

# Revealing POMDPs: Qualitative and Quantitative Analysis for Parity Objectives

Ali Asadi<sup>1</sup>, Krishnendu Chatterjee<sup>1</sup>, David Lurie<sup>2</sup> Raimundo Saona<sup>3</sup> \*

<sup>1</sup>Institute of Science and Technology Austria

<sup>2</sup>Paris Dauphine University, PSL Research University, Paris, France and NyxAir, Paris, France

<sup>3</sup>London School of Economics and Political Science, London, United Kingdom

{ali.asadi, krishnendu.chatterjee}@ista.ac.at, david.lurie@dauphine.eu, raimundo.saona@gmail.com

## Abstract

Partially observable Markov decision processes (POMDPs) are a central model for uncertainty in sequential decision making. The most basic objective is the reachability objective, where a target set must be eventually visited, and the more general parity objectives can model all  $\omega$ -regular specifications. For such objectives, the computational analysis problems are the following: (a) qualitative analysis that asks whether the objective can be satisfied with probability 1 (almost-sure winning) or probability arbitrarily close to 1 (limit-sure winning); and (b) quantitative analysis that asks for the approximation of the optimal probability of satisfying the objective. For general POMDPs, almost-sure analysis for reachability objectives is EXPTIME-complete, but limit-sure and quantitative analyses for reachability objectives are undecidable; almost-sure, limit-sure, and quantitative analyses for parity objectives are all undecidable. A special class of POMDPs, called revealing POMDPs, has been studied recently in several works, and for this subclass the almost-sure analysis for parity objectives was shown to be EXPTIME-complete. In this work, we show that for revealing POMDPs the limit-sure analysis for parity objectives is EXPTIME-complete, and even the quantitative analysis for parity objectives can be achieved in EXPTIME.

## 1 Introduction

**POMDPs** *Partially observable Markov decision processes* (POMDPs) model sequential decision-making under uncertainty (Bertsekas 1976; Papadimitriou and Tsitsiklis 1987; Kaelbling, Littman, and Cassandra 1998). At each step, the environment is in some hidden state. A controller interacts with it by choosing actions. The chosen action and the current state determine a probability distribution over the subsequent state. The controller cannot observe the state directly, but observes a signal that discloses only partial information of the state. POMDPs generalize classic models: *Markov decision processes* (MDPs), in which the state is fully observed (Puterman 2014), and *blind MDPs*, in which no state information is observed and are equivalent to Probabilistic Finite Automata (Rabin 1963; Paz 1971).

\*These authors contributed equally.

Copyright © 2026, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

**Objectives** The controller aims to maximize an *objective function*, which formally captures the desired behaviors of the model. Two main classes of objectives are typically considered: *logical objectives*, e.g., reachability and LTL (linear-temporal logic), and *quantitative objectives*, e.g., discounted sum and limit-average; see (Puterman 2014; Baier and Katoen 2008; Filar and Vrieze 1997) for details. This work focuses on logical objectives.

The most basic example of logical objectives is *reachability*, i.e., given a set of target states, the objective requires that some target state is visited at least once. A more general class of logical objectives is *parity objectives*, which assign to each state a non-negative integer, called *priority*. The objective is satisfied if the smallest priority appearing infinitely often is even. Parity objectives express all commonly used temporal objectives such as liveness, and are a canonical form to express all  $\omega$ -regular objectives (Thomas 1997) and LTL objectives (Baier and Katoen 2008). Hence, the study of POMDPs with parity objectives is a fundamental theoretical problem.

**Applications** POMDPs have been applied across diverse fields, including computational biology (Durbin et al. 1998) and reinforcement learning (Kaelbling, Littman, and Moore 1996). In particular, POMDPs with logical objectives have proven useful in many application areas such as probabilistic planning (Bonet and Geffner 2009); randomized embedded scheduler (de Alfaro et al. 2005); randomized distributed algorithms (Pogosyants, Segala, and Lynch 2000); probabilistic specification languages (Baier, Größer, and Bertrand 2012); and robot planning (Kress-Gazit, Fainekos, and Pappas 2009; Chatterjee et al. 2015).

**Computational Analysis Questions** A policy determines the choice of actions by the controller. The value corresponds to the maximum probability that the controller can guarantee to satisfy an objective. The main computational analysis of POMDPs with reachability and parity objectives considers the following two problems.

- *Qualitative analysis* has two variants: the *almost-sure winning* asks whether there exists a policy that satisfies the objective with probability 1; and the *limit-sure winning* asks whether the objective can be satisfied with probability arbitrarily close to 1.

Problems	Objectives	
	Reachability	Parity
Almost-sure	EXPTIME-complete	Undecidable
Limit-sure	Undecidable	Undecidable
Quantitative	Undecidable	Undecidable

Table 1: Computational complexity for POMDPs with reachability and parity objectives.

- *Quantitative analysis* asks to approximate the maximum or optimal probability with which the objective can be satisfied, up to a given additive error.

**General Undecidability** Most results for POMDPs and the computational analysis are negative (undecidability results). Indeed, the limit-sure analysis for reachability objectives is undecidable (Gimbert and Oualhadj 2010; Chatterjee and Henzinger 2010), which extends to general parity objectives. The almost-sure analysis for reachability objectives is EXPTIME-complete (Baier, Bertrand, and Größer 2008; Chatterjee, Doyen, and Henzinger 2010), but is undecidable for parity objectives even for two priorities (namely coBüchi objectives) (Baier, Bertrand, and Größer 2008; Chatterjee et al. 2010). The quantitative analysis for reachability objectives is undecidable (Madani, Hanks, and Condon 2003), which extends to general parity objectives. Given this wide range of results about non-existence of algorithms a natural direction to explore is the existence of subclasses for which the computational problems are decidable.

**Revealing POMDPs** We study *revealing POMDPs*, a special subclass of POMDPs where each visited state is announced to the controller with a positive probability. Consequently, the controller’s uncertainty about the state, i.e. the probability distribution over states (the *belief*), occasionally collapses into a Dirac distribution on the announced state. Revealing POMDPs have been previously studied in several works, including (Chen and Liew 2023; Belly et al. 2025; Avrachenkov, Dhiman, and Kavitha 2025), which present motivation and applications for this model.

**Previous Results and Open Questions** A key result for revealing POMDPs shows that the almost-sure analysis for parity objectives is EXPTIME-complete (Belly et al. 2025). However, the limit-sure and quantitative analysis problems for reachability and parity objectives remained open for revealing POMDPs.

**Our Contributions** We address the open questions for revealing POMDPs. Our main contributions are the following. For revealing POMDPs with parity objectives, we show that

- The limit-sure analysis coincides with almost-sure analysis, and consequently is EXPTIME-complete.
- The quantitative analysis can be achieved in EXPTIME, i.e., in the same complexity as for qualitative analysis.

The results for POMDPs and revealing POMDPs are summarized in Table 1 and Table 2, respectively.

Problems	Objectives	
	Reachability	Parity
Almost-sure	EXPTIME	EXPTIME-complete
Limit-sure	<b>EXPTIME</b>	<b>EXPTIME-complete</b>
Quantitative	<b>EXPTIME</b>	<b>EXPTIME</b>

Table 2: Computational complexity for revealing POMDPs with reachability and parity objectives. Our contributions are marked in bold.

