

1 Thank you for your excellent feedback. Solving Partially Observable RL is challenging
 2 and, as noted by R3, the proposed method had much better performance than A3C, ACER,
 3 and PPO with LSTM. We address your main questions here and will elaborate in the paper.

4 **How well does Tabu search scale? (R1):** Roughly speaking, Tabu search scales as well as
 5 genetic algorithms (which are commonly used in RL). Its main bottleneck is evaluating the
 6 neighbourhood around the current RM solution, but this step is easily parallelizable. Given
 7 a training set composed of 1 million transitions, a simple Python implementation of Tabu
 8 search took less than 2.5 minutes to learn an RM across all our environments (using 62 workers). In our experiments, the
 9 agent relearned the RM 8.6 times (on average) per run. Note that the size of the neighbourhood depends on the number
 10 of possible abstract observations, and so exhaustively evaluating the neighbourhood may sometimes become impractical.
 11 This well-studied problem has plenty of proposed solutions (known as Large Neighborhood Search methods), though.

12 **What are the limitations of LRM (R1) and when might QRM not work as intended (R2)?:** As R1 mentioned, an
 13 interesting idea from LRM was to optimize over a necessary condition for perfect RMs. This objective favors RMs that
 14 are able to predict possible and impossible future observations at the abstract level given by the labelling function L .
 15 Learning the RM at the abstract level is efficient but requires ignoring (possibly relevant) low-level information. (In the
 16 future, we would like to learn LSTM policies inside the RM to account for any missing information in L .)

17 To the limitations, Figure 1 shows an adversarial example for LRM. The agent receives reward for eating the cookie (☺).
 18 There is an external force pulling the agent down—i.e., the outcome of the “move up” action is actually a downward
 19 movement with high probability. There is a button (●) that the agent can press to turn off (or back on) the external force.
 20 Hence, the optimal policy is to press the button and then eat the cookie. Given $\mathcal{P} = \{\text{☺}, \text{●}\}$, a perfect RM for this
 21 environment is fairly simple (see Figure 1) but LRM might not find it. The reason is that pressing the button changes
 22 the low-level probabilities in the environment but does not change what is possible or impossible at the abstract level.
 23 Moreover, if a perfect RM is found, our heuristic approach to share experiences in QRM would not work as intended
 24 because the experiences collected when the force is on (at u_0) would be used to update the policy for the case where the
 25 force is off (at u_1). From a practical perspective, a simple solution is to add a high-level detector that senses the external
 26 force. Nonetheless, we will discuss the theoretical implications of this interesting adversarial example in the paper.

27 **Question about the evaluation metric and grid sizes (R1):** Indeed, we are reporting the maximum “cumulative”
 28 rewards for the baselines every 10,000 steps. The idea was to highlight that no run of the baselines learned to
 29 consistently solve the tasks. We used small grids in our experiments (it takes 8 steps to go from one room to another).
 30 Hence, a 10-order memory for DDQN should be enough (in principle) to learn to solve our environments.

31 **Relation with Zhang et al.’s work (R1):** This interesting work—which was submitted to arxiv after the NeurIPS
 32 deadline—will definitely be discussed in our paper. In short, the main distinction is that LRM exploits a high-level
 33 abstraction of the observations to understand how to decompose the problem into Markovian subproblems. Exploiting
 34 abstractions has historically helped RL agents to solve challenging domains. In contrast, Zhang et al.’s work has the
 35 merit of learning PSRs at the low-level, making it easier to use as an off-the-shelf tool for Partially Observable RL.

36 **What if the high-level detectors were noisy? (R2):** To handle noise over L (without requiring an explosion of the
 37 size of the RM), it seems necessary to move from deterministic to stochastic RMs. This is to allow the agent to be at
 38 multiple RM states with certain probability. We believe this is an interesting research direction.

39 **Are you planning to release your code? Does LRM+DDQN/DQRN require a stochastic environment? Could
 40 you add more empirical evaluation on why LRM+DQRN converges faster than LRM+DDQN? (R2):** We will
 41 release our code. Our approach works well on the deterministic version of our environments too. We will add further
 42 information, including exploration heatmaps and learned trajectories, to the supplementary materials.

43 **How to prepare an effective labelling function L ? (R3):** Intuitively, any event that might be useful for the agent to
 44 remember is a good candidate to be included in L . We are also interested in investigating methods for learning a suitable
 45 L during environment interaction, and feel this paper demonstrates how one can be exploited if learned (or given).

46 **Is LRM a model-based RL approach and, as such, shouldn’t it be compared with Doshi-Velez et al.? (R3):**
 47 Thanks for pointing this out. LRM+DDQN/DQRN lies somewhere between model-based and model-free as the RM is
 48 learned in a model-based fashion but its policies are learned in a model-free way. This allows our models to leverage
 49 deep RL’s ability to learn policies from low-level inputs (e.g., images). As such, our baselines and part of the related
 50 work discussion was indeed biased towards deep RL approaches for Partially Observable RL. We will partially remedy
 51 this by including a discussion about non-Parametric methods, including Doshi-Velez et al., in the related work section.

52 **Clarification in proof sketch of Theorem 4.3 (R3):** The discrepancy between Def 4.1 and the LRM’s objective value
 53 comes from the fact that LRM is optimizing over a necessary (but not sufficient) condition for finding a perfect RM. If
 54 this does not answer your question, please let us know and we will further elaborate on a revised version of this work.

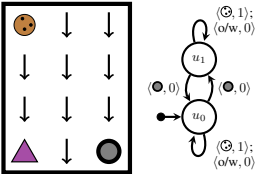


Figure 1: Gravity domain