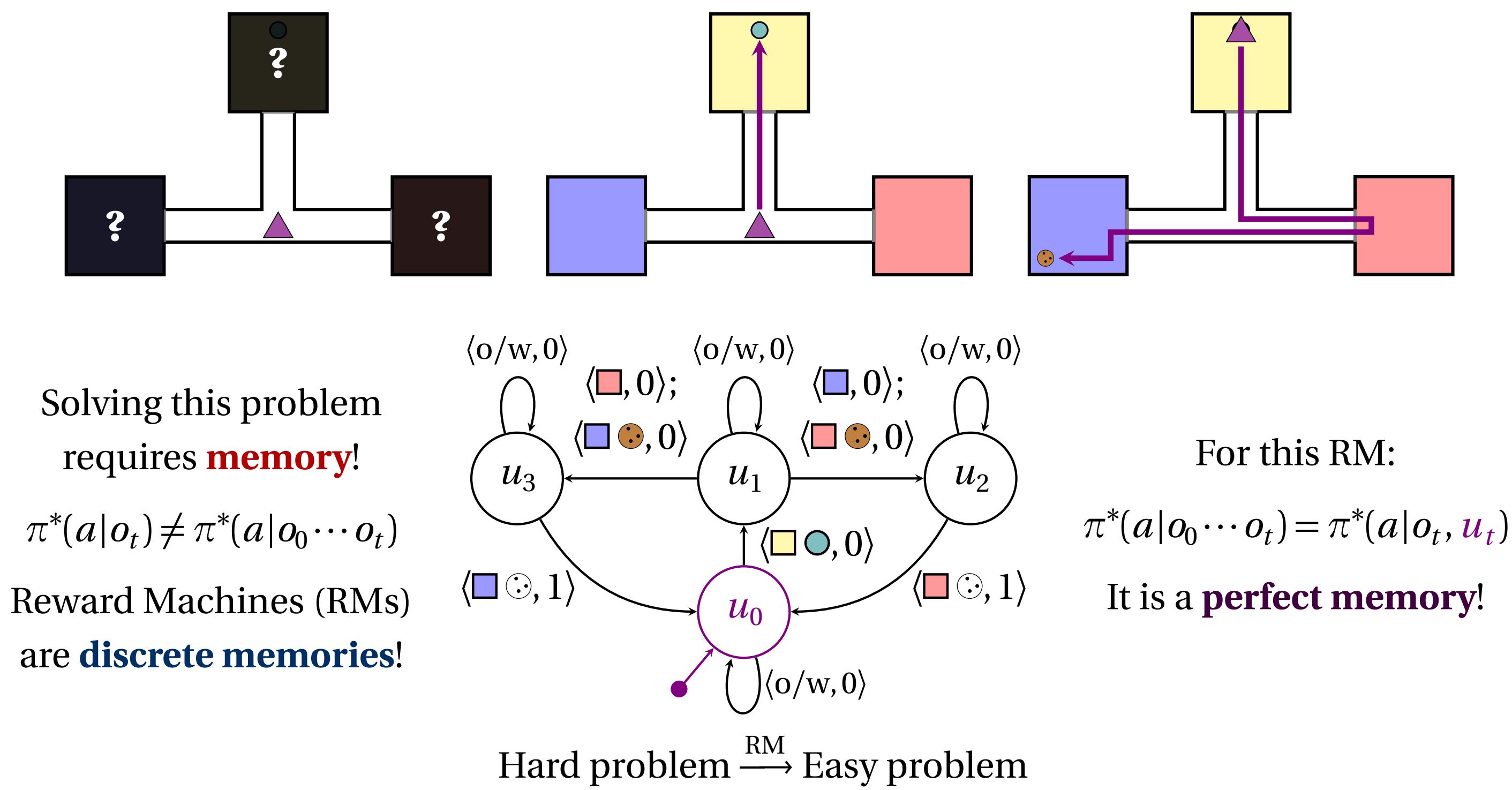


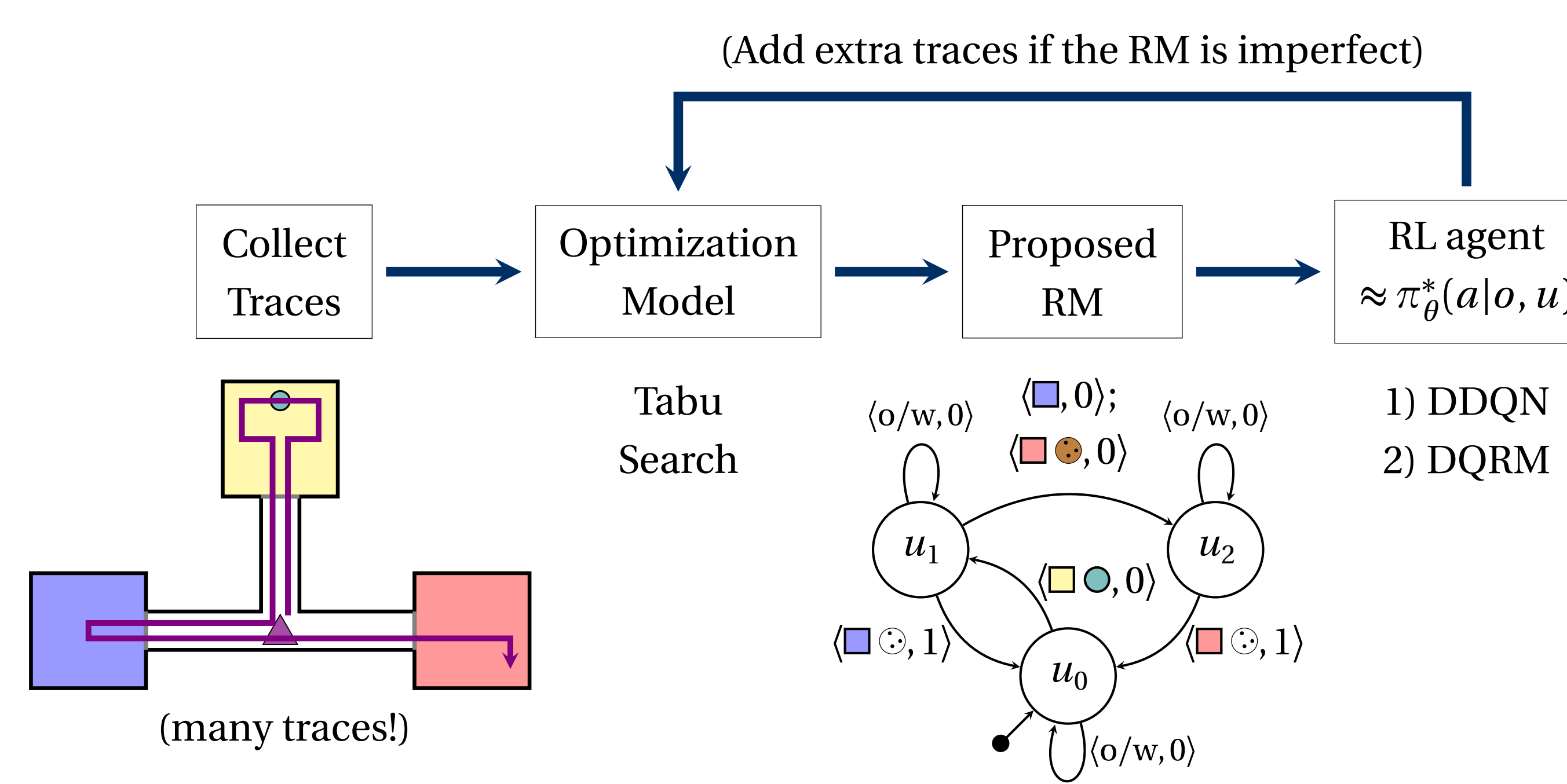


Abstract. Reward Machines (RMs) provide a structured, automata-based representation of a reward function that enables a Reinforcement Learning (RL) agent to decompose an RL problem into structured subproblems that can be efficiently learned via off-policy learning. Here we show that RMs can be learned from experience, instead of being specified by the user, and that the resulting problem decomposition can be used to effectively solve partially observable RL problems. We pose the task of learning RMs as a discrete optimization problem where the objective is to find an RM that decomposes the problem into a set of subproblems such that the combination of their optimal memoryless policies is an optimal policy for the original problem. We show the effectiveness of this approach on three partially observable domains, where it significantly outperforms A3C, PPO, and ACER, and discuss its advantages, limitations, and broader potential.

