

---

# On Evaluating Policies for Robust POMDPs

---

**Merlijn Krale\***  
Radboud University  
Nijmegen, The Netherlands  
merlijn.krale@ru.nl

**Eline M. Bovy\***  
Radboud University  
Nijmegen, The Netherlands  
eline.bovy@ru.nl

**Maris F. L. Galesloot\***  
Radboud University  
Nijmegen, The Netherlands  
maris.galesloot@ru.nl

**Thiago D. Simão**  
Eindhoven University of Technology  
Eindhoven, The Netherlands  
t.simao@tue.nl

**Nils Jansen**  
Ruhr-University Bochum  
Bochum, Germany  
n.jansen@rub.de

## Abstract

Robust partially observable Markov decision processes (RPOMDPs) model partially observable sequential decision-making problems where an agent must be *robust* against a range of dynamics. RPOMDPs can be viewed as two-player games between an agent, who picks actions, and *nature*, who adversarially picks dynamics. Evaluating an agent policy requires finding an adversarial nature policy, which is computationally challenging. In this paper, we advance the evaluation of agent policies for RPOMDPs in three ways. First, we discuss suitable benchmarks. We observe that for some RPOMDPs, an optimal agent policy can be found by considering only subsets of nature policies, making them easier to solve. We formalize this concept of *solvability* and construct three benchmarks that are only solvable for expressive sets of nature policies. Second, we describe a provably sound method to evaluate agent policies for RPOMDPs by solving an equivalent MDP. Third, we lift two well-known POMDP upper value bounds to RPOMDPs, which can be used to efficiently approximate the optimality gap of a policy and serve as baselines. Our experimental evaluation shows that (1) our proposed benchmarks cannot be solved by assuming naive nature policies, (2) our method of evaluating policies is accurate, and (3) the approximations provide solid baselines for evaluation.

## 1 Introduction

Partially observable Markov decision processes [POMDPs; 22] are ubiquitous for representing sequential decision-making problems under partial observability. POMDPs have been used to represent many real-world problems, ranging from robotics [27] to infrastructure maintenance [31, 30] to wildlife conservation [10]. Yet, to model such problems, the dynamics of the problem need to be precisely known, which is often unrealistic [23, 48]. *Robust POMDPs* [RPOMDPs; 37] aim to solve this by capturing *model uncertainty*. More precisely, RPOMDPs model this uncertainty as a two-player game between the agent, who picks actions, and nature, who adversarially picks the dynamics of the model [1]. This makes RPOMDPs a more versatile framework, but also makes them harder to solve. Thus, existing solvers use an extensive range of approximations to find agent policies within a reasonable time [45, 12, 7].

To test such solvers, we require an *evaluation pipeline*, such as shown in Figure 1. To start, we need a set of *benchmarks* that allows us to investigate the limits of the solver. We use the solver to obtain a policy, and then need an *evaluation method* to obtain a value for the policy. We compare this value

---

\*These authors contributed equally to this work

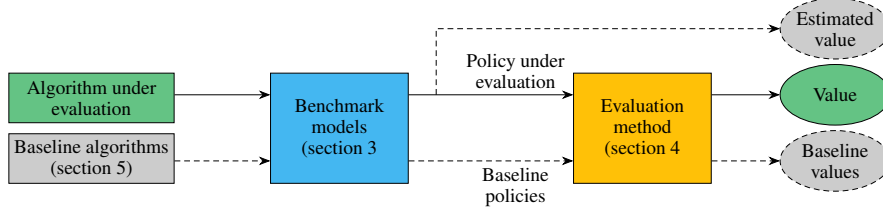


Figure 1: Visualisation of the evaluation pipeline.

with that of suitable *baselines*, and, if available, the estimated value of the policy as provided by the solver. In summary, an evaluation pipeline requires at least these parts: baselines, benchmarks, and an evaluation method.

Most existing RPOMDP literature does not address all these parts. Firstly, benchmarks are used that are based on existing POMDPs with artificial model uncertainty [37, 7, 12, 34, 19], and there has been little research on whether such benchmarks are representative of RPOMDPs in general. Secondly, policy evaluation requires finding a worst-case nature policy, which can be prohibitively expensive. Existing work often either simplifies the evaluation process or reports the estimated value of their approach without explicitly evaluating their policy. Lastly, there are no readily computable baselines for RPOMDPs, hindering the efficient assessment of new solvers. In this paper, we aim to address these gaps and consider all three parts in RPOMDP evaluation with the following contributions:

**(1) We propose novel solvability classes and corresponding benchmark problems (Section 3).**

We discuss the problem of defining benchmarks, and argue that suitable benchmarks should test whether a solver has considered all possible nature policies. Thus, if we can find an optimal agent policy for an RPOMDP by considering only naive subsets of nature policies, then the RPOMDP is not a suitable benchmark. We formalize this intuition using the concept of *solvability*. Based on this concept, we define three small benchmarks for which we prove that the optimal agent policy cannot be found against naive subsets of nature policies.

**(2) We propose a new provably robust policy evaluation method (Section 4).** To evaluate agent policies, we propose to compute a *best-response nature policy*, i.e., the worst-case nature policy against the policy under evaluation. We show that this evaluation method is less computationally expensive than finding an optimal nature policy, and we prove that the aforementioned best-response policy is, in fact, the worst-case evaluation for a fixed agent policy. Lastly, we demonstrate that this evaluation method can be represented as an MDP with continuous state and action spaces, which means we can perform evaluation using off-the-shelf MDP solvers.

**(3) We define two new efficient robust baselines (Section 5).** We lift two existing POMDP approximation algorithms, QMDP [28] and FIB [18], to RPOMDPs. We prove the relative tightness and convergence of these approximations and highlight that, under conventional technical assumptions, both are tractable baselines to compute.

We empirically confirm the solvability classes of our benchmarks, show that our evaluation method is accurate, and see that our approximations are solid baselines for evaluation. For this, we implement all relevant methods in a Julia framework (based on POMDPs.jl [9]), which we will make publicly available so that future research can easily use our methods.

## 2 Robust POMDPs

To start, we give a comprehensive definition of RPOMDPs. We first introduce some basic notation. We denote the set of all possible probability distributions over the (countable) set  $A$  as  $\Delta(A)$ , and all elements with non-zero probability in a distribution  $\mu \in \Delta(A)$  as  $\text{supp}(\mu)$ . For any predicate  $P$ , the *Iverson brackets*  $[P]$  return 1 if  $P$  is true and 0 otherwise. Given a convex set  $C$ , we denote the set of all extreme points as  $\text{Extremes}(C)$ , and the centroid of the set, i.e., the arithmetic mean of all points in  $C$ , as  $\text{Centroid}(C)$ . Lastly, given a function  $f: X \rightarrow \Delta(Y)$  and elements  $x \in X, y \in Y$ ,  $f(y | x)$  denotes the probability of element  $y$  according to  $f(x)$ , and  $y \sim f(x)$  denotes a randomly sampled element  $y$  from  $f(x)$ .

**Model definition.** Robust partially observable Markov decision processes (RPOMDPs) extend POMDPs with model uncertainty. Different definitions of RPOMDPs exist, as discussed in [1]. We focus on a variant of their definition, with some additional assumptions for simplicity.

**Definition 1 (RPOMDP).** An RPOMDP is a tuple  $\mathcal{M} = \langle S, A, \Omega, \mathcal{U}, \mathcal{T}, \mathcal{O}, R, b_0, \gamma \rangle$ , with:

- $S, A$  and  $\Omega$  (finite) sets of states, actions and observations;
- $\mathcal{U} \subseteq \{f: \text{Var} \rightarrow \mathbb{R}\} = \mathbb{R}^{|\text{Var}|}$  the uncertainty set, with  $\text{Var}$  a finite set of decision variables;
- $\mathcal{T}: \mathcal{U} \rightarrow (S \times A \rightarrow \Delta(S))$  the uncertain transition function, which defines a transition function  $T: S \times A \rightarrow \Delta(S)$  for each assignment of decision variables  $u \in \mathcal{U}$ ;
- $\mathcal{O}: \mathcal{U} \rightarrow (S \times A \times S \rightarrow \Delta(\Omega))$  the uncertain observation function which defines an observation function  $O: S \times A \times S \rightarrow \Delta(\Omega)$  for each assignment  $u \in \mathcal{U}$ ;
- $R: S \times A \rightarrow \mathbb{R}$  the reward function;<sup>1</sup>
- $b_0 \in \Delta(S)$  the initial state distribution (or initial belief);
- $\gamma \in [0, 1)$  the discount factor.

For notational convenience, we define  $\mathcal{P}(u)(s', o \mid s, a) = \mathcal{T}(u)(s' \mid s, a)\mathcal{O}(u)(o \mid s, a, s')$  to indicate the joint uncertain transition-observation function. Note that an RPOMDP with a singleton uncertainty set is a POMDP, a fully observable RPOMDP is a robust MDP [RMDP; 20, 35, 46], and an RMDP with a singleton uncertainty set is an MDP [40].

Intuitively, RPOMDPs describe a game between an agent and nature, with policies  $\pi$  and  $\theta$ , respectively [1]. The game starts in a state  $s_0 \in S$  as sampled from the initial state distribution  $b_0$ . At each timestep  $t$ , the agent picks an action  $a_t \in A$  according to its policy  $\pi$ . Then, nature picks a decision variable assignment  $u_t \in \mathcal{U}$  according to policy  $\theta$ . Next, the environment transitions to a state  $s_{t+1} \sim \mathcal{T}(u_t)(s_t, a_t)$ , and emits an observation  $o_t \sim \mathcal{O}(u_t)(s_t, a_t, s_{t+1})$ . Lastly, the agent receives a (non-observable) reward  $r_t = R(s_t, a_t)$ .

**Policies.** Let  $\Pi$  and  $\Theta$  denote sets of agent and nature policies, respectively. We consider randomized agent policies of the form  $\Pi = \{\langle \sigma: X \rightarrow \Delta(A), \tau: X \times A \times \Omega \rightarrow \Delta(X) \rangle\}$  where  $X \subset \mathbb{R}^N$ , with  $N \in \mathbb{N}$ , denotes a *memory space*. This notation generalizes agents where  $X = \Delta(S)$ , i.e., that use beliefs  $b \in \Delta(S)$  [37], or where  $X$  corresponds to memory nodes, i.e., *finite state controllers* [7, 12]. Given a policy  $\pi = \langle \sigma, \tau \rangle \in \Pi$  and *memory element*  $x \in X$ , the agent chooses an action  $a \sim \sigma(x)$  according to the *action selection function*  $\sigma$ , and uses the *memory update function*  $\tau$  to retrieve the next memory element  $x' \sim \tau(x, a, o)$  upon receiving observation  $o \in \Omega$ .

We assume nature has full knowledge of the underlying state of the environment, the last action taken by the agent, and the agent's memory state. We assume history-based deterministic nature policies of the form  $\Theta = \{\theta: (S \times A)^N \times X \rightarrow \mathcal{U} \mid N \in \mathbb{N}\}$ . Furthermore, we assume nature can choose a different decision variable assignment when revisiting a state-action pair, which is known as *dynamic uncertainty* [20] or *zero stickiness* [1] in the literature. We also assume that our uncertainty set is  $\langle s, a \rangle$ -*rectangular*, meaning there exists some partitioning of  $\text{Var}$  into sets  $\text{Var}_{s,a}$ , such that decision variables  $v \in \text{Var}_{s,a}$  only affect transitions and observations in the state-action pair  $\langle s, a \rangle$  [49]. Lastly, we assume the set of joint transition-observation functions in  $\mathcal{P}$  given the uncertainty set  $\mathcal{U}$ , i.e.,  $\{\mathcal{P}(u) \mid u \in \mathcal{U}\}$ , is convex and closed. The former is a common assumption for tractability reasons, while the latter is required to guarantee the existence of an optimal nature policy (as defined below).

**Objective.** The agent's objective is to maximize its *value*, i.e., the infinite-horizon expected cumulative discounted reward. We assume nature is *adversarial*, meaning it aims to minimize the value. Let  $V^{\pi, \theta} := \mathbb{E}_{\pi, \theta}[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 \sim b_0]$  denote the value of the RPOMDP given policies  $\pi \in \Pi$  and  $\theta \in \Theta$ . We denote the optimal value for the agent and nature as  $V^a = \max_{\pi \in \Pi} \min_{\theta \in \Theta} V^{\pi, \theta}$  and  $V^n = \min_{\theta \in \Theta} \max_{\pi \in \Pi} V^{\pi, \theta}$ . If  $V^a = V^n$ , then the underlying game has a *Nash equilibrium* [38], in which case  $\pi^*$  and  $\theta^*$  are called *Nash policies*.

### 3 Solvability Classes and Benchmark Models

In this section, we first propose the concept of *solvability* to analyze the complexity of an RPOMDP. Next, we introduce three small RPOMDP benchmarks that are complex according to our definition.

<sup>1</sup>This definition could trivially be extended to include uncertainty in the reward function [37].

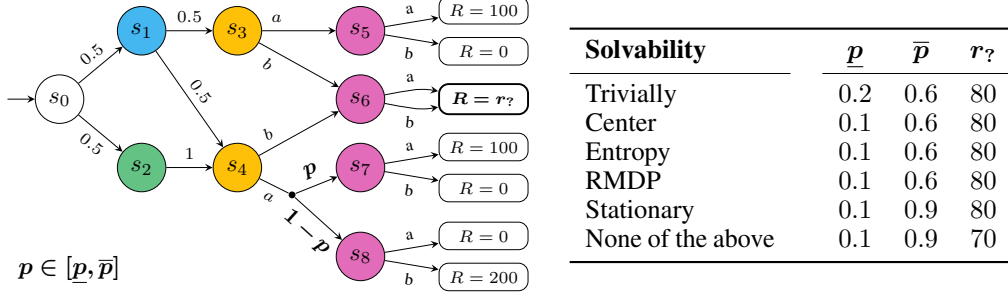


Figure 2: Example RPOMDP with parameter values for solvability classes. We use  $\text{TOY}^*$  to denote the model with the parameter values from the bottom line.

In Section 6, we provide an empirical evaluation of the complexity of these benchmarks, which we compare to standard POMDPs extended with trivial uncertainty sets.

We first introduce the  $\text{TOY}$  environment (Figure 2, left). Throughout this section, we will use variants of this model with different choices for the reward  $r_?$  and uncertainty set  $[p, \bar{p}]$ , as shown on the right of the model. For notational convenience, we assume  $\gamma = 1$ . Intuitively, the agent needs to make two decisions. Firstly, in the  $\bullet$  states, the agent needs to decide whether they want to take the safe action  $b$ , which will always yield a reward  $r_?$ , or take the more risky action  $a$ . In the latter case, the agent needs to make another choice between  $a$  and  $b$  in the  $\bullet$  states. Which action is optimal in the  $\bullet$  and  $\bullet$  states depends both on the history up to that point, as this indicates in which states the agent can be, as well as the nature policy. We highlight two different agent policies for this environment. First, the *safe* policy  $\pi_s^*$  picks the safe action  $b$  in the  $\bullet$  states regardless of the previous observation, which yields a guaranteed value of  $r_?$ . In contrast, the *risky* policy  $\pi_r^*$  picks the riskier action  $a$  in the  $\bullet$  states if the previous observation was  $\bullet$ , and the safe action  $b$  otherwise. Furthermore,  $\pi_r^*$  picks action  $b$  in the  $\bullet$  states. For  $r_? = 80$ ,  $\pi_s^*$  is optimal if  $p \leq 0.6 \leq \bar{p}$ , and  $\pi_r^*$  is optimal if  $p \leq 0.6$  and  $0.2 \leq \bar{p} \leq 0.6$ . Similarly, for  $r_? = 70$ ,  $\pi_s^*$  is optimal if  $p \leq 0.4$  and  $0.65 \leq \bar{p}$ . We prove these policies are optimal for our models in Appendix A.1.

### 3.1 $\Theta$ -solvable

We define the concept of *solvability* as follows:

**Definition 2.** Recall that  $\Theta$  is the set of all nature policies. Define  $\Pi^{\bar{\Theta}}$  as the set of agent policies that are optimal against both  $\Theta$ , as well as a subset of nature policies  $\bar{\Theta} \subseteq \Theta$ :

$$\Pi^{\bar{\Theta}} = \operatorname{argmax}_{\pi \in \Pi} \min_{\bar{\theta} \in \bar{\Theta}} V^{\pi, \bar{\theta}} \cap \operatorname{argmax}_{\pi \in \Pi} \min_{\theta \in \Theta} V^{\pi, \theta}. \quad (1)$$

Then, a model  $\mathcal{M}$  is  $\bar{\Theta}$ -solvable if  $\Pi^{\bar{\Theta}} \neq \emptyset$ .

Intuitively, if a model is  $\bar{\Theta}$ -solvable, then we can find an optimal agent policy against *all* nature policies  $\Theta$  by finding the optimal agent policies against a *subset* of nature policies  $\bar{\Theta} \subseteq \Theta$  and evaluating only those agent policies against *all* nature policies. Thus, a solver that considers only  $\bar{\Theta}$  for a  $\bar{\Theta}$ -solvable model may find optimal solutions more efficiently. Note that  $\bar{\Theta}$ -solvability does not imply that  $\theta^* \in \bar{\Theta}$ , nor that the value against these sets is equal. Below, we define different *solvability classes* based on whether a model is solvable for a particular  $\bar{\Theta}$ .

**Trivial solvability.** We first consider the most extreme case. A model is *trivially solvable* if it is solvable for *any* set  $\bar{\Theta}$ . Such models are unsuitable as benchmarks, since they do not adequately test whether a solver has considered all possible nature policies.

**Example 1.**  $\text{TOY}$  with uncertainty set  $p \in [0.2, 0.6]$  and reward  $r_? = 80$  is trivially solvable, since  $\pi_r^*$  is optimal for any choice of  $\theta$ .

**Naive solvability.** Next, we consider models where the optimal value does depend on the choice of  $\theta$ , but where a naive choice of  $\theta$  suffices. We argue that such models are not adequate to show the capabilities of solvers to be robust against all possible nature policies, and are thus unsuitable as benchmarks. We consider the following policies as naive:

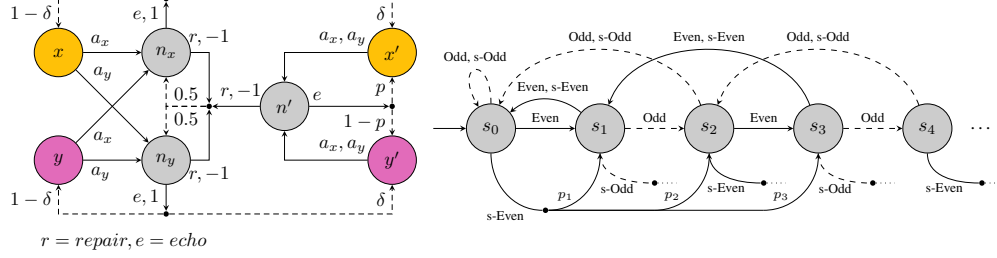


Figure 3: Visualizations of the ECHO environment (left) and PARITY environment (right).

