

Elective AI&R : Reasoning Agents

A.A. 2019/20



SAPIENZA
UNIVERSITÀ DI ROMA

Learning Reward Machines for Partially Observable Reinforcement Learning

Riccardo Gozzovelli

Mario Vetrini

Francesco Caputo

REINFORCEMENT LEARNING

SETTING

Markovian Decision Process (MDP), $M = \{S, A, r, p, \gamma\}$

PROBLEM

Maximize the future expected discounted reward for every state in S

SOLUTION

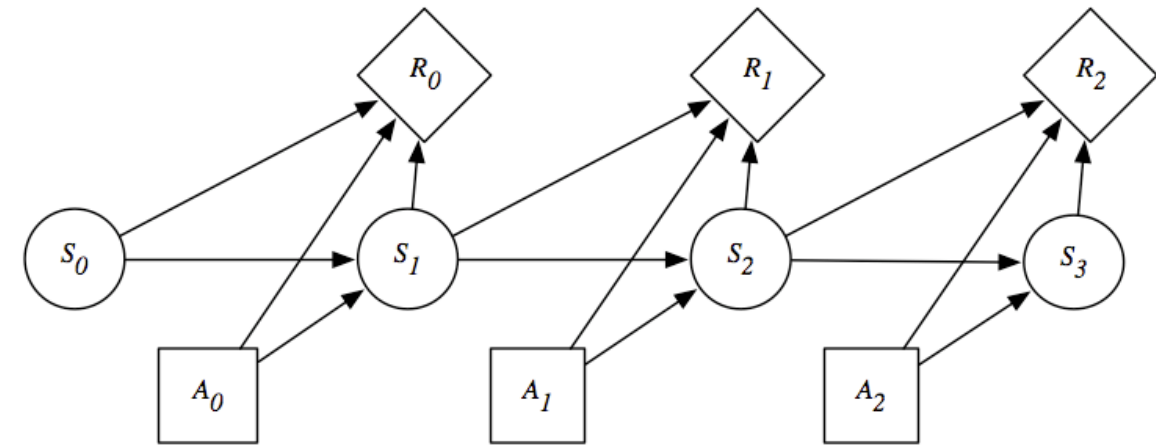
Optimal policy: $\pi^*(a_t|s_t)$

ALGORITHMS

Q-Learning, SARSA, A3C, PPO

DEEP REINFORCEMENT LEARNING

RL algorithms combined with DL principles.



WHY LEARN RMs

SETTING

Partially Observable Markovian Decision Process (POMDP), $M = \{S, O, A, r, p, w, \gamma\}$

PROBLEM

Maximize the future expected discounted reward for every state in S

SOLUTION

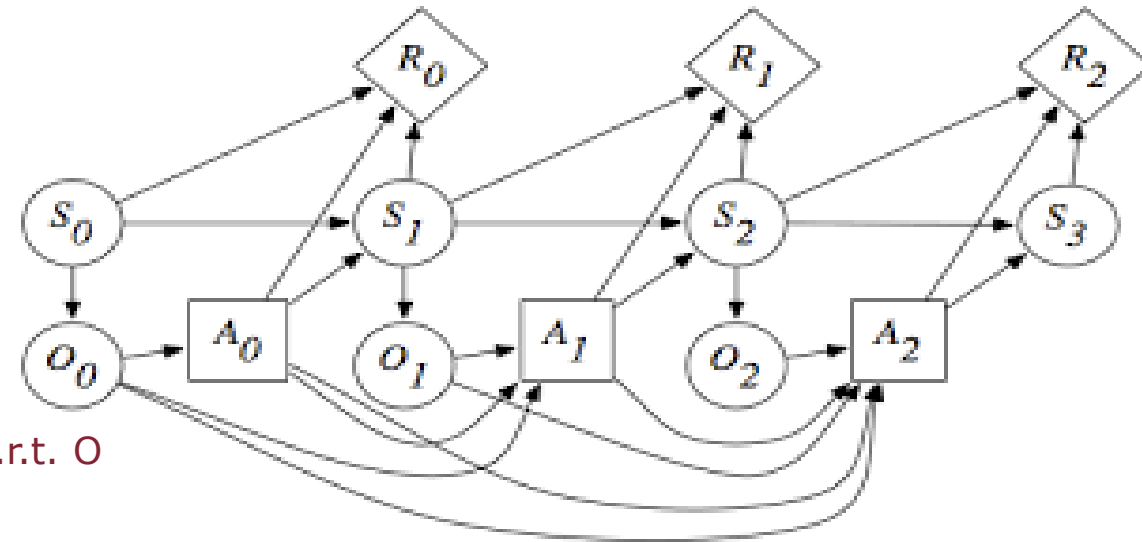
Optimal policy: $\pi^*(a_t|s_t)$

ALGORITHMS

Transition probabilities and reward function not necessarily Markovian w.r.t. O

DEEP REINFORCEMENT LEARNING

RL algorithms combined with DL principles.



PARTIALLY OBSERVABLE DOMAINS

Many real-world applications have partially observable domains

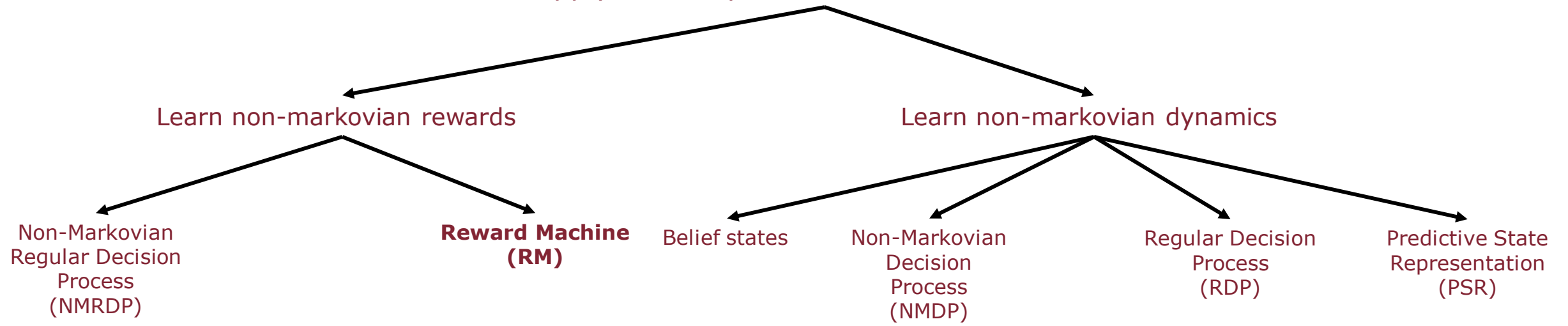


A partially observable domain can have:

- *non-markovian dynamics* or
- *non-markovian rewards* or
- both

RL & POMDPs

How to apply RL to a problem described as a POMDP?