**Technical Contributions** A closely related work established that almost-sure analysis for parity objectives in revealing POMDPs is EXPTIME-complete (Belly et al. 2025). Additionally, (Chatterjee, Chmelík, and Tracol 2016) previously showed that almost-sure analysis for parity objectives in general POMDPs under finite-memory policies is EXPTIME-complete. Therefore, the EXPTIME-completeness result naturally follows once finite memory policies are proven sufficient for almost-sure analysis. However, limit-sure analysis and quantitative analysis for POMDPs remain undecidable in general, even under finite memory policies (Chatterjee, Saona, and Ziliotto 2021). This highlights a sharp contrast with the almost-sure analysis and gives rise to new technical challenges.

First, we consider revealing POMDPs with the belief-reachability objectives. Given a set of target states, the belief-reachability objectives consider the probability of eventually observing such states. We prove that the quantitative analysis of belief-reachability in revealing POMDPs can be achieved in EXPTIME. The argument proceeds in two steps: (i) we prove that this objective can be approximated by their finite-horizon counterparts; and (ii) we show that this finite-horizon objective can be approximated by a point-based approximation.

Then, we consider revealing POMDPs with the parity objectives. We show that the value for parity objectives corresponds to the value for belief-reachability objectives to almost-sure winning states. Finally, we prove the existence of optimal policies for parity objectives in revealing POMDPs.

**Related Works** The study of subclasses of POMDPs with tractable algorithms is a broad topic with many directions. For example, subclasses of POMDPs for qualitative analysis have been explored in (Fijalkow et al. 2015; Chatterjee and Tracol 2012); and for quantitative analysis various subclasses have also been studied such as with ergodicity condition (Chatterjee et al. 2024); multiple environments only (Van Der Vegt, Jansen, and Junges 2023; Chatterjee et al. 2025); or in online learning (Liu et al. 2022; Chen et al. 2023). Our work focuses on such a subclass, namely revealing POMDPs, which has been studied in the literature.

Avrachenkov, Dhiman, and Kavitha (2025) studied strongly-connected revealing POMDPs and restricting attention to the set of belief-stationary policies. They proved that there is an optimal policy within this set, whose induced belief dynamics are contracting and reach a positive recurrent class of beliefs in finite time. They also provided an infinite-

dimensional linear program that characterizes the optimal belief-stationary policy. They considered an application of strongly-connected revealing blind MDPs. (Chen and Liew 2023) introduced the model of revealing blind MDPs. They consider the discounted objective and present algorithms based on finite-horizon approximation to compute the discounted value. Our work addresses the open computational analysis questions for revealing POMDPs.

## 2 Preliminaries

**Notation** For a positive integer  $n$  the set  $\{1, 2, \dots, n\}$  is denoted by  $[n]$ . Sets and correspondences are denoted by calligraphic letters, e.g.,  $\mathcal{S}, \mathcal{A}, \mathcal{Z}$ . Elements of these sets are denoted by lowercase letters, e.g.,  $s, a, z$ . Random elements with values in these sets are denoted by uppercase letters, e.g.,  $S, A, Z$ . The set of probability measures over a finite set  $\mathcal{S}$  is denoted by  $\Delta(\mathcal{S})$ . The Dirac measure at some element  $s \in \mathcal{S}$  is denoted by  $\mathbb{1}[s]$ . The support of a probability measure  $b \in \Delta(\mathcal{S})$  is denoted by  $\text{supp}(b)$ .

**Definition 1 (POMDP).** A POMDP is a tuple  $P = (\mathcal{S}, \mathcal{A}, \mathcal{Z}, \delta, b_0)$ , where:

- $\mathcal{S}$  is a finite set of states;
- $\mathcal{A}$  is a finite set of actions;
- $\mathcal{Z}$  is a finite set of signals;
- $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S} \times \mathcal{Z})$  is a probabilistic transition function that, given a state  $s$  and an action  $a$ , returns the distribution over the successor states and signal;
- $b_0 \in \Delta(\mathcal{S})$  is the initial belief over the states.

Markov Decision Processes (MDPs) are POMDPs in which the signal corresponds to the state (Puterman 2014) and are denoted by  $M = (\mathcal{S}, \mathcal{A}, \delta: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S}), b_0)$ .

**Dynamic** Given a POMDP, a controller knows all defining parameters. At the beginning, nature draws a state  $S_0 \sim b_0$ , which is not informed to the controller. Then, at each step  $t \geq 0$ , the controller chooses an action  $A_t \in \mathcal{A}$ , possibly at random. In response, nature draws the next state and signal  $(S_{t+1}, Z_{t+1}) \sim \delta(S_t, A_t)$ . The signal  $Z_{t+1}$  is revealed to the controller, but the state  $S_{t+1}$  is not announced. A play (or a path) in the POMDP is an infinite sequence  $\rho = (s_0, a_0, z_1, s_1, a_1, z_2, s_2, a_2, \dots)$  of states, actions, and signals such that, for all  $t \geq 0$ , we have  $\delta(s_t, a_t)(s_{t+1}, z_{t+1}) > 0$ . The set of all plays is denoted by  $\Omega$ .

**Belief** At each step  $t \geq 0$ , the controller's (random) belief about the current state can be computed using Bayes' rule and is denoted by  $B_t \in \Delta(\mathcal{S})$ . In the cases where the controller knows the exact state the controller's belief is a Dirac measure  $\mathbb{1}[s]$  at some state  $s \in \mathcal{S}$ .

**Definition 2 (Revealing POMDP).** A POMDP is revealing if, each time a state is visited, the state is also announced to the controller with positive probability. Formally, for each state there is a designated signal, i.e.,  $\mathcal{S} \subseteq \mathcal{Z}$ , and, for all states  $s, s' \in \mathcal{S}$  and actions  $a \in \mathcal{A}$ ,

$$\begin{aligned} \sum_{z \in \mathcal{Z}} \delta(s, a)(s', z) &> 0 \\ \implies \sum_{\tilde{s} \in \mathcal{S}} \delta(s, a)(\tilde{s}, s') &= \delta(s, a)(s', s') > 0. \end{aligned}$$

Revealing POMDPs coincide with the class of strongly revealing defined in (Belly et al. 2025). Indeed, if the controller observes a signal  $s \in \mathcal{Z}$ , then their next belief is  $\mathbb{1}[s]$ , which corresponds to a revelation of the state.

**Remark 1 (Signals).** The signaling structure of POMDPs is modeled in different ways in the literature. We comment on two general cases.

- The controller may receive a set of signals at each step with a transition function  $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S} \times 2^{\mathcal{Z}})$ . In that case, we consider each set of signals as one signal  $\tilde{\mathcal{Z}} := 2^{\mathcal{Z}}$  and reduce to our model. Even with exponentially many signals encoded in a polynomial size input, our EXPTIME upper bound holds.
- In general POMDPs, it is enough to consider transition functions that pair each state with only one signal, known as deterministic observation, see (Chatterjee, Chmélík, and Tracol 2016, Remark 4). For revealing POMDPs, deterministic observations do not capture the general case because they correspond to MDPs.

