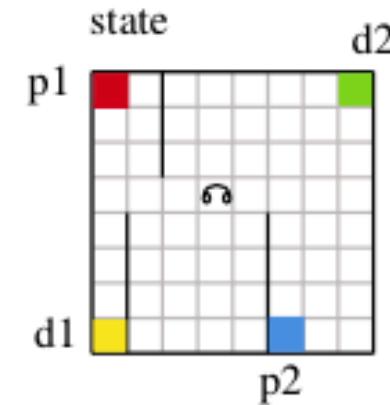
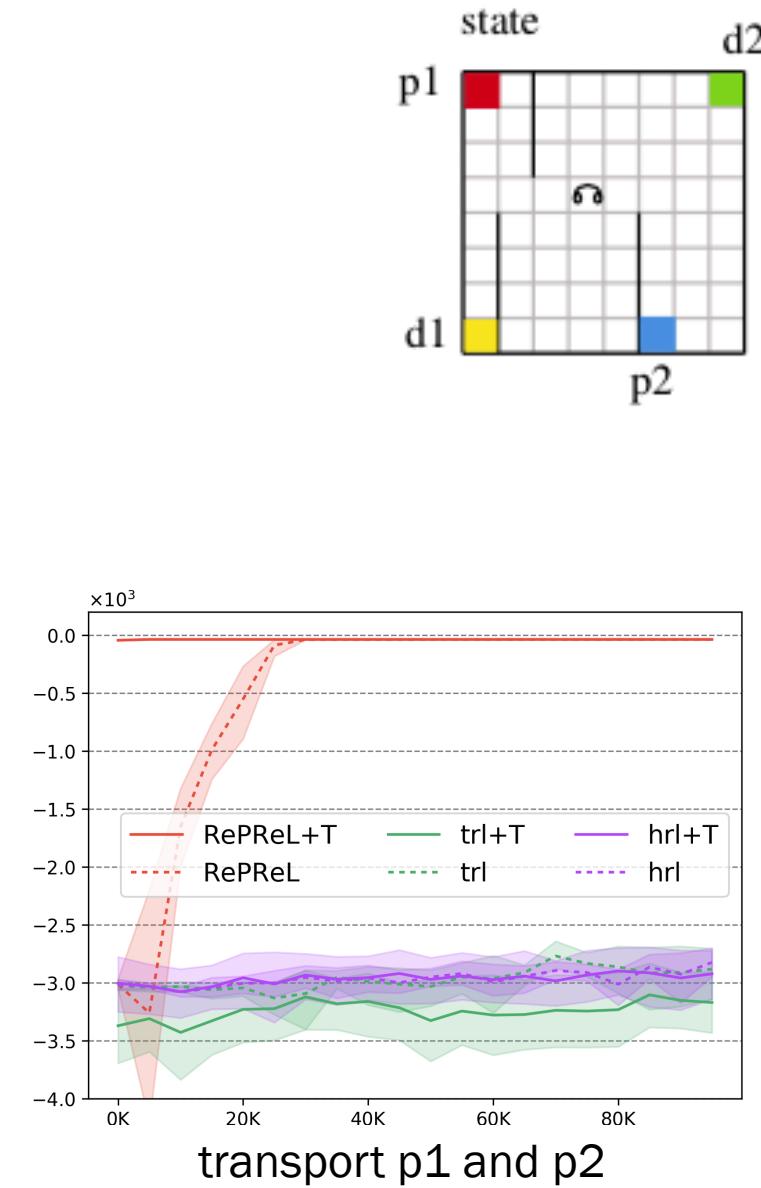
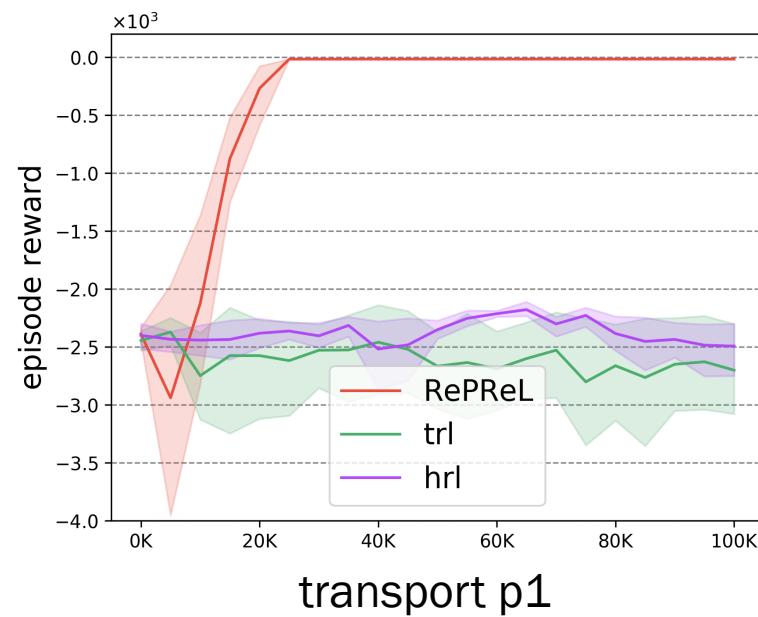
- Sample efficiency
- Transfer across task
- Generalization across objects



