

# 1 The Value Problem for Multiple-Environment 2 MDPs with Parity Objective

3 **Krishnendu Chatterjee** 

4 IST Austria

5 **Laurent Doyen** 

6 CNRS & LMF, ENS Paris-Saclay, France

7 **Jean-François Raskin** 

8 Université Libre de Bruxelles, Belgium

9 **Ocan Sankur** 

10 Université de Rennes, CNRS, Inria, France & Mitsubishi Electric R&D Centre Europe, France

## 11 — Abstract —

12 We consider multiple-environment Markov decision processes (MEMDP), which consist of a finite  
13 set of MDPs over the same state space, representing different scenarios of transition structure and  
14 probability. The value of a strategy is the probability to satisfy the objective, here a parity objective,  
15 in the worst-case scenario, and the value of an MEMDP is the supremum of the values achievable by  
16 a strategy.

17 We show that deciding whether the value is 1 is a PSPACE-complete problem, and even in P  
18 when the number of environments is fixed, along with new insights to the almost-sure winning  
19 problem, which is to decide if there exists a strategy with value 1. Pure strategies are sufficient for  
20 these problems, whereas randomization is necessary in general when the value is smaller than 1. We  
21 present an algorithm to approximate the value, running in double exponential space. Our results are  
22 in contrast to the related model of partially-observable MDPs where all these problems are known  
23 to be undecidable.

24 **2012 ACM Subject Classification** Theory of computation → Logic and verification; Theory of  
25 computation → Probabilistic computation

26 **Keywords and phrases** Markov decision processes, imperfect information, randomized strategies,  
27 limit-sure winning

28 **Digital Object Identifier** 10.4230/LIPIcs...

29 **Funding** *Krishnendu Chatterjee*: ERC CoG 863818 (ForM-SMArt) and Austrian Science Fund  
30 (FWF) 10.55776/COE12;

31 *Jean-François Raskin*: PDR Weave project FORM-LEARN-POMDP funded by FNRS and DFG,  
32 and the support of the Fondation ULB;

33 *Ocan Sankur*: ANR BisoUS (ANR-22-CE48-0012) and ANR EpiRL (ANR-22-CE23-0029).

## 34 **1** Introduction

35 We consider Markov decision processes (MDP), a well-established state-transition model  
36 for decision making in a stochastic environment. The decisions involve choosing an action  
37 from a finite set, which together with the current state determine a probability distribution  
38 over the successor state. The question of constructing a strategy that maximizes the  
39 probability to satisfy a logical specification is a classical synthesis problem with a wide range  
40 of applications [19, 11, 3, 20].

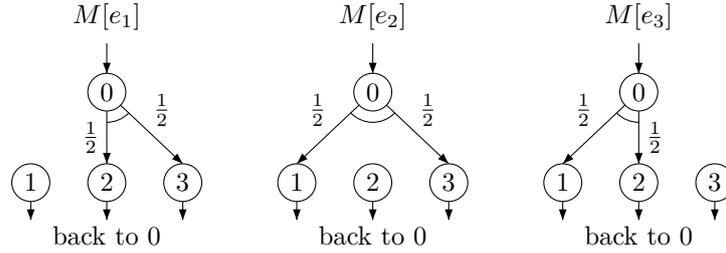
41 The stochastic transitions in MDPs capture the uncertainty in the effect of an action.  
42 Another form of uncertainty arises when the states are (partially) hidden to the decision-  
43 maker, as in the classical model of partially-observable MDPs (POMDP) [15, 18]. Recently,



© Krishnendu Chatterjee and Laurent Doyen and Jean-François Raskin and Ocan Sankur;  
licensed under Creative Commons License CC-BY 4.0

Leibniz International Proceedings in Informatics

**LIPICs** Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



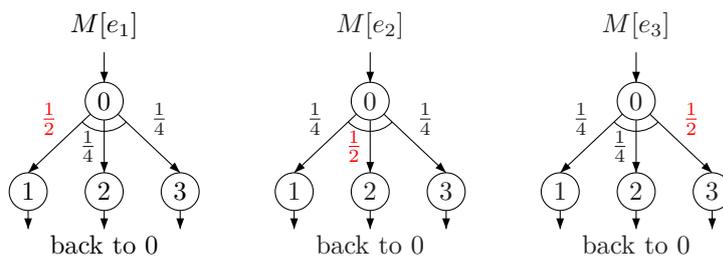
■ **Figure 1** Multiple-environment MDP for the missing card (over 3-card deck). Each  $M[e_i]$  represents the behavior of the MEMDP under environment  $e_i$  where card  $i$  has been removed. The environment can be identified almost-surely (with probability 1).

44 an alternative model of MDPs with partial information has attracted attention, the multiple-  
 45 environment MDPs (MEMDP) [21], which consists of a finite set of MDPs over the same  
 46 state space. Each MDP represents a possible environment, but the decision-maker does not  
 47 know in which environment they are operating. The synthesis problem is then to construct a  
 48 single strategy that can be executed in all environments to ensure the objective be satisfied  
 49 independently of the environment. This model is natural in applications where the structure  
 50 of the transitions and their probability are uncertain such as in robust planning or population  
 51 models with individual variability [4, 6, 1, 23, 22].

52 In contrast to what previous work suggest, the two models of POMDP and MEMDP  
 53 are (syntactically) incomparable: the choice of the environment in MEMDP is adversarial,  
 54 which cannot be expressed in a POMDP, and the partial observability of POMDP can occur  
 55 throughout the execution, whereas the uncertainty in MEMDP is only initial. In particular,  
 56 MEMDP are *not* a subclass of POMDP since pure strategies are sufficient in POMDPs [17, 7]  
 57 while randomization is necessary in general in MEMDPs [21, Lemma 3].

58 The synthesis problem has been considered for traditional  $\omega$ -regular objectives, defined  
 59 as parity [21] or Rabin [22] condition, in three variants: the almost-sure problem is to decide  
 60 whether there exists a strategy that is winning with probability 1 in all environments, the  
 61 limit-sure problem is to decide whether, for every  $\varepsilon > 0$ , there exists a strategy that is  
 62 winning with probability at least  $1 - \varepsilon$  in all environments, and the gap problem, which is  
 63 an approximate version of the quantitative problem to decide, given a threshold  $0 < \lambda \leq 1$ ,  
 64 whether there exists a strategy that is winning with probability at least  $\lambda$  in all environments.  
 65 The limit-sure problem is also called the value-1 problem, where the value of an MEMDP is  
 66 defined as the supremum of the values achievable by a strategy. The value is 1 if and only if  
 67 the answer to the limit-sure problem is Yes.

68 A classical example to illustrate the difference between almost-sure and limit-sure winning  
 69 is to consider an environment consisting of 51 cards, obtained by removing one card from a  
 70 standard 52-card deck (see Figure 1). The decision-maker has two possible actions: the action  
 71 *sample* reveals the top card of the deck and then shuffles the cards (including the top card,  
 72 which remains in the deck); the action *guess*( $x$ ), where  $1 \leq x \leq 52$  is a card, stops the game  
 73 with a win if  $x$  is the missing card, and a lose otherwise. If no guess is ever made, the game  
 74 is also losing. An almost-sure winning strategy is to sample until each of the 51 cards has  
 75 been revealed at least once, then to make a correct guess. It is easy to see that the strategy  
 76 wins with probability 1, even if there exist scenarios (though with probability 0) where some  
 77 of the 51 cards are never revealed and no correct guess is made. Hence the MEMDP is  
 78 almost-sure winning, and we say that it is not sure winning because a losing scenario exists



**Figure 2** Multiple-environment MDP for the duplicate card (over 3-card deck). Each  $M[e_i]$  represents the behavior of the MEMDP under environment  $e_i$  where card  $i$  has been duplicated. The environment can be identified limit-surely (with probability arbitrarily close to 1).

79 in every strategy. Consider now an environment consisting of 53 cards, obtained by adding  
 80 one duplicate card  $c$  to the standard deck, and the same action set and rules of the game,  
 81 except that a correct guess is now the duplicate card  $x = c$  (see Figure 2). The strategy  
 82 that samples for a long time and then makes a guess based on the most frequent card wins  
 83 with probability close to 1 – and closer to 1 as the sampling time is longer – but not equal  
 84 to 1, since no matter how long is the sampling phase there is always a nonzero probability  
 85 that the duplicate card does not have the highest frequency at the time of the guess. In  
 86 this case, the MEMDP is limit-sure winning, but not almost-sure winning. Intuitively, the  
 87 solution of almost-sure winning relies on the analysis of *revealing* transitions, which give a  
 88 sure information allowing to exclude some environment (seeing card  $c$  is a guarantee that we  
 89 are not in the environment where  $c$  is missing); the solution of limit-sure winning involves  
 90 *learning* by sampling, which also allows to exclude some environment, but possibly with a  
 91 nonzero probability to be mistaken.

92 For MEMDPs with two environments, it is known that the almost-sure and limit-sure  
 93 problem for parity objectives are solvable in polynomial time [21, Theorem 33, Theorem 40],  
 94 while the gap problem is decidable in 2-fold exponential space [21, Theorem 30] and is  
 95 NP-hard, even for acyclic MEMDPs with two environments [21, Theorem 26]. With an  
 96 arbitrary number of environments, the almost-sure problem becomes PSPACE-complete [22,  
 97 Theorem 41], even for reachability objectives [23, Lemma 11]. For comparison, in the close  
 98 model of POMDP, the decidability frontier lies between limit-sure winning and almost-sure  
 99 winning: with reachability objectives, the almost-sure problem is decidable (and EXPTIME-  
 100 complete [2]), whereas the limit-sure problem is undecidable [12]. The gap problem is also  
 101 undecidable [16].

102 In this paper, we consider the limit-sure problem and the gap problem for parity objectives  
 103 in MEMDPs with an arbitrary number of environments. Our main result is to show that  
 104 (a) the limit-sure problem is PSPACE-complete and can be solved in polynomial time for a  
 105 fixed number of environments, and (b) the gap problem can be solved in double exponential  
 106 space. Correspondingly, our algorithms significantly extend the solutions that are known for  
 107 two environments, relying on a non-trivial recursive (inductive) analysis.

108 The PSPACE upper bound is obtained by a characterization of limit-sure winning for  
 109 a subclass of MEMDPs, in terms of almost-sure winning conditions (Lemma 14). A pre-  
 110 processing phase transforms general MEMDPs into the subclass. We present a PSPACE  
 111 algorithm to compute the pre-processing and verify the characterization. Since our algorithm  
 112 relies on almost-sure winning, we also give a new characterization of almost-sure winning  
 113 for parity objectives in MEMDPs (Lemma 3), which gives a conceptually simple alternative

114 algorithm to the known solution [22]. The PSPACE lower bound straightforwardly follows  
 115 from the same reduction as for almost-sure winning [22, Theorem 7]. A corollary of our  
 116 characterizations is a refined strategy complexity: pure (non-randomized) strategies are  
 117 sufficient for both limit-sure and almost-sure winning, which was known only for acyclic  
 118 MEMDPs and almost-sure reachability objectives [23, Lemma 12], and exponential memory  
 119 is sufficient. In the last part of the paper, we present an algorithm running in double  
 120 exponential space for solving the gap problem, by computing an approximation of the value  
 121 of the MEMDP. To win with probability at least  $\lambda$  in all environments, randomized strategies  
 122 are more powerful [21, Lemma 3], and thus need to be considered for solving the gap problem.

123 In conclusion, the model of MEMDP is a valuable alternative to POMDPs, from a  
 124 theoretical perspective since the limit-sure problem and gap problem are undecidable for  
 125 POMDPs whereas our results establish decidability for MEMDPs, and from a practical  
 126 perspective since many applications of POMDPs can be expressed by MEMDPs, as was  
 127 observed previously [1, 23].

## 128 2 Definitions

129 A *probability distribution* on a finite set  $Q$  is a function  $d : Q \rightarrow [0, 1]$  such that  $\sum_{q \in Q} d(q) = 1$ .  
 130 The support of  $d$  is  $\text{Supp}(d) = \{q \in Q \mid d(q) > 0\}$ . A Dirac distribution assigns probability 1  
 131 to some  $q \in Q$ . We denote by  $\mathcal{D}(Q)$  the set of all probability distributions on  $Q$ .