- $\theta_{\text{Center}}$ , which picks decision variables such that  $\mathcal{P}(u) = \text{Centroid}(\{\mathcal{P}(u) \mid u \in \mathcal{U}\})$ . For RPOMDPs constructed by extending a POMDP with model uncertainty, this nature policy often corresponds to the original POMDP.
- $\theta_{\text{Ent}}$ , which picks decision variables to maximize the *entropy* of the probability distribution of each transition. This nature policy generally results in high uncertainty for the agent.
- $\theta_{\text{RMDP}}$ , which picks decision variables that are optimal in the underlying RMDP, i. e., optimal against a fully observing agent policy. In particular, this nature policy is good if partial observability has little effect on the optimal agent policy.

Thus, we call a model *naively solvable* if it is either  $\{\theta_{\text{Center}}\}$ -,  $\{\theta_{\text{Ent}}\}$ -, or  $\{\theta_{\text{RMDP}}\}$ -solvable. Note that this definition of naive solvability is not necessarily exhaustive; additional solvability classes can be defined as naive if so desired.

**Example 2.** TOY with  $p \in [0.1, 0.6]$  and  $r_? = 80$  is not trivially solvable, since  $\pi_r^*$  is optimal against  $\Theta$ , but is not optimal for each individual  $\theta$ . For example, if  $\theta(\cdot) = \{p \mapsto 0.1\}$ , the agent can gain a higher value by choosing the risky action  $a$  in the  $\bullet$  states and action  $b$  in the  $\circ$  states, regardless of the observation history. However,  $\pi_r^*$  is optimal against all the naive policies  $\theta_{\text{Center}}(\cdot) = \{p \mapsto 0.35\}$ ,  $\theta_{\text{Ent}}(\cdot) = \{p \mapsto 0.5\}$ , and  $\theta_{\text{RMDP}}(\cdot) = \{p \mapsto 0.6\}$ , which means the model is naively solvable.

**Stationary solvability.** Lastly, a model is *stationary solvable* if it is solvable for the set of stationary nature policies  $\Theta^{\text{Sta}}$ . Intuitively, an RPOMDP with a stationary nature policy induces a POMDP with the same state space as the RPOMDP. In contrast to the previous solvability classes, the set of nature policies remains continuous. We are not aware of a way to exploit the stationary solvability of a model to solve it. Thus, we do not consider stationary solvable models unsuitable benchmarks.

**Example 3.** Consider TOY with  $p \in [0.1, 0.9]$  and  $r_? = 80$ . Using similar logic as above, against all naive nature policies  $\theta \in \{\theta_{\text{Center}}, \theta_{\text{Ent}}, \theta_{\text{RMDP}}\}$ , it is optimal for the agent to take a risky action  $a$  in the  $\bullet$  states for one of the possible observation histories. However, for any agent policy that takes a risky action in  $\bullet$ , there exists a stationary nature policy that yields a lower value than  $r_?$ . Thus,  $\pi_s^*$  is the only optimal policy against both  $\Theta^{\text{Sta}}$  and  $\Theta$ , meaning the model is stationary solvable and not naively solvable. In contrast, if we use the same uncertainty set with  $r_? = 70$ , taking a risky action  $a$  in  $\bullet$  is optimal against all stationary nature policies, while  $\pi_s^*$  is optimal against the set of non-stationary policies. Therefore, the model is not stationary solvable.

### 3.2 Benchmarks

We introduce three novel environments that, according to our classification above, require solving against expressive nature policies. The first, denoted TOY\*, is the variant of the TOY (Figure 2) that does not fall within any of our solvability classes. Next, we provide two more, detailed below.

**ECHO.** The ECHO environment (Figure 3, left) is inspired by predictive maintenance problems. Starting in one of two distinguishable states  $x$  or  $y$ , the agent picks an action  $a_i$  with  $i \in \{x, y\}$ . Based on this action, the agent transitions to the state  $n_i$ , and can then take the *echo* ( $e$ ) action to transition to state  $i$ . However, the machine has a probability  $\delta$  to transition to the broken states  $x'$  or  $y'$  instead, which return the same observations as  $x$  and  $y$ . From these, the agent always transitions to state  $n'$ , where the outcome of the echo action is determined by a parameter  $p \in [\underline{p}, \bar{p}]$ . The agent receives a reward of 1 for taking the echo action in  $n_x$  or  $n_y$ , but none in  $n'$ . Alternatively, they may

take a *repair* ( $r$ ) action in any of these states at a cost of  $-1$ , which returns the agent to  $n_x$  or  $n_y$  with probability 0.5. We make the following claim about this environment:

**Theorem 1.** ECHO, with  $p = 0.01$ ,  $\bar{p} = 0.99$ ,  $\delta = 0.1$  and  $\gamma = 0.95$ , is not in any of the solvability classes defined in Section 3.1.

We give a more general formulation of this theorem, as well as a proof, in Appendix A.2.1. Intuitively, a non-stationary nature policy can pick  $p$  depending on whether the agent previously played  $a_x$  or  $a_y$ , while a stationary nature cannot. Thus, the optimal policy will sometimes repair even if it is uncertain whether the machine is broken, while this is always suboptimal against a stationary nature policy.

**PARITY.** Lastly, we consider an abstract chain environment called  $\text{PARITY}(N)$ , parameterized with  $N \in \mathbb{N}_{>0}$ , as visualized in Figure 3 (right). The environment consists of a chain of indistinguishable states  $s_0$  to  $s_N$ , and the agent receives a reward when reaching  $s_N$ . To do so, the agent can guess the *parity* of its current state  $s_i$ , and can choose actions with either *deterministic* (*Even*, *Odd*) or *stochastic* (*s-Even*, *s-Odd*) outcome. When correctly guessing the parity in state  $s_i$ , the agent moves to state  $s_{i+1}$  when choosing deterministic actions, and to states  $s_{i+1}$ ,  $s_{i+2}$  or  $s_{i+3}$  with probabilities  $p_1 \in P_1, p_2 \in P_2, p_3 \in P_3$  otherwise. However, the agent moves to state  $s_{i-2}$  for incorrect guesses. Stochastic actions take the agent further on average, but make it harder to guess correctly in the future. In addition to finite-length chains, we consider a chain of infinite length with the same dynamics, denoted  $\text{PARITY}(\infty)$ , where the agent receives an immediate reward equal to the number of steps they take. We show this problem can be reduced to a 9-state model in Appendix A.2.2. For this environment, we show the following with regard to solvability:

**Theorem 2.**  $\text{PARITY}(\infty)$ , with  $P_1 = \{0.2\}$ ,  $P_2 = [0.1, 0.7]$ ,  $P_3 = [0.1, 0.7]$ , and  $\gamma \geq 0.7\bar{3}$ , is not naively solvable.

We provide a proof in Appendix A.2.2. Intuitively, our proof shows that optimal policies for naive subsets of nature policies take the riskier stochastic actions, which is suboptimal in the worst case. We empirically show in Section 6 that  $\text{PARITY}(10)$ , with  $P_1 = \{0.1\}$ ,  $P_2 = [0.5, 0.8]$ ,  $P_3 = [0.1, 0.4]$ , and  $\gamma = 0.95$  is also not naively solvable.

## 4 Evaluating Agent Policies in RPOMDPs

In this section, we consider the problem of policy evaluation for RPOMDPs. To evaluate an agent policy  $\pi$ , we must consider against which nature policy to evaluate. One choice is to evaluate against the optimal nature policy  $\theta^* \in \operatorname{argmin}_{\theta \in \Theta} \max_{\pi \in \Pi} V^{\pi, \theta}$ , the worst-case nature policy when considering all possible agent policies. However, reasoning over all possible agent policies is often intractable, and is unnecessary if we will only use  $\theta^*$  to evaluate a single policy. Moreover, there might be alternative nature policies that yield lower values, as the following example illustrates.

**Example 4.** Consider an agent policy in the  $\text{TOY}^*$  environment (Figure 2) that always picks action  $a$  in the  $\bullet$  states, then  $b$  in the  $\circ$  states. This policy achieves a close-to-optimal value of  $66\frac{2}{3}$  against  $\theta^*$ , but is highly exploitable by a nature policy that chooses  $p = 0$ , achieving a value of 0.

Thus, evaluation against  $\theta^*$  does not always give a good indication of the robustness of a policy. Instead, we propose to evaluate the policy  $\pi$  against its own worst-case nature, which we denote as the *best-response* policy  $\theta_\pi^* \in \operatorname{argmin}_{\theta \in \Theta} V^{\pi, \theta}$ . In practice, finding (approximations of)  $\theta_\pi^*$  is computationally cheaper than finding  $\theta^*$ . In particular, we show that for  $\langle s, a \rangle$ -rectangular and dynamic uncertainty models, the policy type of  $\theta_\pi^*$  is typically simpler than that of  $\theta^*$ , which suggests that it is easier to find. Recalling  $\mathcal{P}$  and  $\text{Extremes}(\cdot)$  from the preliminaries, we state this as follows:

**Theorem 3.** For any agent policy  $\pi \in \Pi$  there exists a best-response nature policy  $\theta_\pi^*: X \rightarrow \mathcal{U}$  such that  $\forall x \in X: \mathcal{P}(\theta_\pi^*(x)) \in \text{Extremes}(\{\mathcal{P}(u) \mid u \in \mathcal{U}\})$ .

We provide a full proof in Appendix B. Intuitively, nature does not need to consider the whole history since, due to  $\langle s, a \rangle$ -rectangularity, it can pick variable assignments that are optimal for any reached state-action pair. However, nature’s decision should still be based on what the agent will do in the future, which depends on the agent’s current memory state  $x \in X$ , as the agent policy is stationary over  $X$  and the memory update  $\tau$  is fixed for a policy  $\pi \in \Pi$ . Moreover, since the agent policy is fixed, nature enjoys more expressive power in the variable assignments it chooses. Therefore, nature only needs to consider a subset of the policy class  $\Theta$  to select the best-response actions.

Now, using Theorem 3, we propose a method for finding  $\theta_\pi^*$ . In particular, since  $\theta_\pi^*$  is Markovian in  $X$ , we can represent our problem as a (possibly continuous-space) *nature MDP* with state space  $S_n = X$ , action space  $A_n \subset \mathcal{U}$ , and dynamics that represent both the underlying RPOMDP and the memory update of the agent. To simplify our notation, we use an expanded state space that explicitly includes the current state and agent action. In that case, we define the nature MDP as follows:

**Definition 3** (Nature MDP). *Assume we have an RPOMDP  $\mathcal{M} = \langle S, A, \Omega, \mathcal{U}, \mathcal{T}, \mathcal{O}, R, b_0, \gamma \rangle$  and an agent policy tuple  $\pi = \langle \sigma, \tau \rangle \in \Pi$  that uses a (possibly continuous) memory space  $X$ . Then, the corresponding nature MDP  $M_n^\pi$  is defined as  $M_n^\pi = \langle S_n, A_n, T_n, R_n, \mu_n, \gamma \rangle$ , with:*

- $S_n = S \times X \times A$  the state space;
- $A_n \subset \mathcal{U}$  the action space. The available actions  $A_n(\langle s, x, a \rangle) \subseteq A_n$  for each nature state  $\langle s, x, a \rangle \in S_n$  are defined as:

$$A_n(\langle s, x, a \rangle) = \{u \in \mathcal{U} \mid \mathcal{P}(u) \in \text{Extremes}(\{\mathcal{P}(u) \mid u \in \mathcal{U}\})\}.$$

- $T_n: S_n \times A_n \rightarrow \Delta(S_n)$  the transition function, defined as:

$$T_n(\langle s', x', a' \rangle \mid \langle s, x, a \rangle, a_n) = \sigma(a' \mid x') \sum_{o \in \Omega} \mathcal{P}(a_n)(s', o \mid s, a) \tau(x' \mid x, a, o).$$

- $R_n: S_n \rightarrow \mathbb{R}$  the state-based reward function, with  $R_n(\langle s, x, a \rangle) = R(s, a)$ .
- $\mu_n$  the initial state distribution resulting from  $b_0$  and  $\pi$ .

The objective in the nature MDP is to minimize the expected reward, as the reward function is based on the reward in the original RPOMDP.

**Remark 1.** *Given a policy  $\theta \in \Theta_n^\pi := \{\theta: S_n \rightarrow A_n\}$ , its value  $V_n^{\pi, \theta} := \mathbb{E}_{\pi, \theta} [\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 \sim b_0]$  is equal to  $V^{\pi, \theta}$  by construction. Moreover, via Theorem 3, there exists a policy in  $\Theta_n^\pi$  that is a best-response against  $\pi$ . Thus,  $\theta_n^* \in \arg\min_{\theta \in \Theta_n^\pi} V_n^{\pi, \theta}$  is a best-response nature policy against  $\pi$ .*

Remark 1 implies that we can find  $\theta_\pi^*$  by (approximately) solving the nature MDP using off-the-shelf methods. This generalizes the robust Markov chain construction [12] to arbitrary memorization schemes  $\tau$  of the agent. We note that  $A_n$  is finite if  $\mathcal{P}$  is a convex polytope<sup>2</sup>, and  $S_n$  is discrete if  $X$  is discrete. Both these conditions hold for the setting of [12], which is why their robust policy evaluation scales to relatively large environments. However, the state and action space of  $M_n^\pi$  can be continuous in general, and the value function of  $M_n^\pi$  can be discontinuous and/or non-convex. Thus, in practice, we may have to resort to an off-the-shelf approximation method for solving  $M_n^\pi$ .

## 5 Efficient Approximations for Robust Agent Policies

In this section, we lift two value approximation methods from POMDPs to their robust counterparts in RPOMDP. These approximations are helpful in multiple ways. Many existing POMDP algorithms, both online [50] and offline [44], use approximations as upper bounds to guide exploration and initialize value estimates. While any choice for  $\theta \in \Theta$  is a valid upper bound, using tighter upper bounds often leads to better results more quickly [26]. Furthermore, these approximations are ideal candidates to serve as baselines and provide sanity checks in an evaluation, such as the one we conduct in Section 6, as they provide lower bounds on performance that are efficient to compute.

Recall  $b \in \Delta(S)$  is an agent belief. Let  $\mathbf{b}_s$  be the *unit belief* with  $b(s) = 1$ , and  $\mathcal{B}_S = \{\mathbf{b}_s \mid s \in S\} \subset \Delta(S)$  be the set of all such beliefs. Then, we define robust variants of the QMDP-bound [29] and the *fast informed bound* [FIB; 18] as:

**Definition 4.**  $Q_{\text{RMDP}}$  and  $Q_{\text{RFIB}}$  are the fixed point of the operators  $H_{\text{RMDP}}$  and  $H_{\text{RFIB}}$ , defined as:

$$H_{\text{RMDP}}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{s' \in S} \mathcal{T}(u)(s' \mid s, a) \max_{a' \in A} Q(\mathbf{b}_{s'}, a') \right], \text{ and} \quad (2)$$

$$H_{\text{RFIB}}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o \mid s, a) Q(\mathbf{b}_{s'}, a') \right]. \quad (3)$$

<sup>2</sup>This is the case if the uncertainty set is comprised of probability intervals or the  $\ell_1$  or  $\ell_\infty$  norms, but generally not for sets based on  $\ell_2$  or Kullback-Leibner divergence.



We prove both these operators are contraction mappings in Appendix C of the supplemental material, guaranteeing the existence and uniqueness of these fixed points. Intuitively, RFIB corresponds to the worst-case value under the assumption that the agent observes the current and future states from the next timestep onwards with a one-step delay, while RQMDP corresponds to the same assumption with no delay. This matches their POMDP variants. Let  $Q_{\text{RPOMDP}}$  be the fixed point of the robust Bellman equation for RPOMDPs [37]. We highlight the following property:

**Theorem 4.** *Regarding tightness, the following inequalities on the fixed points hold:*

$$\forall b \in \Delta(S), \forall a \in A: Q_{\text{RMDP}}(b, a) \geq Q_{\text{RFIB}}(b, a) \geq Q_{\text{RPOMDP}}(b, a).$$

We provide a full proof in Appendix C of the supplemental material. Intuitively, the theorem follows from the fact that the operators are contraction mappings and unequal. As for their non-robust variants, both  $Q_{\text{RMDP}}$  and  $Q_{\text{RFIB}}$  depend only on the  $Q$ -values for state-action pairs. Thus, precomputing the (approximate) fixed point for beliefs  $b_s \in \mathcal{B}_S$  allows computing the bounds for any belief  $b$  efficiently. Whether or not this precomputation is tractable depends on the uncertainty sets. Both are convex optimization problems for convex uncertainty sets. In particular, an iteration of Equation (2) is solvable with a linear program, or by using an efficient *bisection method* [35, Section 7.2] if the uncertainty sets are convex polytopes, while Equation (3) is solvable via a linear relaxation of a mixed-integer program, which means both require polynomial time.

## 6 Experiments and Discussion

Next, we empirically evaluate our benchmarks and evaluation pipeline by addressing the following:

- (Q1) **Evaluation pipeline.** How accurate is our pipeline in evaluating policies?
- (Q2) **Benchmarks.** Can our proposed benchmarks be solved using a naive nature heuristic, or one of our efficient approximations? How does this compare to other benchmarks?
- (Q3) **Approximations.** How do our approximations perform as baselines?

We first provide a brief overview of our experimental setup, and then address these questions. We include more details in Appendix D, and will make all code and results publicly available after publication.

### 6.1 Experimental Setup

**Implementation and Algorithms.** We implement our evaluation method in the Julia programming language, using a variant of the POMDPs.jl framework [9] for RPOMDPs with interval uncertainty sets. To solve these RPOMDPs, we use three separate algorithms: a variant of RHSVI [37] (with minor alterations, as described in Appendix D.1), as well as our approximate solvers RQMDP and RFIB (from Section 5). With these solvers, we compute both robust policies on the RPOMDP  $\mathcal{M}$ , as well as non-robust POMDP policies for the simplified models  $M_{\text{Center}}$ ,  $M_{\text{Ent}}$  and  $M_{\text{RMDP}}$ , which correspond to the naive nature policies  $\theta_{\text{Center}}$ ,  $\theta_{\text{Ent}}$ ,  $\theta_{\text{RMDP}}$  in Section 3. Given an agent policy, we construct a nature MDP in the POMDPs.jl framework, and solve it using the native Julia implementation of Monte Carlo tree search (MCTS) for continuous-state MDPs [6]. For evaluation, we run MCTS five times and report the lowest value; unnormalized values and standard deviation are in Appendix D.2.

**Benchmarks.** We test our evaluation pipeline on three sets of benchmarks. Firstly, we use our novel benchmarks as introduced in Section 3: TOY\* and ECHO, as well as the finite- and infinite chain environments PARITY(10) and PARITY( $\infty$ ). Secondly, we lift several POMDPs from the literature into RPOMDPs: TIGER [4], MINIHALLWAY [28], and ALOHA [21], as well as an expanded variant of HEAVENORHELL [2] (also used in [37]). We construct these RPOMDPs such that for any  $\mathcal{T}(u)$ , any transition probability is less than 0.5 times higher or lower than the nominal POMDP, with no alterations to the observation function. Lastly, we add partial observability to two benchmarks from the RMDP literature: HEALTHDETECTION [13] and REPLACEMENT [8]. The former can be interpreted as an RPOMDP with little changes. For the latter, we add partial observability by adding *measuring actions* that reveal the state at the cost of incurring a negative reward [25]. Appendix D.2 provides more detailed descriptions of the environments and their dimensions.