trl: Taskable RL (Illanes et al. ICAPS 2020)

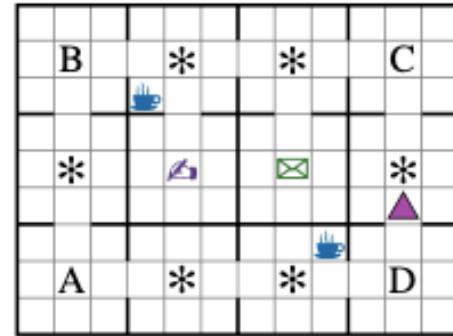
Experiments

- Sample efficiency
- Transfer across task
- Generalization across objects

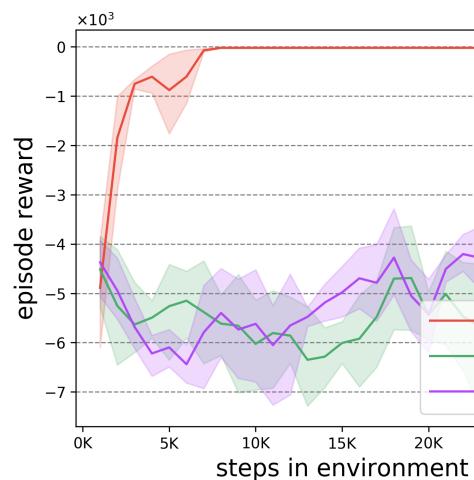


RePReL

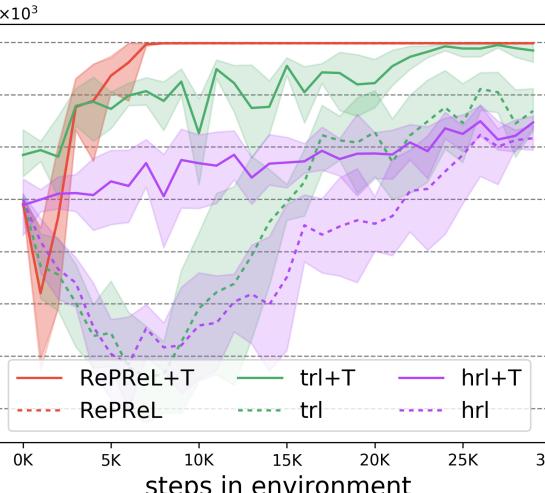
Office World



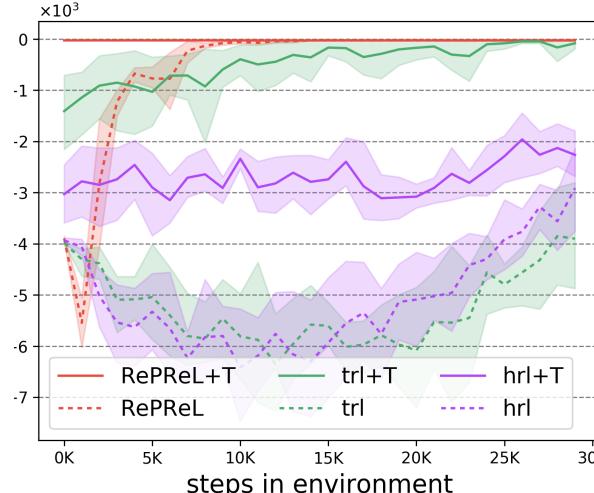
Symbol	Meaning
▲	Agent
*	Furniture
☕	Coffee machine
✉	Mail room
🏢	Office
A, B, C, D Marked locations	



