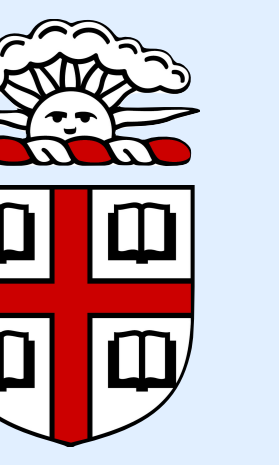


Policy and Value Transfer in Lifelong Reinforcement Learning

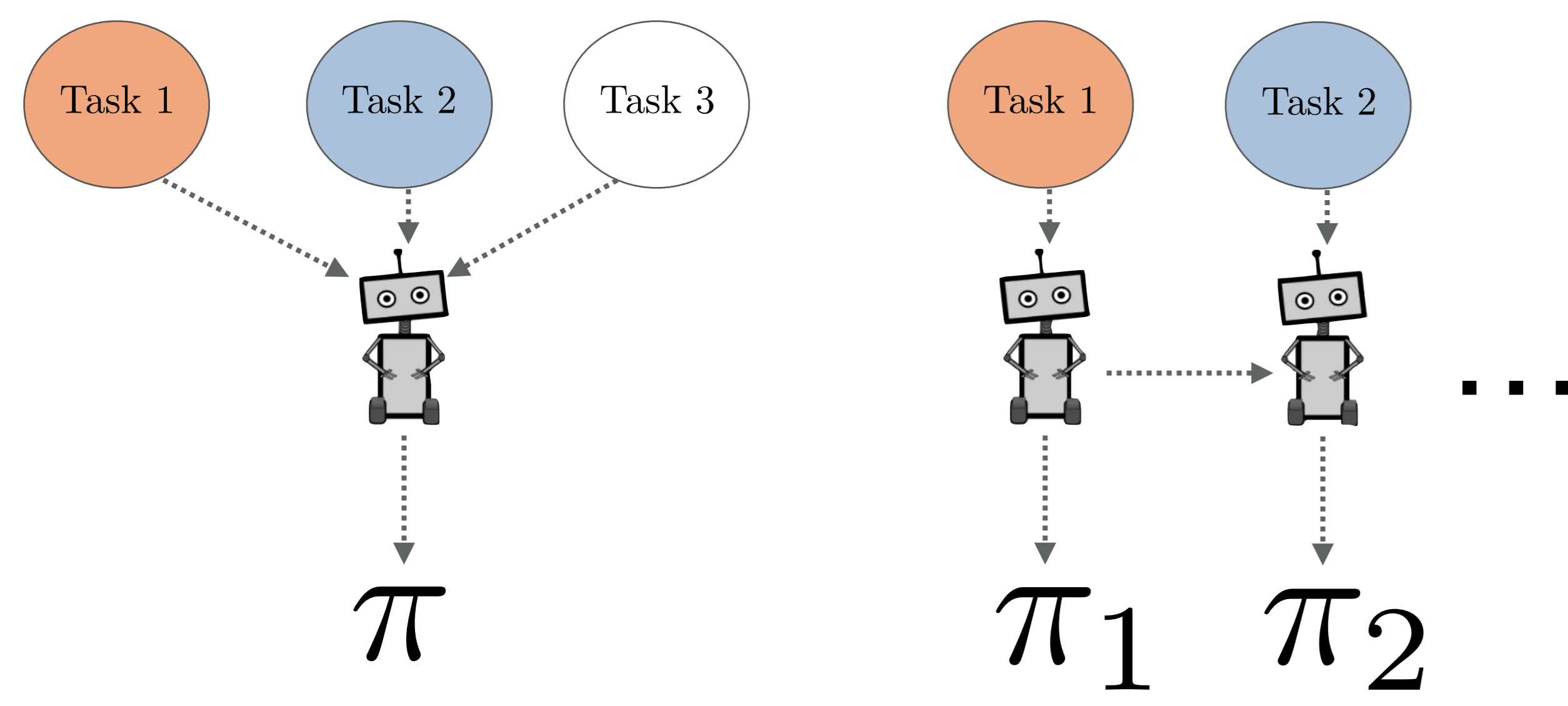


Yuu Jinnai, David Abel, Yue Guo, George Konidaris, Michael L. Littman
{yuu_jinnai, david_abel, yue_guo, george_konidaris, michael_littman}@brown.edu

Department of Computer Science, Brown University

Goal

Understand knowledge transfer in lifelong RL.