### 132 2.1 Markov Decision Processes

133 A *Markov decision process (MDP)* over a finite set  $A$  of actions is a tuple  $M = \langle Q, (A_q)_{q \in Q}, \delta \rangle$   
 134 consisting of a finite set  $Q$  of *states*, a nonempty set  $A_q \subseteq A$  of actions for each state  $q \in Q$ ,  
 135 and a partial probabilistic transition function  $\delta : Q \times A \rightarrow \mathcal{D}(Q)$ . We say that  $(q, a, q')$  is a  
 136 transition if  $\delta(q, a)(q') > 0$ . A state  $q \in Q$  is a *sink* if  $\delta(q, a)(q) = 1$  for all  $a \in A_q$ .

137 A *run* of  $M$  from an initial state  $q_0 \in Q$  is an infinite sequence  $\pi = q_0 a_0 q_1 a_1 \dots$  of  
 138 interleaved states and actions such that  $a_i \in A_{q_i}$  and  $\delta(q_i, a_i)(q_{i+1}) > 0$  for all  $i \geq 0$ . Finite  
 139 prefixes  $\rho = q_0 a_0 \dots q_n$  of runs ending in a state are called *histories* and we denote by  
 140  $\text{last}(\rho) = q_n$  the last state of  $\rho$ . We denote by  $\text{Hist}^\omega(M)$  (resp.,  $\text{Hist}(M)$ ) the set of all runs  
 141 (resp., histories) of  $M$ , and by  $\text{Inf}(\pi)$  the set of states that occur infinitely often along the  
 142 run  $\pi$ .

143 A *sub-MDP* of  $M$  is an MDP  $M' = \langle Q', (A'_q)_{q \in Q'}, \delta' \rangle$  such that  $Q' \subseteq Q$  and  $\text{Supp}(\delta(q, a)) \subseteq$   
 144  $Q'$  for all states  $q \in Q'$  and actions  $a \in A'_q$  (recall the requirement that  $A'_q \neq \emptyset$ ). Consider a  
 145 set  $Q' \subseteq Q$  such that for all  $q \in Q'$ , there exists  $a \in A_q$  with  $\text{Supp}(\delta(q, a)) \subseteq Q'$ . We define  
 146 the *sub-MDP of  $M$  induced by  $Q'$* , denoted by  $M|_{Q'}$ , as the sub-MDP  $M' = \langle Q', (A'_q)_{q \in Q'}, \delta' \rangle$   
 147 where  $A'_q = \{a \in A_q \mid \text{Supp}(\delta(q, a)) \subseteq Q'\}$  for all  $q \in Q'$ .

148 **End-components** An *end-component* of  $M = \langle Q, (A_q)_{q \in Q}, \delta \rangle$  is a pair  $(Q', (A'_q)_{q \in Q'})$   
 149 such that  $(Q', (A'_q)_{q \in Q'}, \delta')$  is a sub-MDP of  $M$ , where  $\delta'$  denotes the restriction of  $\delta$  to  
 150  $\{(q, a) \mid q \in Q', a \in A'_q\}$ , and where the graph  $\langle Q', E' \rangle$  with  $E' = \{(q, q') \in Q' \times Q' \mid \exists a \in$   
 151  $A'_q : \delta(q, a)(q') > 0\}$  is strongly connected [9, 3]. We often identify an end-component as  
 152 the set  $Q' \cup \{(q, a) \mid q \in Q', a \in A_q\}$  of states and state-action pairs, and we say that it is  
 153 *supported* by the set  $Q'$  of states. The (componentwise) union of two end-components with  
 154 nonempty intersection is an end-component, thus one can define *maximal* end-components.  
 155 We denote by  $\text{MEC}(M)$  the set of maximal end-components of  $M$ , which is computable in  
 156 polynomial time [9], and by  $\text{EC}(M)$  the set of all end-components of  $M$ .

157 **Histories and Strategies** A *strategy* is a function  $\sigma : \text{Hist}(M) \rightarrow \mathcal{D}(A)$  such that

158  $\text{Supp}(\sigma(\rho)) \subseteq A_q$  for all histories  $\rho \in \text{Hist}(M)$  ending in  $\text{last}(\rho) = q$ . A strategy is *pure* if  
 159 all histories are mapped to Dirac distributions. A strategy  $\sigma$  is *memoryless* if  $\sigma(\rho) = \sigma(\rho')$   
 160 for all histories  $\rho, \rho'$  such that  $\text{last}(\rho) = \text{last}(\rho')$ . We sometimes view memoryless strategies  
 161 as functions  $\sigma : Q \rightarrow \mathcal{D}(A)$ . A strategy  $\sigma$  uses *finite memory* (of size  $k$ ) if there exists a  
 162 right congruence  $\approx$  over  $\text{Hist}(M)$  (i.e., such that if  $\rho \approx \rho'$ , then  $\rho \cdot a \cdot q \approx \rho' \cdot a \cdot q$  for all  
 163  $\rho, \rho' \in \text{Hist}(M)$  and  $(a, q) \in A \times Q$ ) of finite index  $k$  such that  $\sigma(\rho) = \sigma(\rho')$  for all histories  
 164  $\rho \approx \rho'$  with  $\text{last}(\rho) = \text{last}(\rho')$ .

165 **Objectives** An objective  $\varphi$  is a Borel set of runs. We denote by  $\mathbb{P}_q^\sigma(M, \varphi)$  the standard  
 166 probability measure on the sigma-algebra over the set of (infinite) runs of  $M$  with initial state  $q$ ,  
 167 generated by the cylinder sets spanned by the histories [3]. Given a history  $\rho = q_0 a_0 q_1 \dots q_k$ ,  
 168 the cylinder set  $\text{Cyl}(\rho) = \rho(AQ)^\omega$  has probability  $\mathbb{P}_q^\sigma(M, \text{Cyl}(\rho)) = \prod_{i=0}^{k-1} \sigma(q_0 a_0 q_1 \dots q_i)(a_i) \cdot$   
 169  $\delta(q_i, a_i)(q_{i+1})$  if  $q_0 = q$ , and probability 0 otherwise. We say that a run  $\rho$  is compatible with  
 170 strategy  $\sigma$  if  $\mathbb{P}_q^\sigma(M, \text{Cyl}(\rho)) > 0$ .

171 We consider the following standard objectives for an MDP  $M$ :

- 172 ■ safety objective: given a set  $T \subseteq Q$  of states, let  $\text{Safe}(T) = \{q_0 a_0 q_1 a_1 \dots \in \text{Hist}^\omega(M) \mid$   
 173  $\forall i \geq 0 : q_i \in T\}$ ;
- 174 ■ reachability objective: given a set  $T \subseteq Q$  of states, let  $\text{Reach}(T) = \{q_0 a_0 q_1 a_1 \dots \in$   
 175  $\text{Hist}^\omega(M) \mid \exists i \geq 0 : q_i \in T\}$ ;
- 176 ■ parity objective: given a priority function  $p : Q \rightarrow \mathbb{N}$ , let  $\text{Parity}(p) = \{\pi \in \text{Hist}^\omega(M) \mid$   
 177  $\min\{p(q) \mid q \in \text{Inf}(\pi)\} \text{ is even}\}$ .

178 It is standard to cast safety and reachability objectives as special cases of parity objectives,  
 179 using sink states. Given an objective  $\varphi$ , we denote by  $\neg\varphi = \text{Hist}^\omega(M) \setminus \varphi$  the complement of  
 180  $\varphi$ . We say that a run  $\pi \in \text{Hist}^\omega(M)$  *satisfies*  $\varphi$  if  $\pi \in \varphi$ , and that it *violates*  $\varphi$  otherwise.

181 It is known that under arbitrary strategies, with probability 1 the set  $\text{Inf}(\pi)$  of states  
 182 occurring infinitely often along a run  $\pi$  is the support of an end-component [8, 9].

183 ► **Lemma 1** ([8, 9]). *Given an MDP  $M$ , for all states  $q \in Q$  and all strategies  $\sigma$ , we have*  
 184  $\mathbb{P}_q^\sigma(M, \{\pi \mid \text{Inf}(\pi) \text{ is the support of an end-component}\}) = 1$ .

185 An end-component  $D \in \text{EC}(M)$  is *positive* under strategy  $\sigma$  from  $q$  if  $\mathbb{P}_q^\sigma(M, \{\pi \mid \text{Inf}(\pi) =$   
 186  $D\}) > 0$ . By Lemma 1, we have  $\sum_{D \in \text{EC}(M)} \mathbb{P}_q^\sigma(M, \{\pi \mid \text{Inf}(\pi) = D\}) = 1$ .

187 **Value and qualitative satisfaction** A strategy  $\sigma$  is winning for objective  $\varphi$  from  $q$  with  
 188 probability (at least)  $\alpha$  if  $\mathbb{P}_q^\sigma(M, \varphi) \geq \alpha$ . We denote by  $\text{Val}_q^*(M, \varphi) = \sup_\sigma \mathbb{P}_q^\sigma(M, \varphi)$  the  
 189 *value* of objective  $\varphi$  from state  $q$ . A strategy  $\sigma$  is *optimal* if  $\mathbb{P}_q^\sigma(M, \varphi) = \text{Val}_q^*(M, \varphi)$ .

190 We consider the following classical qualitative modes of winning. Given an objective  $\varphi$ , a  
 191 state  $q$  is:

- 192 ■ *almost-sure winning* if there exists a strategy  $\sigma$  such that is winning with probability 1,  
 193 that is  $\mathbb{P}_q^\sigma(M, \varphi) = 1$ .
- 194 ■ *limit-sure winning* if  $\text{Val}_q^*(M, \varphi) = 1$ , or equivalently for all  $\varepsilon > 0$  there exists a strategy  $\sigma$   
 195 such that  $\mathbb{P}_q^\sigma(M, \varphi) \geq 1 - \varepsilon$ .

196 We denote by  $\text{AS}(M, \varphi)$  and  $\text{LS}(M, \varphi)$  the set of almost-sure and limit-sure winning  
 197 states, respectively. In MDPs, it is known that  $\text{AS}(M, \varphi) = \text{LS}(M, \varphi)$  and pure memoryless  
 198 optimal strategies exist for parity objectives  $\varphi$  [19, 8].

199 We recall that the value of a parity objective  $\varphi = \text{Parity}(p)$  from every state of an  
 200 end-component  $D$  is the same, and is either 0 or 1, which does not depend on the precise  
 201 value of the (non-zero) transition probabilities, but only on the supports  $\text{Supp}(\delta(q, a))$  of the  
 202 transition function at the state-action pairs  $(q, a)$  in  $D$  [9]. When the value 1, there exists a

203 pure memoryless strategy  $\sigma$  such that  $\mathbb{P}_q^\sigma(M, \varphi) = 1$  for all states  $q \in D$ . If such a strategy  
 204 exists, then  $D$  is said to be  $\varphi$ -winning, and otherwise  $\varphi$ -losing.

## 205 2.2 Multiple-Environment MDP

206 A *multiple-environment MDP (MEMDP)* over a finite set  $E$  of environments is a tuple  
 207  $M = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$ , where  $M[e] = \langle Q, (A_q)_{q \in Q}, \delta_e \rangle$  is an MDP that models the  
 208 behaviour of the system in the environment  $e \in E$ . The state space is identical in all  $M[e]$   
 209 ( $e \in E$ ), only the transition probabilities may differ. We sometimes refer to the environments  
 210 of  $M$  as the MDPs  $\{M[e] \mid e \in E\}$  rather than the set  $E$  itself. For  $E' \subset E$ , let  $M[E']$  be  
 211 the MEMDP  $M$  over set  $E'$  of environments. We denote by  $M[-e]$  the MEMDP  $M$  over  
 212 environments  $E \setminus \{e\}$ , and by  $\cup_{e \in E} M[e]$  the MDP  $\langle Q, (A_q)_{q \in Q}, \delta_\cup \rangle$  such that  $\delta_\cup(q, a)$  is the  
 213 uniform distribution over  $\bigcup_{e \in E} \text{Supp}(\delta_e(q, a))$  for all  $q \in Q$  and  $a \in A$ .

214 A transition  $t = (q, a, q')$  is *revealing* in  $M$  if  $K_t = \{e \in E \mid q' \in \text{Supp}(\delta_e(q, a))\}$  is  
 215 a strict subset of  $E$  ( $K_t \subsetneq E$ ). We say that  $K_t$ , which is the set of environments where  
 216 the transition  $t = (q, a, q')$  is possible, is the *knowledge* after observing transition  $t$ . An  
 217 MEMDP is in *revealed form* if for all revealing transitions  $t = (q, a, q')$ , the state  $q'$  is a sink  
 218 in all environments, that is  $\text{Supp}(\delta_e(q', a)) = \{q'\}$  for all environments  $e \in E$  and all actions  
 219  $a \in A_{q'}$ . By extension, we call knowledge after a history  $\rho$  the set of environments in which  
 220 all transitions of  $\rho$  are possible.