**Policies** A policy defines how a controller selects actions based on all the information available up to a given step. Formally, a (history-dependent randomized) policy is a function  $\sigma: (\mathcal{A} \times \mathcal{Z})^* \rightarrow \Delta(\mathcal{A})$ . The set of all policies is denoted by  $\Sigma$ . A policy is pure if it prescribes deterministic actions, i.e., it corresponds to a function  $\sigma: (\mathcal{A} \times \mathcal{Z})^* \rightarrow \mathcal{A}$ .

**Probability Measures** For a finite prefix of a play, an element in  $(\mathcal{S} \times \mathcal{A})^*$ , its cone is the set of plays with it as their prefix. Given a policy  $\sigma$  and an initial belief  $b_0$ , the unique probability measure over Borel sets of infinite plays obtained given  $\sigma$  is denoted by  $\mathbb{P}_{b_0}^\sigma(\cdot)$ , which is defined by Carathéodory's extension theorem by extending the natural definition over cones of plays (Billingsley 2012).

**Objectives** An objective in a POMDP is a Borel set of plays  $\Phi \subseteq \Omega$  in the Cantor topology on  $\Omega$  (Kechris 1995). We consider objectives in the first  $2^{1/2}$  levels of the Borel hierarchy, including the parity objective which expresses all  $\omega$ -regular objectives (Thomas 1997). Denote the state at time  $t$  by  $s_t$  and set of states that occur infinitely often in a play  $\rho$  by  $\mathcal{I}(\rho) := \{s \in \mathcal{S} : \forall t \geq 0 \exists \tilde{t} \geq t \quad s = s_{\tilde{t}} \in \rho\}$ . We consider the following objectives.

- **Reachability:** Given a set  $\mathcal{X} \subseteq \mathcal{S}$  of target states, the reachability objective requires that a target state is visited at least once. Formally, the reachability objective is  $\text{Reach}(\mathcal{X}) := \{\rho \in \Omega : \exists t \geq 0 \quad s_t \in \mathcal{X}\}$ .
- **Parity:** Given a  $d \geq 0$  and a function  $\text{pri}: \mathcal{S} \rightarrow \{0, 1, \dots, d\}$  of priorities, the parity objective requires that the smallest priority that appears infinitely often is even. Formally,  $\text{Parity} := \{\rho \in \Omega : \min\{\text{pri}(s) : s \in \mathcal{I}(\rho)\} \text{ is even}\}$ .

**Value** The value of an objective in a POMDP is the maximum probability a controller can guarantee to satisfy the objective. Formally, given an objective  $\Phi \subseteq \Omega$ , the value is a function of the initial belief  $\text{val}_\Phi: \Delta(\mathcal{S}) \rightarrow [0, 1]$  defined by  $\text{val}_\Phi(b) := \sup_{\sigma \in \Sigma} \mathbb{P}_b^\sigma(\rho \in \Phi)$ . Denote the reachability and parity values by  $\text{val}_{\text{R}(\mathcal{X})}$  and  $\text{val}_{\text{P}}$ , respectively. We may omit  $\mathcal{X}$  in  $\text{val}_{\text{R}(\mathcal{X})}$  if it is clear from the context.

**Approximately Optimal Policies** Given  $\varepsilon \geq 0$  and an objective  $\Phi \subseteq \Omega$ , a policy  $\sigma \in \Sigma$  is  $\varepsilon$ -optimal if it guarantees the value up to an additive error  $\varepsilon$ , i.e., if  $\mathbb{P}_b^\sigma(\rho \in \Phi) \geq \text{val}_\Phi(b) - \varepsilon$ . In particular, we call 0-optimal policies simply optimal.

**Computational Analysis** The computational analysis problems for POMDPs with an objective  $\Phi$  are:

- *Almost-sure* analysis asks whether the objective can be satisfied with probability 1, i.e.,  $\exists \sigma \in \Sigma$  such that  $\mathbb{P}_{b_0}^\sigma(\rho \in \Phi) = 1$ .
- *Limit-sure* analysis asks whether the objective can be satisfied with probability arbitrarily close to 1, i.e.,  $\forall \varepsilon > 0 \exists \sigma \in \Sigma$  such that  $\mathbb{P}_{b_0}^\sigma(\rho \in \Phi) \geq 1 - \varepsilon$  or equivalently whether  $\text{val}_\Phi(b_0) = 1$ .
- *Quantitative* analysis asks to compute an approximation of the optimal value, i.e., for all  $\varepsilon > 0$ , provide  $v \in [0, 1]$  such that  $|v - \text{val}_\Phi(b_0)| \leq \varepsilon$ .

### 3 Overview

We present an overview of our approach and the results.

#### 3.1 Overview of Approach

To solve the qualitative and quantitative analysis of parity objectives, we introduce a new objective called belief-reachability and show that the parity value coincides with the belief-reachability value to the set of almost-sure winning parity states.

**Belief-Reachability** Given a set of target states  $\mathcal{X} \subseteq \mathcal{S}$ , the belief-reachability objective requires that a target state is visited and the controller has complete knowledge of the state at that step at least once. Formally, the set of Dirac beliefs on  $\mathcal{X} \subseteq \mathcal{S}$  is denoted by  $\mathcal{D}_\mathcal{X} := \{b \in \Delta(\mathcal{S}) : \exists s \in \mathcal{X} \text{ such that } b = \mathbb{1}[s]\}$ . Then, the belief-reachability objective is  $\text{Belief-Reach}(\mathcal{X}) := \{\rho \in \Omega : \exists t \geq 0 B_t(\rho) \in \mathcal{D}_\mathcal{X}\}$ . Denote the belief-reachability value by  $\text{val}_{\text{BR}(\mathcal{X})}$ . We may omit  $\mathcal{X}$  in  $\text{val}_{\text{BR}(\mathcal{X})}$  if it is clear from the context.

**Remark 2** (Generality of belief-reachability). *In general POMDPs with reachability objectives, without loss of generality, target states can be considered absorbing, i.e., the dynamic remains in the same state once a target state is visited. Furthermore, one can reduce to the case where there is only one target state. These simplifications affect the dynamic, but neither optimal policies nor the reachability value. Belief-reachability generalizes reachability objectives as follows. Given a POMDP  $P$  with an absorbing target state  $s^*$ , consider a copy of  $P$  but add a new action and two absorbing states  $\top$  and  $\perp$ . After playing the new action, the state moves to  $\top$  if the state was in  $s^*$ , and to  $\perp$  otherwise. Moreover, the controller is announced which of these states is reached. Then, the belief-reachability value to  $\top$  coincides with the original reachability value. Indeed, for an approximately optimal policy of  $P$ , play the new action after sufficiently many steps to obtain approximately the same value in the belief-reachability objective.*

### 3.2 Overview of Results

**Results for Belief-Reachability Objectives** Our main results for belief-reachability objectives are the following, which are proved in Section 4.

**Theorem 1.** *Quantitative analysis for belief-reachability objectives for revealing POMDPs is in EXPTIME.*

By Remark 2, reachability objectives reduce to belief-reachability objectives. Therefore, they have the following consequence.