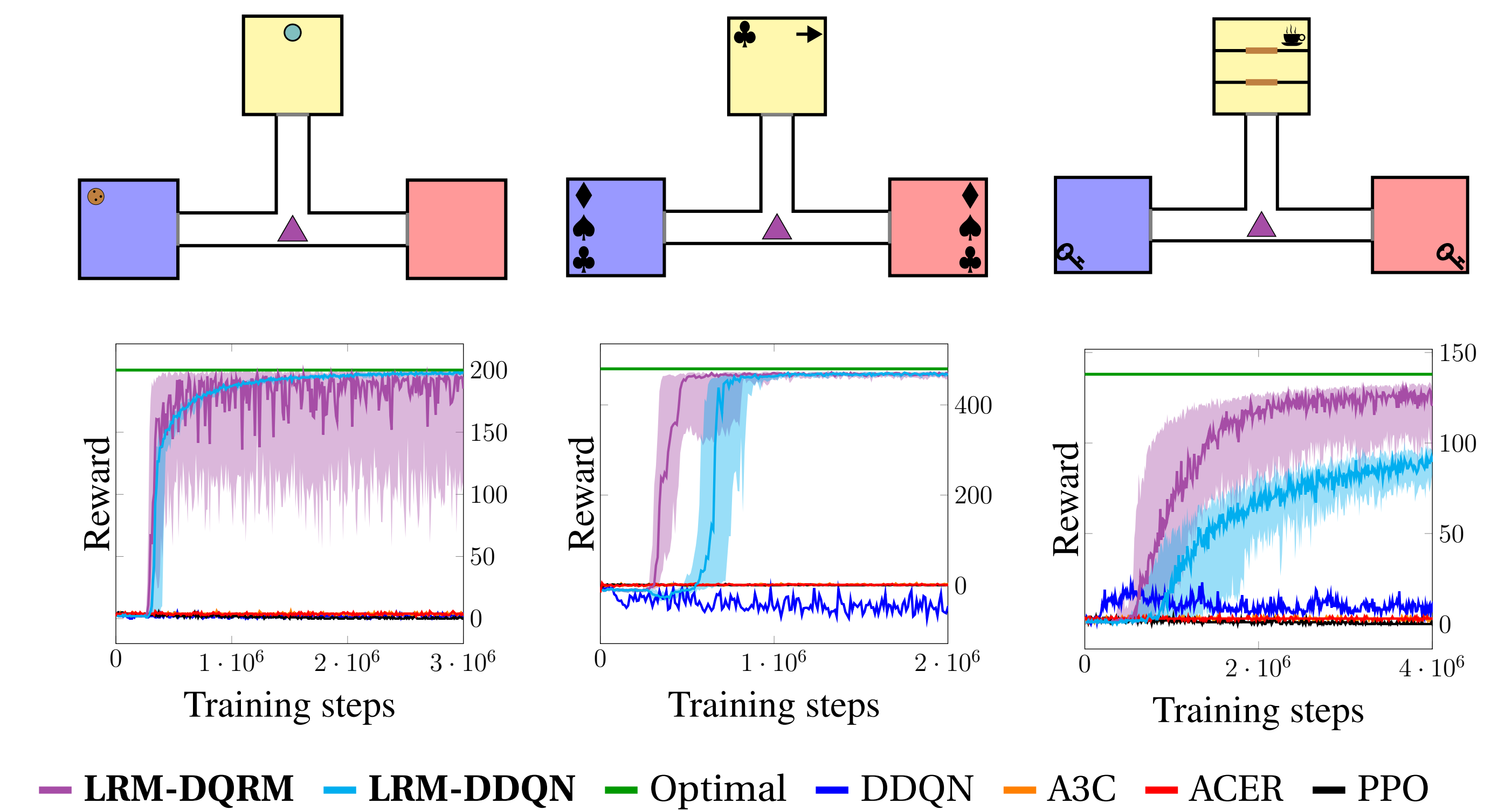
Partially Observable RL



Our Approach (LRM)