**Metric.** We test our evaluation method on the abovementioned algorithms to compute value  $V^{\pi, \theta_{\pi}^*}$ . From these evaluations, we compute a *relative value gap*  $V_{\text{gap}} = V^{\pi, \theta_{\pi}^*} - \tilde{V} / |\tilde{V}|$ , where  $\tilde{V}$  denotes the



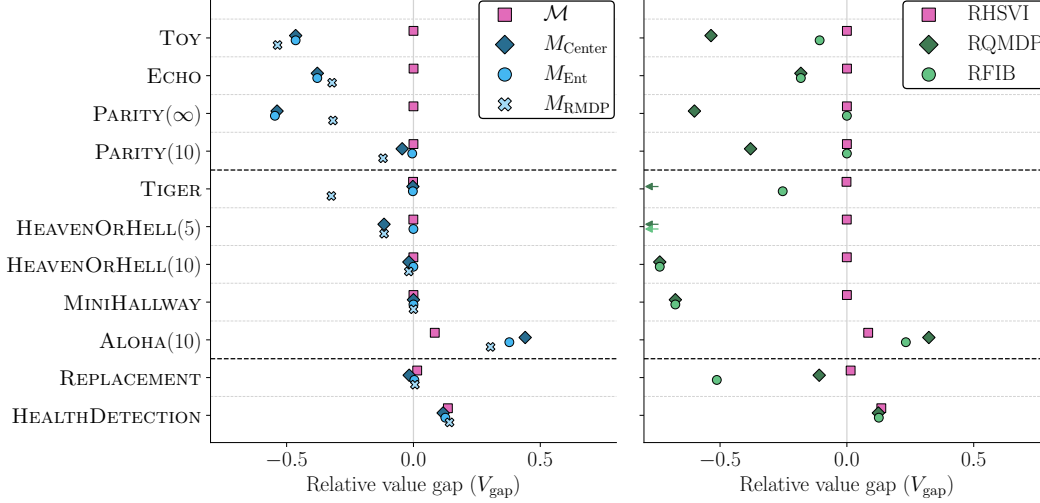


Figure 4: Relative optimality gap  $V_{\text{gap}}$  found by solving the nature MDP. On the left, policies are computed using RHSVI on the RPOMDP model  $\mathcal{M}$  or RHSVI on the naive POMDP models  $M_{\text{Center}}$ ,  $M_{\text{Ent}}$ , and  $M_{\text{RMDP}}$ . On the right, policies are computed using different RPOMDP algorithms. Arrows denote outliers.

(approximate) value of the RPOMDP computed using RHSVI. A negative value gap indicates that our pipeline has determined a policy is suboptimal. In contrast, a positive value gap means our pipeline could not compute the worst-case value, i.e., due to the approximation used to solve the nature MDP.

## 6.2 Results and Discussion

We aggregate our results in two separate plots (Figure 4), which we analyze below:

**(Q1): Our evaluation pipeline is accurate for smaller environments.** For the smaller environments (MINIHALLWAY and above), our evaluation of RHSVI achieves a value gap close to zero, and the value gaps for all approximation methods are less than or equal to zero. Thus, for these environments, our evaluation pipeline is likely accurate. However, for larger environments (ALOHA and HEALTHDETECTION), some policies achieve positive value gaps. In these cases, MCTS fails to find an accurate solution for the nature MDP, leading to inaccurate policy evaluation.

**(Q2): Our novel benchmarks cannot be solved naively.** For TOY\*, ECHO, and PARITY( $\infty$ ), all the policies computed with naive nature heuristics have significant value gaps. Moreover, while RFIB yields an optimal policy for both variants of PARITY, it does not solve TOY\* or ECHO, despite their relatively small state spaces. In contrast, for all other benchmarks, at least one policy computed from a naive nature heuristic performs on par with RHSVI.

**(Q3): RFIB is an adequate baseline.** RFIB outperforms all naive nature heuristics for TOY\*, ECHO, and both variants of PARITY, on which it also is optimal. This confirms RFIB is useful as a computationally inexpensive baseline for RPOMDPs. However, as we see for the other benchmarks, RFIB is not guaranteed to perform well on all benchmarks. In contrast to RFIB, RQMDP performs worse in general, but is still a valuable baseline for models where RFIB is too expensive to compute.

**Limitations.** As shown above, approximation errors in solving the nature MDP may result in incorrect values. Thus, care should be taken in the approximation method and the values reported. Furthermore, our methods and analysis in this paper are restricted to  $\langle s, a \rangle$ -rectangular and convex uncertainty sets. Yet, these are common assumptions for RPOMDPs [37, 45, 7, 12], while studying less conservative cases of rectangularity is still an open problem even in fully observable settings [49, 14].

## 7 Related Work

**RPOMDP Evaluation.** No prior work exists that explicitly tackles the evaluation problem for RPOMDPs. However, methods that aim to solve RPOMDPs often make implicit assumptions about

evaluation. Most notably, Galesloot et al. [12] consider an iterative planning framework that uses an evaluation method similar to ours. However, their work focuses on finding (and evaluating) finite state controllers for stationary nature policies only, while our evaluation method is more generic. More broadly, RPOMDP solvers exist that compute belief-based policies [37], history-based policies [19, 33] or policies represented as finite state controllers [45, 7, 12]. In all these works, the benchmarks consist of POMDPs with  $\epsilon$ -uncertainty around the original transition and observation functions, and include no discussion on why these benchmarks are picked. In particular, such variants of both TIGER and HEAVENORHELL have been used as RPOMDP benchmarks [19, 33, 37], while our work shows that both are naively solvable.

**Related settings.** Next, we discuss a number of settings that are related to RPOMDPs. Burns and Brock [3] study evaluation for robust planning through sampling POMDPs from the uncertainty set. In contrast to our evaluation method, this does not guarantee finding the worst-case value. Other settings include robustifying POMDP policies against observation perturbations [5], robust active measuring [25], non-rectangular but finite sets of POMDPs [11], distributionally robust value iteration with side information [32], and value iteration under varying pessimism levels [41]. However, since these works describe slightly different settings, we cannot directly compare our methods. Lastly, we note that solver evaluation has been studied in different fields, including for MDPs [17] and reinforcement learning, both in general [36, 15] and for RMDPs in particular [51].

## 8 Conclusion

In this paper, we consider three understudied components of the RPOMDP evaluation pipeline: (1) finding suitable benchmarks, (2) robust policy evaluation, and (3) efficient baseline algorithms. We introduce novel methods to tackle all three problems, and empirically confirm that the resulting pipeline is sound. Future work could use our approximations to guide RPOMDP solvers, or introduce more specific approximations of the nature MDP to better scale robust evaluation.

**Practical recommendations.** For future research on RPOMDPs, we have two practical recommendations. Firstly, RPOMDP solvers should be tested on benchmarks that are not naively solvable, which can be tested theoretically (as done in Section 3) or empirically (as done in Section 6). Secondly, RPOMDP solvers should be evaluated using a robust evaluation method that considers the full range of possible nature policies, such as the one proposed in Section 4.

## Acknowledgments

We would like to thank the anonymous reviewers for their useful comments. This work has been partially funded by the ERC Starting Grant DEUCE (101077178).

## References

- [1] Eline M. Bovy, Marnix Suilen, Sebastian Junges, and Nils Jansen. Imprecise probabilities meet partial observability: Game semantics for robust POMDPs. In *IJCAI*, pages 6697–6706. ijcai.org, 2024.
- [2] Darius Braziunas and Craig Boutilier. Stochastic local search for POMDP controllers. In *AAAI*, pages 690–696. AAAI Press / The MIT Press, 2004.
- [3] Brendan Burns and Oliver Brock. Sampling-based motion planning with sensing uncertainty. In *ICRA*, pages 3313–3318. IEEE, 2007.
- [4] Anthony R. Cassandra, Leslie Pack Kaelbling, and Michael L. Littman. Acting optimally in partially observable stochastic domains. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 94, pages 1023–1028, 1994.
- [5] Mahmoud El Chamie and Hala Mostafa. Robust action selection in partially observable Markov decision processes with model uncertainty. In *CDC*, pages 5586–5591. IEEE, 2018.
- [6] Adrien Couëtoux, Jean-Baptiste Hoock, Nataliya Sokolovska, Olivier Teytaud, and Nicolas Bonnard. Continuous upper confidence trees. In *LION*, volume 6683 of *Lecture Notes in Computer Science*, pages 433–445. Springer, 2011.

- [7] Murat Cubuktepe, Nils Jansen, Sebastian Junges, Ahmadreza Marandi, Marnix Suilen, and Ufuk Topcu. Robust finite-state controllers for uncertain POMDPs. In *AAAI*, pages 11792–11800. AAAI Press, 2021.
- [8] Erick Delage and Shie Mannor. Percentile optimization for markov decision processes with parameter uncertainty. *Oper. Res.*, 58(1):203–213, 2010.
- [9] Maxim Egorov, Zachary N. Sunberg, Edward Balaban, Tim A. Wheeler, Jayesh K. Gupta, and Mykel J. Kochenderfer. POMDPs.jl: A framework for sequential decision making under uncertainty. *JMLR*, 18(26):1–5, 2017.
- [10] Paul L. Fackler and Robert G. Haight. Monitoring as a partially observable decision problem. *Resource and Energy Economics*, 37:226–241, 2014.
- [11] Maris F. L. Galesloot, Roman Andriushchenko, Milan Češka, Sebastian Junges, and Nils Jansen. Robust finite-memory policy gradients for hidden-model POMDPs. In *IJCAI*, 2025.
- [12] Maris F. L. Galesloot, Marnix Suilen, Thiago D. Simão, Steven Carr, Matthijs T. J. Spaan, Ufuk Topcu, and Nils Jansen. Pessimistic iterative planning with RNNs for robust POMDPs. In *ECAI*, 2025.
- [13] Joel Goh, Mohsen Bayati, Stefanos A. Zenios, Sundeep Singh, and David Moore. Data uncertainty in markov chains: Application to cost-effectiveness analyses of medical innovations. *Oper. Res.*, 66(3):697–715, 2018.
- [14] Julien Grand-Clément, Nian Si, and Shengbo Wang. Tractable robust Markov decision processes, November 2024. arXiv:2411.08435 [math].
- [15] Shangding Gu, Laixi Shi, Muning Wen, Ming Jin, Eric Mazumdar, Yuejie Chi, Adam Wierman, and Costas J. Spanos. Robust gymnasium: A unified modular benchmark for robust reinforcement learning. In *ICLR*. OpenReview.net, 2025.
- [16] Eric A. Hansen. Cost-effective sensing during plan execution. In *AAAI*, pages 1029–1035. AAAI Press / The MIT Press, 1994.
- [17] Arnd Hartmanns, Sebastian Junges, Tim Quatmann, and Maximilian Weininger. The revised practitioner’s guide to MDP model checking algorithms. *International Journal on Software Tools for Technology Transfer*, 2025.
- [18] Milos Hauskrecht. Value-function approximations for partially observable Markov decision processes. *JAIR*, 13:33–94, 2000.
- [19] Hideaki Itoh and Kiyohiko Nakamura. Partially observable Markov decision processes with imprecise parameters. *Artif. Intell.*, 171(8-9):453–490, 2007.
- [20] Garud N. Iyengar. Robust dynamic programming. *Math. Oper. Res.*, 30(2):257–280, 2005.
- [21] Wha Sook Jeon, Seung Beom Seo, and Dong Geun Jeong. POMDP-based contention resolution for framed slotted-ALOHA protocol in machine-type communications. *IEEE Internet Things J.*, 9(15):13511–13523, 2022.
- [22] Leslie Pack Kaelbling, Michael L. Littman, and Anthony R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1):99–134, 1998.
- [23] Suresh Kalyanasundaram, Edwin K. P. Chong, and Ness B. Shroff. Markov decision processes with uncertain transition rates: sensitivity and robust control. In *CDC*, pages 3799–3804. IEEE, 2002.
- [24] Merlijn Krale, Thiago D. Simão, and Nils Jansen. Act-then-measure: Reinforcement learning for partially observable environments with active measuring. In *ICAPS*, pages 212–220. AAAI Press, 2023.
- [25] Merlijn Krale, Thiago D. Simão, Jana Tumova, and Nils Jansen. Robust active measuring under model uncertainty. In *AAAI*, pages 21276–21284. AAAI Press, 2024.

- [26] Merlijn Krale, Wietze Koops, Sebastian Junges, Thiago D. Simão, and Nils Jansen. Tighter value-function approximations for POMDPs. In *AAMAS*, pages 1200–1208. International Foundation for Autonomous Agents and Multiagent Systems / ACM, 2025.
- [27] Hanna Kurniawati, David Hsu, and Wee Sun Lee. SARSOP: Efficient point-based POMDP planning by approximating optimally reachable belief spaces. In *Robotics: Science and Systems (RSS)*. MIT Press, 2008.
- [28] Michael L. Littman, Anthony R. Cassandra, and Leslie P. Kaelbling. Learning policies for partially observable environments: Scaling up. In *ICML*, pages 362–370. Morgan Kaufmann, 1995.
- [29] Michael L. Littman, Thomas L. Dean, and Leslie Pack Kaelbling. On the complexity of solving Markov decision problems. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 394–402. Morgan Kaufmann, 1995.
- [30] Pablo G. Morato, Konstantinos G. Papakonstantinou, Charalampos P. Andriotis, Jannie Sønderkær Nielsen, and Philippe Rigo. Optimal inspection and maintenance planning for deteriorating structural components through dynamic Bayesian networks and Markov decision processes. *Structural Safety*, 94:102140, 2022.
- [31] Pablo G. Morato, Konstantinos G. Papakonstantinou, Charalampos P. Andriotis, and Philippe Rigo. Managing offshore wind turbines through Markov decision processes and dynamic Bayesian networks. In *13th International Conference on Structural Safety & Reliability (ICOSAR)*, 2022.
- [32] Hideaki Nakao, Ruiwei Jiang, and Siqian Shen. Distributionally robust partially observable Markov decision process with moment-based ambiguity. *SIAM J. Optim.*, 31(1):461–488, 2021.
- [33] Yaodong Ni and Zhi-Qiang Liu. Bounded-parameter partially observable Markov decision processes. In *ICAPS*, pages 240–247. AAAI, 2008.
- [34] Yaodong Ni and Zhi-Qiang Liu. Bounded-parameter partially observable Markov decision processes: Framework and algorithm. *Int. J. Uncertain. Fuzziness Knowl. Based Syst.*, 21(6): 821–864, 2013.
- [35] Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with uncertain transition matrices. *Oper. Res.*, 53(5):780–798, 2005.
- [36] Ian Osband, Yotam Doron, Matteo Hessel, John Aslanides, Eren Sezener, Andre Saraiva, Katrina McKinney, Tor Lattimore, Csaba Szepesvari, Satinder Singh, et al. Behaviour suite for reinforcement learning. In *International Conference on Learning Representations*, 2025.
- [37] Takayuki Osogami. Robust partially observable markov decision process. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 106–115. JMLR.org, 2015.
- [38] Hans Peters. *Game theory: A Multi-leveled approach*. Springer, 2015.
- [39] Joelle Pineau, Geoffrey J. Gordon, and Sebastian Thrun. Point-based value iteration: An anytime algorithm for POMDPs. In *IJCAI*, pages 1025–1032. Morgan Kaufmann, 2003.
- [40] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley Series in Probability and Statistics. Wiley, 1994.
- [41] Soroush Saghafian. Ambiguous partially observable Markov decision processes: Structural results and applications. *J. Econ. Theory*, 178:1–35, 2018.
- [42] Richard D. Smallwood and Edward J. Sondik. The optimal control of partially observable Markov processes over a finite horizon. *Oper. Res.*, 21(5):1071–1088, 1973.
- [43] Trey Smith and Reid G. Simmons. Heuristic search value iteration for POMDPs. In *UAI*, pages 520–527. AUAI Press, 2004.
- [44] Trey Smith and Reid G. Simmons. Point-based POMDP algorithms: Improved analysis and implementation. In *UAI*, pages 542–547, 2005.

- [45] Marnix Suilen, Nils Jansen, Murat Cubuktepe, and Ufuk Topcu. Robust policy synthesis for uncertain POMDPs via convex optimization. In *IJCAI*, pages 4113–4120. ijcai.org, 2020.
- [46] Marnix Suilen, Thom S. Badings, Eline M. Bovy, David Parker, and Nils Jansen. Robust Markov decision processes: A place where AI and formal methods meet. In *Principles of Verification (3)*, volume 15262 of *Lecture Notes in Computer Science*, pages 126–154. Springer, 2024.
- [47] Ole Tange. Gnu parallel 20241222 ('bashar') [stable], December 2024. URL <https://doi.org/10.5281/zenodo.14550073>.
- [48] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. Intelligent robotics and autonomous agents. MIT Press, 2005.
- [49] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov decision processes. *Math. Oper. Res.*, 38(1):153–183, 2013.
- [50] Nan Ye, Adhiraj Somani, David Hsu, and Wee Sun Lee. DESPOT: Online POMDP Planning with Regularization. *Journal of Artificial Intelligence Research*, 58:231–266, January 2017. ISSN 1076-9757. doi: 10.1613/jair.5328.
- [51] Adil Zouitine, David Bertoin, Pierre Clavier, Matthieu Geist, and Emmanuel Rachelson. RRLS: robust reinforcement learning suite. *CoRR*, abs/2406.08406, 2024.

## A Solvability Classes and Benchmark Models

This section contains all proofs related to section 3 of the main paper.

### A.1 Solvability proofs

First, we show the correctness of the policies used in section Section 3 to explain the solvability classes. We restate the RPOMDP used in section 3.1.

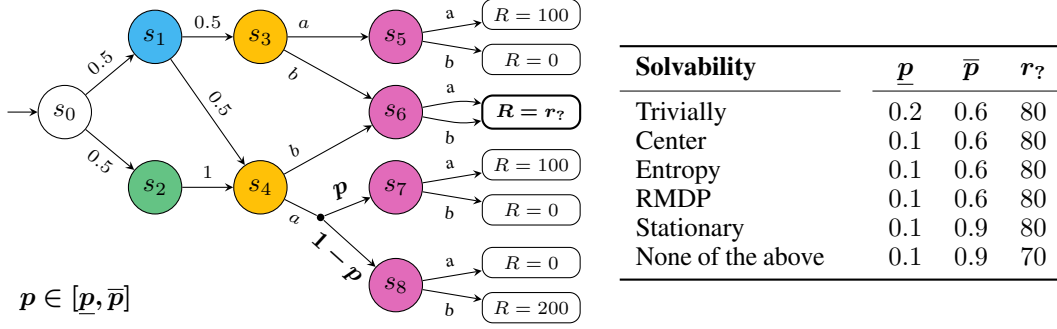


Figure 5: Example RPOMDP with parameter values for solvability classes. We use TOY\* to denote the model with the parameter values from the bottom line.