**Corollary 2.** *Quantitative analysis for reachability objectives for revealing POMDPs is in EXPTIME.*

**Results for Parity Objectives** Our main results for parity objectives are the following, which are proved in Section 5.

**Theorem 3.** *Quantitative analysis for parity objectives for revealing POMDPs is in EXPTIME.*

Theorem 3 generalizes Corollary 2 to parity objectives, and follows from a reduction of parity objectives to belief-reachability to a set that can be computed in EXPTIME, see Lemma 15.

**Theorem 4.** *Optimal policies exist for parity objectives for revealing POMDPs.*

The following results follow from (Belly et al. 2025, Theorem 3) and Theorem 4.

**Corollary 5.** *For revealing POMDPs with parity objectives, limit-sure and almost-sure winning coincide, and limit-sure analysis is in EXPTIME-complete.*

## 4 Belief-Reachability Objectives

In this section, we prove Theorem 1, i.e., that the quantitative analysis for belief-reachability objectives is in EXPTIME. We proceed in five steps:

- We introduce reliable actions, which preserve the current belief-reachability value. Proposition 6 proves that every belief has a reliable action.
- Lemma 7 shows that stopping policies that play only reliable actions are optimal.
- Lemma 8 shows that playing reliable actions uniformly at random is stopping.
- Lemma 10 upper bounds the horizon needed to approximate the belief-reachability value, using stopping optimal policies.
- Lemma 11 presents a point-based dynamic programming algorithm to approximate the finite-horizon belief-reachability value.

**Definition 3** (Reliable Action). *Consider a POMDP with belief-reachability objectives. An action  $a \in \mathcal{A}$  is reliable if it preserves the belief-reachability value. Formally, for each belief  $b \in \Delta(\mathcal{S})$ , define the set of reliable actions  $\mathcal{R}(b)$  by*

$$\mathcal{R}(b) := \{a \in \mathcal{A} : \mathbb{E}_b^a(\text{val}_{\text{BR}(\mathcal{X})}(B_1)) = \text{val}_{\text{BR}(\mathcal{X})}(b)\}.$$

Because the set of actions  $\mathcal{A}$  is finite, we have the following result.

**Proposition 6.** *The set of reliable actions is nonempty.*

**Definition 4** (Terminal State). Consider a POMDP and a set of target states  $\mathcal{X} \subseteq \mathcal{S}$ . A state  $s \in \mathcal{S}$  is called terminal (for  $\mathcal{X}$ ) if  $s \in \mathcal{X}$  or if  $\mathcal{X}$  cannot be reached from  $s$  with positive probability, i.e.,  $\sup_{\sigma \in \Sigma} \mathbb{P}_{\mathbb{1}[s]}^{\sigma}(\exists t \geq 0 \quad \text{supp}(B_t) \cap \mathcal{X} \neq \emptyset) = 0$ . The set of terminal states is denoted by  $\mathcal{T}$ .

**Definition 5** (Stopping Policy). Consider a POMDP, a set of target states  $\mathcal{X} \subseteq \mathcal{S}$ , and a corresponding set of terminal states  $\mathcal{T} \subseteq \mathcal{S}$ . For  $n \in \mathbb{N}$  and  $q > 0$ , a policy  $\sigma \in \Sigma$  is called  $(n, q)$ -stopping (for  $\mathcal{X}$ ) if within  $n$  steps the belief collapses to a Dirac mass on a terminal state with probability at least  $q$ . Formally, from every initial belief  $b \in \Delta(\mathcal{S})$ ,

$$\mathbb{P}_b^{\sigma}(\exists t \leq n \quad B_t \in \mathcal{D}_{\mathcal{T}}) \geq q.$$

We say that a policy is stopping if it is stopping for some  $n$  and  $q > 0$ .

Similar to (Andersson and Miltersen 2009, Lemma 5) that considers finite duration games, it is easy to see that, if a policy is stopping and uses only reliable actions, then it is optimal.

**Lemma 7.** Consider a POMDP with belief-reachability objectives. If a policy is stopping and uses only reliable actions, then it is optimal.

*Proof Sketch.* Consider a stopping policy  $\sigma \in \Sigma$  which uses only reliable actions. Because every action chosen by  $\sigma$  is reliable, the belief-reachability values  $\text{val}_{\text{BR}}(B_t)$  forms a martingale, i.e.,  $\text{val}_{\text{BR}}(b_0) = \mathbb{E}_{b_0}^{\sigma}(\text{val}_{\text{BR}}(B_t))$ . Moreover, since the policy  $\sigma$  is stopping, it reaches a Dirac belief on a terminal state in finite time with probability 1 and the hitting time  $\tau$  of  $\mathcal{D}_{\mathcal{T}}$  is almost-surely finite. Therefore, we have  $\text{val}_{\text{BR}}(b_0) = \mathbb{E}_{b_0}^{\sigma}(\text{val}_{\text{BR}}(B_{\tau}))$ .  $\square$

**Minimal Positive Transition** Denote  $\delta_{\min}$  the minimum non-zero probability in the transition function, i.e.,

$$\delta_{\min} := \min \left\{ \delta(s, a)(s', z) : \forall s, s' \in \mathcal{S}, a \in \mathcal{A}, z \in \mathcal{Z} \right. \\ \left. \delta(s, a)(s', z) > 0 \right\}.$$

**Lemma 8.** Consider a revealing POMDP with belief-reachability objectives. The policy  $\sigma$  that plays reliable actions uniformly at random is  $(n, q)$ -stopping with parameters  $n := |\mathcal{S}| + 2$  and  $q := \delta_{\min}^2 (\delta_{\min}/|\mathcal{A}|)^{|\mathcal{S}|}$ .

*Proof Sketch.* By Proposition 6, every belief has a reliable action, so the policy  $\sigma$  is well-defined. Build the layered sets  $\mathcal{S}_0 := \mathcal{T}$ ,  $\mathcal{S}_{t+1} := \{s : \exists a \in \mathcal{R}(\mathbb{1}[s]) \text{ with a positive-probability move to } \mathcal{S}_t\}$ . The sequence covers all states after at most  $|\mathcal{S}|$  iterations, i.e.,  $\mathcal{S}_{|\mathcal{S}|} = \mathcal{S}$ , because, if there were states outside, then they require to use unreliable actions to get to terminal states, which lowers their value and forms a contradiction. Under the policy  $\sigma$  that plays every reliable action uniformly, (a) the first “reveal” of the Dirac belief happens with probability  $\delta_{\min}$ ; (b) each step toward the next layer then succeeds with probability at least  $\delta_{\min}/|\mathcal{A}|$ ; and (c) the final “reveal” of the Dirac belief happens with probability  $\delta_{\min}$ . Hence, within  $n := |\mathcal{S}| + 2$  steps, a Dirac belief is reached on a terminal state with probability  $q := \delta_{\min}^2 (\delta_{\min}/|\mathcal{A}|)^{|\mathcal{S}|}$ , which proves that  $\sigma$  is  $(n, q)$ -stopping.  $\square$

We deduce directly from Lemma 7 and Lemma 8 the next result.