221 **Decision Problems** We are interested in synthesizing a *single* strategy  $\sigma$  with guarantees  
 222 in *all* environments, without knowing in which environment  $\sigma$  is executing. We consider  
 223 reachability, safety, and parity objectives.

224 A state  $q$  is *almost-sure winning* in  $M$  for objective  $\varphi$  if there exists a strategy  $\sigma$  such  
 225 that in all environments  $e \in E$ , we have  $\mathbb{P}_q^\sigma(M[e], \varphi) = 1$ , and we call such a strategy  $\sigma$   
 226 almost-sure winning. A state  $q$  is *limit-sure winning* in  $M$  for objective  $\varphi$  if for all  $\varepsilon > 0$ ,  
 227 there exists a strategy  $\sigma$  such that in all environments  $e \in E$  we have  $\mathbb{P}_q^\sigma(M[e], \varphi) \geq 1 - \varepsilon$ ,  
 228 and we say that such a strategy  $\sigma$  is  $(1 - \varepsilon)$ -winning.

229 We denote by  $\text{AS}(M, \varphi)$  (resp.,  $\text{LS}(M, \varphi)$ ) the set of all almost-sure (resp., limit-sure)  
 230 winning states in  $M$  for objective  $\varphi$ . We consider the *membership problem for almost-sure*  
 231 (*resp., limit-sure*) *winning*, which asks whether a given state  $q$  is almost-sure (resp., limit-sure)  
 232 winning in  $M$  for objective  $\varphi$ . We refer to these membership problems as *qualitative* problems.

233 We are also interested in the *quantitative* problems. Given MEMDP  $M$ , a parity objective  
 234  $\varphi$ , and probability threshold  $\alpha \geq 0$ , we are interested in the existence of a strategy  $\sigma$  satisfying  
 235  $\mathbb{P}_q^\sigma(M[e], \varphi) \geq \alpha$  for all  $e \in E$ . We present an approximation algorithm for the quantitative  
 236 problem, solving the *gap problem* consisting, given MEMDP  $M$ , state  $q$ , parity objective  $\varphi$ ,  
 237 and thresholds  $0 < \alpha < 1$  and  $\varepsilon > 0$ , in answering

- 238 ■ Yes if there exists a strategy  $\sigma$  such that for all  $e \in E$ , we have  $\mathbb{P}_q^\sigma(M[e], \varphi) \geq \alpha$ ,
- 239 ■ No if for all strategies  $\sigma$ , there exists  $e \in E$  with  $\mathbb{P}_q^\sigma(M[e], \varphi) < \alpha - \varepsilon$ ,
- 240 ■ and arbitrarily otherwise.

241 The gap problem is an instance of promise problems which guarantee a correct answer in  
 242 two disjoint sets of inputs, namely positive and negative instances – which do not necessarily  
 243 cover all inputs, while giving no guarantees in the rest of the input [10, 13].

244 **Results** We solve the membership problem for limit-sure winning with parity objectives  $\varphi$   
 245 (i.e., deciding whether a given state  $q$  is limit-sure winning, that is  $q \in \text{LS}(M, \varphi)$ ), providing a  
 246 PSPACE algorithm with a matching complexity lower bound, and showing that the problem  
 247 is solvable in polynomial time when the number of environments is fixed. Our solution

248 relies on the solution of almost-sure winning, which is known to be PSPACE-complete for  
 249 reachability [23] and Rabin objectives [22]. We revisit the solution of almost-sure winning  
 250 and give a simple characterization for safety objectives (which is also PSPACE-complete),  
 251 that can easily be extended to parity objectives. A corollary of our characterization is that  
 252 pure (non-randomized) strategies are sufficient for both limit-sure and almost-sure winning,  
 253 which was known only for acyclic MEMDPs and reachability objectives [23, Lemma 12].

254 For the gap problem, we present an double exponential-space procedure to approximate  
 255 the value  $\alpha$  that can be achieved in all environments, up to an arbitrary precision  $\varepsilon$ .

### 256 3 Almost-Sure Winning

257 It is known that the membership problem for almost-sure winning in MEMDPs is PSPACE-  
 258 complete with reachability objectives [23] as well as with Rabin objectives [22], an expressively  
 259 equivalent of the parity objectives. We revisit the membership problem for almost-sure  
 260 winning with parity and safety objectives, as it will be instrumental to the solution of limit-  
 261 sure winning. We present a conceptually simple characterization of the winning region for  
 262 almost-sure winning, from which we derive a PSPACE algorithm, thus matching the known  
 263 complexity for almost-sure Rabin objectives. A corollary of our characterization is that pure  
 264 (non-randomized) strategies are sufficient for both limit-sure and almost-sure winning, which  
 265 was known only for acyclic MEMDPs and reachability objectives [23, Lemma 12].

266 ► **Theorem 2** ([23],[22]). *The membership problem for almost-sure winning in MEMDPs*  
 267 *with a reachability, safety, or Rabin objective is PSPACE-complete.*

268 To solve the membership problem for a safety or parity objective  $\varphi$ , we first convert  $M$   
 269 into an MEMDP  $M'$  in revealed form with state space  $Q \uplus \{q_{\text{win}}, q_{\text{lose}}\}$  and each revealing  
 270 transition  $t = (q, a, q')$  in  $M$  is redirected in  $M'$  to  $q_{\text{win}}$  if  $q' \in \text{AS}(M[K_t], \varphi)$  is almost-sure  
 271 winning when the set of environments is the knowledge  $K_t$  after observing transition  $t$ , and  
 272 to  $q_{\text{lose}}$  otherwise. In order to decide if  $q' \in \text{AS}(M[K_t], \varphi)$ , we need to solve the membership  
 273 problem for an MEMDP with strictly fewer environments than in  $M$  as  $K_t \subsetneq E$ , which  
 274 will lead to a recursive algorithm. The base case of the solution is MEMDPs with one  
 275 environment, which is equivalent to plain MDPs.

276 It is easy to see that  $\text{AS}(M, \varphi) \cup \{q_{\text{win}}\} = \text{AS}(M', \varphi \cup \text{Reach}(q_{\text{win}}))$  for all prefix-independent  
 277 objectives  $\varphi$ , and we can transform the objective  $\varphi \cup \text{Reach}(q_{\text{win}})$  into an objective of the  
 278 same type as  $\varphi$  (for example, if  $\varphi$  is a parity objective then assigning the smallest even  
 279 priority to  $q_{\text{win}}$  turns the objective  $\varphi \cup \text{Reach}(q_{\text{win}})$  into a pure parity objective).

280 Hence, the main difficulty is to solve the membership problem for MEMDP in revealed  
 281 form.

#### 282 3.1 Safety

283 Although safety objectives are subsumed by parity objectives which we solve in the next  
 284 section, we give here a simpler algorithm specifically for safety, and also prove PSPACE-  
 285 hardness in this case.

286 The safety objective has the property that almost-sure winning is equivalent to sure  
 287 winning, where a strategy is sure winning if all runs compatible with the strategy satisfy  
 288 the objective. Intuitively, if some runs does not satisfy the safety objective  $\text{Safe}(T)$ , then it  
 289 contains a state outside  $T$  after a finite prefix, thus with positive probability (the probability  
 290 of the finite prefix). In the sure-winning mode, we can consider the probabilistic choices to

291 be adversarial, which entails that only the support of the probability distributions in the  
 292 transition function is relevant.

293 It follows that, as long as the knowledge remains the set  $E$  of all environments a winning  
 294 strategy for a safety objective can play all actions that are safe (i.e., that ensure the successor  
 295 state remains in the winning region) in all environments. We obtain the following property:  
 296 almost-sure winning for a safety objective in a MEMDP  $M$  in revealed form is equivalent to  
 297 almost-sure winning in the MDP  $\cup_{e \in E} M[e]$ .

298 An algorithm for solving almost-sure safety is as follows: (1) for each revealing transition  
 299  $t = (q, a, q')$  in  $M$ , decide if  $q' \in \text{AS}(M[K_t], \text{Safe}(T))$  (using a recursive call), and redirect  
 300 the transition  $t$  to  $q_{\text{win}}$  or  $q_{\text{lose}}$  accordingly, transforming  $M$  into revealed form; (2) assuming  
 301  $M$  is in revealed form, compute the almost-sure winning states  $W = \text{AS}(M_{\cup}, \text{Safe}(T))$  where  
 302  $M_{\cup} = \cup_{e \in E} M[e]$  is an MDP. Return  $W \setminus \{q_{\text{win}}\}$ . The depth of recursive calls is bounded  
 303 by the number of environments, and the almost-sure safety in MDPs can be solved in  
 304 polynomial time, namely, in time  $O(|Q|^2|A|)$ . It follows that almost-sure safety in MEMDPs  
 305 can be solved in PSPACE, and in time  $O(|Q|^2 \cdot |A| \cdot 2^{|E|})$ . A PSPACE lower bound can be  
 306 established by a similar reduction from QBF as for reachability, the constructed MEMDP  
 307 being acyclic [23].

308 Note that for a fixed number of environments, the membership problem for almost-sure  
 309 safety in MEMDPs is solvable in polynomial time by our algorithm since the depth of the  
 310 recursion is then constant. This is also the case in Theorem 2 as shown in [23].

### 311 3.2 Parity

312 By definition, the almost-sure winning region  $W = \text{AS}(M, \text{Parity}(p))$  for a parity objective  
 313 in an MEMDP  $M$  is such that there exists a strategy  $\sigma$  that is almost-sure winning for the  
 314 parity objective from every state  $q \in W$  in every MDP  $M[e]$  (where  $e$  is an environment of  
 315  $M$ ). In contrast, we show the following characterization (note the order of the quantifiers).

316 ► **Lemma 3.** *Given an MEMDP  $M$  in revealed form with state space  $Q$ , if  $W \subseteq Q$  is such  
 317 that in every environment  $e$ , from every state  $q \in W$ , there exists a strategy  $\sigma_e$  that is almost-  
 318 sure winning for the parity objective  $\text{Parity}(p)$  in  $M|_W[e]$  from  $q$ , then  $W \subseteq \text{AS}(M, \text{Parity}(p))$ .  
 319 Moreover, for all  $q \in W$ , there exists a pure  $(|Q| \cdot |E|)$ -memory strategy ensuring  $\text{Parity}(p)$   
 320 from  $q$  in  $M$ .*

321 **Proof.** For each environment  $M[e]$ , consider a memoryless strategy  $\sigma_e$  almost-surely winning  
 322 for the objective  $\text{Parity}(p)$  in  $M|_W[e]$  from every state of  $W$ . Recall that almost-sure winning  
 323 strategies can be assumed to be memoryless in MDPs with single environments; and that  
 324 one can build a single memoryless strategy that is almost-surely winning from all winning  
 325 states. Let  $\text{EC}(\sigma_e) = \{D \in \text{EC}(M[e]) \mid \exists q \in W : \mathbb{P}_q^{\sigma_e}(M[e], \text{Inf} = D) > 0\}$  be the set of  
 326 positive end-components under strategy  $\sigma_e$ . Note that the least priority in an end-component  
 327  $D \in \text{EC}(\sigma_e)$  is even since the parity objective is satisfied with probability 1.

328 Let  $E = \{1, \dots, k\}$  be the set of environments of  $M$ . We construct a pure almost-sure  
 329 winning strategy  $\sigma$  for the MEMDP  $M$  as follows, where initially  $e = 1$ :

- 330 (1) play according to  $\sigma_e$  for  $|W|$  steps;
- 331 (2) if the current state is  $q_{\text{win}}$  or belongs to a positive end-component  $D \in \text{EC}(\sigma_e)$ , keep  
 332 playing according to  $\sigma_e$  forever. Otherwise, increment  $e$  (modulo  $k$ ) and go to (1).

333 The strategy  $\sigma$  uses memory of size at most  $|Q| \cdot |E|$  since  $W \subseteq Q$ .

334 Fix environment  $f \in E$ . We show that strategy  $\sigma$  is almost-sure winning in  $M[f]$ . Because  
 335 all strategies  $\sigma_e$  are defined in  $M|_W$ , the region  $W$  is never left while playing  $\sigma$ , and during

336 phase (1) of the strategy there is a lower-bounded probability to reach an end-component  
 337  $D \in \text{EC}(\sigma_e)$  when  $e = f$ .