Let  $h^{\bullet\bullet} = \circ\bullet\bullet$ ,  $h^{\bullet\bullet\bullet} = \circ\bullet\bullet\bullet$ ,  $h^{\bullet\bullet\bullet\bullet} = \circ\bullet\bullet\bullet\bullet$ , and  $h^{\bullet\bullet\bullet\bullet\bullet} = \circ\bullet\bullet\bullet\bullet\bullet$  denote the histories for which the agent or nature need to make non-singleton choices. Given agent and nature policies  $\pi$  and  $\theta$ , the value function of the RPOMDP in fig. 5 can be expressed by the following function:

$$\begin{aligned}
 V^{\pi, \theta} = & 0.5 \cdot 0.5 \cdot \pi(h^{\bullet\bullet})(a) \cdot \pi(h^{\bullet\bullet\bullet\bullet})(a) \cdot 100 \\
 & + 0.5 \cdot 0.5 \cdot \pi(h^{\bullet\bullet})(b) \cdot r_? \\
 & + 0.5 \cdot 0.5 \cdot \pi(h^{\bullet\bullet\bullet})(b) \cdot r_? \\
 & + 0.5 \cdot 0.5 \cdot \pi(h^{\bullet\bullet})(a) \cdot \theta(h^{\bullet\bullet}, a)(p) \cdot \pi(h^{\bullet\bullet\bullet\bullet})(a) \cdot 100 \\
 & + 0.5 \cdot 0.5 \cdot \pi(h^{\bullet\bullet})(a) \cdot (1 - \theta(h^{\bullet\bullet}, a)(p)) \cdot \pi(h^{\bullet\bullet\bullet\bullet\bullet})(b) \cdot 200 \\
 & + 0.5 \cdot \pi(h^{\bullet\bullet\bullet})(b) \cdot r_? \\
 & + 0.5 \cdot \pi(h^{\bullet\bullet\bullet})(a) \cdot \theta(h^{\bullet\bullet\bullet}, a)(p) \cdot \pi(h^{\bullet\bullet\bullet\bullet\bullet})(a) \cdot 100 \\
 & + 0.5 \cdot \pi(h^{\bullet\bullet\bullet})(a) \cdot (1 - \theta(h^{\bullet\bullet\bullet}, a)(p)) \cdot \pi(h^{\bullet\bullet\bullet\bullet\bullet})(b) \cdot 200
 \end{aligned}$$

We simplify and rewrite this function to remove the  $b$  terms:

$$\begin{aligned}
 = & 25 \cdot \pi(h^{\bullet\bullet})(a) \cdot \pi(h^{\bullet\bullet\bullet\bullet})(a) \\
 & + 0.5 \cdot r_? \cdot (1 - \pi(h^{\bullet\bullet})(a)) \\
 & + 25 \cdot \pi(h^{\bullet\bullet})(a) \cdot \theta(h^{\bullet\bullet}, a)(p) \cdot \pi(h^{\bullet\bullet\bullet\bullet})(a) \\
 & + 50 \cdot \pi(h^{\bullet\bullet})(a) \cdot (1 - \theta(h^{\bullet\bullet}, a)(p)) \cdot (1 - \pi(h^{\bullet\bullet\bullet\bullet})(a)) \\
 & + 0.5 \cdot r_? \cdot (1 - \pi(h^{\bullet\bullet\bullet})(a)) \\
 & + 50 \cdot \pi(h^{\bullet\bullet\bullet})(a) \cdot \theta(h^{\bullet\bullet\bullet}, a)(p) \cdot \pi(h^{\bullet\bullet\bullet\bullet\bullet})(a) \\
 & + 100 \cdot \pi(h^{\bullet\bullet\bullet})(a) \cdot (1 - \theta(h^{\bullet\bullet\bullet}, a)(p)) \cdot (1 - \pi(h^{\bullet\bullet\bullet\bullet\bullet})(a))
 \end{aligned}$$

We can now further simplify the notation with  $a^\bullet = \pi(h^{\bullet\bullet})(a)$ ,  $a^\circ = \pi(h^{\bullet\circ})(a)$ ,  $a^{\bullet\circ} = \pi(h^{\bullet\circ\bullet})(a)$ ,  $a^{\circ\circ} = \pi(h^{\circ\circ\bullet})(a)$  and  $p^\bullet = \theta(h^{\bullet\bullet}, a)(p)$ ,  $p^\circ = \theta(h^{\bullet\circ}, a)(p)$ .

$$\begin{aligned}
&= 25a^\bullet a^{\bullet\circ} \\
&\quad + 0.5r_\gamma(1 - a^\bullet) \\
&\quad + 25a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + 50a^\bullet(1 - p^\bullet)(1 - a^{\bullet\circ}) \\
&\quad + 0.5r_\gamma(1 - a^\circ) \\
&\quad + 50a^\circ p^\circ a^{\circ\circ} \\
&\quad + 100a^\circ(1 - p^\circ)(1 - a^{\circ\circ}) \\
&= 25a^\bullet a^{\bullet\circ} \\
&\quad + 0.5r_\gamma - 0.5r_\gamma a^\bullet \\
&\quad + 25a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + 50a^\bullet - 50a^\bullet p^\bullet - 50a^\bullet a^{\bullet\circ} + 50a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + 0.5r_\gamma - 0.5r_\gamma a^\circ \\
&\quad + 50a^\circ p^\circ a^{\circ\circ} \\
&\quad + 100a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 100a^\circ p^\circ a^{\circ\circ} \\
&= r_\gamma + (50 - 0.5r_\gamma)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_\gamma)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ}
\end{aligned}$$

Recall the optimal value for the agent  $V^a = \max_{\pi \in \Pi} \min_{\theta \in \Theta} V^{\pi, \theta}$  and nature  $V^n = \min_{\theta \in \Theta} \max_{\pi \in \Pi} V^{\pi, \theta}$  with corresponding optimal agent and nature policies  $\pi^*$  and  $\theta^*$ . If  $V^a = V^n$ , then the underlying game has a *Nash equilibrium* [38] in which case  $\pi^*$  and  $\theta^*$  are called *Nash policies*. We can therefore show that an agent policy is optimal if it can guarantee the same value as a nature policy can guarantee, and vice versa. Using this logic, we consider the following two agent policies.

First, the *safe* policy  $\pi_s^*$  always picks the safe action  $b$  in the  $\bullet$  states, so regardless of the previous observation, which yields a guaranteed value of  $r_\gamma$ . The action in the  $\circ$  states does not matter after the safe action  $b$ .

$$\begin{aligned}
\forall \theta \in \Theta : V^{\pi_s^*, \theta} &= r_\gamma + (50 - 0.5r_\gamma) \cdot 0 - 50 \cdot 0 \cdot p^\bullet - 25 \cdot 0 \cdot a^{\bullet\circ} + 75 \cdot 0 \cdot p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_\gamma) \cdot 0 - 100 \cdot 0 \cdot p^\circ - 100 \cdot 0 \cdot a^{\circ\circ} + 150 \cdot 0 \cdot p^\circ a^{\circ\circ} \\
&= r_\gamma + 0 - 0 - 0 + 0 + 0 - 0 - 0 + 0 \\
&= r_\gamma
\end{aligned}$$

Next, the *risky* policy  $\pi_r^*$  picks the riskier action  $a$  in these  $\bullet$  states if it's previous observation was  $\circ$ , and the safe action  $b$  otherwise. Furthermore,  $\pi_r^*$  picks action  $b$  in the  $\circ$  states. The value that  $\pi_r^*$  can guarantee depends on the bounds on  $p$ .

$$\begin{aligned}
\forall \theta \in \Theta : V^{\pi_r^*, \theta} &= r_\gamma + (50 - 0.5r_\gamma) \cdot 0 - 50 \cdot 0 \cdot p^\bullet - 25 \cdot 0 \cdot a^{\bullet\circ} + 75 \cdot 0 \cdot p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_\gamma) \cdot 1 - 100 \cdot 1 \cdot p^\circ - 100 \cdot 1 \cdot 0 + 150 \cdot 1 \cdot p^\circ \cdot 0 \\
&= r_\gamma + 0 - 0 - 0 + 0 + 100 - 0.5r_\gamma - 100p^\circ - 0 + 0 \\
&= 0.5r_\gamma + 100 - 100p^\circ
\end{aligned}$$

Similarly, we can look at the value certain nature policies can guarantee.

We find that nature can ensure a value of at most  $r_\gamma$  when choosing  $p^\bullet$  and  $p^\circ$  within certain intervals.

$$\begin{aligned}
\forall \pi \in \Pi : V^{\pi, \theta} &= r_\gamma + (50 - 0.5r_\gamma)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_\gamma)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ}
\end{aligned}$$

$a^\bullet$ ,  $a^{\bullet\circ}$ , and  $p^\bullet$  are independent of  $a^\circ$ ,  $a^{\circ\circ}$ , and  $p^\circ$ , so  $\forall \pi \in \Pi : V^{\pi, \theta} \leq r_\gamma$  iff:

$$\forall \pi \in \Pi : (50 - 0.5r_\gamma)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \leq 0$$



and

$$\forall \pi \in \Pi : (100 - 0.5r_\gamma)a^\bullet - 100a^\bullet p^\bullet - 100a^\bullet a^{\bullet\bullet} + 150a^\bullet p^\bullet a^{\bullet\bullet} \leq 0$$

Given arbitrary  $a^\bullet, a^{\bullet\bullet} \in [0, 1]$ , we first compute the interval for  $p^\bullet$ :

$$\begin{aligned} (50 - 0.5r_\gamma)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} &\leq 0 \\ a^\bullet(50 - 0.5r_\gamma - 50p^\bullet - 25a^{\bullet\bullet} + 75p^\bullet a^{\bullet\bullet}) &\leq 0 \end{aligned}$$

$a^\bullet \in [0, 1]$ , so the inequality above holds for all  $\pi \in \Pi$  iff:

$$\begin{aligned} 50 - 0.5r_\gamma - 50p^\bullet - 25a^{\bullet\bullet} + 75p^\bullet a^{\bullet\bullet} &\leq 0 \\ 50 - 0.5r_\gamma - 50p^\bullet &\leq 25a^{\bullet\bullet} - 75p^\bullet a^{\bullet\bullet} \\ 50 - 0.5r_\gamma - 50p^\bullet &\leq a^{\bullet\bullet}(25 - 75p^\bullet) \end{aligned}$$

$a^{\bullet\bullet} \in [0, 1]$ , so the inequality above holds for all  $\pi \in \Pi$  iff:

$$\begin{aligned} 50 - 0.5r_\gamma - 50p^\bullet &\leq 0 \quad \wedge \quad 50 - 0.5r_\gamma - 50p^\bullet \leq 25 - 75p^\bullet \\ 50 - 0.5r_\gamma &\leq 50p^\bullet \quad \wedge \quad 25p^\bullet \leq 0.5r_\gamma - 25 \\ 1 - 0.01r_\gamma &\leq p^\bullet \quad \wedge \quad p^\bullet \leq 0.02r_\gamma - 1 \end{aligned}$$

So we get our first interval  $p^\bullet \in [1 - 0.01r_\gamma, 0.02r_\gamma - 1]$ . Given arbitrary  $a^\bullet, a^{\bullet\bullet} \in [0, 1]$ , we continue with the interval for  $p^\bullet$ :

$$\begin{aligned} (100 - 0.5r_\gamma)a^\bullet - 100a^\bullet p^\bullet - 100a^\bullet a^{\bullet\bullet} + 150a^\bullet p^\bullet a^{\bullet\bullet} &\leq 0 \\ a^\bullet(100 - 0.5r_\gamma - 100p^\bullet - 100a^{\bullet\bullet} + 150p^\bullet a^{\bullet\bullet}) &\leq 0 \end{aligned}$$

$a^\bullet \in [0, 1]$ , so the inequality above holds for all  $\pi \in \Pi$  iff:

$$\begin{aligned} 100 - 0.5r_\gamma - 100p^\bullet - 100a^{\bullet\bullet} + 150p^\bullet a^{\bullet\bullet} &\leq 0 \\ 100 - 0.5r_\gamma - 100p^\bullet &\leq 100a^{\bullet\bullet} - 150p^\bullet a^{\bullet\bullet} \\ 100 - 0.5r_\gamma - 100p^\bullet &\leq a^{\bullet\bullet}(100 - 150p^\bullet) \end{aligned}$$

$a^{\bullet\bullet} \in [0, 1]$ , so the inequality above holds for all  $\pi \in \Pi$  iff:

$$\begin{aligned} 100 - 0.5r_\gamma - 100p^\bullet &\leq 0 \quad \wedge \quad 100 - 0.5r_\gamma - 100p^\bullet \leq 100 - 150p^\bullet \\ 100 - 0.5r_\gamma &\leq 100p^\bullet \quad \wedge \quad 50p^\bullet \leq 0.5r_\gamma \\ 1 - 0.005r_\gamma &\leq p^\bullet \quad \wedge \quad p^\bullet \leq 0.01r_\gamma \end{aligned}$$

So we get our second interval  $p^\bullet \in [1 - 0.005r_\gamma, 0.01r_\gamma]$ . We can conclude that  $\forall \pi \in \Pi : V^{\pi, \theta} \leq r_\gamma \iff p^\bullet \in [1 - 0.01r_\gamma, 0.02r_\gamma - 1] \wedge p^\bullet \in [1 - 0.005r_\gamma, 0.01r_\gamma]$ . We consider two values for  $r_\gamma$  in our example, i. e., 70 and 80, so we get the following intervals for  $p^\bullet$  and  $p^\bullet$ :

$$\begin{aligned} r_\gamma = 70 : p^\bullet &\in [0.3, 0.4], p^\bullet \in [0.65, 0.7] \\ r_\gamma = 80 : p^\bullet &\in [0.2, 0.6], p^\bullet \in [0.6, 0.8] \end{aligned}$$

### A.1.1 Trivially solvable

Since the trivially solvable model in fig. 5 has  $r_\gamma = 80$  and bounds  $p \in [0.2, 0.6]$ , we know that nature can guarantee a value of 80 by playing  $p^\bullet = p^\bullet = 0.6$ . We therefore know that any agent policy that can guarantee a value of 80 is optimal as well. In particular, we know that  $\pi_r^*$  is optimal, since:

$$\begin{aligned} \forall \theta \in \Theta : V^{\pi_r^*, \theta} &= 0.5r_\gamma + 100 - 100p^\bullet \\ &= 0.5 \cdot 80 + 100 - 100p^\bullet \\ &= 140 - 100p^\bullet \\ &\geq 140 - 100 \cdot 0.6 \\ &= 80 \end{aligned}$$

To show that this model is trivially solvable, we show that  $\pi_r^*$  is optimal for any  $\theta \in \Theta$ . In other words, for all nature policies in the set of nature policies, there is no agent policy  $\pi'$  that achieves a higher reward than  $\pi_r^*$  against that particular nature policies.

Given an arbitrary  $\theta \in \Theta$ , we construct the optimal agent policy  $\pi'$  by maximizing over the value function:

$$\begin{aligned}
\pi' &= \operatorname{argmax}_{\pi \in \Pi} V^{\pi, \theta} \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( r_? + (50 - 0.5r_?)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + (100 - 0.5r_?)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + (50 - 0.5 \cdot 80)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + (100 - 0.5 \cdot 80)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 10a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + 60a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 10a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + 60a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right)
\end{aligned}$$

$a^\bullet, a^{\bullet\bullet}$ , and  $p^\bullet$  are independent of  $a^\circ, a^{\circ\circ}$ , and  $p^\circ$ , so we can compute two parts of the agent policy separately:

$$\begin{aligned}
&= \operatorname{argmax}_{a^\bullet, a^{\bullet\bullet}} \left( 10a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right) \\
&\quad \times \operatorname{argmax}_{a^\circ, a^{\circ\circ}} \left( 60a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right)
\end{aligned}$$

We begin with  $a^\bullet$  and  $a^{\bullet\bullet}$ :

$$\operatorname{argmax}_{a^\bullet, a^{\bullet\bullet}} \left( 10a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right) = \operatorname{argmax}_{a^\bullet, a^{\bullet\bullet}} \left( a^\bullet (10 - 50p^\bullet - a^{\bullet\bullet} (25 - 75p^\bullet)) \right)$$

As we are maximizing, we know that  $a^\bullet$  should be 1 if  $10 - 50p^\bullet - a^{\bullet\bullet} (25 - 75p^\bullet) > 0$  and we can set it to 0 otherwise. We show  $10 - 50p^\bullet - a^{\bullet\bullet} (25 - 75p^\bullet) \leq 0$  regardless of  $a^{\bullet\bullet}$ :

$$\begin{aligned}
10 - 50p^\bullet - a^{\bullet\bullet} (25 - 75p^\bullet) &\leq 0 \\
10 - 50p^\bullet &\leq a^{\bullet\bullet} (25 - 75p^\bullet)
\end{aligned}$$

$a^{\bullet\bullet} \in [0, 1]$ , so the inequality above holds for all  $a^{\bullet\bullet}$  iff:

$$\begin{aligned}
10 - 50p^\bullet &\leq 0 \quad \wedge \quad 10 - 50p^\bullet \leq 25 - 75p^\bullet \\
10 &\leq 50p^\bullet \quad \wedge \quad 25p^\bullet \leq 15 \\
0.2 &\leq p^\bullet \quad \wedge \quad p^\bullet \leq 0.6
\end{aligned}$$

Since  $\forall \theta \in \Theta : p^\bullet \in [0.2, 0.6]$ , we have that choosing  $a^\bullet = 0$  is optimal.

We continue with  $a^\circ$  and  $a^{\circ\circ}$ :

$$\operatorname{argmax}_{a^\circ, a^{\circ\circ}} \left( 60a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) = \operatorname{argmax}_{a^\circ, a^{\circ\circ}} \left( a^\circ (60 - 100p^\circ - a^{\circ\circ} (100 - 150p^\circ)) \right)$$

As we are maximizing, we know that  $a^\circ$  should be 1 if  $60 - 100p^\circ - a^{\circ\circ} (100 - 150p^\circ) \geq 0$  and we can set it to 0 otherwise. We know  $60 - 100p^\circ \geq 0$  regardless of  $p^\circ$ :

$$\begin{aligned}
60 - 100p^\circ &\geq 60 - 100 \cdot 0.6 \\
&= 0
\end{aligned}$$

Since we can set  $a^{\text{blue}}_{\text{red}}^{\text{blue}}$  to 0, we already know that  $60 - 100p^{\text{blue}} - a^{\text{blue}}_{\text{red}}^{\text{blue}}(100 - 150p^{\text{blue}}) \geq 0$ , so we set  $a^{\text{blue}}_{\text{red}}^{\text{blue}}$  to 1. Finally, we determine the optimal  $a^{\text{blue}}_{\text{red}}^{\text{blue}}$ .

$$\begin{aligned} a^{\text{blue}}_{\text{red}}^{\text{blue}}(100 - 150p^{\text{blue}}) &\geq a^{\text{blue}}_{\text{red}}^{\text{blue}}(100 - 150 \cdot 0.6) \\ &= a^{\text{blue}}_{\text{red}}^{\text{blue}} \cdot 0 \\ &= 0 \end{aligned}$$

Since we subtract with  $a^{\text{blue}}_{\text{red}}^{\text{blue}}(100 - 150p^{\text{blue}}) \geq 0$ , it is optimal to set  $a^{\text{blue}}_{\text{red}}^{\text{blue}}$  to 0.