**Corollary 9.** Every revealing POMDP with belief-reachability objectives has an optimal policy.

We turn to the quantitative analysis for the belief-reachability value. Note that the existence of optimal stopping policies implies that this value can be approximated using a finite horizon.

**Lemma 10.** Consider a POMDP with belief-reachability objectives and an  $(n, q)$ -stopping optimal policy  $\sigma$ . For all  $\varepsilon > 0$ , the finite horizon  $T := n \left\lceil \frac{\log(\varepsilon)}{\log(1-q)} \right\rceil$  is such that

$$\mathbb{P}_{b_0}^{\sigma}(\exists t \leq T \quad B_t \in \mathcal{D}_{\mathcal{T}}) \geq \text{val}_{\text{BR}}(b_0) - \varepsilon.$$

*Proof Sketch.* Split the play into blocks of  $n$  steps. Because the optimal policy  $\sigma$  is  $(n, q)$ -stopping, each block reaches a Dirac belief on a terminal state with probability at least  $q$ . Thus for the hitting time  $\tau$  of  $\mathcal{D}_{\mathcal{T}}$  can be bounded by  $\mathbb{P}_{b_0}^{\sigma}(\tau > kn) \leq (1-q)^k$ . Therefore,  $\tau$  has a geometric tail and is almost surely finite. Choose  $T = n \left\lceil \frac{\log(\varepsilon)}{\log(1-q)} \right\rceil$  so that  $\mathbb{P}_{b_0}^{\sigma}(\tau > T) \leq \varepsilon$ , showing that a horizon of  $T$  steps suffices to approximate the belief-reachability value.  $\square$

**Reduction to Finite Horizon** To solve the quantitative analysis for revealing POMDPs with belief-reachability objectives, it suffices to compute the finite horizon value for a horizon that is exponentially large on the input size. Indeed, by Lemma 8, there exists an  $(n, q)$ -stopping optimal policy with  $n = |\mathcal{S}| + 2$  and  $q = (\delta_{\min})^2 (\delta_{\min}/|\mathcal{A}|)^{|\mathcal{S}|}$ . Hence, by Lemma 10, for every  $\varepsilon > 0$ , we can consider  $T = n \left\lceil \frac{\log(\varepsilon)}{\log(1-q)} \right\rceil = O\left(|\mathcal{S}||\mathcal{A}|^{|\mathcal{S}|} \log(1/\varepsilon) / \delta_{\min}^{|\mathcal{S}|+2}\right)$ , which is at most exponentially large in the input size.

**Previous Approaches to Finite Horizon** Finite horizon objectives have been studied by a long time. A naive approach to the quantitative analysis for belief-reachability objectives would take the time bound  $T$  in Lemma 10 and solve a (fully observable) MDP with  $|\mathcal{A} \times \mathcal{S}|^T$  states. This approach implies a 2EXPTIME complexity upper bound. A fundamental result states that, for horizons that are polynomially large with respect to the input, computing the finite horizon value of POMDPs is PSPACE-complete (Papadimitriou and Tsitsiklis 1987, Theorem 6, page 448). The technique, instead of listing explicitly all exponentially many histories, compactly represents them using nondeterminism and space proportional to the horizon. The conclusion follows from PSPACE being closed under nondeterminism by Savitch’s theorem. For exponentially large horizons, this technique leads to an EXSPACE upper bound. Instead, we prove an EXPTIME upper bound.

Point-based algorithms were introduced for POMDPs by (Pineau, Gordon, and Thrun 2003) as an alternative to approximating the value of POMDPs by considering a fixed subset of beliefs and projected belief updates. For a survey on point-based algorithms see (Shani, Pineau, and Kaplow 2013; Walraven and Spaan 2019), which focuses on finite horizon objectives but does not provide the complexity upper bound we require.

**Finite-Horizon Belief-Reachability** Let  $T \in \mathbb{N}$  be a finite horizon. The  $T$ -step belief-reachability value to  $\mathcal{X} \subseteq \mathcal{S}$  is defined by  $\text{val}_{\text{BR}(\mathcal{X})}^T(b) := \max_{\sigma \in \Sigma} \mathbb{P}_b^\sigma(\exists t \leq T \ B_t \in \mathcal{D}_{\mathcal{X}})$ .

**Lemma 11.** *Approximating the  $T$ -step belief-reachability value of revealing POMDPs up to an additive error of  $\varepsilon > 0$  can be computed in exponential time, formally, in time  $O(T^{|\mathcal{S}|} |\mathcal{S}|^{|\mathcal{S}+1|} |\mathcal{A}| |\mathcal{Z}| / \varepsilon^{(|\mathcal{S}|-1)})$ .*

*Proof Sketch.* We approximate the  $T$ -step belief-reachability value by discretizing the belief simplex and applying projected Bellman updates over this grid. Define a  $k$ -uniform grid  $G_k \subseteq \Delta(\mathcal{S})$  so that every belief  $b \in \Delta(\mathcal{S})$  is at  $L_1$ -distance at most  $\varepsilon/(T+1)$  from some grid point  $\Pi_k(b)$ , where  $k = \lceil (T+1)|\mathcal{S}|/\varepsilon \rceil$ . Then, compute the Bellman updates on grid points, carefully projecting back to  $G_k$  after each update, for  $T$  steps. By induction on the horizon, the error at each step accumulates by at most  $\varepsilon/(T+1)$ , yielding a total error at most  $\varepsilon$ . The grid has size  $O((T|\mathcal{S}|/\varepsilon)^{|\mathcal{S}|-1})$ , and each Bellman update takes  $O(|\mathcal{S}|^2 |\mathcal{A}| |\mathcal{Z}|)$  time, so the total running time is  $O(T^{|\mathcal{S}|} |\mathcal{S}|^{|\mathcal{S}+1|} |\mathcal{A}| |\mathcal{Z}| \varepsilon^{-(|\mathcal{S}|-1)})$ , which is exponential in the input size.  $\square$

The following result follows from Lemmas 10 and 11.

**Theorem 1.** *Quantitative analysis for belief-reachability objectives for revealing POMDPs is in EXPTIME.*

## 5 Parity Objectives

### 5.1 Reduction to Belief-Reachability Objectives

In the sequel, we show that, for revealing POMDPs, the parity objective value coincides with the belief-reachability to the set of Dirac on the almost-sure winning parity states. This proof uses the standard notion of *end components* in MDPs, introduced in (Alfaro 1998); see also (Baier and Katoen 2008, Section 10.6.3).

**Underlying MDPs of POMDPs** For a POMDP  $P = (\mathcal{S}, \mathcal{A}, \mathcal{Z}, \delta, b_0)$ , we define its underlying MDP as  $M = (\mathcal{S} \times (\mathcal{Z} \cup \{\square\}), \mathcal{A}, \tilde{\delta}, \tilde{b}_0)$ , where the transition function is defined as, for all  $s, s' \in \mathcal{S}, z \in \mathcal{Z} \cup \{\square\}, z' \in \mathcal{Z}$ ,

$$\tilde{\delta}((s, z), a)((s', z')) := \delta(s, a)(s', z'),$$

and the initial belief is defined as  $\tilde{b}_0((s, \square)) := b_0(s)$  for all states  $s \in \mathcal{S}$ . This construction captures the entire dynamic of  $P$ . We use this notion to analyze the policies derived from the original POMDP.