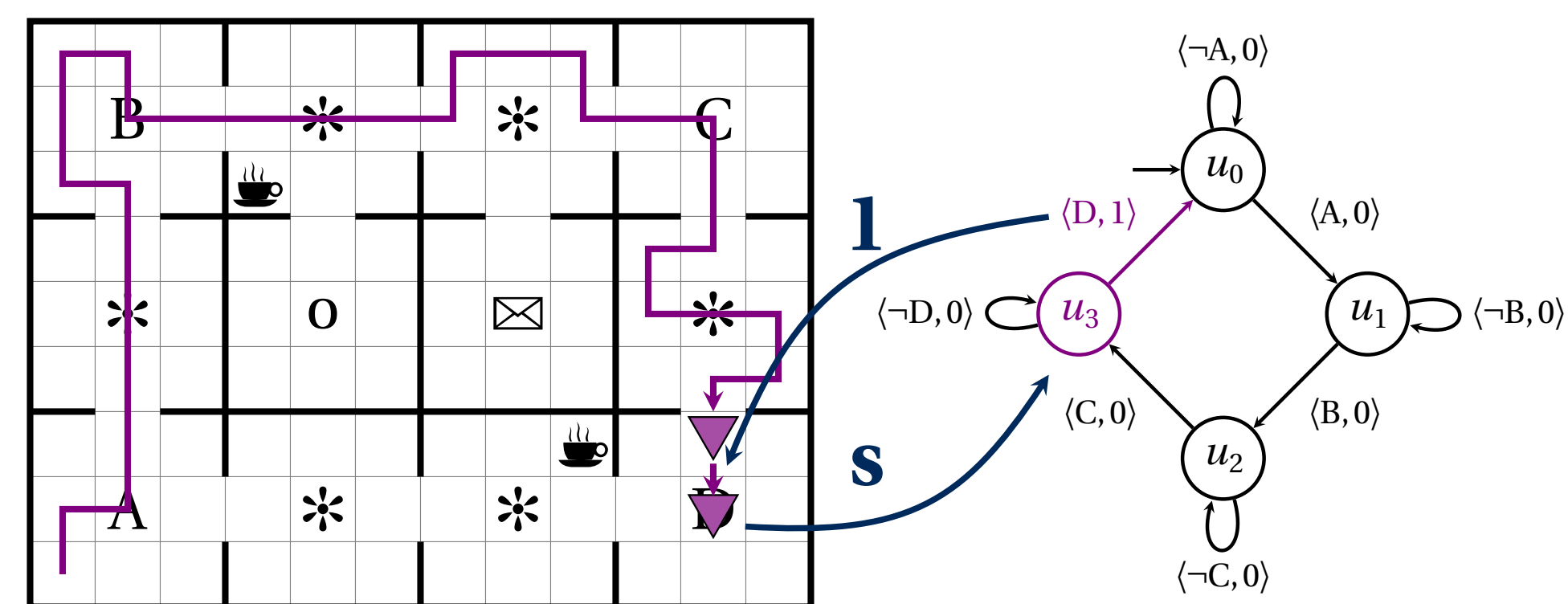
Results



Summary ↑ / ↓ Details

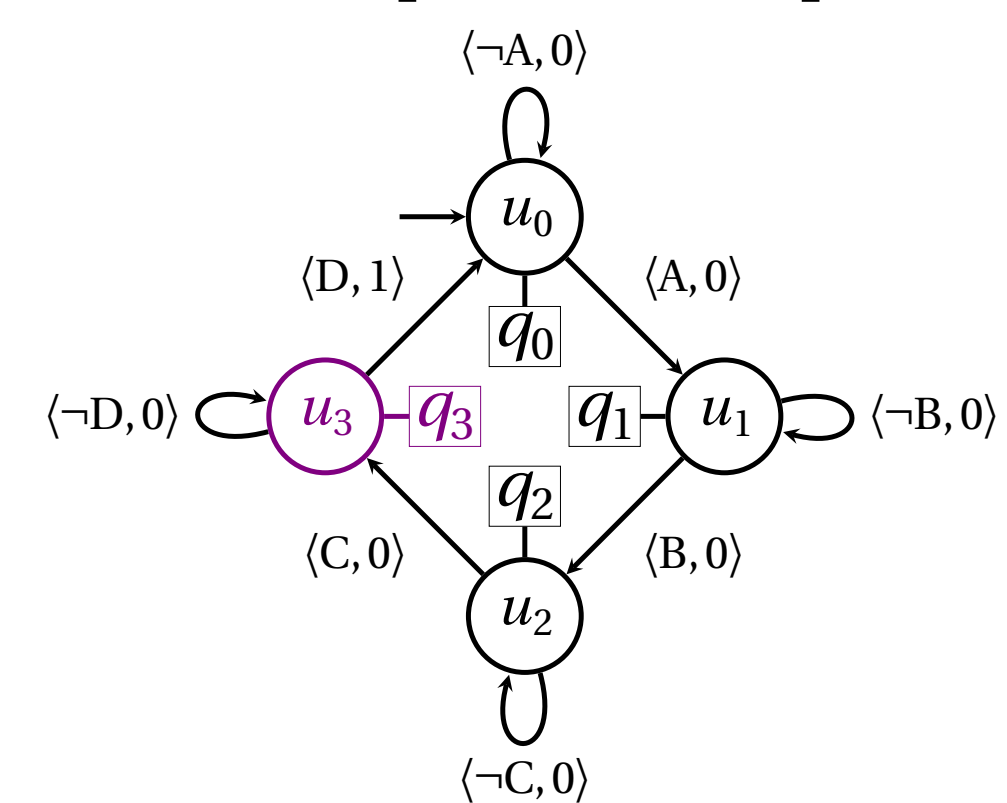
Reward Machines (RMs)

RMs are automata-based representations of a reward function that allow RL agents to learn policies faster [1].



Q-Learning for Reward Machines (QRM)

- Learn one policy (q-function) per state in the RM.
- Select actions using the policy of the current RM state.
- Reuse experience to update all the q-values at the same time.