338 We show that eventually phase (2) is executed forever with probability 1, that is, some  
 339 end-component  $D \in \text{EC}(\sigma_e)$  for some  $e$  is reached with probability 1. Towards contradiction,  
 340 assume that phase (1) of the strategy  $\sigma$  is executed infinitely often with positive probability  
 341  $p$ . Then phase (1) for  $e = f$  and  $\sigma_f$  is also executed infinitely often and it follows that,  
 342 conditioned on phase (1) being executed infinitely often, a positive end-component  $D \in \text{EC}(\sigma_f)$   
 343 is reached with probability 1; hence phase (2) is executed forever from that point on. Thus  
 344 with probability  $1 - p + p = 1$  phase (1) is executed only finitely often, contradicting our  
 345 assumption.

346 As phase (2) of the strategy  $\sigma$  is eventually executed forever with probability 1, let  $e$  be  
 347 the corresponding environment (i.e., such that  $\sigma$  plays according to  $\sigma_e$ ) and let  $D \neq \{q_{\text{win}}\}$   
 348 be the reached end-component of  $M[e]$  (the other case where  $q_{\text{win}}$  is reached is trivial). If  
 349 some transition of  $D$  is not present in  $f$ , then it must be a revealing transition in  $e$ , thus  
 350 leading in  $M[e]$  to  $q_{\text{win}}$  outside  $D$ , which is impossible since  $D$  is an end-component in  $M[e]$ .  
 351 Hence all transitions of  $D$  are present in all environments.

352 We show that  $\sigma$  is almost-sure winning in  $f$ . The result is immediate if  $D$  is an end-  
 353 component of  $M[f]$  (in particular if  $f = e$ ). If  $D$  is not an end-component of  $M[f]$ , then in  
 354  $M[f]$  the strategy would leave  $D$  and reach  $q_{\text{win}}$ , thus  $\sigma$  is almost-sure winning as well in  
 355 that case. ◀

356 The characterization in the first part of Lemma 3 holds simply because parity objectives  
 357 are prefix-independent (runs that differ by a finite prefix are either both winning or both  
 358 losing), and thus the characterization holds for all prefix-independent objectives.

359 The converse of Lemma 3 is immediate, which entails that the almost-sure winning region  
 360  $W = \text{AS}(M, \text{Parity}(p))$  is the largest set of states satisfying the condition in Lemma 3. We  
 361 exploit this characterization in Algorithm 1 to compute the winning region for almost-sure  
 362 parity. After transforming the MEMDP into revealed form (through recursive calls to the  
 363 algorithm), we compute the winning region for almost-sure parity in each environment  
 364 (line 11), and then the set  $P'$  of states from which we can remain in the intersection  $P$  of all  
 365 these winning regions (line 12). We iterate this process on  $M|_{P'}$  until a fixpoint  $P = P'$  is  
 366 reached.

367 It is easy to see that the fixpoint satisfies the characterization of Lemma 3, and thus  
 368  $P' \subseteq \text{AS}(M, \text{Parity}(p)) \cup \{q_{\text{win}}\}$ . Also by the proof of Lemma 3, we can construct a pure  
 369 almost-sure winning ( $|Q| \cdot |E|$ )-memory strategy from all states in  $P'$ , and define (recursively,  
 370 for each subset of the environments) a pure almost-sure winning strategy from the states that  
 371 were replaced by  $q_{\text{win}}$  in the revealed form, with a total memory size at most  $|Q| \cdot |E| \cdot 2^{|E|}$ ,  
 372 corresponding to the memory bound from Lemma 3 for each subset  $K \subseteq E$  of environments  
 373 (representing the belief, i.e., the set of environments where the current history is possible).

374 To show the converse inclusion, we show the invariant that every state  $q \in Q \setminus P'$  is not  
 375 almost-sure winning in  $M$ : for all strategies  $\sigma$  from  $q$ , in some environment  $M[e]$  the set  $P$   
 376 is left with positive probability (along some history  $\rho$ ). Given a state  $q' \in Q \setminus P$  reached  
 377 in  $M[e]$ , there is an environment  $f \in E$  where the parity objective is violated with positive  
 378 probability under  $\sigma$  from  $q'$ . The crux is to show that the state  $q'$  is reached with positive  
 379 probability in  $M[f]$  as well. Towards contradiction, assume that the history  $\rho$  from  $q$  to  
 380  $q'$  (in  $M[e]$ ) is not possible in  $M[f]$ . Then  $\rho$  contains a revealing transition in  $M[e]$ , and  
 381  $q' = q_{\text{win}} \in P$ , which is a contradiction since  $q' \in Q \setminus P$ . Hence, in  $M[f]$  with strategy  $\sigma$  the  
 382 parity objective is violated with positive probability.

---

**Algorithm 1** AS\_Parity( $M, p$ )

---

**Input** :  $M = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$  an MEMDP,  $p : Q \rightarrow \mathbb{N}$  a priority function.  
**Output**: The winning region  $\text{AS}(M, \text{Parity}(p))$  for almost-sure parity.

**begin**

/\* pre-processing \*/

1  $M' \leftarrow M$

2 add two sink states  $q_{\text{win}}, q_{\text{lose}}$  to  $M'$

3 define  $p(q_{\text{win}}) = 0$  and  $p(q_{\text{lose}}) = 1$

4 **foreach** revealing transition  $t = (q, a, q')$  in  $M$  **do**

/\*  $K_t \subsetneq E$  \*/

5     **if**  $q' \in \text{AS\_Parity}(M[K_t], p)$  **then**

6          $\perp$  replace  $t$  by  $(q, a, q_{\text{win}})$  in  $M'$

7         **else**

8              $\perp$  replace  $t$  by  $(q, a, q_{\text{lose}})$  in  $M'$

9  $M \leftarrow M'$

/\*  $M$  is in revealed form \*/

10  $P \leftarrow \emptyset; P' \leftarrow \emptyset$

11 **repeat**

12      $P \leftarrow \bigcap_{e \in E} \text{AS}(M[e], \text{Parity}(p))$      /\*  $M[e]$  and  $\bigcup_{e \in E} M[e]$

13      $P' \leftarrow \text{AS}(\bigcup_{e \in E} M[e], \text{Safe}(P))$      are MDPs \*/

14      $M \leftarrow M|_{P'}$

**until**  $P'$  is unchanged

15 **return**  $P' \setminus \{q_{\text{win}}\}$

**end**

---

383     Algorithm 1 can be implemented in PSPACE by a similar argument as for almost-sure  
384 safety: the depth of recursive calls is bounded by the number of environments, both almost-  
385 sure safety and almost-sure parity can be solved in polynomial time in MDPs, and the  
386 repeat-loop runs at most  $|Q|$  times. The algorithm runs in polynomial time if the number of  
387 environments is fixed. The PSPACE-hardness follows from Theorem 2.

388 ► **Theorem 4.** *The membership problem for almost-sure parity in MEMDPs is PSPACE-*  
389 *complete. Pure exponential-memory strategies are sufficient for almost-sure winning in*  
390 *MEMDPs with parity (thus also reachability and safety) objectives. When the number of*  
391 *environments is fixed, the problem is solvable in polynomial time.*

392     The time complexity of Algorithm 1 is established as follows. Each recursive call,  
393 corresponds to a subset of the initial environment set  $E$  that we can compute once and  
394 tabulate. In each call, the second loop runs at most  $|Q|$  times, and the set of almost-sure  
395 winning states for parity conditions (that is, the set  $\text{AS}(M[e], \text{Parity}(p))$ ) can be computed in  
396 time  $O(|Q| \cdot |\delta|)$  [3]. Since  $|\delta|$  is in  $O(|Q|^2 \cdot |A|)$ , each recursive call takes  $O(|Q|^4 \cdot |E| \cdot |A|)$   
397 time, and overall, this is  $O(|Q|^4 \cdot |E| \cdot |A| \cdot 2^{|E|})$ .

398     Note that pure exponential-memory strategies for almost-sure parity in MEMDPs are  
399 provided by Lemma 3. The algorithm for almost-sure parity can be used to solve almost-sure  
400 safety with optimal PSPACE complexity, although the specific algorithm for safety is slightly  
401 simpler (the repeat-loop can be replaced by just line 12 where  $P = T$  is the set of states

402 defining the safety objective  $\text{Safe}(T)$ .

403 The PSPACE procedure can be implemented in exponential time by solving all subprob-  
404 lems and storing their solutions. Moreover, for large numbers of environments, the exponent  
405 in the complexity can be made to depend only on the size of  $M$ . In fact, intuitively, two  
406 environments with identical supports yield the same result so one can derive a dynamic  
407 programming solution where at most one environment per support is solved.

408 Define the *support* of a probabilistic transition relation  $\delta : Q \times A \rightarrow \mathcal{D}(Q)$  as the family  
409 of supports of its transitions, that is,  $\text{Supp}(\delta) = (\text{Supp}(\delta(q, a)))_{(q, a) \in Q \times A}$ . Define the support  
410 of a family of transition relations as  $\text{Supp}((\delta_e)_{e \in E}) = \{\text{Supp}(\delta_e) \mid e \in E\}$ .

411 Two environments  $\delta_e$  and  $\delta_f$  are said to be *equivalent* if they have the same support. One  
412 can check whether two environments are equivalent in polynomial time, by going through all  
413 triples  $(q, a, q')$  and verifying that  $\delta_e(q, a, q') = 0$  iff  $\delta_f(q, a, q') = 0$ .

414 Almost sure parity in MEMDPs does not depend on the precise probability values in the  
415 given environments in  $M$  but only on their supports.

416 In addition to Theorem 4, we can obtain a complexity bound whose exponent is inde-  
417 pendent of the number of environments (Theorem 6), using the following result: if in two  
418 environments, the support of the transition relation is the same, we can discard one of the  
419 environment (all strategies that are almost-sure winning in one are also almost-sure winning  
420 in the other one, as shown in Lemma 5) and thus consider at most one environment for each  
421 support. Here, we denote by  $\text{Supp}((\delta_e)_{e \in E}) = (\text{Supp}(\delta_e))_{e \in E}$  where  $\text{Supp}(\delta_e)$  denotes the set  
422 of transitions with positive probability under  $\delta_e$ .

423 **► Lemma 5.** *Consider two MEMDPs  $M_i = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E_i} \rangle$  for  $i = 1, 2$ , with the same  
424 state and action sets, and with the same supports of their transition relation,  $\text{Supp}((\delta_e)_{e \in E_1}) =$   
425  $\text{Supp}((\delta_e)_{e \in E_2})$ . Given a parity condition  $\text{Parity}(p)$ , for all states  $q$  and all finite-memory  
426 strategies  $\sigma$ , the following equivalence holds:  $\mathbb{P}_q^\sigma[M_1[e], \text{Parity}(p)] = 1$  for all  $e \in E_1$  if  
427 and only if  $\mathbb{P}_q^\sigma[M_2[e], \text{Parity}(p)] = 1$  for all  $e \in E_2$ . In particular,  $\text{AS}(M_1, \text{Parity}(p)) =$   
428  $\text{AS}(M_2, \text{Parity}(p))$ .*

429 **Proof.** Given state  $q$  and finite-memory strategy  $\sigma$ , assume that  $\mathbb{P}_q^\sigma[M_1[e_1], \text{Parity}(p)] = 1$   
430 for all  $e_1 \in E_1$ . Consider any  $e_2 \in E_2$ , and let  $e_1 \in E_1$  be such that  $\text{Supp}(\delta_{e_1}) = \text{Supp}(\delta_{e_2})$ ;  
431 such a  $e_2$  exists by the hypothesis  $\text{Supp}((\delta_e)_{e \in E_1}) = \text{Supp}((\delta_e)_{e \in E_2})$ . Consider the Markov  
432 chain obtained as the product of the MDP  $M_1[e_1]$  with the Moore machine describing the  
433 finite-memory strategy  $\sigma$ . Because  $\mathbb{P}_q^\sigma[M_1[e_1], \text{Parity}(p)] = 1$ , all bottom strongly connected  
434 components (BSCC) in this product are winning for  $\text{Parity}(p)$  (i.e., the smallest priority of  
435 their states is even). But the product of  $M_2[e_2]$  and the Moore machine for  $\sigma$  have the same  
436 set of BSCCs since the supports are identical. It follows that  $\mathbb{P}_q^\sigma[M_2[e_2], \text{Parity}(p)] = 1$ . By  
437 symmetry, this proves the first statement.

438 It follows that  $\text{AS}(M_1, \text{Parity}(p)) = \text{AS}(M_2, \text{Parity}(p))$  since finite-memory strategies suffice  
439 for almost-sure parity in MEMDPs by Theorem 4.