**End Components** Let  $M = (\mathcal{S}, \mathcal{A}, \delta, b_0)$  be an MDP. An end component is a pair  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$  where  $\mathcal{Q} \subseteq \mathcal{S}$  is a subset of states and  $\mathcal{E}: \mathcal{Q} \rightrightarrows \mathcal{A}$  assigns each state to a set of nonempty actions such that

- *Closedness.* For all states  $q \in \mathcal{Q}$  and actions  $a \in \mathcal{E}(q)$ , we have  $\text{supp}(\delta(q, a)) \subseteq \mathcal{Q}$ ; and
- *Strong connectivity.* The directed graph with states  $\mathcal{Q}$  and edges  $(q, q')$  where there exists  $a \in \mathcal{E}(q)$  such that  $\delta(q, a)(q') > 0$  is strongly connected.

A play  $\rho$  visits infinitely often an end component  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$ , denoted by  $\mathcal{I}(\rho) = (\mathcal{Q}, \mathcal{E})$ , if

- The set of states visited infinitely often in  $\rho$  is  $\mathcal{Q}$ ; and
- For all state  $q \in \mathcal{Q}$ , the set of actions played infinitely often when in  $q$  is  $\mathcal{E}(q)$ .

We say  $\mathcal{I}(\rho) \subseteq (\mathcal{Q}, \mathcal{E})$  if  $\mathcal{I}(\rho) = (\tilde{\mathcal{Q}}, \tilde{\mathcal{E}})$ ,  $\tilde{\mathcal{Q}} \subseteq \mathcal{Q}$  and  $\tilde{\mathcal{E}} \subseteq \mathcal{E}$ .

The following statement is a fundamental result of end components in MDPs.

**Lemma 12** ((Alfaro 1998, Theorems 3.1 and 3.2)). *For MDPs, the following assertions hold.*

- *For every end component  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$  and state  $q \in \mathcal{Q}$ , there exists a policy  $\sigma$  such that*

$$\mathbb{P}_{\mathbb{1}[q]}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) = 1.$$

- *For all policies  $\sigma$ , we have*

$$\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) \text{ is an end component}) = 1.$$

The following result relates end components in MDPs to policies in general POMDPs.

**Lemma 13.** *Consider a POMDP  $P$  and its underlying MDP  $M$ . For every reachable end component  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$  of  $M$ , i.e., there exists a policy  $\sigma$  on  $P$  such that  $\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) > 0$ , we have that, for all states  $q \in \mathcal{Q}$ , actions  $a \in \mathcal{E}(q)$ , and signals  $z \in \mathcal{Z}$ , there exists a policy  $\sigma_{q,a,z}$  on  $P$  such that*

$$\mathbb{P}_b^{\sigma_{q,a,z}}(\mathcal{I}(\rho) \subseteq \mathcal{U}) = 1,$$

where  $b$  is the belief after starting with belief  $\mathbb{1}[q]$ , playing action  $a$ , and receiving signal  $z$ .

*Proof Sketch.* Consider a reachable end component  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$  of  $M$  with corresponding policy  $\sigma$ , i.e.,  $\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) > 0$ . On the event  $\{\mathcal{I}(\rho) = \mathcal{U}\}$ , every state-action pair  $(q, a)$  with  $q \in \mathcal{Q}$  and  $a \in \mathcal{E}(q)$  is taken infinitely often, and each such occurrence results in at least one signal  $z$  that occurs with positive probability. Fix  $q \in \mathcal{Q}$ ,  $a \in \mathcal{E}(q)$  and a signal  $z$  that can follow  $(q, a)$ . By contradiction, assume that from the belief  $b$  obtained after starting with  $\mathbb{1}[q]$ , playing action  $a$ , and receiving signal  $z$ , no policy can stay inside  $\mathcal{U}$  almost surely. Then, starting from  $b$ , any policy has a positive, state-independent strictly positive lower bound on the probability of leaving  $\mathcal{Q}$  within a bounded number of steps. Because  $(q, a)$  is visited infinitely often on the event  $\{\mathcal{I}(\rho) = \mathcal{U}\}$ , these lower-bounded leaving probabilities accumulate, driving the probability of ever leaving  $\mathcal{Q}$  to 1. This contradicts  $\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) > 0$ . Hence, such a leaving bound cannot exist: there must be a policy  $\sigma_{q,a,z}$  that, from  $b$ , keeps the play inside  $\mathcal{Q}$  (and therefore inside  $\mathcal{U}$ ) with probability 1.  $\square$

The following example demonstrates that the above result is tight in the sense that it can not guarantee a policy  $\sigma_{q,a,z}$  such that  $\mathbb{P}_b^{\sigma_{q,a,z}}(\mathcal{I}(\rho) = \mathcal{U}) = 1$ .

**Example 1.** *Consider the example presented in (Chatterjee, Saona, and Ziliotto 2021, Example 4.3) of a POMDP with only one possible signal in Figure 1, commonly known as a blind MDP. It has four states  $\mathcal{S} = \{s_0, s_1, \perp, \top\}$  and two*

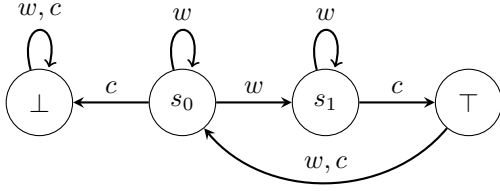


Figure 1: Example of a classic general POMDP that requires infinite memory policies for parity objectives. Edges represent a positive probability transition between states when the corresponding action in its label is used.

actions: wait ( $w$ ) and commit ( $c$ ), i.e.,  $\mathcal{A} = \{w, c\}$ . The priority function is defined as

$$\text{pri}(s_0) := 1, \text{pri}(s_1) := 1, \text{pri}(\perp) := 1, \text{pri}(\top) = 0.$$

The transitions are defined as follows. (a) Under action  $w$ , state  $s_0$  moves either to  $s_1$  or loops, both with probability  $1/2$ ; states  $s_1$  and  $\perp$  loop; and state  $\top$  moves to  $s_0$ . (b) Under action  $c$ , state  $s_0$  moves to  $\perp$ ; state  $s_1$  moves to  $\top$ ; state  $\perp$  loops; and state  $\top$  moves to  $s_0$ . The only end component of the underlying MDP that satisfies the parity condition is

$$\mathcal{U} = (\{s_1, \top\}, \{(s_0, w), (s_1, w), (s_1, c), (\top, w), (\top, c)\}).$$

Note that, for every  $\varepsilon$ -optimal policy  $\sigma$ , we have that  $\mathbb{P}_{\mathbb{1}[s_0]}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) \geq 1 - \varepsilon$ . This is achieved only by policies that wait long enough before committing, which requires unbounded memory. Conversely, no policy can guarantee  $\mathbb{P}_{\mathbb{1}[s_0]}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) = 1$ . Indeed, whenever action  $c$  is taken, the belief on state  $s_0$  is positive and therefore the state absorbs at  $\perp$  with positive probability. The best that can be achieved almost surely is to stay forever in the end component  $\tilde{\mathcal{U}} = (\{s_1\}, \{(s_1, w)\})$ , whose minimal priority is 1 and therefore violates the parity condition.  $\square$

We now show that in revealing POMDPs, the above result can be extended to guarantee a policy  $\sigma_q$  such that  $\mathbb{P}_{\mathbb{1}[q]}^{\sigma_q}(\mathcal{I}(\rho) = \mathcal{U}) = 1$ .