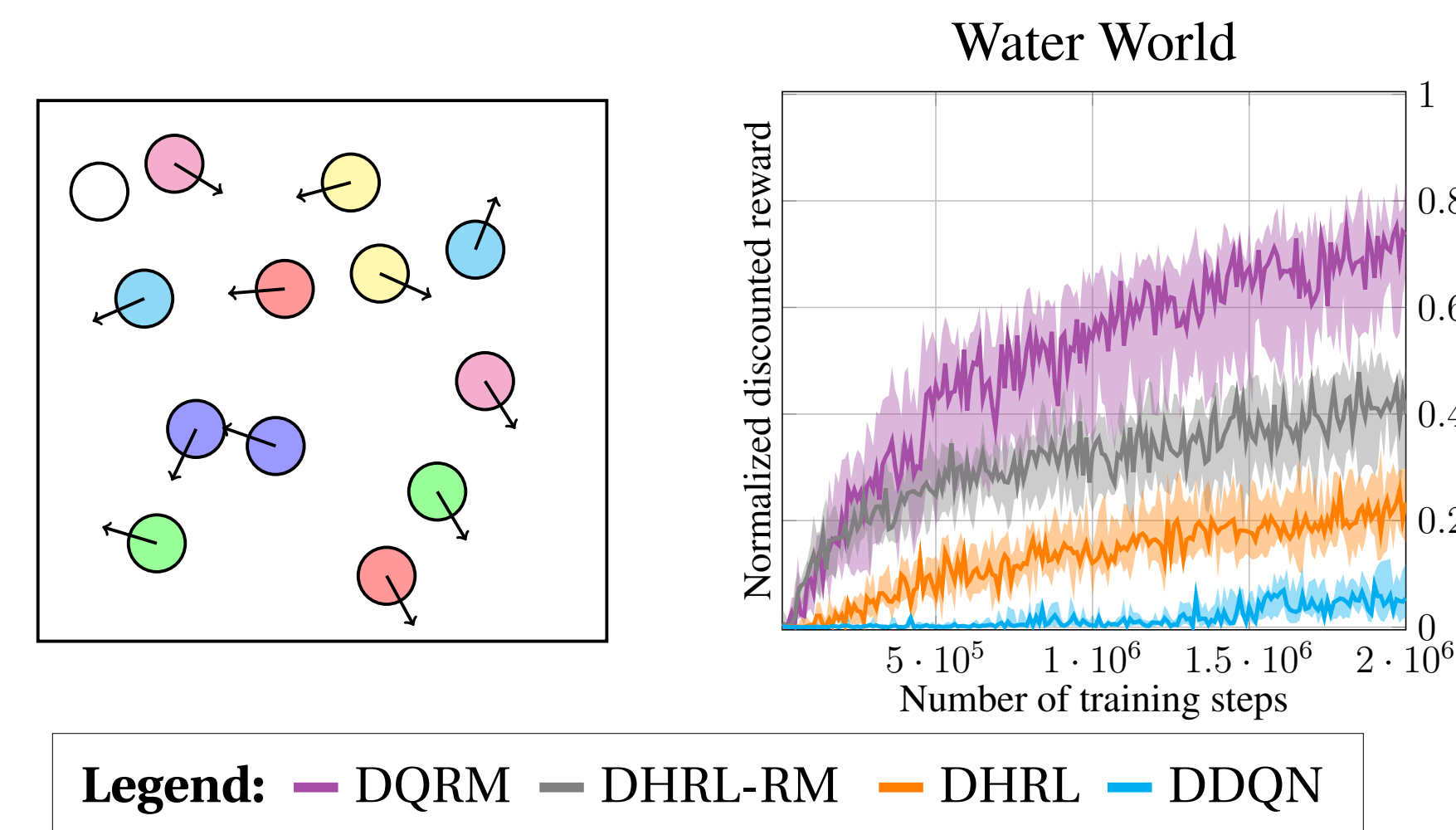
Q-updates

$$\begin{aligned} q_0(s, a) &\leftarrow \frac{\alpha}{\gamma} 0 + \gamma \max_{a'} q_0(s', a') \\ q_1(s, a) &\leftarrow \frac{\alpha}{\gamma} 0 + \gamma \max_{a'} q_1(s', a') \\ q_2(s, a) &\leftarrow \frac{\alpha}{\gamma} 0 + \gamma \max_{a'} q_2(s', a') \\ q_3(s, a) &\leftarrow \frac{\alpha}{\gamma} 1 + \gamma \max_{a'} q_0(s', a') \end{aligned}$$

RMs in Previous Work

- RMs were specified by the user (directly [1] or via LTL [2]).
- QRM did not work under partial observability.

Our 2018 RM Results



Main Contributions

- First approach for learning RMs from experience.
- Extended RMs and QRM to work under partial observability.
- Developed a theory for 1 and 2.

RMs under Partial Observability

RMs are defined over a set of propositional symbols \mathcal{P} that correspond to high-level events (e.g., $\mathcal{P} = \{\square, \square, \square, \square, \square, \square\}$) that the agent can detect using a labelling function $L: O \times A \times O \rightarrow 2^{\mathcal{P}}$.