Deliver mail



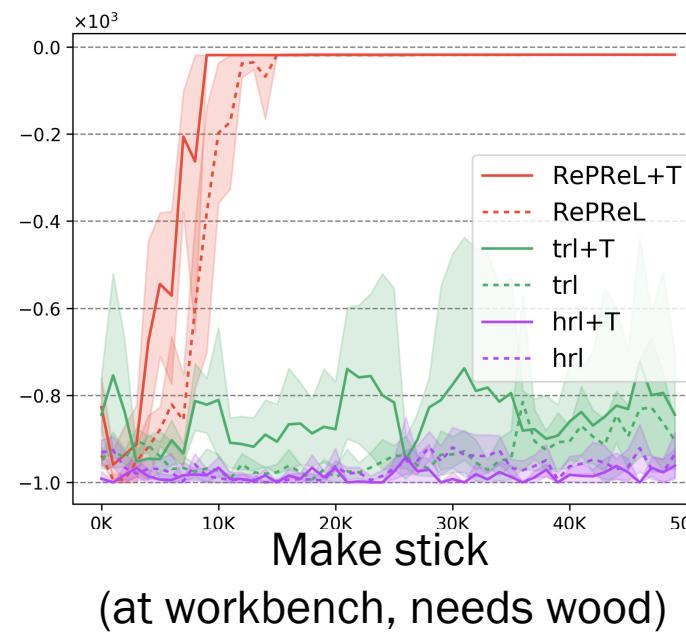
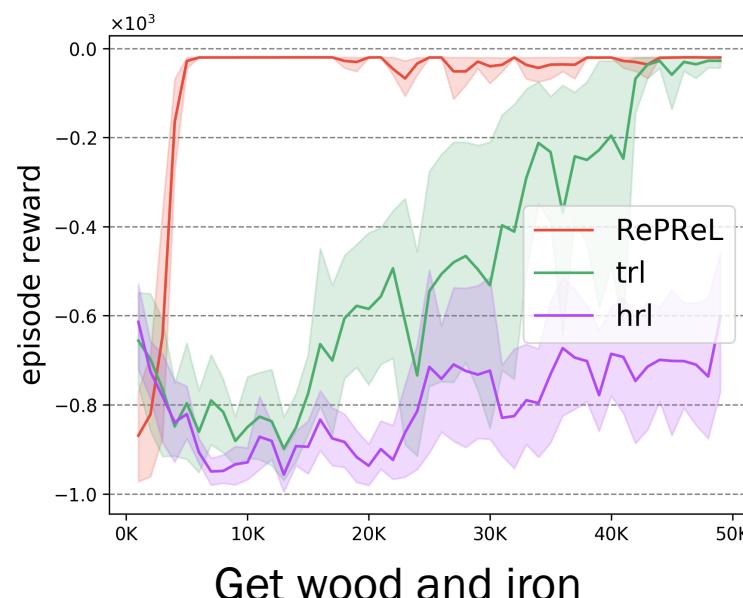
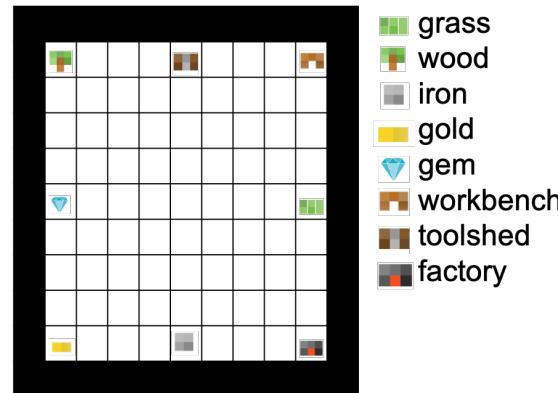
Deliver coffee