The agent policy that we constructed  $\pi'$  with  $a^{\text{blue}} = 0$ ,  $a^{\text{green}} = 1$ , and  $a^{\text{blue}}_{\text{red}}^{\text{blue}} = 0$  is optimal against all  $\theta \in \Theta$ , and this policy is exactly  $\pi_r^*$ . Since the agent policy we compute against any nature policy  $\theta \in \Theta$  is optimal against the entire set of nature policies  $\Theta$ , we conclude that our trivially solvable model in fig. 5 is indeed trivially solvable.

### A.1.2 Naively solvable

Like in the trivially solvable model in fig. 5, nature can guarantee a value of 80 in the center, entropy, and RMDP solvable models in fig. 5 by playing  $p^{\text{blue}} = p^{\text{green}} = 0.6$ . We therefore also know  $\pi_r^*$  is again optimal, since it can guarantee a value of 80.

The naive nature policies in the center, entropy, and RMDP solvable models are:

- $\theta_{\text{Center}}$  assigns  $p^{\text{blue}} = p^{\text{green}} = 0.35$ , as Centroid ( $p$ ) = 0.35.
- $\theta_{\text{Ent}}$  assigns  $p^{\text{blue}} = p^{\text{green}} = 0.5$ , as this creates the maximum entropy between states  $s_7$  and  $s_8$ , namely  $\{s_7 \mapsto 0.5, s_8 \mapsto 0.5\}$ .
- $\theta_{\text{RMDP}}$  assigns  $p^{\text{blue}} = p^{\text{green}} = 0.6$ , as in the fully observable model, it is always optimal for nature to minimize the chance of reaching state  $s_8$ , and therefore to maximize  $p$ .

All three of these naive nature policies are contained in the set of policies of the trivially solvable model, for which we have shown that  $\pi_r^*$  is optimal.

Since the agent policies we compute against  $\theta_{\text{Center}}$ ,  $\theta_{\text{Ent}}$ , and  $\theta_{\text{RMDP}}$  are optimal against the entire set of nature policies  $\Theta$ , we conclude that our center, entropy and RMDP solvable models in fig. 5 are indeed center, entropy, and RMDP solvable.

We also note that the center, entropy, and RMDP solvable models in fig. 5 are not trivially solvable. For example,  $\pi_r^*$  is not optimal against the nature policy  $\theta'$  with  $p^{\text{blue}} = p^{\text{green}} = 0.1$ . We again construct the optimal agent policy  $\pi'$  by maximizing over the value function:

$$\begin{aligned} \pi' &= \operatorname{argmax}_{\pi \in \Pi} V^{\pi, \theta'} \\ &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 10a^{\text{blue}} - 50a^{\text{blue}} \cdot 0.1 - 25a^{\text{blue}} a^{\text{blue}}_{\text{red}}^{\text{blue}} + 75a^{\text{blue}} \cdot 0.1 \cdot a^{\text{blue}}_{\text{red}}^{\text{blue}} \right. \\ &\quad \left. + 60a^{\text{green}} - 100a^{\text{green}} \cdot 0.1 - 100a^{\text{green}} a^{\text{blue}}_{\text{red}}^{\text{blue}} + 150a^{\text{green}} \cdot 0.1 \cdot a^{\text{blue}}_{\text{red}}^{\text{blue}} \right) \\ &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 10a^{\text{blue}} - 5a^{\text{blue}} - 25a^{\text{blue}} a^{\text{blue}}_{\text{red}}^{\text{blue}} + 7.5a^{\text{blue}} a^{\text{blue}}_{\text{red}}^{\text{blue}} \right. \\ &\quad \left. + 60a^{\text{green}} - 10a^{\text{green}} - 100a^{\text{green}} a^{\text{blue}}_{\text{red}}^{\text{blue}} + 15a^{\text{green}} a^{\text{blue}}_{\text{red}}^{\text{blue}} \right) \\ &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 5a^{\text{blue}} - 17.5a^{\text{blue}} a^{\text{blue}}_{\text{red}}^{\text{blue}} + 50a^{\text{green}} - 85a^{\text{green}} a^{\text{blue}}_{\text{red}}^{\text{blue}} \right) \end{aligned}$$

Since we are maximizing, it is optimal to set  $a^{\text{blue}}$  and  $a^{\text{green}}$  to 1 and  $a^{\text{blue}}_{\text{red}}^{\text{blue}}$  and  $a^{\text{green}}_{\text{red}}^{\text{blue}}$  to 0. This  $\pi'$  results in value  $V^{\pi', \theta'} = 80 + 5 + 50 = 135$ , whereas  $\pi_r^*$  results in value  $V^{\pi_r^*, \theta'} = 80 + 50 = 130$ .

Since the optimal agent policy  $\pi'$  for nature policy  $\theta' \in \Theta$  is not optimal against the entire set of nature policies  $\Theta$ , the center, entropy, and RMDP solvable models in fig. 5 are not trivially solvable.

### A.1.3 Stationary solvable

Since the stationary solvable model in fig. 5 has  $r_{\gamma} = 80$  and bounds  $p \in [0.1, 0.9]$ , we know that nature can guarantee a value of 80 by playing any policy with  $p^{\text{blue}} \in [0.2, 0.6]$  and  $p^{\text{green}} \in [0.6, 0.8]$ . We therefore know that any agent policy that can guarantee a value of 80 is optimal as well. In particular, we know that  $\pi_s^*$  is optimal, since this policy always results in a value of  $r_{\gamma} = 80$ .

To show that this model is stationary solvable, we show that  $\pi_s^*$  is also optimal against the set of stationary nature policies  $\Theta^{\text{Sta}}$ . Since  $\pi_s^*$  guarantees a value of 80, this agent policy is optimal if there is stationary nature policy that can also guarantee a value of 80. We know that the stationary nature policy  $p^\bullet = p^\circ = 0.6$  guarantees a value of 80, so we can conclude that the stationary solvable model in fig. 5 is indeed stationary solvable.

We also note that the stationary solvable model in fig. 5 is not trivially or naively solvable, as the optimal agent policies against the nature policies  $\theta_{\text{Center}}$ ,  $\theta_{\text{Ent}}$ , and  $\theta_{\text{RMDP}}$  cannot guarantee a value of 80 against the entire set of nature policies.

The naive policies in the stationary solvable models are:

- $\theta_{\text{Center}}$  assigns  $p^\bullet = p^\circ = 0.5$ , as  $\text{Centroid}(p) = 0.5$ .
- $\theta_{\text{Ent}}$  assigns  $p^\bullet = p^\circ = 0.5$ , as this creates the maximum entropy between states  $s_7$  and  $s_8$ , namely  $\{s_7 \mapsto 0.5, s_8 \mapsto 0.5\}$ .
- $\theta_{\text{RMDP}}$  assigns  $p^\bullet = p^\circ = 0.9$ , as in the fully observable model, it is always optimal for nature to minimize the chance of reaching state  $s_8$ , and therefore to maximize  $p$ .

As shown in appendix A.1.2, we know that  $\pi_r^*$  is optimal against  $\theta_{\text{Center}}$  and  $\theta_{\text{Ent}}$ . However,  $\pi_r^*$  is not optimal against the entire set of nature policies, as nature can achieve a value  $< 80$  when playing a nature policy  $\theta'$  with  $p^\circ > 0.6$ :

$$\begin{aligned}
 V^{\pi_r^*, \theta'} &= 0.5 \cdot r_? + 100 - 100p^\circ \\
 &= 0.5 \cdot 80 + 100 - 100p^\circ \\
 &= 140 - 100p^\circ \\
 &< 140 - 100 \cdot 0.6 \\
 &= 140 - 60 \\
 &= 80
 \end{aligned}$$

Next, we show the stationary solvable model is not RMDP solvable by constructing the optimal agent policy  $\pi'$  against  $\theta_{\text{RMDP}}$  and showing that this agent policy cannot guarantee a value of 80 against the entire set of nature policies  $\Theta$ .

$$\begin{aligned}
 \pi' &= \operatorname{argmax}_{\pi \in \Pi} V^{\pi, \theta_{\text{RMDP}}} \\
 &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 10a^\bullet - 50a^\bullet \cdot 0.9 - 25a^\bullet a^{\bullet\circ} + 75a^\bullet \cdot 0.9 \cdot a^{\circ\circ} \right. \\
 &\quad \left. + 60a^\circ - 100a^\circ \cdot 0.9 - 100a^\circ a^{\bullet\circ} + 150a^\circ \cdot 0.9 \cdot a^{\circ\circ} \right) \\
 &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 + 10a^\bullet - 45a^\bullet - 25a^\bullet a^{\bullet\circ} + 67.5a^\bullet a^{\circ\circ} \right. \\
 &\quad \left. + 60a^\circ - 90a^\circ - 100a^\circ a^{\bullet\circ} + 135a^\circ a^{\circ\circ} \right) \\
 &= \operatorname{argmax}_{\pi \in \Pi} \left( 80 - 35a^\bullet + 42.5a^\bullet a^{\bullet\circ} + -30a^\circ + 35a^\circ a^{\circ\circ} \right)
 \end{aligned}$$

Since we are maximizing, it is optimal to set  $a^\bullet, a^\circ, a^{\bullet\circ}$ , and  $a^{\circ\circ}$  all to 1. This  $\pi'$  results in value  $V^{\pi', \theta_{\text{RMDP}}} = 80 - 35 + 42.5 - 30 + 35 = 92.5$ , whereas  $\pi_s^*$  results in value  $V^{\pi_s^*, \theta_{\text{RMDP}}} = 80$ . However,  $\pi'$  is not optimal against the entire set of nature policies  $\Theta$ , as nature can achieve a value of  $< 80$

when playing a nature policy  $\theta'$  with  $p^\bullet + 2p^\circ < 1.8$ :

$$\begin{aligned}
V^{\pi', \theta'} &= 80 + 10a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + 60a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 80 + 10 \cdot 1 - 50 \cdot 1 \cdot p^\bullet - 25 \cdot 1 \cdot 1 + 75 \cdot 1 \cdot p^\bullet \cdot 1 \\
&\quad + 60 \cdot 1 - 100 \cdot 1 \cdot p^\circ - 100 \cdot 1 \cdot 1 + 150 \cdot 1 \cdot p^\circ \cdot 1 \\
&= 80 + 10 - 50p^\bullet - 25 + 75p^\bullet + 60 - 100p^\circ - 100 + 150p^\circ \\
&= 35 + 25p^\bullet + 50p^\circ \\
&= 35 + 25(p^\bullet + 2p^\circ) \\
&< 35 + 25 \cdot 1.8 \\
&= 35 + 45 \\
&= 80
\end{aligned}$$

Since the stationary solvable model in fig. 5 is not center, entropy, or RMDP solvable, we can conclude it is not trivially or naively solvable.

#### A.1.4 Not stationary solvable

Finally, we show that when we change  $r_?$  to 70 and keep the bounds on  $p \in [0.1, 0.9]$  the same for the stationary solvable model in fig. 5, we end up with a model (the *None of the above* model in fig. 5) that is neither stationary, nor naively, nor trivially solvable. We know that nature can guarantee a value of  $r_? = 70$  by playing  $p^\bullet \in [0.3, 0.4]$  and  $p^\circ \in [0.65, 0.7]$ . We therefore know that any agent policy that can guarantee a value of 70 is optimal as well. In particular, we know that  $\pi_s^*$  is optimal, since this policy always results in a value of  $r_? = 70$ .

We again identify the three naive nature policies:

- $\theta_{\text{Center}}$  assigns  $p^\bullet = p^\circ = 0.5$ , as  $\text{Centroid}(p) = 0.5$ .
- $\theta_{\text{Ent}}$  assigns  $p^\bullet = p^\circ = 0.5$ , as this creates the maximum entropy between states  $s_7$  and  $s_8$ , namely  $\{s_7 \mapsto 0.5, s_8 \mapsto 0.5\}$ .
- $\theta_{\text{RMDP}}$  assigns  $p^\bullet = p^\circ = 0.9$ , as in the fully observable model, it is always optimal for nature to minimize the chance of reaching state  $s_8$ , and therefore to maximize  $p$ .

We first construct the optimal agent policy  $\pi'$  against the center and entropy nature policies, which are the same:

$$\begin{aligned}
\pi' &= \operatorname{argmax}_{\pi \in \Pi} V^{\pi, \theta_{\text{Center}}} \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( r_? + (50 - 0.5r_?)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \right. \\
&\quad \left. + (100 - 0.5r_?)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 + (50 - 0.5 \cdot 70)a^\bullet - 50a^\bullet \cdot 0.5 - 25a^\bullet a^{\bullet\circ} + 75a^\bullet \cdot 0.5 \cdot a^{\bullet\circ} \right. \\
&\quad \left. + (100 - 0.5 \cdot 70)a^\circ - 100a^\circ \cdot 0.5 - 100a^\circ a^{\circ\circ} + 150a^\circ \cdot 0.5 \cdot a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 + 15a^\bullet - 25a^\bullet - 25a^\bullet a^{\bullet\circ} + 37.5a^\bullet a^{\bullet\circ} \right. \\
&\quad \left. + 65a^\circ - 50a^\circ - 100a^\circ a^{\circ\circ} + 75a^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 - 10a^\bullet + 12.5a^\bullet a^{\bullet\circ} + 15a^\circ - 25a^\circ a^{\circ\circ} \right)
\end{aligned}$$

Since we are maximizing, it is optimal to set  $a^\bullet, a^\circ$ , and  $a^{\bullet\circ}$  to 1 and  $a^{\circ\circ}$  to 0. This  $\pi'$  results in value  $V^{\pi', \theta_{\text{Center}}} = 70 - 10 + 12.5 + 15 - 0 = 87.5$ , whereas  $\pi_s^*$  results in value  $V^{\pi_s^*, \theta_{\text{Center}}} = 70$ . However,  $\pi'$  is not optimal against the entire set of nature policies  $\Theta$ , as nature can achieve a value of  $< 70$

when playing a nature policy  $\theta'$  with  $4p^\bullet - p^\circ > 2.2$ :

$$\begin{aligned}
V^{\pi', \theta'} &= 70 + 15a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \\
&\quad + 65a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + 15 \cdot 1 - 50 \cdot 1 \cdot p^\bullet - 25 \cdot 1 \cdot 1 + 75 \cdot 1 \cdot p^\bullet \cdot 1 \\
&\quad + 65 \cdot 1 - 100 \cdot 1 \cdot p^\circ - 100 \cdot 1 \cdot 0 + 150 \cdot 1 \cdot p^\circ \cdot 0 \\
&= 70 + 15 - 50p^\bullet - 25 + 75p^\bullet + 65 - 100p^\circ - 0 + 0 \\
&= 125 + 25p^\bullet - 100p^\circ \\
&= 125 - 25(4p^\bullet - p^\circ) \\
&< 125 - 25 \cdot 2.2 \\
&= 125 - 55 \\
&= 70
\end{aligned}$$

The optimal agent policy  $\pi'$  against the center and entropy policies cannot guarantee a value of 70 against the entire set of nature policies  $\Theta$ , hence we can conclude that  $\pi'$  is not an optimal agent policy in the *None of the above* model in fig. 5 and that this model is not center or entropy solvable, nor trivially solvable.

We continue with the RMDP nature policy  $\theta_{\text{RMDP}}$ . We again construct the optimal agent policy  $\pi'$ :

$$\begin{aligned}
\pi' &= \operatorname{argmax}_{\pi \in \Pi} V^{\pi, \theta_{\text{RMDP}}} \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( r_? + (50 - 0.5r_?)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + (100 - 0.5r_?)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 + (50 - 0.5 \cdot 70)a^\bullet - 50a^\bullet \cdot 0.9 - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet \cdot 0.9 \cdot a^{\bullet\bullet} \right. \\
&\quad \left. + (100 - 0.5 \cdot 70)a^\circ - 100a^\circ \cdot 0.9 - 100a^\circ a^{\circ\circ} + 150a^\circ \cdot 0.9 \cdot a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 + 15a^\bullet - 45a^\bullet - 25a^\bullet a^{\bullet\bullet} + 67.5a^\bullet a^{\bullet\bullet} \right. \\
&\quad \left. + 65a^\circ - 90a^\circ - 100a^\circ a^{\circ\circ} + 135a^\circ a^{\circ\circ} \right) \\
&= \operatorname{argmax}_{\pi \in \Pi} \left( 70 - 30a^\bullet + 42.5a^\bullet a^{\bullet\bullet} - 25a^\circ + 35a^\circ a^{\circ\circ} \right)
\end{aligned}$$

Since we are maximizing, it is optimal to set  $a^\bullet, a^\circ, a^{\bullet\bullet}$ , and  $a^{\circ\circ}$  all to 1. This  $\pi'$  results in value  $V^{\pi', \theta_{\text{Center}}} = 70 - 30 + 42.5 - 25 + 35 = 92.5$ , whereas  $\pi_s^*$  results in value  $V^{\pi_s^*, \theta_{\text{Center}}} = 70$ . However,  $\pi'$  is not optimal against the entire set of nature policies  $\Theta$ , as nature can achieve a value of  $< 70$  when playing a nature policy  $\theta'$  with  $p^\bullet + 2p^\circ < 1.8$ :

$$\begin{aligned}
V^{\pi', \theta'} &= 70 + 15a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\bullet} + 75a^\bullet p^\bullet a^{\bullet\bullet} \\
&\quad + 65a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + 15 \cdot 1 - 50 \cdot 1 \cdot p^\bullet - 25 \cdot 1 \cdot 1 + 75 \cdot 1 \cdot p^\bullet \cdot 1 \\
&\quad + 65 \cdot 1 - 100 \cdot 1 \cdot p^\circ - 100 \cdot 1 \cdot 1 + 150 \cdot 1 \cdot p^\circ \cdot 1 \\
&= 70 + 15 - 50p^\bullet - 25 + 75p^\bullet + 65 - 100p^\circ - 100 + 150p^\circ \\
&= 25 + 25p^\bullet + 50p^\circ \\
&= 25 + 25(p^\bullet + 2p^\circ) \\
&< 25 + 25 \cdot 1.8 \\
&= 25 + 45 \\
&= 70
\end{aligned}$$

The optimal agent policy  $\pi'$  against the RMDP nature policies cannot guarantee a value of 70 against the entire set of nature policies  $\Theta$ , hence we can conclude that  $\pi'$  is not an optimal agent policy in the *None of the above* model in fig. 5 and that this model is RMDP solvable.

Finally, we show the *None of the above* model in fig. 5 is not stationary solvable. We therefore construct an agent policy  $\pi'$  that guarantees a value of 75 against the set of stationary nature policies, which is better than the value  $\pi_s^*$  achieves, therefore we can conclude that  $\pi_s^*$  is not optimal against the set of stationary nature policies.

Let  $\pi'$  assign 1 to  $a^\bullet$ ,  $a^\circ$ , and  $a^{\bullet\circ}$ , and assign 0.5 to  $a^{\circ\circ}$ . We get the following value:

$$\begin{aligned}
\forall \theta \in \Theta^{\text{Sta}} : V^{\pi', \theta} &= r_? + (50 - 0.5r_?)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_?)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + (50 - 0.5 \cdot 70) \cdot 1 - 50 \cdot 1 \cdot p^\bullet - 25 \cdot 1 \cdot 1 + 75 \cdot 1 \cdot p^\bullet \cdot 1 \\
&\quad + (100 - 0.5 \cdot 70) \cdot 1 - 100 \cdot 1 \cdot p^\circ - 100 \cdot 1 \cdot 0.5 + 150 \cdot 1 \cdot p^\circ \cdot 0.5 \\
&= 70 + 15 - 50p^\bullet - 25 + 75 \cdot p^\bullet \\
&\quad + 65 - 100 \cdot p^\circ - 50 + 75p^\circ \\
&= 75 + 25p^\bullet - 25p^\circ \\
&= 75
\end{aligned}$$

Where the last step follows from the restriction to stationary nature policies.  $\pi'$  hence guarantees a higher value than  $\pi_s^*$  against the set of stationary nature policies, so  $\pi_s^*$  is not optimal against the set of stationary nature policies.