We learn policies and automata!

REWARD MACHINES

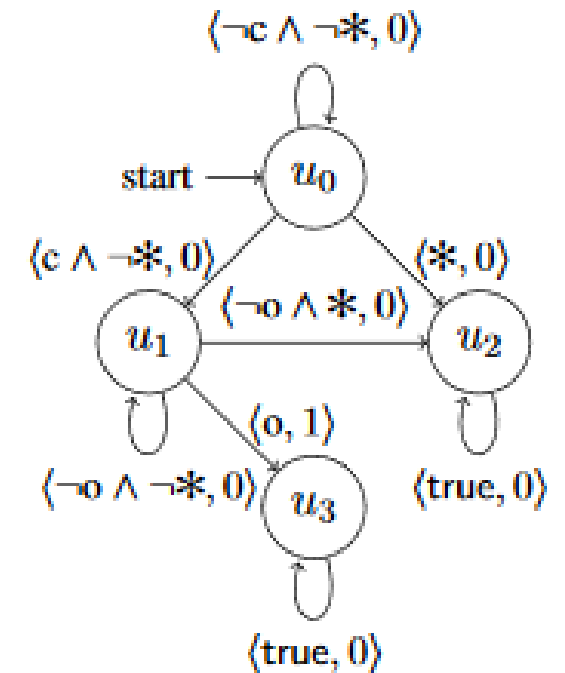
Structured Finite State Automata (FSA) representation.

Expose the reward function and speed up learning.

Provide «memory» to a learning agent.

Variant of Q-Learning algorithm used (QRM).

Initially handcrafted and suitable only for fully observable environments.



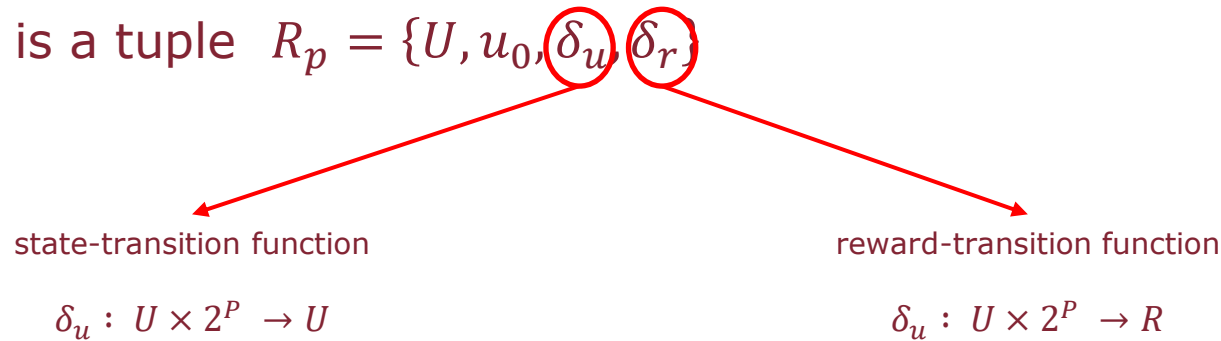
RMs AND POMDPs

RMs defined over set of propositional symbols P (binary properties).

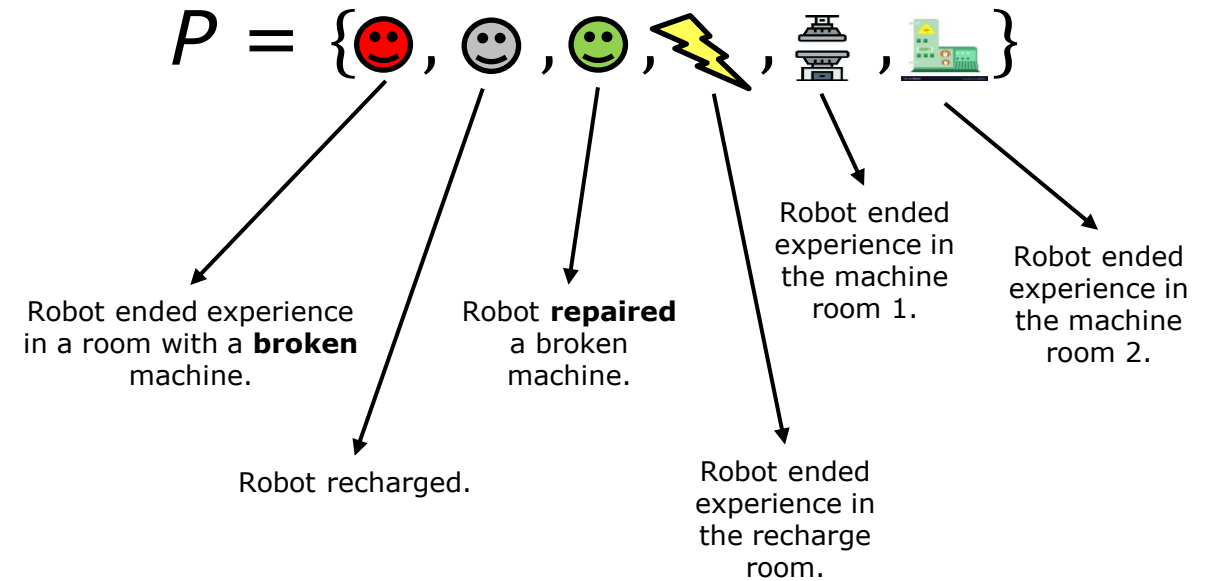
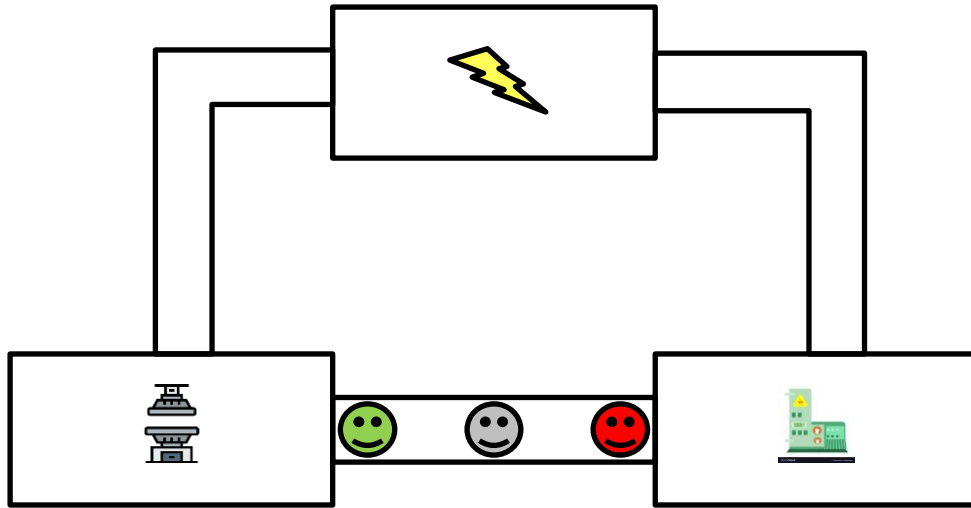
Elements of P are detected by a labelling function $L : O_\phi \times A_\phi \times O \rightarrow 2^P$