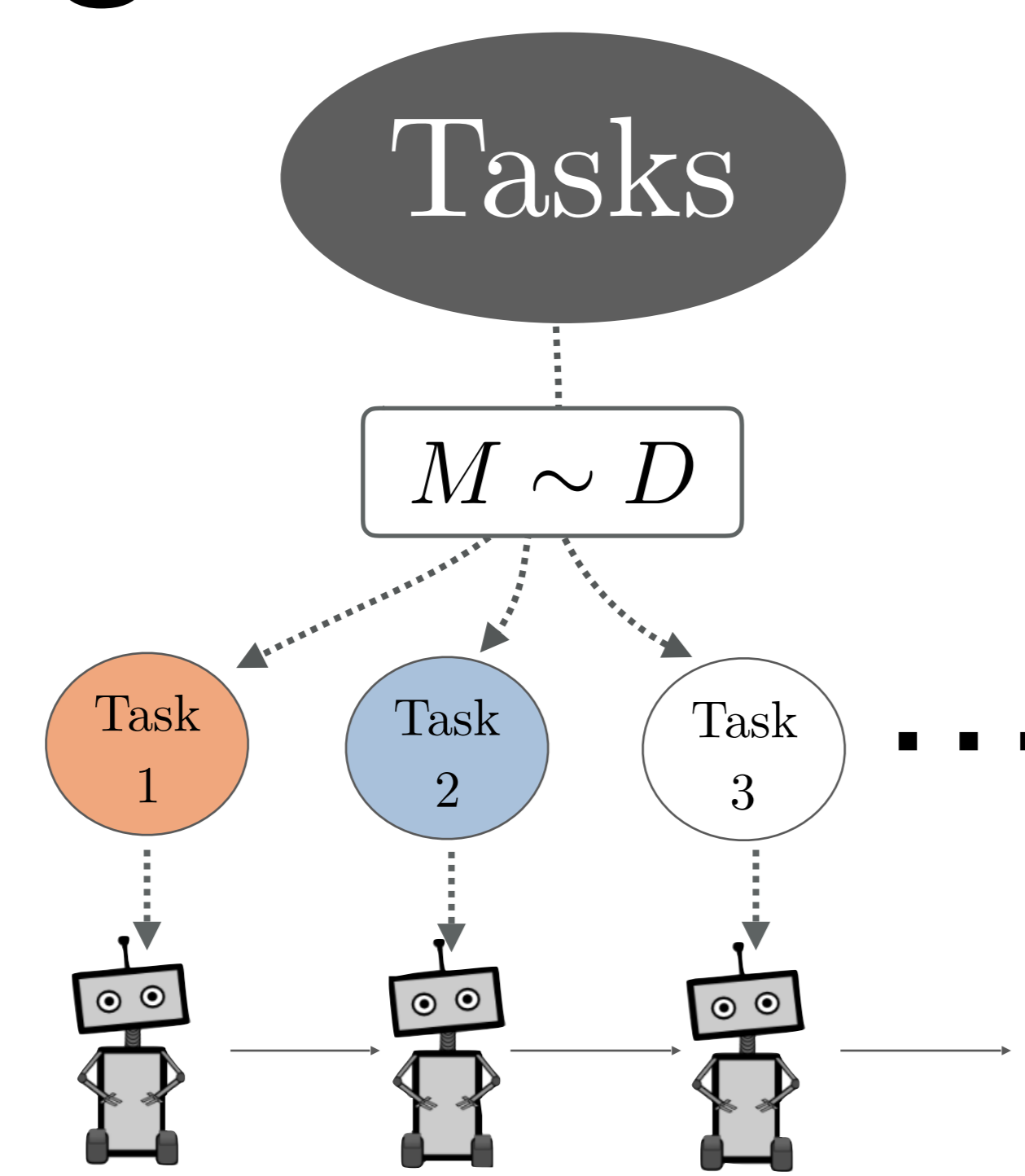


1. Optimal Fixed Policy 2. Informed Transfer for RL

Lifelong Reinforcement Learning [1]

repeat:

1. Sample a task from a distribution: $M \sim D$.
2. Solve the sampled task M .



Objective 1: Find the policy that maximizes the expected performance over the distribution:

$$\arg \max_{\pi \in \Pi} \mathbb{E}_{M \sim D} [V_M^\pi(s_0)]$$

Performance of π on task M .