**Lemma 14.** Consider a revealing POMDP  $P$  and its underlying MDP  $M$ . For every policy  $\sigma$  on  $P$  and end component  $\mathcal{U} = (\mathcal{Q}, \mathcal{E})$  of  $M$  such that  $\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) > 0$ , we have, for all states  $q \in \mathcal{Q}$ , there exists a policy  $\sigma_q$  on  $P$  such that

$$\mathbb{P}_{\mathbb{1}[q]}^{\sigma_q}(\mathcal{I}(\rho) = \mathcal{U}) = 1.$$

*Proof Sketch.* Consider a policy  $\sigma$  that  $\mathbb{P}_{b_0}^\sigma(\mathcal{I}(\rho) = \mathcal{U}) > 0$ . Fix a state  $q \in \mathcal{Q}$ . We construct a policy  $\sigma_q$  that proceeds in phases, one for each ordered pair  $(\tilde{q}, \tilde{a})$  where  $\tilde{q} \in \mathcal{Q}$  and  $\tilde{a} \in \mathcal{E}(\tilde{q})$ . The goal of a phase is to reach  $\tilde{q}$  and to play action  $\tilde{a}$ . When the current belief is non-Dirac, the policy follows the almost-sure safety policy from Lemma 13 that stays inside  $\mathcal{U}$  until a Dirac belief  $\mathbb{1}[q']$  is reached, which is guaranteed in finite time because  $P$  is revealing. From  $q'$ , the policy moves inside  $\mathcal{U}$  along a fixed finite path of “safe” actions to the target state  $\tilde{q}$ ; if the belief becomes non-Dirac, the policy returns to the safety policy and retries. Once the belief is  $\mathbb{1}[\tilde{q}]$ , the policy plays the action  $\tilde{a}$ ; upon observing

the resulting signal  $z$ , switches to the safety policy and starts the next phase. Because every ordered pair  $(q, a) \in \mathcal{E}$  is visited infinitely often, each edge of  $\mathcal{U}$  is taken infinitely often. The safety policies ensure the play never leaves  $\mathcal{Q}$ . Consequently, the event  $\{\mathcal{I}(\rho) = \mathcal{U}\}$  is satisfied with probability 1 under  $\sigma_q$  when starting from the belief  $\mathbb{1}[q]$ .  $\square$

We are ready to present the reduction of parity to belief-reachability of the almost-sure winning states for parity.

**Lemma 15.** Consider a revealing POMDP with parity objectives. Denote the set of states for which, if the initial belief were a Dirac on that state, then the parity condition can be satisfied almost-surely by  $\mathcal{X} := \{s \in \mathcal{S} : \exists \sigma \in \Sigma \ \mathbb{P}_{\mathbb{1}[s]}^\sigma(\text{Parity}) = 1\}$ . Then, the parity value coincides with the belief-reachability value to  $\mathcal{X}$ , i.e., for all beliefs  $b \in \Delta(\mathcal{S})$ ,

$$\text{val}_P(b) = \text{val}_{\text{BR}(\mathcal{X})}(b).$$

*Proof sketch.* Consider a policy  $\sigma \in \Sigma$ . Consider its end components in the underlying MDP on  $\mathcal{S} \times \mathcal{Z}$ . Note that each end component either does satisfy or does not satisfy the parity condition. Moreover, if they satisfy the parity condition, then  $\sigma$  is an almost-sure winning policy starting from any state inside the end component. Therefore, the probability of satisfying the parity condition under  $\sigma$  corresponds to the belief-reachability to the end components where the parity condition is satisfied. In particular, in these end components, parity condition is satisfied almost-surely. We conclude because the policy is arbitrary.  $\square$

## 5.2 Proofs of Main Results

Building on Lemma 15, which reduces parity to belief-reachability, and the EXPTIME procedure in Theorem 1 for approximating the belief-reachability value, we now (a) obtain an EXPTIME algorithm for approximating parity values (Theorem 3); and (b) show that limit-sure and almost-sure winning coincide (Theorem 4).

**Theorem 3.** Quantitative analysis for parity objectives for revealing POMDPs is in EXPTIME.

*Proof Sketch.* By (Belly et al. 2025, Theorem 4), computing almost-sure winning parity states in the revealing POMDP is EXPTIME. Therefore, this result is a direct implication of Lemma 15 and Theorem 1.  $\square$

**Theorem 4.** Optimal policies exist for parity objectives for revealing POMDPs.

*Proof Sketch.* It is a direct implication of Corollary 9 and Lemma 15.  $\square$

**Concluding Remarks** In this work, we consider revealing POMDPs which have been studied in the literature and provide decidability results for the fundamental computational problems for this model. Interesting directions for future work include exploring the practical applicability of our algorithms and extending our decidability results to other classes of POMDPs.

## Acknowledgements

This work was partially supported by the ANRT under the French CIFRE Ph.D program in collaboration between NyxAir and Paris-Dauphine University (Contract: CIFRE N° 2022/0513), by the French Agence Nationale de la Recherche (ANR) under reference ANR-21-CE40-0020 (CONVERGENCE project), by Austrian Science Fund (FWF) 10.55776/COE12, and by the ERC CoG 863818 (ForM-SMArt) grant.