L assigns truth values to symbols in P given an experience $e = (o, a, o')$

A reward machine is a tuple $R_p = \{U, u_0, \delta_u, \delta_r\}$

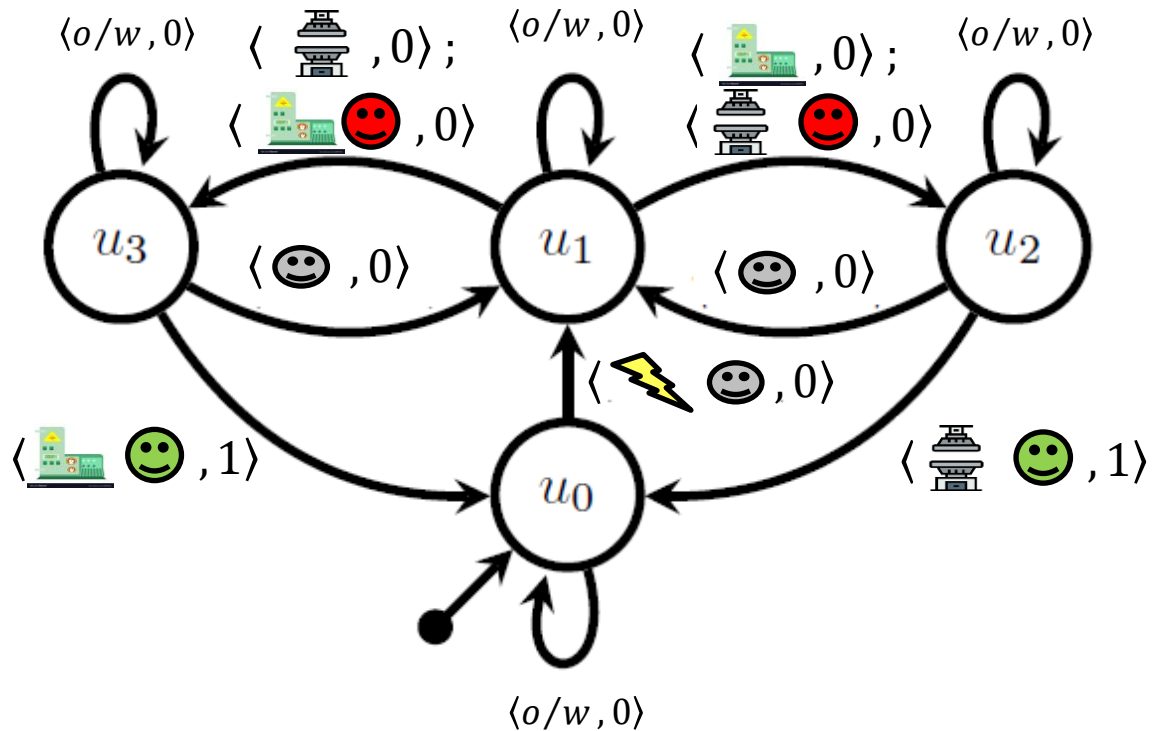


RM – MACHINE MAINTENANCE EXAMPLE



Assumptions: only one machine is broken, the robot can only repair it if it's charged

RM – MACHINE MAINTENANCE EXAMPLE



- Each RM starts in the initial state u_0
- Edge label provides a visual representation of δ_u and δ_r
- Multiple labels separated by a semicolon used to describe different conditions for transitioning between the RM states
- The label “otherwise” on an edge means that that transition will be made if none of the other transitions from u can be taken