440 ◀

441 **► Theorem 6.** *The membership problem for almost-sure parity for an MEMDP  $M =$   
442  $\langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$  can be solved in time  $O((|E|^2 + |Q|^4 \cdot |E| \cdot |A|) \cdot 2^{\min(|E|, 2^{|Q|^2 \cdot |A|})})$ .*

443 **Proof.** Consider an MEMDP  $M = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$  and a parity objective  $\varphi$ . If  
444  $|E| \leq 2^{|Q|^2 \cdot |A|}$ , then we apply the PSPACE procedure from Theorem 4. The number of  
445 recursive calls is then bounded by  $2^{|E|}$ , and each call itself takes polynomial time, so the  
446 result follows.

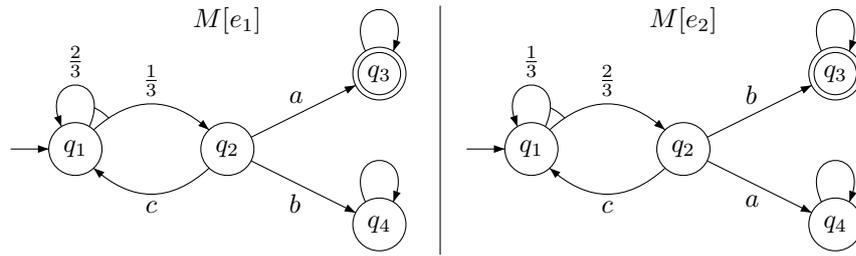


Figure 3 An end-component  $\{q_1, q_2\}$  with different transition probabilities in environments  $e_1$  and  $e_2$ .

447 Otherwise, we scan the set of environments given as input, and store a subset  $E'$  of these:  
 448 we include an environment  $e$  in  $E'$  if and only if none of the previously stored environments  
 449 is equivalent to  $e$ . This takes  $O(|E|^2)$  time. This yields a subset with at most  $2^{|Q|^2|A|}$   
 450 environments, with at most one representative for each possible support. We then apply the  
 451 recursive algorithm on the MEMDP  $M[E']$ , which yields the same result as if it was applied  
 452 to  $M = M[E]$  by Lemma 5. ◀

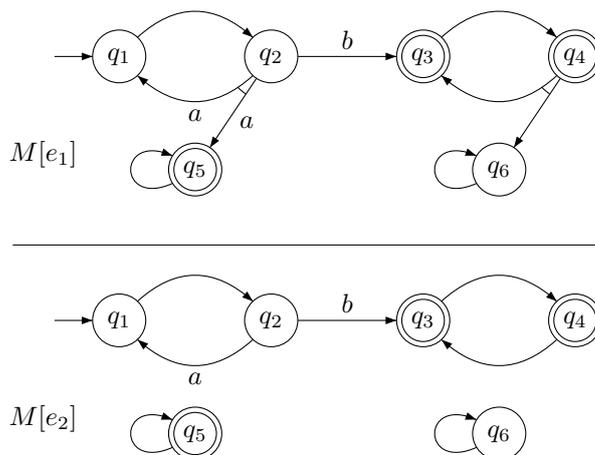
453 **4 Limit-Sure Winning**

454 We refer to the examples of the duplicate card and the missing card in Section 1 to illustrate  
 455 the difference between limit-sure and almost-sure winning. We present in Section 4.1 two  
 456 other scenarios where limit-sure winning and almost-sure winning do not coincide, which will  
 457 be useful to illustrate the key ideas in the algorithmic solution.

458 **4.1 Examples**

459 In the example of Figure 3, the set  $D = \{q_1, q_2\}$  is an end-component in both environments  
 460  $e_1$  and  $e_2$  (the actions are shown in the figures only when relevant, that is in  $q_2$ ). However,  
 461 the transition probabilities from  $q_1$  are different in the two environments  $e_1$  and  $e_2$ , and  
 462 intuitively we can learn (with high probability) in which environment we are by playing  $c$   
 463 for a long enough (but finite) time and collecting the frequency of the visits to  $q_1$  and  $q_2$ .  
 464 Then, in order to reach the target  $q_3$ , if there are more  $q_1$ 's than  $q_2$ 's in the history we play  
 465  $a$  in  $q_2$ , otherwise  $b$ . The intuition is that the histories with more  $q_1$ 's than  $q_2$ 's have a high  
 466 probability (more than  $1 - \varepsilon$ ) in  $M[e_1]$  and a small probability (less than  $\varepsilon$ ) in  $M[e_2]$ , where  
 467  $\varepsilon$  can be made arbitrarily small (however not 0) by playing  $c$  for sufficiently long. Hence  $q_1$   
 468 is limit-sure winning, but not almost-sure winning.

469 In the second scenario (Figure 4), the transition probabilities do not matter. The objective  
 470 is to visit some state in  $\{q_3, q_4, q_5\}$  infinitely often (those states have priority 0, the other  
 471 states have priority 1). The state  $q_1$  is limit-sure winning, but not almost-sure winning. To  
 472 win with probability  $1 - \varepsilon$ , a strategy can play  $a$  (in  $q_2$ ) for a sufficiently long time, then  
 473 switch to playing  $b$  (unless  $q_5$  was reached before that). The crux is that playing  $a$  does  
 474 not harm, as it does not leave the limit-sure winning region, but ensures in at least one  
 475 environment (namely,  $e_1$ ) that the objective is satisfied with probability 1 (by reaching  $q_5$ ).  
 476 This allows to “discard” the environment  $e_1$  if  $q_5$  was not reached, and to switch to a strategy  
 477 that is winning with probability at least  $1 - \varepsilon$  in  $e_2$ , namely by playing  $b$ . With an arbitrary



■ **Figure 4** The set  $\{q_1, q_2\}$  is an end-component in  $e_2$ , not in  $e_1$ .

478 number of environments, the difficulty is to determine in which order the environments can  
479 be “discarded”.

480 Note that the transition  $(q_2, a, q_1)$  is not revealing, since it is present in both environments.  
481 However, after crossing this transition a large number of times, we can still learn that the  
482 environment is  $e_2$  (and be mistaken with arbitrarily small probability). In contrast, the  
483 transition  $(q_2, a, q_5)$  is revealing and the environment is  $e_1$  with certainty upon crossing that  
484 transition.

485 To solve the membership problem for limit-sure parity, we first convert  $M$  into a revealed-  
486 form MEMDP  $M'$ , similar to the case of almost-sure winning, with the obvious difference that  
487 revealing transitions  $t = (q, a, q')$  of  $M[e]$  are redirected in  $M'[e]$  to  $q_{\text{win}}$  if  $q' \in \text{LS}(M[K_t], \varphi)$   
488 is limit-sure winning when the set of environments is the knowledge  $K_t$  after observing  
489 transition  $t$ . Thus, we aim for a recursive algorithm, where the base case is limit-sure winning  
490 in MEMDPs with one environment, which are equivalent to plain MDPs, for which limit-sure  
491 and almost-sure parity coincide. Note that the examples of Figure 3 and Figure 4 are in  
492 revealed form.

## 493 4.2 Common End-Components and Learning

494 A *common end-component (CEC)* of an MEMDP  $M = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$  is a pair  $(Q', A')$   
495 that is an end-component in  $M[e]$  for all environments  $e \in E$ . A CEC  $D$  is *trivial* if it  
496 contains a single state.  $D$  is said *winning* for a parity condition  $\text{Parity}(p)$ , if for all  $e \in E$ ,  
497 there is a strategy in  $M[e]$  which, when started inside  $D$ , ensures  $\text{Parity}(p)$  with probability 1.  
498 Notice that since  $D$  is a common end-component, such a strategy ensures  $\text{Parity}(p)$  with  
499 probability 1 in  $M[e]$  iff it does in  $M[e']$ .

500 We note that the common end-components of an MEMDP are the end-components of the  
501 MDP  $\cup_{e \in E} M[e]$  assuming  $M$  is in revealed form, and thus can be computed using standard  
502 algorithm for end-components [9].

503 ► **Lemma 7.** Consider an MEMDP  $M$  in revealed form. The common end-components of  $M$   
504 are exactly the end-components of  $\cup_{e \in E} M[e]$ .

505 **Proof.** Consider a common end-component  $D$  of  $M$ . Because in each  $M[e]$ , all state-action

## XX:14 The Value Problem for Multiple-Environment MDPs with Parity Objective

506 pairs in  $D$  stay inside  $D$ , and  $D$  is strongly connected, this is also the case in  $\cup_{e \in E} M[e]$ ;  
 507 thus  $D$  is an end-component of the latter.

508 Conversely, consider an end-component  $D$  of  $\cup_{e \in E} M[e]$ . If  $D$  consists of a single sink  
 509 state, then it is indeed a common end-component. Otherwise  $D$  contains more than one  
 510 state. We show that all state-action pairs  $(q, a)$  of  $D$  must have the same support in all  
 511 environments, and it follows that  $D$  is an end-component in every environment, thus a  
 512 common end-component. By contradiction, if a transition  $(q, a, q')$  with  $(q, a) \in D$  exists in  
 513  $M[e]$  but not in  $M[f]$ , then it is revealing and  $q'$  is a sink state. Hence  $D$  is not strongly  
 514 connected in  $\cup_{e \in E} M[e]$  because  $D$  does not consist of a single sink state. ◀

515 A CEC may have different transition probabilities in different environments. We call a  
 516 CEC *distinguishing* if it contains a transition  $(q, a, q')$  (called a distinguishing transition)  
 517 such that  $\delta_e(q, a)(q') \neq \delta_f(q, a)(q')$  for some environments  $e, f \in E$ . Given a distinguishing  
 518 transition  $(q, a, q')$  and environment  $e$ , define  $K_1 = \{f \in E \mid \delta_f(q, a)(q') = \delta_e(q, a)(q')\}$  and  
 519  $K_2 = E \setminus K_1$ . We say that  $(K_1, K_2)$  is a *distinguishing partition* of  $D$  that is *induced* by the  
 520 distinguishing transition  $(q, a, q')$  and environment  $e$ .

521 Distinguishing transitions can be used to learn the partition  $(K_1, K_2)$ , that is to guess  
 522 (correctly with high probability) whether the current environment is in  $K_1$  or  $K_2$ , as in the  
 523 example of Figure 3, where the set  $D = \{q_1, q_2\}$  is a distinguishing end-component with  
 524 distinguishing transition  $(q_1, \cdot, q_2)$  and partition  $(\{e_1\}, \{e_2\})$ . A distinguishing CEC may  
 525 have several distinguishing transitions and induced partitions.

526 We formalize how a strategy can distinguish between  $K_1$  and  $K_2$  with high probability  
 527 inside a distinguishing CEC. First let us recall Hoeffding's inequality.

528 ▶ **Theorem 8** (Hoeffding's Inequality [14]). *Let  $X_1, X_2, \dots, X_n$  be a sequence of independent*  
 529 *and identical Bernoulli variables with  $\mathbb{P}[X_i] = p$ , and write  $S_n = X_1 + \dots + X_n$ . For all  $t > 0$ ,*  
 530  *$\mathbb{P}[S_n - \mathbb{E}[S_n] \geq t] \leq e^{-2t^2/n}$ , and  $\mathbb{P}[\mathbb{E}[S_n] - S_n \geq t] \leq e^{-2t^2/n}$ .*

531 Given a distinguishing CEC with distinguishing partition  $(K_1, K_2)$  induced by a transition  
 532  $(q, a, q')$ , a strategy can sample the distribution  $\delta_e(q, a)$  by repeating the following two phases:  
 533 first, use a pure memoryless strategy to almost-surely visit  $q$ , then play action  $a$ ; by repeating  
 534 this long enough (the precise bound depends on a given  $\varepsilon$  and is derived from Theorem 8)  
 535 while storing the frequency of visits to  $q'$  in the second phase, we can learn and guess in  
 536 which block  $K_i$  belongs the environment, with sufficiently small probability of mistake to  
 537 ensure winning with probability  $1 - \varepsilon$ .