Next we show that  $\pi_s^*$  is the only agent policy that can guarantee a value of 70 against the entire set of nature policies. Let  $\pi''$  be an agent policy with  $a^\bullet, a^\circ \in (0, 1]$ , then we get:

$$\begin{aligned}
\forall \theta \in \Theta : V^{\pi'', \theta} &= r_? + (50 - 0.5r_?)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5r_?)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + (50 - 0.5 \cdot 70)a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + (100 - 0.5 \cdot 70)a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + 15a^\bullet - 50a^\bullet p^\bullet - 25a^\bullet a^{\bullet\circ} + 75a^\bullet p^\bullet a^{\bullet\circ} \\
&\quad + 65a^\circ - 100a^\circ p^\circ - 100a^\circ a^{\circ\circ} + 150a^\circ p^\circ a^{\circ\circ} \\
&= 70 + a^\bullet(15 - 50p^\bullet - 25a^{\bullet\circ} + 75p^\bullet a^{\bullet\circ}) \\
&\quad + a^\circ(65 - 100p^\circ - 100a^{\circ\circ} + 150p^\circ a^{\circ\circ})
\end{aligned}$$

Since  $a^\bullet, a^\circ > 0$ , the agent can only guarantee a of 70 if at least  $\forall \theta \in \Theta$ :

$$15 - 50p^\bullet - 25a^{\bullet\circ} + 75p^\bullet a^{\bullet\circ} \geq 0 \vee 65 - 100p^\circ - 100a^{\circ\circ} + 150p^\circ a^{\circ\circ} \geq 0$$

However, when  $p^\circ \in (0.3, 0.4)$  and  $p^\bullet \in (0.65, 0.7)$  this requirement does not hold, for example when  $p^\bullet = 0.35$  and  $p^\circ = \frac{2}{3}$ :

$$\begin{aligned}
15 - 50p^\bullet - 25a^{\bullet\circ} + 75p^\bullet a^{\bullet\circ} &= 15 - 50 \cdot 0.35 - 25a^{\bullet\circ} + 75 \cdot 0.35 \cdot a^{\bullet\circ} \\
&= 15 - 17.5 - 25a^{\bullet\circ} + 26.25 \cdot a^{\bullet\circ} \\
&= -2.5 + 1.25a^{\bullet\circ} \\
&< 0
\end{aligned}$$

Where the last step follows from  $a^{\bullet\circ} \in [0, 1]$ . And similarly:

$$\begin{aligned}
65 - 100p^\circ - 100a^{\circ\circ} + 150p^\circ a^{\circ\circ} &= 65 - 100 \cdot \frac{2}{3} - 100a^{\circ\circ} + 150 \cdot \frac{2}{3} \cdot a^{\circ\circ} \\
&= 65 - 66\frac{2}{3} - 100a^{\circ\circ} + 100 \cdot a^{\circ\circ} \\
&= -1\frac{2}{3} \\
&< 0
\end{aligned}$$

Hence, an agent policy can only guarantee a value of 70 by playing  $a^\bullet, a^\circ = 0$

Since  $\pi_s^*$  cannot be found by considering the subset of stationary nature policies, while it is the only optimal policy against the entire set of nature policies, we can conclude that the *None of the above* model in fig. 5 is not stationary solvable.



## A.2 Benchmarks

Next, we provide proofs for the theorems related to our proposed benchmarks.

### A.2.1 Echo machine

Let us first restate the theorem in the main text:

**Theorem 1.** ECHO, with  $p = 0.01$ ,  $\bar{p} = 0.99$ ,  $\delta = 0.1$  and  $\gamma = 0.95$ , is not in any of the solvability classes defined in Section 3.1.

To prove the theorem in the main text, we first prove two helping lemmas:

**Lemma 1.** If  $0.5 \in \mathcal{P}$ , then the optimal adversarial stationary policy  $\theta_S^*$  picks  $p = 0.5$ .

*Proof.* Let  $\Theta_{\text{delayed}}$  denote the set of nature policies that is not stationary, but for which the choice of  $p$  only applies with a delay of 1 timestep, which means each  $\theta \in \Theta_{\text{delayed}}$  must pick  $p$  for step  $t$  based on history  $h_{t-1} = (\dots, a_{t-2}, o_t)$  and previous state  $s_{1-t}$ . We show that the Nash-optimal policy in this class always picks 0.5, which means it is stationary. Moreover, since  $\Theta_{\text{delayed}}$  contains all stationary policies, this proves this policy is also the Nash-optimal stationary policy.

To do so, we first notice that the choice of any policy  $\theta \in \Theta_{\text{delayed}}$  only has an impact at histories where the agent is in state  $x'$  but the agent has some belief over  $x$  and  $x'$ , or similarly with  $y'$  and  $y$ . Let  $h_t$  be such a history, and assume  $s_t = x'$ . In that case, denote  $\theta(h_t, x) = p$ ,  $h_{a \in \{x, y\}} = (h_t, a, \perp)$  and  $h_{a \in \{x, y\}, i \in \{x, y\}} = (h_t, a, \perp, g, o_i)$ . Furthermore, let  $V^\pi(h, s)$  denote the value of policy  $\pi$  against  $\theta$  given history  $h$  and state  $s$ . In that case, the agent has the following value function for history  $h_t$

$$\begin{aligned} V^\pi(h_t, x') &= \frac{1}{2}\pi(x | h_t)\pi(r | h_x)\left(V^\pi((h_x, r, \perp), n_x) + V^\pi((h_x, r, \perp), n_y)\right) \\ &\quad + \frac{1}{2}\pi(y | h_t)\pi(r | h_y)\left(V^\pi((h_y, r, \perp), n_x) + V^\pi((h_y, r, \perp), x_y)\right) \\ &\quad + \pi(x | h)\pi(g | h_x)\left[pV^\pi(h_{xx}, x') + (1-p)V^\pi(h_{xy}, y')\right] \\ &\quad + \pi(y | h)\pi(g | h_y)\left[pV^\pi(h_{yx}, x') + (1-p)V^\pi(h_{yy}, y')\right] \end{aligned}$$

We simplify this formula in three ways. Firstly, we collect all values that do not directly depend on the choice of  $p$  (i.e., the top two lines) into a constant  $C^\pi$ . Secondly, using symmetry of the model, we define  $V(h_{xx}, x') = V(h_{yy}, y') := V_{\text{same}}$  and  $V(h_{xy}, y') = V(h_{yx}, x') := V_{\text{diff}}$ . Lastly, we denote  $\pi(x | h_t) = \pi_x$ , and since no new information is contained in  $h_x$  or  $h_y$  we may assume  $\pi(g | h_x) = \pi(g | h_y) = \pi_g$ . Then, our formula simplifies as follows:

$$V^\pi(h_t, x') = C^\pi + \pi_g\left(\pi_x(pV_{\text{same}} + (1-p)V_{\text{diff}}) + (1-\pi_x)((1-p)V_{\text{same}} + pV_{\text{diff}})\right)$$

We compute the partial derivatives of this function (with constant factor  $\pi_g$  remove for readability) as:

$$\begin{aligned} \frac{\partial V^\pi(h_t, x')}{\partial \pi_x} &= pV_{\text{same}} + (1-p)V_{\text{diff}} - (1-p)V_{\text{same}} - pV_{\text{diff}} \\ &= (1-2p)(V_{\text{diff}} - V_{\text{same}}) \\ \frac{\partial V^\pi(h_t, x')}{\partial p} &= \pi_x V_{\text{same}} - \pi_x V_{\text{diff}} - (1-\pi_x)(V_{\text{diff}} - V_{\text{same}}) \\ &= (1-2\pi_x)(V_{\text{diff}} - V_{\text{same}}) \end{aligned}$$

We find that there exists a saddle point at  $\pi_x = p = 0.5$ . Thus, for any history  $h_t$  leading to  $x'$ , we find that  $\theta_{\text{delayed}}^*(h_t, x) = 0.5$ , and the same holds for  $y'$  via symmetry of the model. Lastly, we note that the choice of  $p$  is only relevant in these two states, which means the choice of  $p = 0.5$  (or any other arbitrary choice) is optimal for other states. Thus, always picking  $p = 0.5$  is Nash-optimal for our policy class  $\Theta_{\text{delayed}}$ , which proves our lemma.  $\square$

**Lemma 2.** For any history  $h_t$ , let  $\tau_a, a \in A$  denote the highest  $t$  at which the agent has picked action  $a$ . Then, the following history-based nature policy is optimal:

$$\theta^*(h_t, s) = \begin{cases} \sup(\mathcal{P}) & \text{if } \tau_x > \tau_y \\ \inf(\mathcal{P}) & \text{otherwise.} \end{cases} \quad (4)$$

**Corollary.** Let  $h_t$  denote some history such that the  $s_t \in \{x, y, x', y'\}$ . Then, the optimal agent policy is as follows:

$$\pi^*(h_t) = \begin{cases} x & \text{if } 1 - \sup(\mathcal{P}) > \inf(\mathcal{P}) \\ y & \text{otherwise} \end{cases}$$

*Proof.* The choice of  $\theta$  only matters if  $s = n'$ , in which case either  $\tau_x = t - 1$  or  $\tau_y = t - 1$ . We can denote the two value function for these cases as follows:

$$\begin{aligned} V^\pi(h_x, n') &= \pi(r \mid h_x)C^\pi + \gamma\pi(g \mid h_x)(pV^\pi(h_{x,x}, x') + (1-p)V^\pi(h_{x,y}, y')) \\ V^\pi(h_y, n') &= \pi(r \mid h_y)C^\pi + \gamma\pi(g \mid h_y)(pV^\pi(h_{y,x}, x') + (1-p)V^\pi(h_{y,y}, y')) \end{aligned}$$

Since  $x'$  and  $y'$  have the same successor states and immediate rewards, the difference between  $V^\pi(h_{x,x}, x')$  and  $V^\pi(h_{x,y}, y')$  depends only on the history. We see that history  $h_{x,y}$  (and  $h_{y,x}$ ) give the agent strictly more information than  $h_{x,x}$  (and  $h_{y,y}$ ): for the former we know we are in  $x'$  (and  $y'$ ), while for the latter we could be in either  $x$  or  $x'$  (and  $y$  or  $y'$ ) with non-zero probability. Since more information can only allow the agent to pick better actions, we conclude  $V^\pi(h_{x,x}, x') \leq V^\pi(h_{x,y}, y')$  (and  $V^\pi(h_{y,y}, y') \leq V^\pi(h_{y,x}, x')$ ), and thus that the value is minimized for  $p = 1 - \sup(\mathcal{P})$  for history  $h_x$  (and for  $p = \inf(\mathcal{P})$  for history  $h_y$ ). The corollary holds via the same argument.  $\square$

Next, we show that the optimal policy is suboptimal against all stationary policies. More precisely, we give conditions for which, against the worst-case nature policy, it is only optimal to repair if the agent has detected that the machine is broken, i.e., if it takes action  $a_x$  but observes  $y$  two timesteps later, or similarly for  $s_y$  and  $x$ . In contrast, we give a similar condition such that, for the worst-case non-stationary policy, repair is optimal without detecting that the machine is broken. We formalize this logic as follows:

**Theorem 5.** In the maintenance benchmark, let  $\theta$  be any stationary nature policy, and  $h_t$  be any history such that the agent has some non-zero probability  $< 1$  to be in state  $n'$ . Then, the optimal policy against  $\theta$  chooses action  $\pi(h_t) = g$  as long as:

$$\frac{(1-\delta)}{1-0.5\gamma^2} + \frac{1-\gamma^2\delta}{1-\gamma^2(1+\gamma\delta-\delta)} \left[ \frac{0.5\gamma^2}{1-0.5\gamma^2} - 1 \right] > 0. \quad (5)$$

*Proof.* Via our lemma above, we may assume  $p = 0.5$ , in which case our problem is a standard POMDP for which we can talk about beliefs. (Note that if  $p \notin \mathcal{P}$ , then the agent can always get strictly more information, and thus has less incentive to repair without seeing a malfunction.) We first note that if the agent has detected the machine is broken since it's last repair, then no history exists such that the probability of being in state  $n'$  lies between 0 and 1. In that case, we may denote the two most 'extreme' beliefs possible as follows.  $b_\top$  denotes any belief such that  $b(n_x) + b(n_y) = 1$  (due to symmetry, these states are interchangeable for both nature and the agent), and  $b_\perp$  denotes the belief such that  $b(n')$  is maximized. In particular, let  $n$  denote the number of number of times the agent has observed  $x$  or  $y$  since the last repair action, or the beginning of the episode if no repair action has yet occurred. Then, we calculate the probability of being in state  $n'$  as follows:

$$b_\perp(n') = \delta \sum_{n=1} 0.5^n \leq \frac{0.5\delta}{1-0.5} = \delta, \quad (6)$$

We remark that if repairing is optimal for any belief  $h_t$ , then it must be optimal in  $b_\perp$ . Thus, we consider the values for both actions  $g$  and  $r$  in  $b_\perp$ :

$$\begin{aligned} Q(b_\perp, g) &:= Q_g = (1-\delta) + \gamma^2[0.5Q_g + 0.5Q_r] \\ &= \frac{(1-\delta) + 0.5\gamma^2Q_r}{1-0.5\gamma^2} \\ Q(b_\perp, r) &:= Q_r = -1 + \gamma Q(b_\top, g) \end{aligned}$$

To determine when measuring is optimal, we may look at the difference between these two values. In particular, action  $g$  is optimal as long as the following holds:

$$Q_g - Q_r = \frac{(1-\delta)}{1-0.5\gamma^2} + Q_r \left[ \frac{0.5\gamma^2}{1-0.5\gamma^2} - 1 \right] \geq 0 \quad (7)$$

(8)

We note that  $\frac{0.5\gamma^2}{1-0.5\gamma^2} < 1$  for all  $\gamma$ . Thus, we can use an upper bound for  $Q_r$  to find an overapproximation of when measuring is optimal. For this, we use the expected value given the fully observable setting, which we denote as  $V_{\text{RMDP}}$ . We note that our theorem is trivially true if repairing is not worth it in the fully observable case, thus we may assume that  $V_{\text{RMDP}}(n') = \gamma V_{\text{RMDP}}(n_x) - 1$ . Then, we get:

$$\begin{aligned} V_{\text{RMDP}}(n_x) &= V_{\text{RMDP}}(n_y) = 1 + \gamma^2 \left( \delta V_{\text{RMDP}}(n') + (1-\delta) V_{\text{RMDP}}(n_x) \right) \\ &= 1 + \gamma^2 \left( \delta(\gamma V_{\text{RMDP}}(n_x) - 1) + (1-\delta) V_{\text{RMDP}}(n_x) \right) \\ &= 1 - \gamma^2 \delta + \gamma^2 V_{\text{RMDP}}(n_x) (1 + \delta\gamma - \delta) \\ &= \frac{1 - \gamma^2 \delta}{1 - \gamma^2(1 + \gamma\delta - \delta)} \\ &\geq Q_r \end{aligned}$$

Filling this in for  $Q_r$ , we find that:

$$Q_g - Q_r \leq \frac{(1-\delta)}{1-0.5\gamma^2} + \frac{1 - \gamma^2 \delta}{1 - \gamma^2(1 + \gamma\delta - \delta)} \left[ \frac{0.5\gamma^2}{1-0.5\gamma^2} - 1 \right]$$

Action  $g$  is only optimal as long as this value is  $\geq 0$ , which gives us our bound.  $\square$

**Theorem 6.** *In the maintenance benchmark, let  $\theta_{\text{Nash}}$  be the Nash-optimal history-based nature policy, and define  $q := 1 - \min(\inf(\mathcal{P}), 1 - \sup(\mathcal{P}))$ . Then, there exists a history  $h_t$  such that the agent has a non-zero, but  $< 1$ , probability of being in state  $n'$  but  $\pi(h_t) = r$ , as long as:*

$$\frac{1}{1 - \gamma(1 - \delta)} \left[ 1 - \frac{(1-q)\gamma^2}{1 - q\gamma^2} \right] - \frac{(1-\delta)}{1 - q\gamma^2} \geq 0 \quad (9)$$

*Proof.* Without loss of generality, assume  $1 - \sup(\mathcal{P}) \geq \inf(\mathcal{P})$ , in which case lemma 2 and its corollary state  $\pi^*$  always picks action  $x$  and  $\theta^*$  chooses  $p = \sup(\mathcal{P}) := q$ . (The proof for  $1 - \sup(\mathcal{P}) \leq \inf(\mathcal{P})$  follows via symmetry of the model.) In that case, following the same logic used in the proof of theorem 5, we define the belief with the highest probability to be in state  $n'$  as

$$b_{\perp}(n') = \delta \sum_{n=1} q^n \leq \frac{q\delta}{1-q}$$

We define  $Q_g$  and  $Q_r$  as before, which yields the following condition:

$$Q_r - Q_g = Q_r \left[ 1 - \frac{(1-q)\gamma^2}{1 - q\gamma^2} \right] - \frac{(1-\delta)}{1 - q\gamma^2} \geq 0$$

Since  $\forall \gamma, \gamma \in (0, 1)$ :  $\frac{(1-q)\gamma^2}{1 - q\gamma^2} < 1$ , we can use a lower bound on  $Q_r$  to find a sufficient condition. One such lower bound is given by taking an agent policy that never measures, in which case

$$Q_r \geq \sum_{n=0} (\gamma(1-\delta))^{2n} = \frac{1}{1 - \gamma(1-\delta)},$$

which yields the condition given in the theorem.  $\square$

Lastly, the proof of theorem 1 follows from the fact that the parameters satisfy the conditions of both theorem 5 and theorem 6.

### A.2.2 Parity

We start by restating our theorem for convenience.

**Theorem 2.**  $\text{PARITY}(\infty)$ , with  $P_1 = \{0.2\}$ ,  $P_2 = [0.1, 0.7]$ ,  $P_3 = [0.1, 0.7]$ , and  $\gamma \geq 0.7\bar{3}$ , is not naively solvable.

*Proof.* We assume that the agent always picks an action according to its current most likely parity, which is clearly optimal. Thus, we can summarize the agent’s uncertainty as the using the *even-odd ratio*  $e = \max(\Pr(\text{even}), \Pr(\text{odd}))$ . Since the agent never receives any information in  $\text{PARITY}$ , it follows that any optimal policy must be *cyclic*. More precisely, let  $\pi_n$  denote a policy that repeatedly takes  $n$  s-actions, then a normal action which resets  $e$  to 1. Then, any optimal policy must be representable as  $\pi_n$ , for some  $n \in \mathbb{N}$ .