SEARCHING FOR AN RM – CASE 1

$\langle \text{🏠} \text{😊}, 1 \rangle;$

$\langle \text{🏢} \text{😊}, 1 \rangle;$

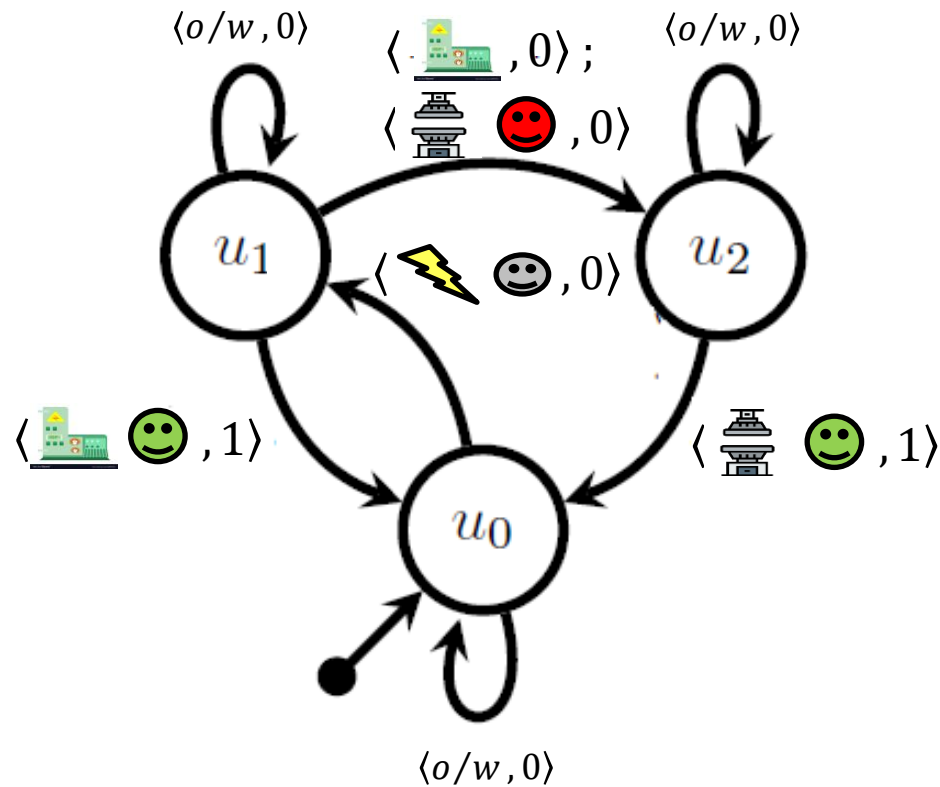
$\langle o/w, 0 \rangle$



NAIVE RM

- Learning the smallest RM that correctly mimics the external reward signal given by the environment
- This naive RM correctly predicts reward in the domain but provides no memory in support of solving the task

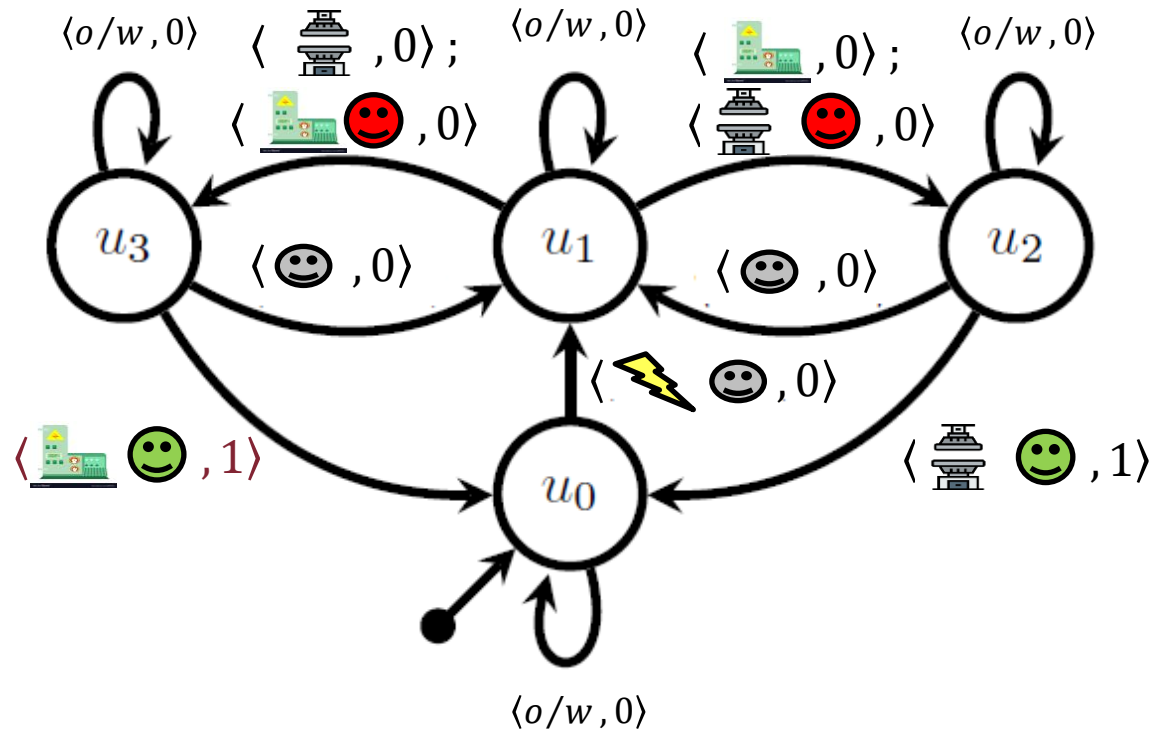
SEARCHING FOR AN RM – CASE 2



OPTIMAL RM

- Looking for the RM whose optimal policy receives the most reward
- Requires computing optimal policies in order to compare the relative quality of RMs, which seems prohibitively expensive

SEARCHING FOR AN RM - CASE 3



PERFECT RM

- Learning the RM that remembers sufficient information about the history to make accurate Markovian predictions about the next observation (w.r.t. $O \times U$)
- Since keeping track of more information will not result in better predictions, this RM is perfect

PERFECT RM - DEFINITIONS AND THEOREMS

DEFINITION: An RM R_p is considered perfect for a POMDP P_o with respect to a labelling function L if and only if for every trace generated by any policy over P_o , the following holds:

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

where $x_0 = u_0$ and $x_t = \delta_u(x_{t-1}, L(o_{t-1}, a_{t-1}, o_t))$

THEOREM: Given any POMDP P_o with a finite reachable belief space, there will always exist at least one perfect RM R_p for P_o with respect to some labelling function L .

THEOREM: Let R_p be a perfect RM for a POMDP P_o w.r.t. a labelling function L , then any optimal policy for R_p w.r.t. the environmental reward is also optimal for P_o .