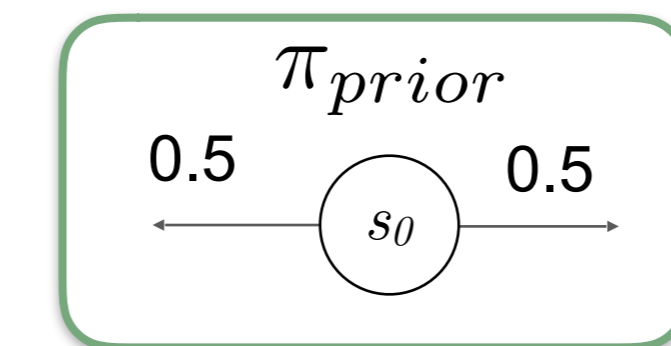
Jumpstart Policies

Π	$\mathcal{R} \sim D$	$G \sim D$
(Deterministic) $\Pi_d : \mathcal{S} \mapsto \mathcal{A}$	[2] Avg. MDP	Avg. V [3]
(Stochastic) $\Pi_s : \mathcal{S} \mapsto \Pr(\mathcal{A})$	Avg. MDP	Avg. V
(Any policy) $\Pi_b : \mathcal{S} \times \Pr(\mathcal{M}) \mapsto \mathcal{A}$	Belief MDP	Belief MDP [4]

Action Prior [5]

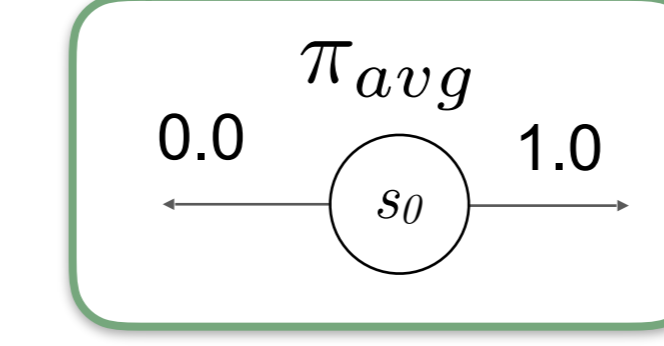
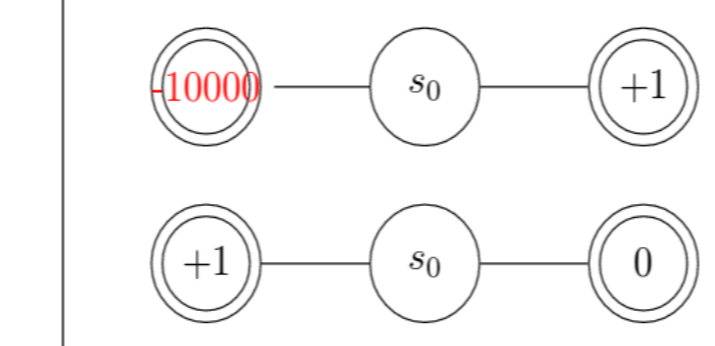
Probability the action is optimal:

$$\pi_{prior}(a | s) := \Pr_{M \sim D} (a = \arg \max_{a'} Q_M^*(s, a'))$$



Average MDP

Choose optimal action in averaged MDP:



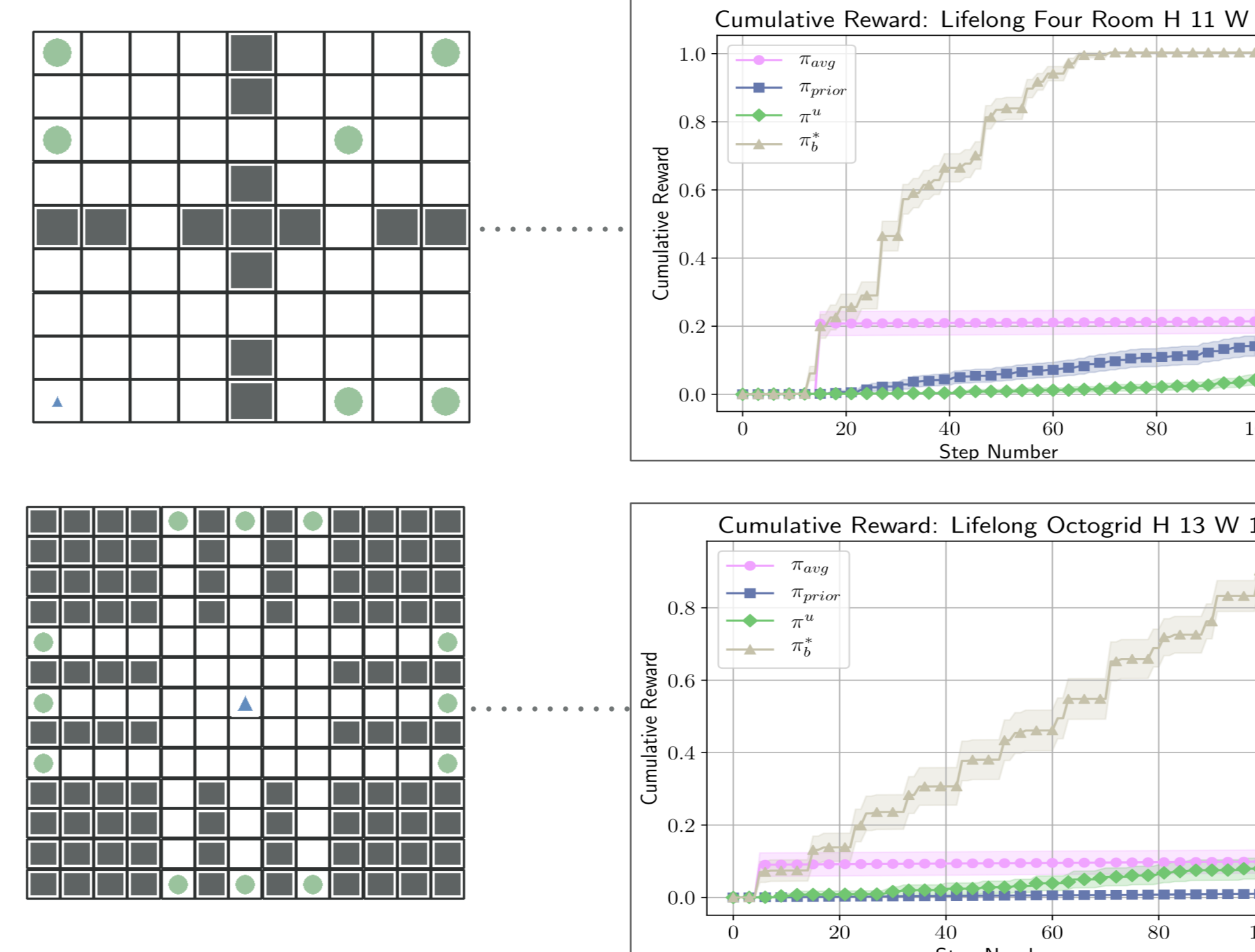
Theorem 1. For a distribution of MDPs with $R \sim D$, Average MDP is an optimal fixed policy. [2]

Theorem 2. For a distribution of MDPs with $R \sim D$, the performance of Average MDP has a lower bound:

$$\mathbb{E}_{M \sim D} [V_M^{\pi_{avg}}(s)] \geq \max_{M \in \mathcal{M}} \Pr(M) V_M^*(s)$$

Jumpstart Experiments

- π_{avg} (pink circle)
- π_{prior} (blue square)
- π^u (green diamond)
- π_b^* (brown triangle)



Informed Transfer for RL

PAC-MDP [6]

Sample complexity of many PAC-MDP algorithms depends on the overestimation of the initial value function:

$$\tilde{O} \left(\sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \max \{Q_{init}(s,a) - V^*(s), 0\} \right)$$

(Objective) Minimize the **overestimate**

(Constraint) Do not **underestimate**: $Q_{init}(s,a) \geq Q^*(s,a)$

(Our Solution) Initialize with: $Q_{max}(s,a) := \max_{M \in \mathcal{M}} Q_M^*(s,a)$

Proposed Algorithm: MaxQinit

Approximate $Q_{max}(s,a)$ from n sampled MDPs:

$$\text{MAXQINIT}(s,a) := \max_{M \in \{M_1, \dots, M_n\}} \hat{Q}_M(s,a)$$

Theorem 3. For a given δ , after $n \geq \frac{\ln(\delta)}{\ln(1-p_{min})}$ sampled MDPs, MaxQinit will retain optimism with probability $1 - \delta$:

$$\forall_{s,a} : \text{MAXQINIT}(s,a) \geq \max_{M \in \mathcal{M}} Q_M^*(s,a)$$

Transfer Experiments