We show that the value of  $\pi_0$  is higher than that of any other  $\pi_n$ . First, the value of  $\pi_0$  is given as:

$$V^{\pi_0} = \sum_{n=0}^{\infty} \gamma^n = \frac{1}{1-\gamma} \quad (10)$$

If we can find any choice of probabilities for which any policy  $\pi_{n \neq 0}$  has a lower value, then we are done. We start with  $\pi_1$ , for which we pick  $p_1 = 0.2, p_2 = 0.5$  and  $p_3 = 0.3$ . In that case:

$$\begin{aligned} V^{\pi_1} &\leq \left( p_1(1+\gamma) + p_2(2-2\gamma) + p_3(3+\gamma) \right) \sum_{n=0}^{\infty} \gamma^{2n} \\ &= \frac{2.1 - 0.5\gamma}{1 - \gamma^2} \end{aligned} \quad (11)$$

This value is smaller than that of  $\pi_0$  for  $\gamma > \frac{11}{15} = 0.7\bar{3}$ . For larger values of  $n$ , we assume that the adversary starts off with the choice above, then picks  $p'_1 = 0.1, p'_2 = 0.1, p'_3 = 0.7$  as long as the agent keeps taking stochastic actions. We claim that under the second set of dynamics,  $e$  converges to the value of  $\frac{10}{19} \approx 0.526$ . To prove this, we invoke Banach’s fixed point theorem [40]. Denote  $e'$  as the even-odd ratio after a single step under the dynamics above, then  $e' = -0.9e + 1$ . Using the  $L^\infty$  distance, we find the following:

$$\left| e' - \frac{10}{19} \right| = \left| -0.9e + 1 - \frac{10}{19} \right| = 0.9 \left| e - \frac{10}{19} \right| \leq \left| e - \frac{10}{19} \right| \quad (12)$$

Thus, using Banachs theorem,  $e$  converges to  $\frac{10}{19}$ . In particular, this means that the maximum  $e$  that will be reached is given as  $e_{\max} = 0.5 + \left| 0.5 - \frac{11}{19} \right| = \frac{11}{19} \approx 0.579$ . Next, we note that the expected immediate return for the stochastic actions can be given as follows:

$$\mathbb{E}[r|e] = 2.1e - 2(1 - e) = 4.1e - 2. \quad (13)$$

For our maximum value of  $e$ , this yields an immediate return of  $\frac{71}{190} \approx 0.37$ . Thus, for any  $n > 1$ , we write the value function as follows:

$$V^{\pi_n} = \left[ \sum_{t=0}^{\infty} \gamma^{nt+t} \cdot 2.1 + \left[ \sum_{t'=1}^{n-1} \gamma^{nt+t'} \frac{71}{190} \right] + \gamma^{(n+1)t} (3e - 2) \right] < V^{\pi_0} \quad (14)$$

Thus,  $\pi_0$  is optimal.

Next, we show that the expected value for  $\pi_1$  is higher than that of  $\pi_0$  for all naive nature policies, which implies that  $\pi_0$  is not the policy with the highest value. Starting with  $\theta_{\text{Center}}$  and  $\theta_{\text{Ent}}$ , we find that both pick parameters  $p_1 = 0.2, p_2 = p_3 = 0.4$ , in which case:

$$V^{\pi_1, \theta_{\text{Center}}} = V^{\pi_1, \theta_{\text{Ent}}} = \frac{2.2 - 0.2\gamma}{1 - \gamma} \quad (15)$$

For  $\theta_{\text{RMDP}}$ , we find the parameters  $p_1 = 0.2, p_2 = 0.7, p_3 = 0.1$ . In contrast to the other policies, this means the most likely parity does not change after a stochastic step, and the expected value is given as:

$$V^{\pi_1, \theta_{\text{RMDP}}} = \frac{2.1 + 0.2\gamma}{1 - \gamma} \quad (16)$$

Both values are strictly larger than  $\frac{1}{1-\gamma}$  for any  $\gamma \in [0, 1]$ , which proves the model is not naively solvable.  $\square$

## B Evaluating Agent Policies in RPOMDPs

In this appendix, we provide proof for theorem 3, which we restate for convenience:

**Theorem 3.** *For any agent policy  $\pi \in \Pi$  there exists a best-response nature policy  $\theta_\pi^*: X \rightarrow \mathcal{U}$  such that  $\forall x \in X: \mathcal{P}(\theta_\pi^*(x)) \in \text{Extremes}(\{\mathcal{P}(u) \mid u \in \mathcal{U}\})$ .*

*Proof.* We first write out the value function for a policy  $\pi$  against an optimal nature policy  $\theta_\pi^*$ :

$$V^{\pi, \theta_\pi^*}(h_t, s, x) = \sum_{a \in A} \sigma(a \mid x) \inf_{u \in \mathcal{U}} \left[ R(s, a) + \gamma \left( \sum_{s' \in S} \sum_{o \in \Omega} \mathcal{P}(u)(s', o \mid s, a) \sum_{x' \in X} \tau(x' \mid x, o, a) V^{\pi, \theta_\pi^*}((h_t, a, o), s', x') \right) \right]$$

We notice that none of the terms in this formula depend on  $h_t$  (except recursively via  $V^{\pi, \theta_\pi^*}$ ), which means we can remove this dependency. With that, we rewrite the value function with simplified notation as follows:

$$\begin{aligned} V^{\pi, \theta_\pi^*}(s, x) &= \sum_{a \in A} \pi(a \mid x) Q^{\pi, \theta_\pi^*}(s, a, x, \theta_\pi^*(s, a, x)) \\ Q^{\pi, \theta_\pi^*}(s, a, x, u) &= R(s, a) + \gamma \sum_{s' \in S} \sum_{o \in \Omega} \mathcal{P}(u)(s', o \mid s, a) \sum_{x' \in X} \tau(x' \mid x, o, a) V^{\pi, \theta_\pi^*}(s', x') \\ \theta_\pi^*(s, a, x) &= \arg \inf_{u \in \mathcal{U}} Q^{\pi, \theta_\pi^*}(s, a, x, u) \end{aligned}$$

Given this formula, we first show that an optimal nature policy exists with signature  $\theta: X \rightarrow \mathcal{U}$ . We start by defining the following nature policy:

$$\theta_X(x) = \arg \inf_{u \in \mathcal{U}} \sum_{s \in S} \sum_{a \in A} Q^{\pi, \theta_X}(s, a, x, u) \quad (17)$$

Recall that  $\mathcal{U}$  is constructed such that every  $(s, a)$  pair has dedicated decision variables. Thus, there exists an  $u \in \mathcal{U}$  that minimizes eq. (17) for each  $(s, a)$  independently.

We prove that this choice is optimal, using proof by contradiction. If  $\theta_X(x)$  is suboptimal for  $x \in X$ , then there must exist some history  $h$  where its choice of  $u$  is suboptimal. In particular, let  $x$  be the first memory state reached where  $\theta_X$  makes a suboptimal choice, i.e. where choosing decision variables according to  $\theta_X$  and then following  $\theta_\pi^*$  would lead to a higher value. In that case, there must be at least one state-action pair  $s_{\text{diff}}, a_{\text{diff}} \in S \times A$  where the following holds:

$$Q^{\pi, \theta_\pi^*}(s_{\text{diff}}, a_{\text{diff}}, x, \theta_\pi^*(s_{\text{diff}}, a_{\text{diff}}, x)) - Q^{\pi, \theta_X}(s_{\text{diff}}, a_{\text{diff}}, x, \theta_X(x)) < 0$$

Since our model is  $(s, a)$ -rectangular, there exists a distinct set of decision variables that affect any state-action pair  $\mathcal{P}(\cdot, \cdot \mid s, a)$ , which we denote as  $\text{Var}_{s, a}$ . Then, denoting  $p \in \text{Var}$  as a decision variable, we define the following memory-based nature policy:

$$\theta'_X(x)(p) = \begin{cases} \theta_\pi^*(x)(p) & \text{if } p \in \text{Var}_{s_{\text{diff}}, a_{\text{diff}}} \\ \theta_X(x)(p) & \text{otherwise.} \end{cases}$$

Looking at eq. (17), we see that  $\theta'_X$  achieves the same same  $Q$ -values as  $\theta_X$  for all state-action tuples  $(s, a) \neq (s_{\text{diff}}, a_{\text{diff}})$ , but achieves a lower value for  $(s_{\text{diff}}, a_{\text{diff}})$ . However, we had defined  $\theta_X$  as the function that minimizes eq. (17), so we have a contradiction. This means no state-action tuple can exist where  $\theta_X$  is suboptimal, in which case  $\theta_X$  can never make a first suboptimal choice as compared to  $\theta_\pi^*$ . Thus,  $\theta_X$  is optimal.

Next, we need only show that an optimal policy exists such that  $\forall x: \theta_\pi^*(x) \in \text{Extremes}(\mathcal{P})$ . This immediately follows from the definition of  $Q^{\pi, \theta_\pi^*}$  with the observation that, if  $\pi$  is fixed,  $V^{\pi, \theta_\pi^*}$  is only dependent on the current choice  $u$  via it's arguments  $s'$  and  $\tau(x, o, a)$ . Thus, a variable assignment  $u$  that greedily maximizes the probabilities of reaching tuples  $(s', o)$  with low expected values is both optimal and complies with our condition.  $\square$

## C Efficient Approximations for Robust Agent Policies

Here, we provide extended analysis and proofs of the theoretical results in section 5. Let  $\mathcal{B} \subset \Delta(S)$  be the finite set of reachable beliefs.

First, we introduce the  $H_{RPOMDP}$  operator [37] in our notation. It will be used later in this appendix.

**Definition 5.**  $Q_{RPOMDP}$  is the fixed point of the operator  $H_{RPOMDP}$ , which is defined as:

$$H_{RPOMDP}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \sum_{s' \in S} \mathcal{P}(u)(s', o \mid s, a) \max_{a' \in A} Q(\tau_u(b, a, o), a') \right], \quad (18)$$

where  $\tau_u$  is the belief update under variable assignment  $u \in \mathcal{U}$  defined as:

$$\tau_u(b, a, o)(s') = b'(s') \propto \sum_{s \in S} b(s) \mathcal{P}(u)(s', o \mid s, a) \quad (19)$$

Then, let us restate the definitions of the upper bounds we introduced in the main body of the paper.

**Definition 4.**  $Q_{RMDP}$  and  $Q_{RFIB}$  are the fixed point of the operators  $H_{RQMDP}$  and  $H_{RFIB}$ , defined as:

$$H_{RQMDP}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{s' \in S} \mathcal{T}(u)(s' \mid s, a) \max_{a' \in A} Q(\mathbf{b}_{s'}, a') \right], \text{ and} \quad (2)$$

$$H_{RFIB}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o \mid s, a) Q(\mathbf{b}_{s'}, a') \right]. \quad (3)$$

Below, we provide the reasoning for the uniqueness and existence of fixed points through the Banach fixed point theorem [40]. Then, it suffices to show that the operators are contraction mappings.

**Lemma 3.** The fixed point  $Q_{RMDP}$  is synonymous with the fixed point of robust dynamic programming on the fully observable RMDP and lifting the resulting values into belief space.

**Corollary.** Consequently, the operator  $H_{RQMDP}$  is a contraction mapping, and the existence and uniqueness of the fixed point  $Q_{RMDP}$  are guaranteed through the Banach fixed point theorem [20, 35].

To prove that  $H_{RFIB}$  is a contraction, we introduce the following two well-known lemmas.

**Lemma 4.** Let  $X$  be a compact set and  $f, g$  be functions of type  $X \rightarrow \mathbb{R}$ . Then:

$$\left| \sup_{x \in X} f(x) - \sup_{x \in X} g(x) \right| \leq \sup_{x \in X} |f(x) - g(x)|, \text{ and, } \left| \inf_{x \in X} f(x) - \inf_{x \in X} g(x) \right| \leq \sup_{x \in X} |f(x) - g(x)|.$$

It is a well-established lemma that occurs relatively often. For a proof, see for instance [Lemma B.2; 26].

**Lemma 5** (Triangle Inequality). The triangle inequality states that, for any two real numbers  $u, v \in \mathbb{R}$  the following inequality holds:

$$|u + v| \leq |u| + |v|.$$

Now, we are set to state the main theorem to prove that  $H_{RFIB}$  is indeed a contraction.

**Theorem 7.** The operator  $H_{RFIB}: (\mathcal{B} \times A \rightarrow \mathbb{R}) \rightarrow (\mathcal{B} \times A \rightarrow \mathbb{R})$  is a contraction mapping in terms of the infinity norm  $\|\cdot\|_\infty$  and the discount factor  $0 \leq \gamma < 1$  as Lipschitz constant.

*Proof.* Let  $b \in \mathcal{B}$  and  $a \in A$  be any belief and action, and let  $Q: \mathcal{B} \times A \rightarrow \mathbb{R}$  and  $Q': \mathcal{B} \times A \rightarrow \mathbb{R}$  be any two Q-functions. Then:

$$\begin{aligned}
|H_{\text{RFIB}}Q(b, a) - H_{\text{RFIB}}Q'(b, a)| &= \left| \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right] \right. \\
&\quad \left. - \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q'(\mathbf{b}_{s'}, a') \right] \right| \\
&\leq \gamma \sum_{s \in S} b(s) \left| \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right. \\
&\quad \left. - \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q'(\mathbf{b}_{s'}, a') \right| \\
&\leq \gamma \sum_{s \in S} b(s) \sup_{u \in \mathcal{U}} \left| \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right. \\
&\quad \left. - \sum_{o \in \Omega} \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q'(\mathbf{b}_{s'}, a') \right| \\
&\leq \gamma \sum_{s \in S} b(s) \sup_{u \in \mathcal{U}} \sum_{o \in \Omega} \left| \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right. \\
&\quad \left. - \max_{a' \in A} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q'(\mathbf{b}_{s'}, a') \right| \\
&\leq \gamma \sum_{s \in S} b(s) \sup_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \left| \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right. \\
&\quad \left. - \sum_{s' \in S} \mathcal{P}(s', o | u, s, a) Q'(\mathbf{b}_{s'}, a') \right| \\
&\leq \gamma \sum_{s \in S} b(s) \sup_{u \in \mathcal{U}} \left[ \sum_{o \in \Omega} \max_{a' \in A} \left[ \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) \left| Q(\mathbf{b}_{s'}, a') - Q'(\mathbf{b}_{s'}, a') \right| \right] \right] \\
&\leq \gamma \sum_{s \in S} b(s) \sup_{u \in \mathcal{U}} \left[ \sum_{o \in \Omega} \max_{a' \in A} \left[ \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) \|Q - Q'\|_{\infty} \right] \right] \\
&= \gamma \|Q - Q'\|_{\infty}
\end{aligned}$$

Thus, it follows that, making use of the definition of the infinity norm:

$$\|H_{\text{RFIB}}Q(b, a) - H_{\text{RFIB}}Q'(b, a)\|_{\infty} \leq \max_{(b, a) \in \mathcal{B} \times A} |H_{\text{RFIB}}Q(b, a) - H_{\text{RFIB}}Q'(b, a)| \leq \gamma \|Q - Q'\|_{\infty}.$$

□

Lastly, we note that for  $Q_{\text{RFIB}}$  the set of reachable beliefs  $\mathcal{B}_S$  considered by the operator is finite, as it contains only  $|\mathcal{B}_S| = |S|$  unit beliefs. That is, the variable assignments  $u$  chosen by nature do not lead to an explosion of the set of reachable beliefs. Therefore, computing the fixed point only requires computing iterations of the operator over  $\mathcal{B}_S \times A$ .

The following definition and theorem help establish the tightness of the heuristics.

**Definition 6** (Monotone mapping). *A mapping  $H: (\mathcal{B} \times A \rightarrow \mathbb{R}) \rightarrow (\mathcal{B} \times A \rightarrow \mathbb{R})$  is monotone if for any two  $Q, Q'$  and for all  $(b, a) \in \mathcal{B} \times A$ , we have that  $Q(b, a) \leq Q'(b, a) \rightarrow HQ(b, a) \leq HQ'(b, a)$ .*

**Theorem 8** (Theorem 6, [18]). *Let  $H_1: \mathcal{B} \times A \rightarrow \mathcal{B} \times A$  and  $H_2: \mathcal{B} \times A \rightarrow \mathcal{B} \times A$  be two mappings defined on  $\mathcal{Q}_1$  and  $\mathcal{Q}_2$ . If:*



- $H_1$  and  $H_2$  are contractions with fixed points  $Q_1^*$  and  $Q_2^*$ ,
- $Q_1^* \in \mathcal{Q}_2$  and  $H_2 Q_1^* \geq H_1 Q_1^* = Q_1^*$ ,
- $H_2$  is a monotone mapping,

then  $Q_2^* \geq Q_1^*$ .

Note that we may have  $\mathcal{Q}_1 \subset \mathcal{Q}_2$ , i.e.,  $\mathcal{Q}_1$  may cover a smaller space of Q-value functions.

Now, we are set to prove the following theorem of the main paper:

**Theorem 4.** *Regarding tightness, the following inequalities on the fixed points hold:*

$$\forall b \in \Delta(S), \forall a \in A: Q_{RMDP}(b, a) \geq Q_{RFIB}(b, a) \geq Q_{RPOMDP}(b, a).$$

*Proof.* Let us first restate that it is known that  $H_{RPOMDP}$  is a contraction mapping [37]. Furthermore, note that it can be shown that  $H \in \{H_{RQMDP}, H_{RFIB}, H_{RPOMDP}\}$  are monotone mappings, see for instance [Appendix B.1.4; 26]. Then, it follows from the following observation [18]:

$$\begin{aligned} H_{RQMDP}Q(b, a) &= \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{s' \in S} \mathcal{T}(u)(s' | s, a) \max_{a' \in A} Q(\mathbf{b}_{s'}, a') \right] \\ &\geq H_{RFIB}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \max_{a' \in A} \sum_{s \in S} \mathcal{P}(u)(s', o | s, a) Q(\mathbf{b}_{s'}, a') \right] \\ &\geq H_{RPOMDP}Q(b, a) = \sum_{s \in S} b(s) \left[ R(s, a) + \gamma \inf_{u \in \mathcal{U}} \sum_{o \in \Omega} \sum_{s' \in S} \mathcal{P}(u)(s', o | s, a) \max_{a' \in A} Q(\tau_u(b, a, o), a') \right]. \end{aligned}$$

□

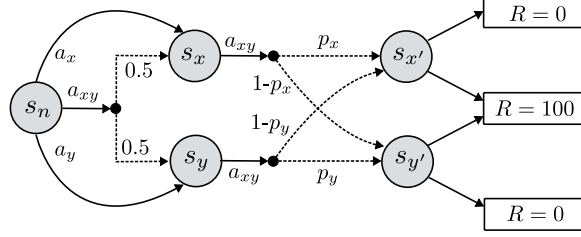


Figure 6: An example RPOMDP. We assume  $p_x \in [0.7, 0.9]$  and  $p_y \in [0.7, 0.9]$ . We assume all states give the same observation.

## D Experiments

### D.1 RHSVI

In this section, we give a brief overview of the RHSVI solver used in our experiments. We start with a brief introduction on value iteration for POMDPs, then repeat the initial formulation of RHSVI from Osogami [37]. Lastly, we describe several alterations made to RHSVI.

We provide an explanation of our alterations here for reproducibility. However, since we only use RHSVI as a baseline RPOMDP solver, we leave a complete description of the correctness and/or necessity of our alterations for future work.