PERFECT RM - LEARNING

- Learning a perfect RM from traces, assuming one exists w.r.t. the given labelling function L
- One solution can be to fit a predictive model for the previously probability and picking the RM that makes better predictions but it's very expensive
- Alternative that focuses on a necessary condition for a perfect RM: the RM must predict what is possible and impossible in the environment at the abstract level

PERFECT RM - LEARNING

- Let $\mathcal{T} = \{\mathcal{T}_0, \dots, \mathcal{T}_n\}$ be a set of traces with $\mathcal{T}_i = (o_{i,0}, a_{i,0}, r_{i,0}, \dots, a_{i,t_i-1}, r_{i,t_i-1}, o_{i,t_i})$
- Look for an RM $\{U, u_0, \delta_u, \delta_r\}$ that predicts $L(e_{i,t+1})$ from $L(e_{i,t})$ and $x_{i,t}$
- Model parameters:
 - \mathcal{T} (set of traces)
 - P (set of propositional symbols)
 - L (labelling function)
 - u_{\max} (maximum number of states)
 - $I = \{0, \dots, n\}$ (the index of the traces)
 - $T_i = \{0, \dots, t_i - 1\}$ (time steps of trace \mathcal{T}_i)
 - $N_{u,l} \subseteq 2^{2^p}$ (set of all the next abstract observations seen from the RM state u and the abstract observations l at some point in T)

LRM OPTIMIZATION PROBLEM

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

$$s.t. \langle U, u_0, \delta_u, \delta_r \rangle \in \mathcal{RP} \quad (3)$$

$$|U| \leq u_{\max} \quad (4)$$

$$x_{i,t} \in U \quad \forall i \in I, t \in T_i \cup \{t_i\} \quad (5)$$

$$x_{i,0} = u_0 \quad \forall i \in I \quad (6)$$

$$x_{i,t+1} = \delta_u(x_{i,t}, L(e_{i,t+1})) \quad \forall i \in I, t \in T_i \quad (7)$$

$$N_{u,l} \subseteq 2^{2^{\mathcal{P}}} \quad \forall u \in U, l \in 2^{\mathcal{P}} \quad (8)$$

$$L(e_{i,t+1}) \in N_{x_{i,t}, L(e_{i,t})} \quad \forall i \in I, t \in T_i \quad (9)$$

PERFECT RM – SOLVING THE MINIMIZATION PROBLEM

- Local search was the most effective category of methods
- Tabu search guarantees convergence
- Start from a random RM and iteratively evaluates *neighbouring* RMs
- The new RM will be the one with the minimum value of the objective function
- Pruning and methods to avoid local minima also used

PERFECT RM – SIMULTANEOUS LEARNING

In order to learn both an RM and a policy:

1. Collect a training set of traces T generated by a policy during t_w 'warmup' steps.
2. Use T to find an initial RM using tabu search.
3. Initialize policy π , set the initial state of RM to u_o and set the current label for the initial observation $L(\emptyset, \emptyset, o)$.
4. Repeat until convergence:
 - 4.1 Select action a following $\pi(o, u)$
 - 4.2 Get observation o' and reward r
 - 4.3 Update the RM state $u' = \delta_u(u, L(o, a, o'))$
 - 4.4 Update π by using the last experience $(\langle o, u \rangle, a, r, \langle o', u' \rangle)$
5. If in any step of 4 there is an evidence that the current RM might not be the best one, a new one is learned.

PERFECT RM – LEARNING A NEW RM

What can be an evidence that a new RM must be learned?

RM R is selected in order to:
$$\underset{\langle U, u_0, \delta_u, \delta_r \rangle}{\text{minimize}} \sum_{i \in I} \sum_{t \in T_i} \log(|N_{x_{i,t}, L(e_{i,t})}|)$$

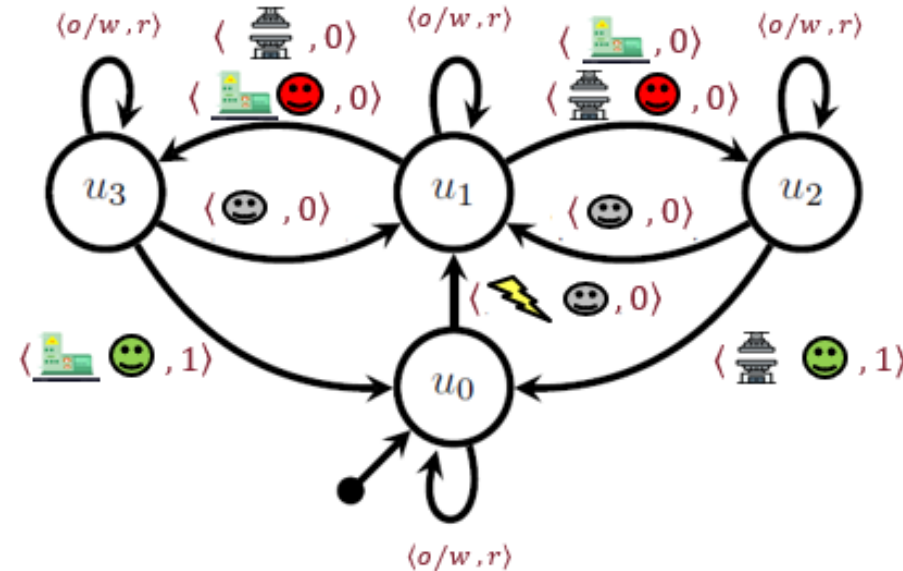
If the current abstract observation I' is not in $N_{u,I}$ then the current trace will increase the size of $N_{u,I}$ and the cost of R .

So the trace will be add to T and a new RM is found using the Tabu search.

If the new RM is better than R then a new policy must be learned for scratch.

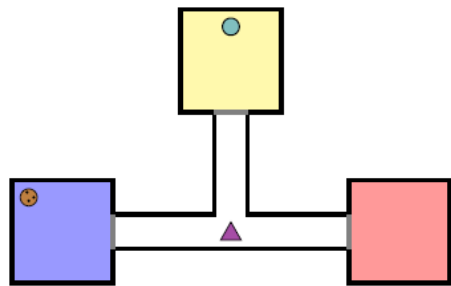
PERFECT RM – DQRM UNDER PARTIAL OBSERVABILITY

An experience $e=(o,a,o')$ can be more or less likely depending on the RM state where the experience was collected.