Def. A Reward Machine is a tuple $R_{\mathcal{P}} = \langle U, u_0, \delta_u, \delta_r \rangle$ where U is a finite set of states, $u_0 \in U$ is an initial state, δ_u is the state-transition function, $\delta_u: U \times 2^{\mathcal{P}} \rightarrow U$, and δ_r is the reward-transition function, $\delta_r: U \times 2^{\mathcal{P}} \rightarrow \mathbb{R}$.

Perfect Reward Machines

Perfect RMs make the environment Markovian w.r.t. $O \times U$, i.e.:

$$\Pr(o_{t+1}, r_{t+1} | o_0, a_0, \dots, o_t, a_t) = \Pr(o_{t+1}, r_{t+1} | o_t, u_t, a_t)$$

for every possible trace $o_0, a_0, \dots, o_t, a_t$ generated by any policy.

Properties: Given an environment \mathcal{E} : (i) If the set of belief states of \mathcal{E} is finite, then there exists a perfect RM for \mathcal{E} , and (ii) Optimal policies over $O \times U$ for perfect RMs are also optimal for \mathcal{E} .

Optimization Model

Inputs:

- A set of traces: $\mathcal{T} = \{\mathcal{T}_0, \dots, \mathcal{T}_n\}$ where each trace \mathcal{T}_i is:

$$\mathcal{T}_i = (o_{i,0}, a_{i,0}, r_{i,0}, \dots, a_{i,t-1}, r_{i,t-1}, o_{i,t}).$$

- A labelling function $L: O \times A \times O \rightarrow 2^{\mathcal{P}}$.