## References

- Alfaro, L. 1998. *Formal Verification of Probabilistic Systems*. Ph.D. diss., Stanford University, Stanford, CA, USA.
- Andersson, D.; and Miltersen, P. B. 2009. The Complexity of Solving Stochastic Games on Graphs. In *Algorithms and Computation*, volume 5878, 112–121. Berlin, Heidelberg: Springer.
- Avrachenkov, K.; Dhiman, M.; and Kavitha, V. 2025. Constrained Average-Reward Intermittently Observable MDPs. arXiv:2504.13823.
- Baier, C.; Bertrand, N.; and Größer, M. 2008. On Decision Problems for Probabilistic Büchi Automata. In *Foundations of Software Science and Computational Structures*, volume 4962, 287–301. Berlin, Heidelberg: Springer.
- Baier, C.; Größer, M.; and Bertrand, N. 2012. Probabilistic  $\omega$ -automata. *Journal of the ACM (JACM)*, 59(1): 1–52.
- Baier, C.; and Katoen, J.-P. 2008. *Principles of Model Checking*. Cambridge, MA, USA: MIT press.
- Belly, M.; Fijalkow, N.; Gimbert, H.; Horn, F.; Pérez, G. A.; and Vandenhoove, P. 2025. Revelations: A Decidable Class of POMDPs with Omega-Regular Objectives. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(25): 26454–26462.
- Bertsekas, D. P. 1976. *Dynamic Programming and Stochastic Control*. New York: Academic Press.
- Billingsley, P. 2012. *Probability and Measure*. Hoboken, NJ, USA: Wiley.
- Bonet, B.; and Geffner, H. 2009. Solving POMDPs: RTDP-Bel vs. Point-based Algorithms. In *Proceedings of the 21st International Joint Conference on Artificial Intelligence, IJCAI'09*, 1641–1646. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Chatterjee, K.; Chmelik, M.; Gupta, R.; and Kanodia, A. 2015. Qualitative Analysis of POMDPs with Temporal Logic Specifications for Robotics Applications. In *IEEE International Conference on Robotics and Automation (ICRA)*, 325–330. Seattle, WA, USA: IEEE.
- Chatterjee, K.; Chmelik, M.; and Tracol, M. 2016. What is Decidable about Partially Observable Markov Decision Processes with  $\omega$ -Regular Objectives. *Journal of Computer and System Sciences*, 82(5): 878–911.
- Chatterjee, K.; Doyen, L.; Gimbert, H.; and Henzinger, T. A. 2010. Randomness for Free. In *Proceedings of the 35th International Conference on Mathematical Foundations of Computer Science*, 246–257. Berlin, Heidelberg: Springer.
- Chatterjee, K.; Doyen, L.; and Henzinger, T. A. 2010. Qualitative Analysis of Partially-Observable Markov Decision Processes. In *International Symposium on Mathematical Foundations of Computer Science*, 258–269. Berlin, Heidelberg: Springer.
- Chatterjee, K.; Doyen, L.; Raskin, J.-F.; and Sankur, O. 2025. The Value Problem for Multiple-Environment MDPs with Parity Objective. In *52nd International Colloquium on Automata, Languages, and Programming (ICALP 2025)*, volume 334 of *Leibniz International Proceedings in Informatics (LIPIcs)*, 150:1–150:17. Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik.
- Chatterjee, K.; and Henzinger, T. A. 2010. Probabilistic Automata on Infinite Words: Decidability and Undecidability Results. In *International Symposium on Automated Technology for Verification and Analysis*, 1–16. Berlin, Heidelberg: Springer.
- Chatterjee, K.; Lurie, D.; Saona, R.; and Ziliotto, B. 2024. Ergodic Unobservable MDPs: Decidability of Approximation. arXiv:2405.12583.
- Chatterjee, K.; Saona, R.; and Ziliotto, B. 2021. Finite-Memory Strategies in POMDPs with Long-Run Average Objectives. *Mathematics of Operations Research*, 47(1): 100–119.
- Chatterjee, K.; and Tracol, M. 2012. Decidable Problems for Probabilistic Automata on Infinite Words. In *27th Annual IEEE Symposium on Logic in Computer Science*, 185–194. Dubrovnik, Croatia: IEEE.
- Chen, F.; Wang, H.; Xiong, C.; Mei, S.; and Bai, Y. 2023. Lower Bounds for Learning in Revealing POMDPs. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *ICML'23*, 5104–5161. Honolulu, Hawaii, USA: JMLR.org.
- Chen, G.; and Liew, S.-C. 2023. Intermittently Observable Markov Decision Processes. arXiv:2302.11761.
- de Alfaro, L.; Faella, M.; Majumdar, R.; and Raman, V. 2005. Code Aware Resource Management. In *Proceedings of the 5th ACM International Conference on Embedded Software, EMSOFT '05*, 191–202. New York, NY, USA: Association for Computing Machinery.
- Durbin, R.; Eddy, S. R.; Krogh, A.; and Mitchison, G. 1998. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge, UK: Cambridge University Press.
- Fijalkow, N.; Gimbert, H.; Kelmendi, E.; and Oualhadj, Y. 2015. Deciding the Value 1 Problem for Probabilistic Leaktight Automata. *Logical Methods in Computer Science*, 11.
- Filar, J.; and Vrieze, K. 1997. *Competitive Markov Decision Processes*. New York, NY, USA: Springer.
- Gimbert, H.; and Oualhadj, Y. 2010. Probabilistic Automata on Finite Words: Decidable and Undecidable Problems. In *Automata, Languages and Programming*, volume 6199, 527–538. Berlin, Heidelberg: Springer.
- Kaelbling, L. P.; Littman, M. L.; and Cassandra, A. R. 1998. Planning and Acting in Partially Observable Stochastic Domains. *Artificial Intelligence*, 101(1-2): 99–134.

- Kaelbling, L. P.; Littman, M. L.; and Moore, A. W. 1996. Reinforcement Learning: A Survey. *Journal of Artificial Intelligence Research*, 4: 237–285.
- Kechris, A. S. 1995. *Classical Descriptive Set Theory*. New York, NY, USA: Springer.
- Kress-Gazit, H.; Fainekos, G. E.; and Pappas, G. J. 2009. Temporal-Logic-Based Reactive Mission and Motion Planning. *IEEE Transactions on Robotics*, 25(6): 1370–1381.
- Liu, Q.; Chung, A.; Szepesvari, C.; and Jin, C. 2022. When Is Partially Observable Reinforcement Learning Not Scary? In *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, 5175–5220. PMLR.
- Madani, O.; Hanks, S.; and Condon, A. 2003. On the Undecidability of Probabilistic Planning and Related Stochastic Optimization Problems. *Artificial Intelligence*, 147(1-2): 5–34.
- Papadimitriou, C. H.; and Tsitsiklis, J. N. 1987. The Complexity of Markov Decision Processes. *Mathematics of Operations Research*, 12(3): 441–450.
- Paz, A. 1971. *Introduction to Probabilistic Automata*. New York, NY, USA: Academic Press.
- Pineau, J.; Gordon, G.; and Thrun, S. 2003. Point-Based Value Iteration: An Anytime Algorithm for POMDPs. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence*, 1025–1030. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
- Pogosyants, A.; Segala, R.; and Lynch, N. 2000. Verification of the Randomized Consensus Algorithm of Aspnes and Herlihy: a Case Study. *Distributed Computing*, 13(3): 155–186.
- Puterman, M. L. 2014. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Hoboken, NJ, USA: Wiley.
- Rabin, M. O. 1963. Probabilistic Automata. *Information and control*, 6(3): 230–245.
- Shani, G.; Pineau, J.; and Kaplow, R. 2013. A Survey of Point-Based POMDP Solvers. *Autonomous Agents and Multi-Agent Systems*, 27(1): 1–51.
- Thomas, W. 1997. Languages, Automata, and Logic. In *Handbook of Formal Languages: Volume 3 Beyond Words*, 389–455. Berlin, Heidelberg: Springer.
- Van Der Vegt, M.; Jansen, N.; and Junges, S. 2023. Robust Almost-Sure Reachability in Multi-Environment MDPs. In *Tools and Algorithms for the Construction and Analysis of Systems*, volume 13993, 508–526. Berlin, Heidelberg: Springer.
- Walraven, E.; and Spaan, M. T. J. 2019. Point-Based Value Iteration for Finite-Horizon POMDPs. *Journal of Artificial Intelligence Research*, 65: 307–341.