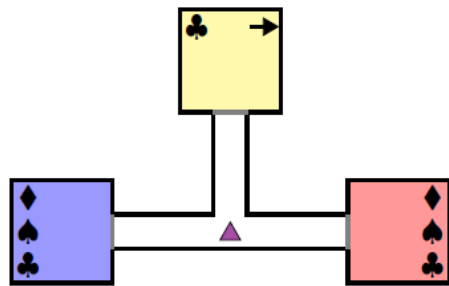


Update q_u using (o,a,o') if and only if $L(o,a,o') \in N_{u,l}$ with l abstract observation that generated e .

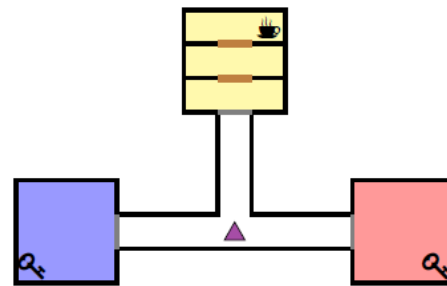
EVALUATION



(a) Cookie domain.



(b) Symbol domain.



(c) 2-keys domain.

Stochastic partially observable domains.

Tested two versions of LRM:

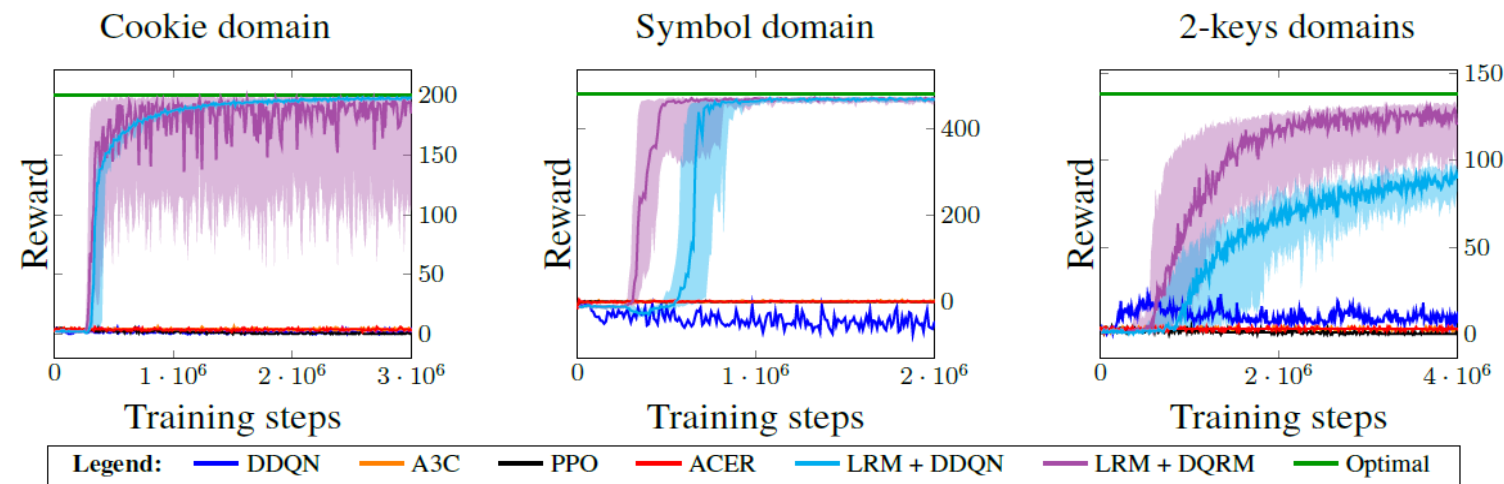
- LRM + DDQN
- LRM + DQRM

Cumulative reward every 10,000 training steps.

Median of 30 training runs per domain.

The LRM versions outperform every baseline.

LRM-DQRM faster but more unstable than LRM-DDQN.



ADVANTAGES AND LIMITATIONS

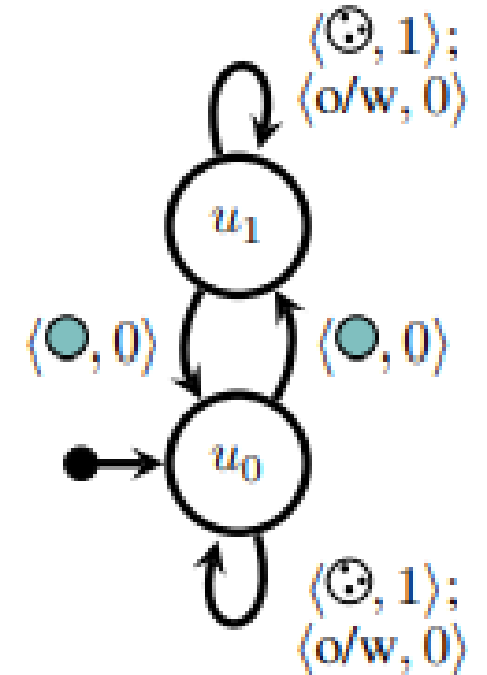
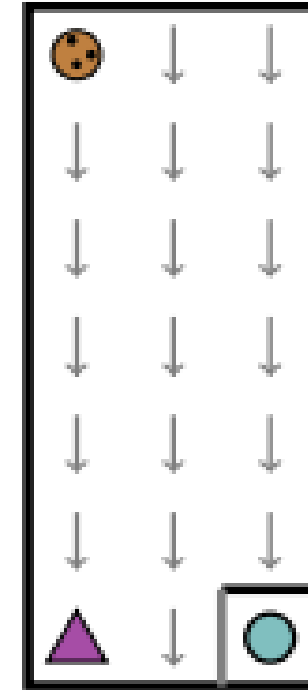
- RMs proved to be able to predict possible and impossible future observations.
- Defining high-level properties might not be trivial for more complex environments.
- Tabu search is the main bottleneck of the entire procedure.
- RMs might ignore (relevant) low-level informations.
- Unclear how to handle noise over the L detectors and how to apply transfer learning from previously learned policies when a new RM is learned.