538 ▶ **Lemma 9.** *Given an MEMDP  $M$  containing a distinguishing common end-component  $D$*   
 539 *with partition  $(K_1, K_2)$  induced by a distinguishing transition, and parity objective  $\varphi$ , for all*  
 540 *states  $q_0$  in  $D$ , all pairs of strategies  $\sigma_1, \sigma_2$ , and all  $\varepsilon > 0$ , there exists a strategy  $\sigma$  such that:*

$$541 \quad \mathbb{P}_{q_0}^\sigma(M[e], \varphi) \geq (1 - \varepsilon) \mathbb{P}_{q_0}^{\sigma_1}(M[e], \varphi) \text{ for all } e \in K_1,$$

$$542 \quad \mathbb{P}_{q_0}^\sigma(M[e], \varphi) \geq (1 - \varepsilon) \mathbb{P}_{q_0}^{\sigma_2}(M[e], \varphi) \text{ for all } e \in K_2.$$

543 *Moreover, the strategy  $\sigma$  is pure if both  $\sigma_i$  are pure; and if each strategy  $\sigma_i$  uses*  
 544 *a memory of size  $m_i$ , then  $\sigma$  uses finite memory of size  $m_1 + m_2 + \lceil 8 \frac{\log(1/\varepsilon)}{\eta^2} \rceil$  where*  
 545  *$\eta = \min(\{|\delta_e(q, a)(q') - \delta_f(q, a)(q')| \mid e, f \in E, q, q' \in Q, a \in A\} \setminus \{0\})$ .*

546 **Proof.** Consider  $M = \langle Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E} \rangle$  and  $D = \langle Q', (A'_q)_{q \in Q'} \rangle$  as in the statement of  
 547 the lemma, and let  $q_0 \in D$ .

548 Consider a distinguishing transition  $(q, a, q')$  and environment  $e_0$  that induces the distin-  
 549 guishing partition  $(K_1, K_2)$ . Consider  $\varepsilon > 0$ , and define  $N = \lceil \frac{2 \log(1/\varepsilon)}{\eta^2} \rceil$ ,

550 The strategy  $\sigma$  runs in two phases. In the first phase, the goal is to estimate the  
 551 distribution of  $(q, a, q')$ . For this, it executes a pure memoryless strategy which has a nonzero  
 552 probability of reaching  $q$  while staying in  $D$  (such a strategy can be defined based on the  
 553 supports of state-action pairs of  $D$ ) and keeps two counters:  $c_{q,a}$  that counts the number of  
 554 times the state-action pair  $(q, a)$  is selected; and  $c_{q,a,q'}$  the number of times the transition  
 555  $(q, a, q')$  is observed. The second round of the strategy starts when  $c_{q,a} = N$ . Note that this  
 556 happens with probability 1. Then, we go back to  $q_0$  (with probability 1), and we switch to

- 557 ■  $\sigma_1$  if  $\left| \frac{c_{q,a,q'}}{c_{q,a}} - \delta_{e_0}(q, a)(q') \right| < \eta/2$ ,
- 558 ■  $\sigma_2$  otherwise.

559 We now analyze this strategy and show that because  $N$  is sufficiently large, the estimation  
 560 error is bounded, so that we obtain the desired result.

561 In each environment  $e$ , at each visit at  $q$  and choice of  $a$ , we have a Bernoulli trial with  
 562 mean  $\delta_e(q, a)(q')$ , and  $c_{q,a,q'}$  is the number of successful trials. By Hoeffding's inequality  
 563 (Theorem 8), we have

$$564 \quad \mathbb{P}_{q_0}^{\sigma} \left( M[e], |c_{q,a,q'}/c_{q,a} - \delta_e(q, a)(q')| \geq \eta/2 \mid c_{q,a} = N \right) \leq e^{-2N(\frac{\eta}{2})^2} \leq \varepsilon.$$

565 Thus, in  $M[e]$  with  $e \in K_i$ , the probability of not switching to  $\sigma_i$  is at most  $\varepsilon$ . It follows  
 566 that  $\mathbb{P}_{q_0}^{\sigma} (M[e], \varphi) \geq (1 - \varepsilon) \mathbb{P}_{q_0}^{\sigma_i} (M[e], \varphi)$ .

567 The memory requirement comes from the fact that  $\sigma$  must store two counters up to  $N$   
 568 values, and it has two modes (before and after reaching  $c_{q,a} = N$ ). ◀

569 It follows that the membership problem for limit-sure winning can be decomposed into  
 570 subproblems where the set of environments is one of the blocks  $K_i$  in the partition.

571 ▶ **Lemma 10.** *Given an MEMDP  $M$  containing a distinguishing common end-component*  
 572  *$D$  with a partition  $(K_1, K_2)$  induced by a distinguishing transition, and a parity objective*  
 573  *$\varphi$  the following equivalence holds:  $D \subseteq \text{LS}(M, \varphi)$  if and only if  $D \subseteq \text{LS}(M[K_1], \varphi)$  and*  
 574  *$D \subseteq \text{LS}(M[K_2], \varphi)$ .*

575 **Proof.** Immediate consequence of Lemma 9. ◀

### 576 4.3 Characterization and Algorithm

577 Here, we assume that MEMDPs are in revealed form with sink states  $q_{\text{win}}$  and  $q_{\text{lose}}$ .

578 We show that the winning region  $W = \text{LS}(M, \varphi)$  for limit-sure parity is a closed set: from  
 579 every state  $q \in W$ , there exists an action  $a$  ensuring in all environments that all successors  
 580 of  $q$  are in  $W$ . We call such actions *limit-sure safe* for  $q$ . We show in Lemma 11 that a  
 581 limit-sure safe action always exists in limit-sure winning states. Note that playing actions  
 582 that are *not* limit-sure safe may be useful for limit-sure winning, as in the example of Figure 3  
 583 where action  $a$  is limit-sure safe, but action  $b$  is not (from  $q_2$ ).

584 By definition of limit-sure winning, if a state  $q$  is not limit-sure winning, there exists  $\varepsilon_q > 0$   
 585 such that for all strategies  $\sigma$ , there exists an environment  $e \in E$  such that  $\mathbb{P}_q^{\sigma} (M[e], \varphi) < 1 - \varepsilon_q$ .  
 586 We denote by  $\varepsilon_0 = \min\{\varepsilon_q \mid q \in Q \setminus \text{LS}(M, \varphi)\}$  a uniform bound.

587 ▶ **Lemma 11.** *Given an MEMDP  $M$  (in revealed form) over environments  $E$ , a parity*  
 588 *objective  $\varphi$ , and a state  $q$ , if  $q \in \text{LS}(M, \varphi)$  is limit-sure winning, then there exists an*  
 589 *action  $a$  such that for all environments  $e \in E$ , all successors of  $q$  are limit-sure winning, i.e*  
 590  *$\text{Supp}(\delta_e(q, a)) \subseteq \text{LS}(M, \varphi)$ .*

## XX:16 The Value Problem for Multiple-Environment MDPs with Parity Objective

591 **Proof.** Consider  $q \in \text{LS}(M, \varphi)$  and let  $0 < \varepsilon < \frac{\nu \varepsilon_0}{|A|}$ , where  $A$  is the set of actions in  $M$ , and  $\nu$   
 592 is a lower bound on the smallest nonzero transition probability (in all environments), and  $\varepsilon_0$   
 593 is the uniform bound defined above. Let  $\sigma$  be a strategy ensuring  $\varphi$  from  $q$  with probability  
 594 at least  $1 - \varepsilon$  in all environments.

595 Towards contradiction, assume that there is no limit-sure safe action from state  $q$ . Let  
 596  $a$  be the action chosen by  $\sigma$  with the highest probability at the history  $q$ , that is  $a =$   
 597  $\arg \max_a \sigma(q)(a)$ , and thus  $\sigma(q)(a) \geq \frac{1}{|A|}$ . By our assumption, there exists an environment  
 598  $e \in E$  and a state  $t \notin \text{LS}(M, \varphi)$  (in particular  $t \neq q_{\text{win}}$ ) such that  $\delta_e(q, a)(t) > 0$ , hence  
 599  $\delta_e(q, a)(t) \geq \nu$ . It is immediate that  $t \neq q_{\text{lose}}$  as otherwise the strategy  $\sigma$  would ensure  $\varphi$  with  
 600 probability at most  $1 - \nu \leq 1 - \varepsilon$  from  $q$ . So  $t \notin \{q_{\text{win}}, q_{\text{lose}}\}$  and therefore  $\delta_e(q, a)(t) \geq \nu$  in  
 601 all environments  $e$ . By definition of the uniform bound  $\varepsilon_0$ , there exists an environment  $e$  such  
 602 that  $\mathbb{P}_t^\sigma(M[e], \varphi) \leq 1 - \varepsilon_0$ , hence from  $q$  we have  $\mathbb{P}_q^\sigma(M[e], \neg \varphi) \geq \frac{\nu \varepsilon_0}{|A|} > \varepsilon$ , in contradiction  
 603 to  $\sigma$  ensuring  $\varphi$  with probability at least  $1 - \varepsilon$  from  $q$ . We conclude that there exists a  
 604 limit-sure safe action from  $q$ . ◀

605 Given an MEMDP  $M$ , consider the limit-sure winning region  $W = \text{LS}(M, \varphi)$  for  $\varphi =$   
 606 Parity( $p$ ). For the purpose of the analysis, consider the (memoryless) randomized strategy  
 607  $\sigma_{\text{LS}}$  that plays uniformly at random all limit-sure safe actions in every state  $q \in W$ , which is  
 608 well-defined by Lemma 11.

609 Consider an arbitrary environment  $e$ , and an end-component  $D$  in  $M[e]$  that is positive  
 610 under  $\sigma_{\text{LS}}$  (recall Lemma 1 and the definition afterward). There are three possibilities:

- 611 1.  $D$  is not a common end-component (as in the example of Figure 4, for  $D = \{q_1, q_2\}$   
 612 in  $M[e_2]$ ), that is,  $D$  is not an end-component in some environment  $e'$  (in the example  
 613  $e' = e_1$ ), then we can learn (and be mistaken with arbitrarily small probability) that we  
 614 are not in  $e'$ , reducing the problem to an MEMDP with fewer environments (namely,  
 615  $M[\neg e']$ );
- 616 2.  $D$  is a common end-component and is distinguishing (as in the example of Figure 3, for  
 617  $D = \{q_1, q_2\}$ ), then we can also learn a distinguishing partition  $(K_1, K_2)$  and reduce the  
 618 problem to MEMDPs with fewer environments (namely,  $M[K_1]$  and  $M[K_2]$ );
- 619 3.  $D$  is a common end-component and is non-distinguishing, then we show in Lemma 12  
 620 below that  $D$  is almost-sure winning ( $D \subseteq \text{AS}(M, \varphi)$ ), obviously in all environments.

621 ► **Lemma 12.** *Given an MEMDP  $M$  over environments  $E$  (in revealed form), a parity objec-*  
 622 *tive  $\varphi$ , and a state  $q$ , if  $q \in \text{LS}(M, \varphi)$ , then all non-distinguishing common end-components*  
 623  *$D$  that are positive under strategy  $\sigma_{\text{LS}}$  from  $q$  in  $M[e]$  (for some  $e \in E$ ) are almost-sure*  
 624 *winning for  $\varphi$  (that is  $D \subseteq \text{AS}(M, \varphi)$ ).*

625 **Proof.** Consider a positive non-distinguishing common end-component  $D$  as in the statement  
 626 of the lemma. Using Lemma 11, note that  $D \subseteq \text{LS}(M, \varphi)$  since  $\sigma_{\text{LS}}$  plays only limit-sure safe  
 627 actions and  $D$  is a common end-component.

628 Assume towards contradiction that  $D$  is not almost-sure winning for the parity objective  
 629  $\varphi$ . It follows that in  $M$ , all strategies that play only limit-sure safe actions ensure the parity  
 630 objective  $\varphi$  with probability 0 from all states in  $D$  (in all environments since  $D$  is a common  
 631 end-component).

632 Denote by  $\Omega_{\text{safe}}$  the set of all runs that contain only limit-sure safe actions. For all  
 633 strategies  $\sigma$  (in  $M$ ), and  $q \in D$  we have  $\mathbb{P}_q^\sigma(M[e], \varphi \mid \Omega_{\text{safe}}) = 0$  (for all  $e \in E$ ) and therefore:

$$\begin{aligned}
634 \quad \mathbb{P}_q^\sigma(M[e], \varphi) &= \mathbb{P}_q^\sigma(M[e], \varphi \mid \Omega_{safe}) \cdot \mathbb{P}_q^\sigma(M[e], \Omega_{safe}) \\
635 \quad &\quad + \mathbb{P}_q^\sigma(M[e], \varphi \mid \neg\Omega_{safe}) \cdot \mathbb{P}_q^\sigma(M[e], \neg\Omega_{safe}) \\
636 \quad &= \mathbb{P}_q^\sigma(M[e], \varphi \mid \neg\Omega_{safe}) \cdot \mathbb{P}_q^\sigma(M[e], \neg\Omega_{safe}) \\
637 \quad &\leq 1 - \mathbb{P}_q^\sigma(M[e], \neg\varphi \mid \neg\Omega_{safe})
\end{aligned}$$

638 Given  $\varepsilon < \frac{\varepsilon_0 \cdot \nu}{|E|}$  where  $\nu$  is the smallest positive probability in  $M$ , we show that there  
639 exists an environment  $e \in E$  such that  $\mathbb{P}_q^\sigma(M[e], \varphi) < 1 - \varepsilon$ , which entails that  $q$  is not  
640 limit-sure winning for  $\varphi$ , establishing a contradiction since  $q \in D \subseteq \text{LS}(M, \varphi)$ . It will follow  
641 that  $D$  is almost-sure winning for  $\varphi$  and conclude the proof.