**HSVI.** Since the value function for POMDPs is piecewise-linear [42], it can be represented as a (possibly infinite) set of linear lower bounds  $\Gamma = \{\alpha: S \rightarrow \mathbb{R}\}$  known as  $\alpha$ -vectors. As shown [39], this set can be approximated by iteratively performing a *backup operation* on a finite set of beliefs  $\mathcal{B}$ . To find such a belief set, *heuristic search value iteration* [43, 44] builds a belief tree via sampling, guided by upper bounds on the value function. Thus, the belief set is restricted to reachable beliefs, which increases tractability. Moreover, by keeping track of an upper bound, the algorithm can determine when the policy it has found is  $\epsilon$ -optimal.

**Robust HSVI.** Osogami [37] generalizes HSVI to RPOMDPs by changing the robust backup operator and belief update function to robust variants. We repeat both here in this paper’s notation. Given an uncertainty assignment  $u$  and belief-action pair  $(b, a)$ , the backup operator is defined as follows:

$$\alpha_\Gamma(b, a, u)(s) = \sum_{o \in \mathcal{O}} \alpha_{\Gamma, o}(b, a, u)(s) \quad (20)$$

$$\alpha_{\Gamma, o}(b, a, u) \in \operatorname{argmax}_{\alpha \in \Gamma} \sum_{s, s' \in S} b(s) \mathcal{P}(u)(s', o | s, a) \alpha(s'). \quad (21)$$

Given this, we can define the robust backup function and the corresponding worst-case nature policy (which directly defines the belief update function) as follows:

$$\theta_\Gamma(b, a) \in \operatorname{argmin}_{u \in \mathcal{U}} \sum_{s \in S} b(s) \alpha_\Gamma(b, a, u)(s), \quad (22)$$

$$\alpha_\Gamma(b, a) = \alpha_\Gamma(b, a, \theta_\Gamma(b, a)). \quad (23)$$

If our uncertainty set is given by intervals, then both can be computed using a linear program with at most  $\mathcal{O}(|\mathcal{O}||S|^2)$  variables and  $\mathcal{O}(|S|^2|\mathcal{O}| + |\mathcal{O}||\Gamma|)$  constraints.<sup>3</sup> However, the support of both the belief and possible successor beliefs are often much smaller than  $|S|$ , which means the complexity mostly depends on  $|\mathcal{O}|$  and  $|\Gamma|$  in practice.

Our implementation of RHSVI makes a number of changes from the description of [37], which we describe below:

**Robust backup.** First, we fix a problem with the backup procedure described above. Equation (21) implies that for each belief-action-observation tuple  $(b, a, o)$ , we can pick *any* alpha-vector that yields an optimal value against  $\theta_\Gamma(b, a)$  to compute our backup. This yields alpha-vectors that give the

<sup>3</sup>For eq. (21), we implement the  $\operatorname{argmax}$  operator via the constraint that  $\sum_{s \in S} b(s) \alpha_\Gamma(b, a, \theta_\Gamma(b, a))(s)$  must be at least as high as any choice  $\alpha \in \Gamma$ .

correct value for the current belief, but may lead to problems if we use the alpha-vector for different beliefs.

To illustrate this problem, consider the RPOMDP shown in fig. 6. Working backwards, the following  $\alpha$ -vectors can be found using the backup in states  $s_{x'}$  and  $s_{y'}$  and are thus valid:

$$\begin{aligned}\alpha_{x'}(s) &= 100 \cdot [s = s_{x'}], \\ \alpha_{y'}(s) &= 100 \cdot [s = s_{y'}].\end{aligned}$$

Next, we use these alpha-vectors to perform backups in the beliefs  $b_x$  and  $b_{xy}$ , which denote the beliefs reached from  $s_n$  after action  $a_x$  and  $a_{xy}$ , respectively. For the first, we can quickly see that the optimal nature policy picks  $\{p_x = 0.7, p_y = 0.3\}$ , which yields a value of 70. For belief  $b_{xy}$ , there are multiple optimal variable assignments for nature, including  $\{p_x \mapsto 0.9, p_y \mapsto 0.9\}$ . We use this assignment, in which case  $\alpha_{x'}$  is a valid solution to eq. (21). In that case, our backup returns the following  $\alpha$ -vector:

$$\alpha_{xy}(s) = \begin{cases} 90 & \text{if } s = s_x \\ 10 & \text{if } s = s_y. \end{cases}$$

This gives us the correct value of 0.5 for  $b_{xy}$ , but yields a value of 0.9 for  $b_x$ , which is higher than the actual value of 0.7 we computed before. Thus, even though  $\alpha_{xy}$  can be found using the robust backup, it is not a valid underapproximation of the value function.

We discovered and empirically confirmed the problem with the backup function described above, but were unable to fully address the problem. Instead, we implement an ad-hoc fix that aims to ensure  $\alpha_\Gamma(b, a)$  has the following properties:

1. *Indifference to state.*  $\alpha_\Gamma(b, a)$  should have roughly equal values for each state in the support of the current belief.
2. *Indifference to nature.*  $\alpha_\Gamma(b, a)$  should yield at least the same value against suboptimal nature policies, as compared to the optimal one.

To achieve both, we first perform the robust backup described above, which yields a nature policy  $\theta^*$  and robust value  $V$ . We then solve a second LP to find an  $\alpha$ -vector  $\alpha_r^*$  that yields the same value (up to a small error bound) but also has the properties above. To encode (1), we define the *state exploitability* of an  $\alpha$ -vector as follows:

**Definition 7.** Given an  $\alpha$ -vector  $\alpha$ , belief  $b$ , and corresponding value  $V_\alpha(b) = \sum_{s \in S} b(s)\alpha(s)$ , the state exploitability of  $\alpha$  with respect to  $b$  is defined as:

$$Expl(\alpha, b) = \sum_{s \in S} b(s)|\alpha(s) - V_\alpha(b)| \quad (24)$$

Intuitively,  $Expl(\alpha, b)$  is zero if the expected value of  $\alpha$  is independent of the actual state of the environment, while  $Expl(\alpha, b) > Expl(\alpha', b)$  implies  $\alpha$  is more robust against changes in the underlying state than  $\alpha'$ . We use  $Expl(\alpha, b)$  as a *penalty term* for our LP, i.e., we aim to maximize  $\left(\sum_{s \in S} b(s)\alpha_r^*(s)\right) - \delta e(\alpha_r^*, b)$  for some small value  $\delta > 0$ . To encode (2), we specify a number of nature policies that differ slightly from  $\theta^*$ , and add constraints so that  $\alpha_r^*$  achieves at least value  $V$  against all of these.

We empirically find that the approach above yields accurate value approximations, which is sufficient for the purposes of this paper. However, our testing has been limited to the experiments described in this paper, which were not specifically aimed at testing our backup. Apart from that, we have no theoretical basis to claim the approach is correct. We leave a more systematic analysis and solution to future work.

**Policy randomization.** Osogami [37] focuses on finding the value function for an RPOMDP, but does not consider the problem of constructing a policy which matches this value. In contrast to POMDPs, this is not a trivial problem, for similar reasons to the ones described above.

To illustrate this point, consider the environment of fig. 7, with  $\delta \in [0, 1]$ , and with  $\gamma = 1$  for simplicity. Here, the agent's only meaningful action is to guess whether they are in state  $x$  or  $y$ . Assuming  $p \in [1, 0]$ ,  $\theta_{\text{Nash}}$  should pick  $p$  such that the value given for both actions  $x$  and  $y$  is equal. In this case, the agent is ambivalent about what action to pick, which means  $\Gamma = \{\alpha_x\}$ , with  $\alpha_x(s) = [s = s_x]$ , is a valid representation of the value function. However, a policy that always picks

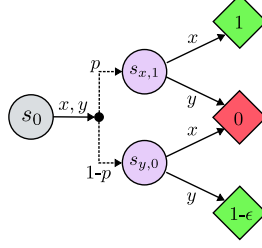


Figure 7: An RPOMDP showing the necessity of probabilities for robust policies.

action  $x$  is clearly suboptimal, since it performs poorly against  $\{p \mapsto 0\}$ . Thus, it is not possible to create an optimal policy based on this  $\Gamma$ . From this example, we can conclude the following:

**Remark 2.** To construct optimal RPOMDP policies,  $\Gamma$  should generally include *all*  $\alpha$ -vectors corresponding to optimal actions for each reachable belief. In contrast, for POMDPs, *any* such  $\alpha$ -vector is sufficient.

Note that we need not consider  $\alpha$ -vectors corresponding to *suboptimal* actions for reachable beliefs. With this condition satisfied, our policy should pick probabilities for each optimal action such that the expected value is robust against different nature policies. More precisely, we want to be robust against nature choosing different dynamics in previous timesteps, which would have led us to a different belief. Luckily, we have already defined the concept of state exploitability above, which aims to solve exactly this problem. Thus, we define our policy to aim and minimize state exploitability, which we formally define as follows:

$$\begin{aligned} d_{\Gamma}^{Expl} &= \operatorname{argmin}_{d \in \Delta(\Gamma)} \operatorname{Expl}\left(\left(\sum_{\alpha \in \Gamma^*} d(\alpha)\alpha\right), b\right) \\ \alpha_{\Gamma}^{Expl}(s) &= \sum_{\Gamma \in \Gamma^*} d_{\Gamma}^{Expl} \alpha(s) \\ \pi_{\Gamma}^{Expl}(a | b) &= \sum_{\alpha \in \Gamma} d_{\Gamma}^{Expl}(\alpha) [\alpha \in \Gamma_a] \end{aligned}$$

Although we have no proof that this is optimal, we show below that this choice is equivalent to minimizing a particular upper bound on the value function:

**Lemma 6.** Let  $\pi^*$  be an optimal policy and  $\pi$  be a policy that picks actions with different probabilities for some belief  $b$  but is otherwise identical. Then:

$$V^{\pi^*} - V^{\pi} \leq \operatorname{Expl}(b, \pi). \quad (25)$$

In particular, if  $\forall b: \operatorname{Expl}(b, \pi) = 0$ , then  $\pi$  is optimal.

*Proof.* Assuming our model is graph-preserving,  $\operatorname{supp}(b) = \operatorname{supp}(b')$ . Let  $s_+, s_-$  denote the states with the biggest difference in expected value, i.e.,

$$(s_+, s_-) \in \operatorname{argmax}_{s, s' \in \operatorname{supp}(b)} \sum_{\alpha \in \Lambda^*(b)} \pi(a) [\alpha(s) - \alpha(s')]. \quad (26)$$

Recall that  $\mathbf{b}_s$  is the unit belief with  $b(s) = 1$ . Let  $b_- = \mathbf{b}_{s_-}$ , and  $b_+ = \mathbf{b}_{s_+}$  in which case  $V^{\pi}(b_-) \leq V^{\pi}(b') \leq V^{\pi}(b) = V^{\pi^*}(b) \leq V^{\pi}(b_+)$ . Thus, we rewrite eq. (25) as follows:

$$\begin{aligned} V^{\pi^*} - V^{\pi} &\leq V^{\pi}(b_+) - V^{\pi}(b_-) \\ &= \sum_{\alpha \in \Lambda^*(b)} \pi(a) [\alpha(s_+) - \alpha(s_-)] \\ &\leq \operatorname{Exploit}(b, \pi), \end{aligned}$$

which proves our lemma.  $\square$

### Computational optimizations.

Next, we highlight a number of significant alterations to RHSVI that we make to improve performance:

*Belief Tree.* Equation (23) shows that the worst-case nature policy, and thus the belief update, depends on the current set of  $\alpha$ -vectors  $\Gamma$ . Thus, for any found belief  $b$ , the possible successor beliefs may change over time. This problem is not addressed by Osogami [37], which suggest they do not keep track of the belief tree at all. This is a valid approach, but it does not allow the reuse of belief nodes, which makes the method computationally expensive in practice. To deal with this, we use a belief tree, but periodically reset it at exponentially increasing intervals. This way, we guarantee that RHSVI finds all reachable beliefs (though this may take many iterations and resets in practice), while we can still use a tree structure to find new beliefs and compute tighter upper bounds more efficiently.

*Vector Pruning.* As in HSVI, we prune  $\alpha$ -vectors using point-wise domination with two changes. Firstly, since the complexity of both the backup- and belief update functions depend strongly on  $|\Gamma|$ , we try to keep this set as small as possible by only adding new  $\alpha$ -vectors if they are not dominated by any  $\alpha$ -vectors in  $\Gamma$ . This requires additional overhead but drastically decreases the cost of backups. Secondly, to allow us to compute random policies, we only consider domination between  $\alpha$ -vectors that correspond to the same action. Thus, if  $\alpha$  dominates  $\alpha'$  but corresponds to a different action, then  $\alpha'$  does not get pruned.

*Upper bounds.* We initialize the upper bounds with the robust upper bounds introduced in section 5 of the main body of the paper.

## D.2 Benchmarks & Infrastructure

**Infrastructure.** All experiments were conducted in Julia (version 1.11.5) on the same Ubuntu machine (version 22.04.5 LTS), which has an Intel(R) Core(TM) i9-10980XE CPU @ 3.00GHz and 256GB RAM (8 x 32GB DDR4-3200). We parallelize experiments across three workers [47].

**Benchmark descriptions.** Our new benchmarks, TOY\* and ECHO, as well as the finite and infinite chain environments PARITY(10) and PARITY( $\infty$ ), are described in section 3.

For the second set of benchmarks, we lift several (classic) POMDPs from the literature into RPOMDPs: TIGER [4] (also used in [19, 33]), MINIHALLWAY [28], and ALOHA [21], as well as an expanded variant of HEAVENORHELL [2] (also used in [37]).

- **TIGER:** The classic problem where an agent has to decide between two doors to open. One door has a tiger behind it, with a high negative reward associated, and one door has a prize with a high reward. The agent only knows the initial state distribution; the tiger is initialized to one of the doors with uniform probability. It therefore has to find a balance between noisy listening actions to gather information and deciding to open one of the doors. Therefore, it is intuitive that the nature policy heuristics that maximizes the entropy of the agent’s belief leads to a good approximation of the worst-case.
- **MINIHALLWAY:** A small POMDP problem where the agent must navigate a corridor with noisy observations of the goal location, where a reward is located.
- **ALOHA( $N$ ):** The agent must set the parameters of a communication protocol formulated as a POMDP, parameterized by  $N$ . A detailed description can be found in Jeon et al. [21].
- **HEAVENORHELL( $N$ ):** An agent is positioned to approach a corridor parameterized by  $N$  with, at each end, either a reward (heaven) or a penalty (hell). The location of the reward is observable only in the starting state, so the agent must remember it while traversing the branching point of the corridor.

We construct these RPOMDPs such that for any  $\mathcal{T}(u)$ , any transition probability is less than 0.5 times higher or lower than the nominal POMDP, with no alterations to the observation function.

For the last set, we added partial observability to two benchmarks from the RMDP literature: HEALTHDETECTION [13] and REPLACEMENT [8].

- **HEALTHDETECTION:** The agent must schedule different methods of medical screening for colorectal cancer (or CRC) for a patient. The dynamics of the model are based on real medical data. In the original paper, the authors only consider a number of existing

Dimensions	$ S $	$ \Omega $	$ A $
TOY*	9	2	2
ECHO	8	2	2
PARITY( $\infty$ )	9	1	4
PARITY (10)	12	1	4
TIGER	2	2	3
HEAVENORHELL(5)	28	11	4
HEAVENORHELL(10)	48	21	4
MINIHALLWAY	13	9	3
ALOHA(10)	30	30	9
REPLACEMENT	7	7	4
HEALTHDETECTION	378	9	3

Table 1: Dimensions of the RPOMDP benchmarks used in the experimental evaluation.

Algorithm	RHSVI( $M_{\text{Center}}$ )		RHSVI( $M_{\text{Ent}}$ )		RHSVI( $M_{\text{RMDP}}$ )	
Metric	min.	std.	min.	std.	min.	std.
TOY*	37.48	35.76	37.48	0.01	32.49	17.88
ECHO	19.31	0.01	19.30	0.01	21.12	0.00
PARITY( $\infty$ )	9.25	0.74	9.07	1.05	13.65	0.00
PARITY (10)	59.92	0.10	62.40	0.04	55.19	0.07
TIGER	19.35	0.01	19.35	0.01	13.12	0.02
HEAVENORHELL(5)	-24.04	0.01	-21.55	0.01	-24.04	0.00
HEAVENORHELL(10)	-37.35	0.01	-36.71	0.01	-37.36	0.02
MINIHALLWAY	0.76	0.00	0.76	0.00	0.76	0.00
ALOHA(10)	62.41	0.15	59.70	0.11	56.49	0.10
REPLACEMENT	-47.64	0.76	-46.69	0.82	-46.58	0.50
HEALTHDETECTION	-5718.42	36.29	-5661.49	67.73	-5554.34	32.23

Table 2: Detailed statistics for the evaluation of the naive nature policies. We report the worst value (min.) out of 5 runs as computed on the nature MDP using MCTS and the standard deviation (std.) of the values found by the 5 MCTS runs.

screening protocols, which they combine with their model to construct robust Markov chains. However, the model can also be interpreted as an RPOMDP.

- **REPLACEMENT:** the agent must schedule costly repairs for a machine, which is represented by a chain environment. To transform this model into an RPOMDP, we assume the agent cannot observe the state of the machine unless they pay a measurement cost. Such *active measure* environments have been considered both for POMDPs [16, 24] and RPOMDPs [25].

We use discount factor  $\gamma = 1$  for TOY\*, of  $\gamma = 0.99$  for ECHO en HEAVENORHELL, and of  $\gamma = 0.95$  for all other environments.

**Benchmark dimensions.** Table 1 shows the dimensions of the benchmarks used in the experimental evaluation.

### D.3 Error margins

In the plot in the main body of the paper, we normalize and plot the worst result among the 5 MCTS runs on the nature MDP. In tables 2 and 3, we provide the raw value (not normalized) of the worst evaluation out of the 5 runs, together with the standard deviation of the set of values found in the 5 MCTS runs, for all the algorithms tested.

Algorithm Metric	RHSVI		RQMDP		RFIB	
	min.	std.	min.	std.	min.	std.
TOY*	69.99	0.00	32.49	28.27	62.47	0.01
ECHO	31.10	0.00	25.45	0.01	25.44	0.01
PARITY( $\infty$ )	20.00	0.00	7.98	0.02	20.00	0.00
PARITY (10)	62.71	0.00	38.89	1.60	62.71	0.00
TIGER	19.36	0.01	-20.04	7.33	14.49	1.27
HEAVENORHELL(5)	-21.55	0.01	-63.76	0.00	-63.76	0.00
HEAVENORHELL(10)	-36.71	0.01	-63.76	0.00	-63.76	0.00
MINIHALLWAY	0.76	0.00	0.25	0.00	0.25	0.00
ALOHA(10)	46.96	1.98	57.32	0.27	53.39	0.90
REPLACEMENT	-46.16	0.15	-51.96	0.37	-70.87	2.02
HEALTHDETECTION	-5596.72	46.17	-5671.29	36.57	-5660.60	64.77

Table 3: Detailed statistics for the evaluation of the new baselines and RHSVI. We report the worst value (min.) out of 5 runs as computed on the nature MDP using MCTS and the standard deviation (std.) of the values found by the 5 MCTS runs.