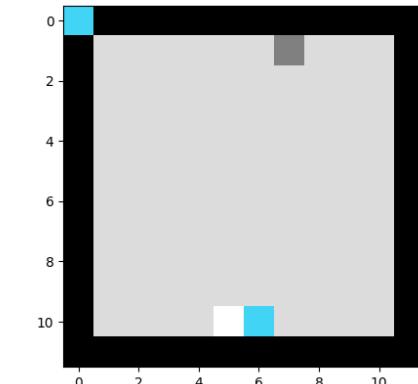
Deliver mail and coffee

RePReLU

CRAFT WORLD

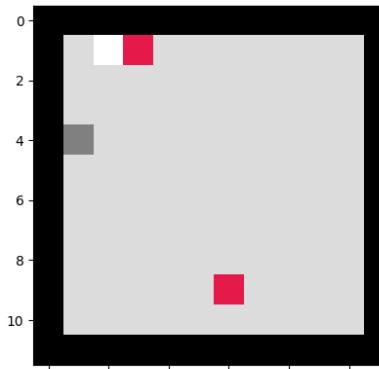


RePReLU



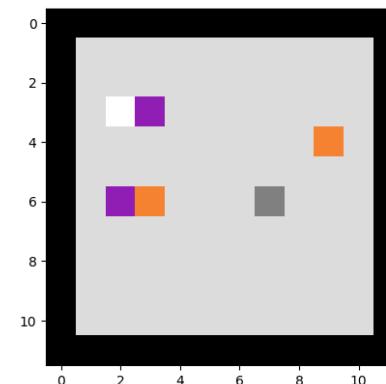
Open lock

Relational Box World

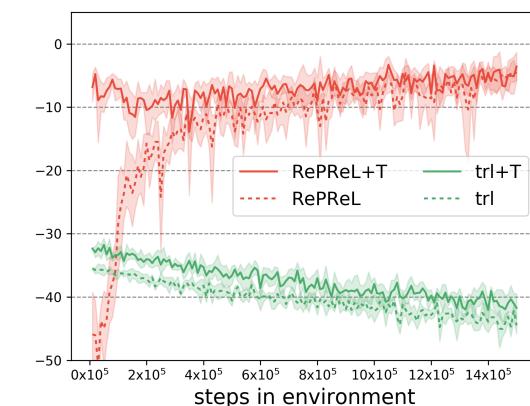
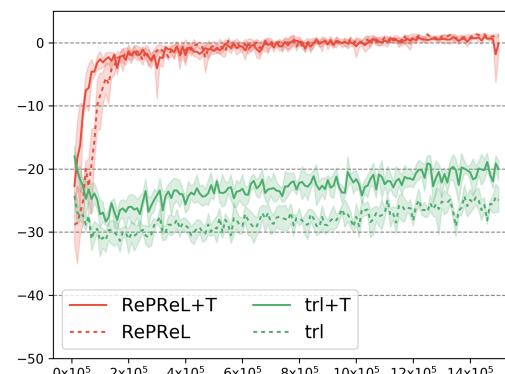
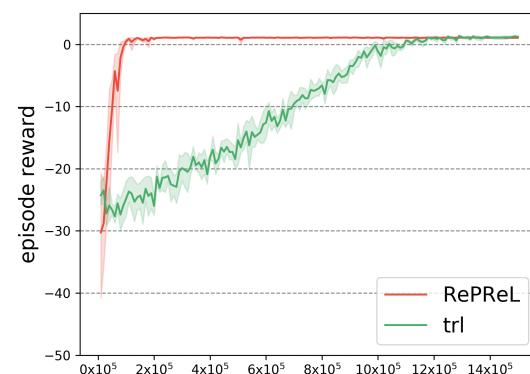


Collect key and open lock

Generalization across objects



Collect key and
open 2 locks



For human-level general intelligence, the ability to detect compositional structure in the domain and form task-specific abstractions are necessary.

Other relevant work

- Learning the high-level planner [*Ludovico et al IJCLR 2021*]
- Modify the plan based on RL agents capability [*Lyu et al AAAI 2019*]
- Automating task termination condition [*Lee et al 2021*]
- Learning task-specific state representation [*Abdulhai et al. 2021*]
- Learning to plan and act simultaneously [*Patra et al AI 2021*]
- Improving Robot Navigation [*Wöhlke et al. ICRA 2021*]
- Extending the RePReL framework to Deep RL setting (under prep)



QUESTIONS?