By definition of limit-sure safe actions, to every pair  $(q, a)$  such that  $a \in A_q$  is not limit-sure safe in  $q$ , we can associate an environment  $e$  such that:

$$\text{Supp}(\delta_e(q, a)) \cap (Q \setminus \text{LS}(M, \varphi)) \neq \emptyset,$$

642 and thus from some state  $q' \in \text{Supp}(\delta_e(q, a))$ , we have  $\mathbb{P}_{q'}^\sigma(M[e], \varphi) \leq 1 - \varepsilon_0$  where  $\varepsilon_0$  is the  
643 uniform bound for non-limit-sure winning states. Assuming that a non-limit-sure safe action  
644 is played by  $\sigma$ , since there are finitely many environments, by the pigeonhole principle there  
645 is an environment  $e$  such that with probability at least  $\frac{1}{|E|}$  an action that is not limit-sure  
646 safe and associated with  $e$  is played, which leads with probability at least  $\nu$  to a state outside  
647  $\text{LS}(M, \varphi)$ . It follows that  $\mathbb{P}_q^\sigma(M[e], \neg\varphi \mid \neg\Omega_{safe}) \geq \varepsilon_0 \cdot \frac{\nu}{|E|} > \varepsilon$  and thus  $\mathbb{P}_q^\sigma(M[e], \varphi) < 1 - \varepsilon$ ,  
648 which concludes the proof.  $\blacktriangleleft$

649 Our approach to compute the limit-sure winning states is to first identify the distinguishing  
650 CECs that are limit-sure winning. We can compute the maximal CECs using Lemma 7, and  
651 note that a maximal CEC containing a distinguishing CEC is itself distinguishing, so it is  
652 sufficient to consider maximal CECs. By Lemma 10, we can decide if a given distinguishing  
653 CEC is limit-sure winning using a recursive procedure on MEMDPs with fewer environments.  
654 We show in Lemma 13 below that we can replace the limit-sure CECs by a sink state  $q_{\text{win}}$ .

655 **► Lemma 13.** *Given an MEMDP  $M$  with parity objective  $\varphi$  and a set  $T \subseteq \text{LS}(M, \varphi)$  of*  
656 *limit-sure winning states, we have  $\text{LS}(M, \varphi) = \text{LS}(M, \varphi \cup \text{Reach}(T))$ .*

657 **Proof.** The inclusion  $\text{LS}(M, \varphi) \subseteq \text{LS}(M, \varphi \cup \text{Reach}(T))$  is immediate since  $\varphi \subseteq \varphi \cup \text{Reach}(T)$ .

658 To show the converse inclusion, consider  $q \in \text{LS}(M, \varphi \cup \text{Reach}(T))$  and show that  $q \in$   
659  $\text{LS}(M, \varphi)$ . Given  $\varepsilon > 0$ , let  $\varepsilon_1 = \frac{\varepsilon}{2}$  and let  $\sigma$  be a strategy such that  $\mathbb{P}_q^\sigma(M, \varphi \cup \text{Reach}(T)) \geq$   
660  $1 - \varepsilon_1$ . We construct a strategy  $\tau$  that satisfies the objective  $\varphi$  with probability at least  
661  $1 - \varepsilon$  as follows: for all histories  $\rho$ , if  $\rho$  does not visit  $T$ , then let  $\tau(\rho) = \sigma(\rho)$ ; otherwise,  
662 consider the suffix  $\rho'$  of  $\rho$  after the first visit to a state  $t \in T$ , and let  $\sigma_t$  be strategy  
663 that ensures  $\varphi$  is satisfied with probability at least  $1 - \varepsilon_1$  from  $t$  (such a strategy exists  
664 since  $T \subseteq \text{LS}(M, \varphi)$ ). Define  $\tau(\rho) = \sigma_t(\rho')$ . We easily show below that  $\mathbb{P}_q^\tau(M, \varphi) \geq 1 - \varepsilon$ ,  
665 establishing that  $q \in \text{LS}(M, \varphi)$ :

$$\begin{aligned}
666 \quad \mathbb{P}_q^\tau(M, \varphi) &= \mathbb{P}_q^\tau(M, \varphi \cap \text{Reach}(T)) + \mathbb{P}_q^\tau(M, \varphi \cap \neg\text{Reach}(T)) \\
667 \quad &= \mathbb{P}_q^\tau(M, \varphi \mid \text{Reach}(T)) \cdot \mathbb{P}_q^\tau(M, \text{Reach}(T)) + \mathbb{P}_q^\tau(M, \varphi \mid \neg\text{Reach}(T)) \\
668 \quad &= \mathbb{P}_q^\tau(M, \varphi \mid \text{Reach}(T)) \cdot \mathbb{P}_q^\sigma(M, \text{Reach}(T)) + \mathbb{P}_q^\sigma(M, \varphi \mid \neg\text{Reach}(T)) \\
669 \quad &\quad \text{(since } \tau \text{ agrees with } \sigma \text{ as long as } T \text{ is not reached)} \\
670 \quad &\geq (1 - \varepsilon_1) \cdot \mathbb{P}_q^\sigma(M, \text{Reach}(T)) + \mathbb{P}_q^\sigma(M, \varphi \mid \neg\text{Reach}(T)) \\
671 \quad &\geq (1 - \varepsilon_1) \cdot \mathbb{P}_q^\sigma(M, \text{Reach}(T)) + (1 - \varepsilon_1) \cdot \mathbb{P}_q^\sigma(M, \varphi \mid \neg\text{Reach}(T)) \\
672 \quad &\geq (1 - \varepsilon_1) \cdot \mathbb{P}_q^\sigma(M, \varphi \cup \text{Reach}(T)) \geq (1 - \varepsilon_1)^2 \geq 1 - \varepsilon.
\end{aligned}$$

673



674 We can now assume that MEMDPs contain no limit-sure winning distinguishing CEC,  
 675 and present a characterization for the remaining possibility, illustrated by the scenario of  
 676 Figure 4, where playing the action  $a$  (in  $q_2$ , forever) ensures, in some environment (namely,  
 677  $e_1$ ), almost-sure satisfaction of the parity objective while remaining inside the limit-sure  
 678 winning region in all other environments.

679 ► **Lemma 14.** *Consider an MEMDP  $M$  (in revealed form) over environments  $E$  with*  
 680  *$|E| \geq 2$ , that contains no limit-sure winning distinguishing common end-component, and*  
 681 *a parity objective  $\varphi$ . Writing  $T_e = \text{LS}(M[\neg e], \varphi)$ , we have the following:  $\text{LS}(M, \varphi) =$*   
 682  *$\text{AS}\left(M, \text{Reach}\left(\bigcup_{e \in E} \text{AS}(M[e], \varphi \cap \text{Safe}(T_e))\right)\right)$ .*

683 **Proof.** First we show the inclusion

$$684 \quad \text{LS}(M, \varphi) \subseteq \text{AS}\left(M, \text{Reach}\left(\bigcup_{e \in E} \text{AS}(M[e], \varphi \cap \text{Safe}(T_e))\right)\right).$$

685 Consider the (memoryless) strategy  $\sigma_{\text{LS}}$  that plays all limit-sure safe actions uniformly at  
 686 random from every state in  $\text{LS}(M, \varphi)$ . The strategy  $\sigma_{\text{LS}}$  is well-defined by Lemma 11 and  
 687 to establish the inclusion, we show that, from every state  $q \in \text{LS}(M, \varphi)$ , it is almost-sure  
 688 winning (in all environments  $e' \in E$ ) for the objective  $\text{Reach}\left(\bigcup_{e \in E} \text{AS}(M[e], \varphi \cap \text{Safe}(T_e))\right)$ .

689 Consider an arbitrary environment  $e' \in E$  and an arbitrary end-component  $D$  that is  
 690 positive under  $\sigma_{\text{LS}}$  in  $M[e']$ . Since positive end-components are reached with probability 1  
 691 (Lemma 1), it is sufficient to show that for all such  $D$ , there exists an environment  $e \in E$   
 692 such that every state in  $D$  is almost-sure winning for the objective  $\varphi \cap \text{Safe}(T_e)$  in  $M[e]$ . We  
 693 consider two cases:

694 ■ if  $D$  is a common end-component, then we show that  $D$  is non-distinguishing. Note that  
 695  $D$  must be limit-sure winning, by definition of limit-sure safe actions (played by  $\sigma_{\text{LS}}$ ).  
 696 It follows by the assumption of the lemma that  $D$  is non-distinguishing and therefore  
 697 almost-sure winning for  $\varphi$  (in all environments) by Lemma 12. We take  $e = e'$  and it is  
 698 easy to see that there exists an almost-sure winning strategy for  $\varphi$  from  $D$  (that stays in  
 699  $D$ ), which is also almost-sure winning for  $\varphi \cap \text{Safe}(T_e)$ .

700 ■ otherwise  $D$  is not a common end-component, and there exists an environment  $e$  where  
 701  $D$  is not an end-component. We first show that all transitions of  $D$  are present in  $M[e]$ ,  
 702 since otherwise  $D$  would contain a revealing transition, thus leading to a state that is a  
 703 sink in all environments (revealed form). Then  $D$  being strongly connected would not  
 704 contain another state, and thus in particular all transitions in  $D$  would be present in  
 705  $M[e]$ .

706 It follows that playing  $\sigma_{\text{LS}}$  from  $D$  in  $M[e]$  ensures with probability 1 that a (revealing)  
 707 transition not present in  $M[e']$  is executed, which leads to  $q_{\text{win}}$  since  $\sigma_{\text{LS}}$  never leaves the  
 708 limit-sure winning region (by definition of limit-sure safe actions). Hence  $\varphi$  is satisfied  
 709 with probability 1 in  $M[e]$  while playing only limit-sure safe actions, thus remaining in  
 710 the limit-sure winning region  $\text{LS}(M, \varphi) \subseteq \text{LS}(M[\neg e], \varphi) = T_e$ , thereby satisfying  $\text{Safe}(T_e)$   
 711 as well. This shows that in  $M[e]$ , the states in  $D$  are almost-sure winning for the objective  
 712  $\varphi \cap \text{Safe}(T_e)$ .

713 For the converse inclusion, given a state  $q$  and a pure<sup>1</sup> strategy  $\sigma$  that is almost-sure  
 714 winning for objective  $\text{Reach}\left(\bigcup_{e \in E} \text{AS}(M[e], \varphi \cap \text{Safe}(T_e))\right)$  (in all environments), we show

---

<sup>1</sup> By Theorem 4, pure strategies are sufficient for almost-sure winning in MEMDPs.

715 that for all  $\varepsilon > 0$  there is a pure strategy  $\tau$  that ensures that  $\varphi$  is satisfied with probability  
716 at least  $1 - \varepsilon$  (from  $q$  in all environments).

717 Given  $\varepsilon > 0$ , let  $\tau$  be the strategy that plays as follows:

- 718 (1) play like  $\sigma$  until a state  $t \in \bigcup_{e \in E} \text{AS}(M[e], \varphi \cap \text{Safe}(T_e))$  is reached, and let  $e \in E$  be an  
719 environment such that from  $t$  there is a (pure memoryless) strategy  $\sigma_t$  that is almost-sure  
720 winning in  $M[e]$  for the objective  $\varphi \cap \text{Safe}(T_e)$ ;  
721 (2) play like  $\sigma_t$  for  $k \cdot |Q|$  steps, where  $k$  is such that  $(1 - \nu^{|Q|})^k \leq \varepsilon$  (where  $\nu$  is the smallest  
722 positive probability in  $M$ );  
723 (3) if the current state belongs to a positive end-component  $D_t$  of  $\sigma_t$  (in  $M[e]$ ), then keep  
724 playing like  $\sigma_t$  (forever); otherwise switch to a strategy that ensures that  $\varphi$  is satisfied  
725 with probability at least  $1 - \varepsilon$  from the current state in all environments of  $E \setminus \{e\}$  – such  
726 a strategy exists because from  $t$  the strategy  $\sigma_t$  ensures the objective  $\text{Safe}(T_e)$  is satisfied  
727 almost-surely (and thus surely as well).