Model:

$$\text{minimize}_{\langle U, u_0, \delta_u, \delta_r \rangle} \sum_{i \in I} \sum_{t \in \mathcal{T}_i} \log(|N_{x_{i,t}, L(e_{i,t})}|) \quad (\text{LRM})$$

$$s.t. \langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{R}_{\mathcal{P}} \quad (1)$$

$$|U| \leq u_{\max} \quad (2)$$

$$x_{i,t} \in U \quad \forall i \in I, t \in \mathcal{T}_i \cup \{t_i\} \quad (3)$$

$$x_{i,0} = u_0 \quad \forall i \in I \quad (4)$$

$$x_{i,t+1} = \delta_u(x_{i,t}, L(e_{i,t})) \quad \forall i \in I, t \in \mathcal{T}_i \quad (5)$$

$$N_{u,l} \subseteq 2^{\mathcal{P}} \quad \forall u \in U, l \in 2^{\mathcal{P}} \quad (6)$$

$$L(e_{i,t+1}) \in N_{x_{i,t}, L(e_{i,t})} \quad \forall i \in I, t \in \mathcal{T}_i \quad (7)$$

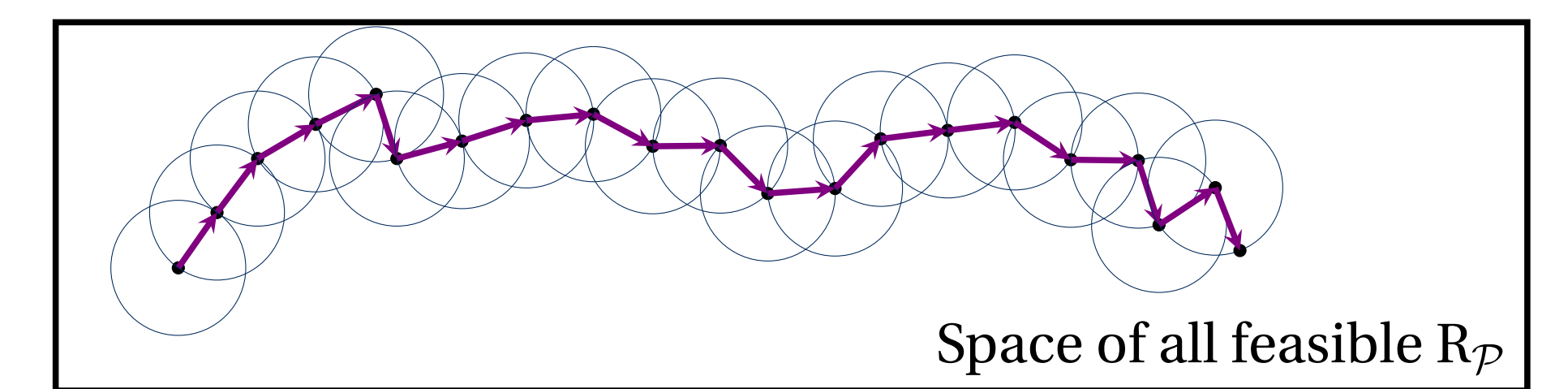
Outputs:

- A Reward Machine: $\langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{R}_{\mathcal{P}}$.
- A set of possible next high-level observations $N_{x_{i,t}, L(e_{i,t})}$.

Theorem (necessary conditions):

When $\mathcal{T} \rightarrow \infty$, any perfect RM is an optimal solution of LRM.

Tabu Search



Experimental Remarks

- The binary classifiers were used by all the approaches.
- LRM works because Tabu search finds high-quality RMs.
- DQRM pays off in domains with sparse rewards.
- Code:** bitbucket.org/RToroIcarte/lrm.

References

- R. Toro Icarte, T. Klassen, R. Valenzano, and S. McIlraith. Using reward machines for high-level task specification and decomposition in reinforcement learning. In *ICML18*.
- A. Camacho, R. Toro Icarte, T. Klassen, R. Valenzano, and S. McIlraith. LTL and beyond: Formal languages for reward function specification in reinforcement learning. In *IJCAI19*.
- R. Toro Icarte, E. Waldie, T. Klassen, R. Valenzano, M. Castro, and S. McIlraith. Learning reward machines for partially observable reinforcement learning. In *NeurIPS19*.