RM LIMITATION - EXAMPLE

Given $\mathcal{P} = \{\odot, \bullet\}$ a perfect RM is very simple.

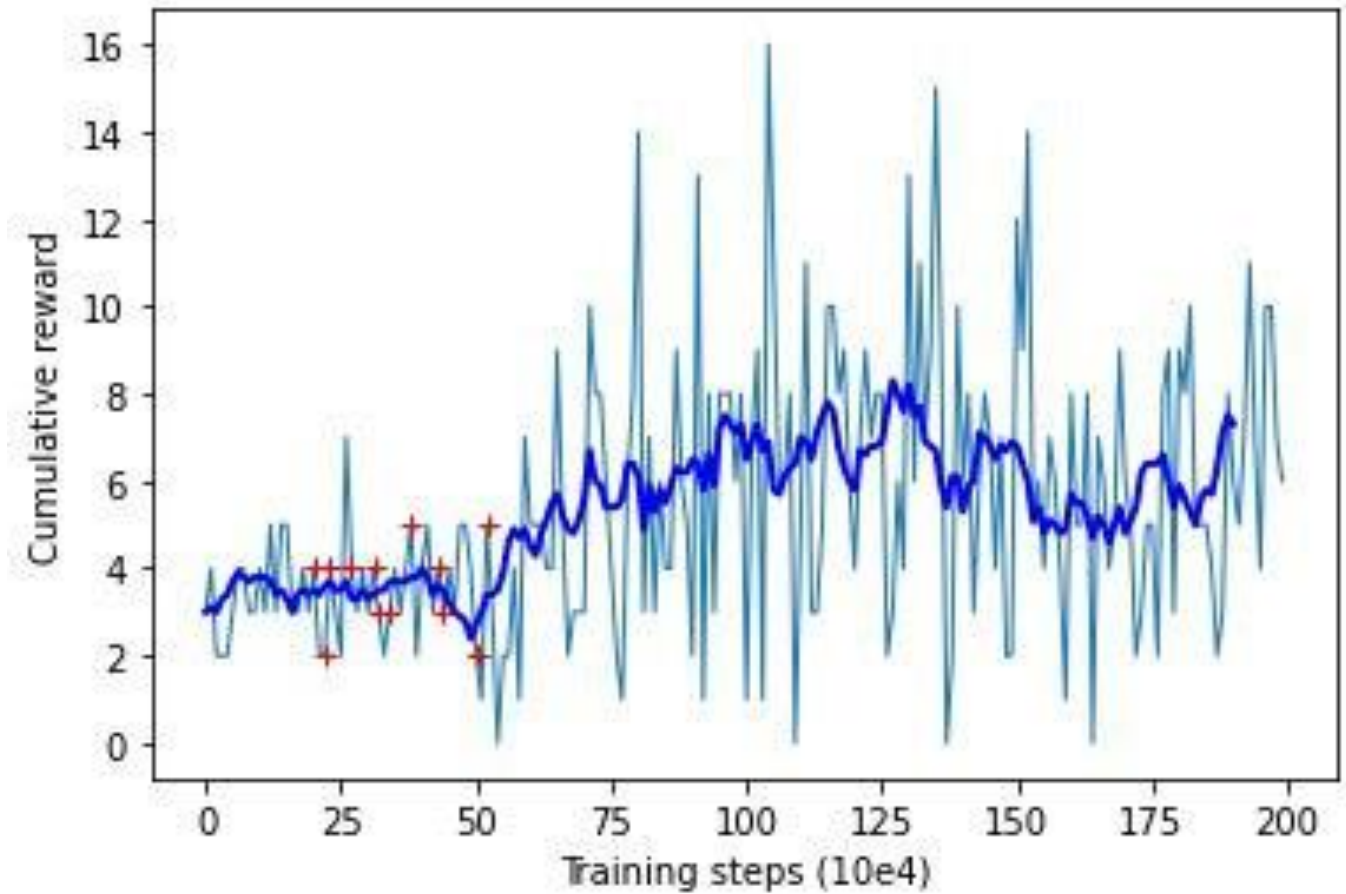
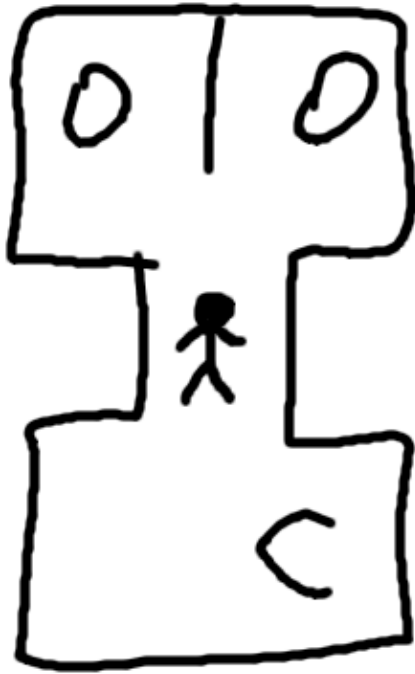
LRM might not find it because the button changes only the low level probability and nothing on the abstract level.

The heuristic in QRM would not work since the experience with the force on will be used to learn a policy in both RM states.