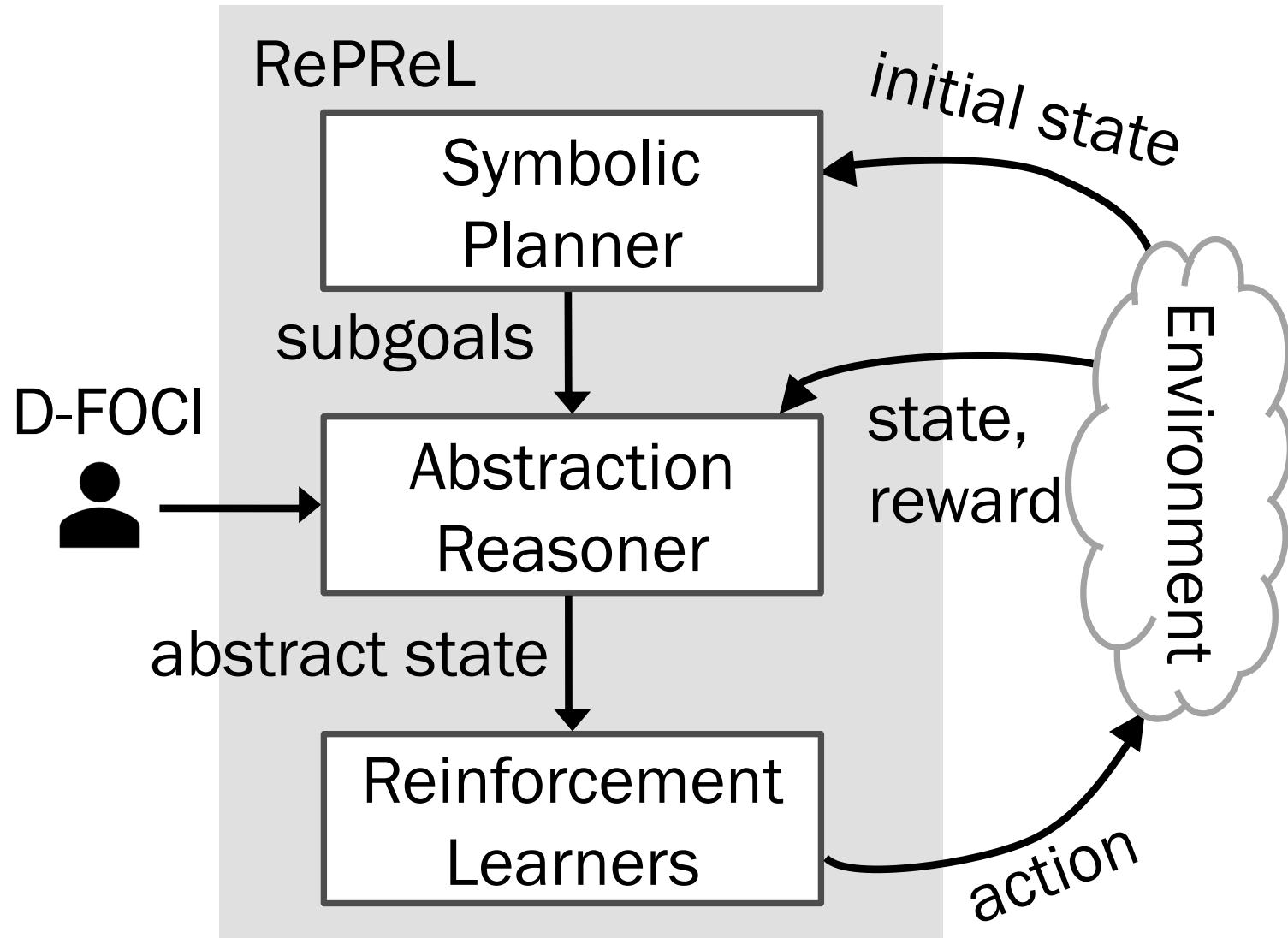
THANKS



starling.utdallas.edu
harshakokel.com

@starling_lab
@harsha_kokel





Abstraction