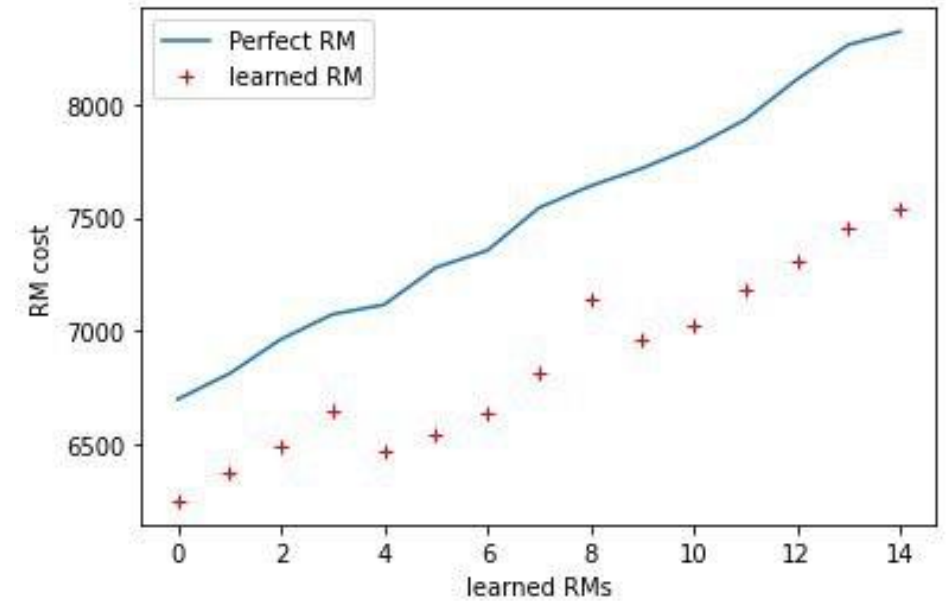
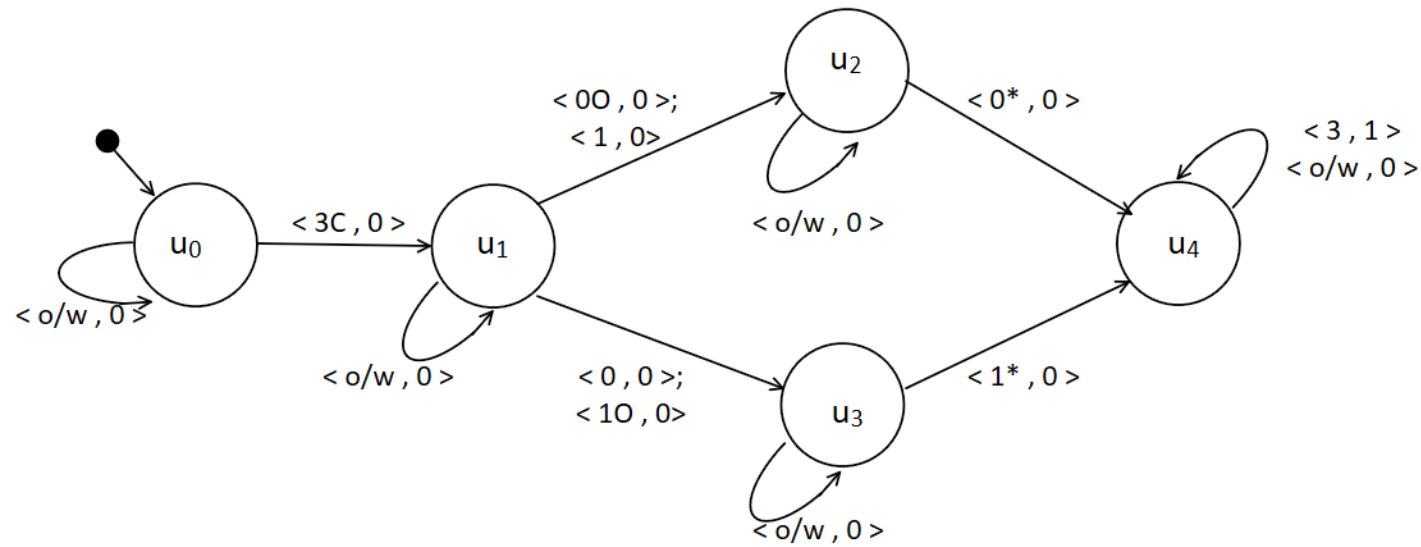
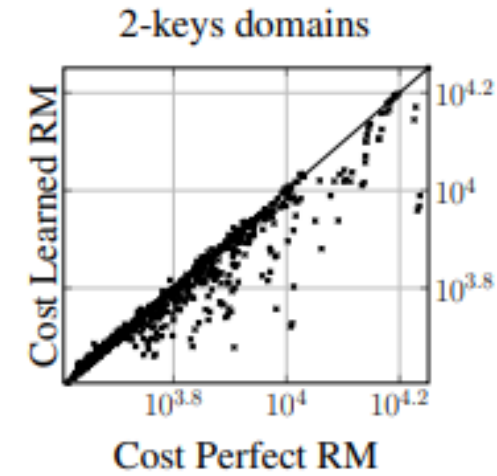
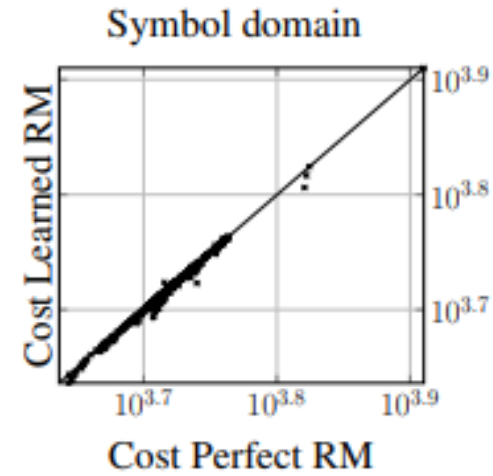
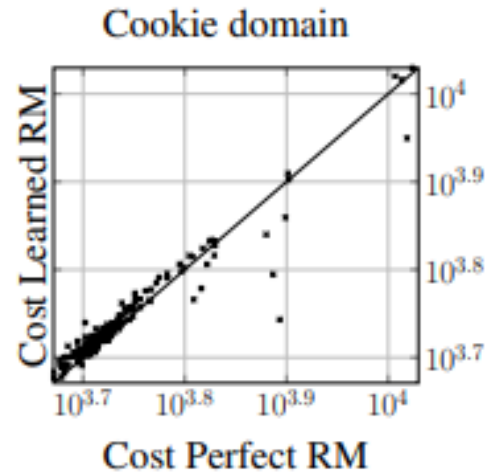


OUR OWN WORLD

- Waiter world: $P = \langle 0, 1, 2, 3, *, C, O \rangle$



RESULTS



CONCLUSIONS

This work showed how RL agents can solve cognitively challenging partially observable tasks.

High level abstraction with RM and low-level policy using deep RL.

The study showed some obstacles regarding:

- Abstraction
- Observability
- Properties of language used for RM construction