728 Consider an arbitrary environment  $e \in E$ , and show that  $\mathbb{P}_q^\tau(M[e], \varphi) \geq 1 - \varepsilon$ , which  
729 establishes that  $q$  is limit-sure winning,  $q \in \text{LS}(M, \varphi)$ .

730 First note that phase (2) (and thus also phase (3)) is reached with probability 1, and  
731 let  $e_t$  be the environment corresponding to the state  $t$  reached at the end of phase (1). We  
732 consider two cases:

- 733 ■ if  $e_t = e$ , then by standard analysis the probability that after phase (2) a positive  
734 end-component of  $\sigma_t$  is *not yet* reached is at most  $(1 - \nu^{|Q|})^k \leq \varepsilon$  since within  $|Q|$  steps a  
735 positive end-component is reached with probability at least  $\nu^{|Q|}$ . Hence with probability  
736 at least  $1 - \varepsilon$ , a positive (winning since  $\sigma_t$  almost-sure winning in  $M[e]$  for the objective  
737  $\varphi$ ) end-component of  $\sigma_t$  is reached and the strategy  $\sigma_t$  is played forever in phase (3), thus  
738 winning with probability at least  $1 - \varepsilon$ .  
739 ■ otherwise  $e_t \neq e$  and we consider the following cases in phase (3):  
740 (a) if the strategy  $\sigma_t$  is played forever, then either the set  $D_t$  (which is an end-component  
741 in  $M[e_t]$ ) is never left, or it is left (via a revealing transition, as  $D_t$  is not left in  $M[e_t]$ )  
742 and since  $\sigma_t$  ensures  $\text{Safe}(T_{e_t})$  the sink  $q_{\text{win}}$  is reached in  $M[e]$ , thus in both cases the  
743 objective  $\varphi$  is satisfied (with probability 1);  
744 (b) otherwise, by construction the strategy  $\tau$  switches to a strategy that ensures  $\varphi$  is  
745 satisfied with probability at least  $1 - \varepsilon$ .

746 In all cases, the objective  $\varphi$  holds with probability at least  $1 - \varepsilon$ , showing that  $\mathbb{P}_q^\tau(M[e], \varphi) \geq$   
747  $1 - \varepsilon$  as claimed.

748 ◀

749 **Algorithm Overview** Given a MEMDP  $M = (Q, (A_q)_{q \in Q}, (\delta_e)_{e \in E})$ , the algorithm proceeds  
750 by recursion on the size of the environment set  $E$  (Algorithm 2). The base case is that of a  
751 singleton set  $E$  where  $\text{LS}(M, \varphi) = \text{AS}(M, \varphi)$  and this can be computed in polynomial time.

752 Assume  $|E| \geq 2$ . We first convert  $M$  into an MEMDP  $M'$  in revealed form with state space  
753  $Q \uplus \{q_{\text{win}}, q_{\text{lose}}\}$  and each revealing transition  $t = (q, a, q')$  in  $M$  is redirected in  $M'$  to  $q_{\text{win}}$  if  
754  $q' \in \text{LS}(M[K_t], \varphi)$  is limit-sure winning when the set of environments is the knowledge  $K_t$   
755 after observing transition  $t$ , and to  $q_{\text{lose}}$  otherwise. Notice that each query  $q' \in \text{LS}(M[K_t], \varphi)$   
756 uses a set  $K_t$  that is strictly smaller than  $E$ .

757 We now assume that  $M$  is in revealed form and we compute the maximal end-components  
758 of the MDP  $\bigcup_{e \in E} M[e]$ ; these are maximal common end-components of  $M$  by Lemma 7.  
759 For each distinguishing maximal CEC  $D$ , we determine whether it is limit-sure winning



783 *environments is fixed, the problem is solvable in polynomial time.*

784 The time complexity of Algorithm 2 is established as follows. Let us consider a single  
 785 recursive call. The maximal end-components of  $\cup_{e \in E} M'[e]$  can be computed in  $O(|Q| \cdot |\delta|)$   
 786 where  $|\delta|$  denotes the number of transitions. Then, determining whether each MEC is  
 787 distinguishing, and replacing them with sink states can be done in time  $O(|\delta| \cdot |E|)$  since  
 788 one needs to go over each transition and check whether their probability differs in two  
 789 environments. The last step requires solving almost-sure parity and safety for MDPs defined  
 790 for each  $e \in E$ , which can be done in time  $O(|E| \cdot |Q| \cdot |\delta|)$  (similarly as in the discussion  
 791 following Theorem 4). The most costly operation is almost-sure reachability for the MEMDP  
 792  $M$ , which by Theorem 4 takes  $O(|Q|^4 \cdot |E| \cdot |A| \cdot 2^{|E|})$ . There are  $2^{|E|}$  recursive calls (the  
 793 algorithm can be run once for each subset of  $E$  using memoization), so overall we get  
 794  $O(|Q|^4 \cdot |E| \cdot |A| \cdot 2^{2|E|})$ .

795 We do not know if a technique similar to that of Theorem 6 can be used for the limit-sure  
 796 case to obtain an exponent independent of  $|E|$ .

## 797 **5 The Gap Problem**

798 The goal of this section is to give a procedure that solves the gap problem for parity objectives.  
 799 For this, we show that an arbitrary strategy in  $M$  can be imitated by a finite-memory one  
 800 (with a computable bound on the memory size) while achieving the same probability of  
 801 winning up to  $\varepsilon$  in all environments. Once this is established, we show how to guess a  
 802 finite-memory strategy of the appropriate size in order to solve the gap problem.

803 To establish the memory bound for such an  $\varepsilon$ -approximation, we need a few intermediate  
 804 lemmas. First, we define a transformation on MEMDPs consisting in collapsing non-  
 805 distinguishing maximal CECs (MCECs) of the MEMDP  $M$ ; the resulting MEMDP is denoted  
 806  $\text{purge}(M)$ . We show that  $M$  and  $\text{purge}(M)$  have the same probabilities of satisfaction of the  
 807 considered parity objective under all environments.

808 Intuitively, removing non-distinguishing MCECs ensures that in  $\text{purge}(M)$ , under all  
 809 strategies, with high probability, within a fixed number of steps, either a maximal CEC  
 810 is reached (which is either distinguishing, or non-distinguishing but trivial – recall that a  
 811 trivial CEC contains a single absorbing state.) or enough samples are gathered to improve  
 812 the knowledge about the current environment, as shown in Section 5.2 This observation will  
 813 help us constructing the finite-memory strategy inductively since in each case the knowledge  
 814 can be improved correctly with high probability: in trivial MCECs, the strategy is extended  
 815 arbitrarily; inside distinguishing MCECs, the strategy can be extended so that it stays  
 816 inside the MCEC while sampling distinguishing transitions with any desired precision as in  
 817 Lemma 9; last, if no MCECs are reached but enough samples are gathered along the way, we  
 818 prove that the knowledge can also be improved with high probability. The final strategy is  
 819 obtained by combining finite-memory strategies constructed inductively for smaller sets of  
 820 environments. This is done in Section 5.3

### 821 **Maximal Common End-Components Revisited**

822 We extend the definition of common end-components (CEC) which, in Section 4.2, were  
 823 defined assuming MEMDPs are in revealed form. In this section, MEMDPs are **not** assumed  
 824 to be in revealed form: in fact, upon observing a revealing transition, we cannot conclude  
 825 recursively since we cannot determine which value vector must be achieved in the recursive  
 826 call. Here, we define a CEC for MEMDP  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$  as a pair  $(Q', A')$  such that

827 for all  $e \in E$ ,  $\langle Q', A', \delta_e \rangle$  is an end-component of  $M[e]$ . A maximal CEC (MCEC) is a CEC  
 828 which does not contain a smaller CEC.

829 There are two types of MCECs:

- 830 ■ MCEC  $(Q', A')$  is non-distinguishing if for all  $q \in Q'$ , and  $a \in A'(q)$ , the distributions  
 831  $\delta_e(q, a)$  and  $\delta_{e'}(q, a)$  are identical for all  $e, e' \in E$ ;
- 832 ■ MCEC  $(Q', A')$  is distinguishing otherwise.

833 While non-distinguishing MCECs have state-action pairs with identical supports in all  
 834 environments, a distinguishing MCEC may contain revealing transitions, that is, state-action  
 835 pairs  $(q, a)$  with different supports in different environments. This is the difference with  
 836 Section 4. The only result we need from Section 4.2 is Lemma 9 which holds for the new  
 837 definition of distinguishing MCECs: in fact, we do require that  $(Q', A')$  is an end-component  
 838 (i.e., closed and strongly connected) in all environments, so revealing transitions are simply  
 839 seen as distinguishing transitions, and thanks to the strong connectivity of  $(Q', A')$  in all  
 840 environments, one can define a strategy that samples a distinguishing transition a desired  
 841 number of times.

842 As previously, a MCEC  $D$  is *trivial* if it contains a single state.

843 In terms of computability, we cannot use Lemma 7 to compute MCECs since this is  
 844 only valid for MEMDPs in revealed form. The  $\varepsilon$ -gap procedure given in this section does  
 845 not actually compute MCECs; these are only used in the proof of the existence of a finite-  
 846 memory strategy (Lemma 24). Nevertheless, for completeness, let us describe how MCECs  
 847 can be computed in polynomial time. For  $|E| = 1$ , the MCECs are exactly the maximal  
 848 end-components (MECs) of  $M[e]$  where  $E = \{e\}$ . For  $|E| \geq 2$ , we pick an environment  
 849  $e \in E$ , and compute the MECs of  $M[e]$ . For each MEC  $D$  of  $M[e]$ , we recursively compute  
 850 the MCECs of  $D$  in the MEMDP  $M[E \setminus \{e\}]$ . This is sound because a MCEC, being an  
 851 end-component in all environments, is necessarily a subset of some MEC in each  $M[e]$ ; so by  
 852 restricting the search for MCECs to MECs of some  $M[e]$ , we do not discard any MCECs.  
 853 Furthermore, each recursive call splits the state space to disjoint sets, so we get an overall  
 854 polynomial-time complexity.

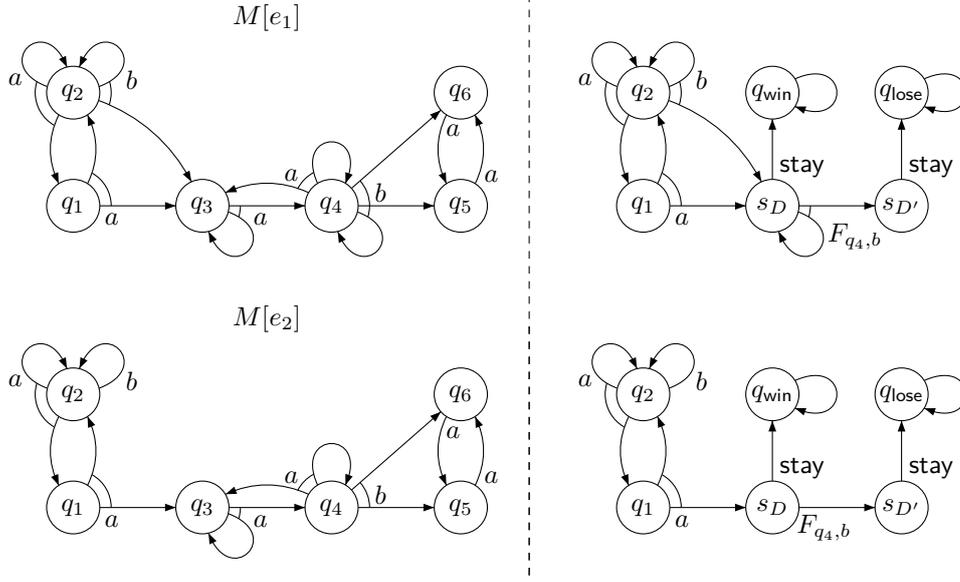
855 Given an MEMDP  $M$  over environments  $E$ , the notation  $\mathbb{P}_q^\sigma(M, \varphi)$  refers to the vector  
 856 of probability values  $(\mathbb{P}_q^\sigma(M[e], \varphi))_{e \in E}$ .

## 857 5.1 Purge: Removing Non-Distinguishing MCECs

858 We first describe a transformation that collapses non-distinguishing MCECs, and keeps  
 859 only trivial ones. Since all trivial MCECs can be classified into winning and losing for the  
 860 objective  $\varphi$ , we assume that the only non-distinguishing MCECs in the resulting MEMDP are  
 861 called  $q_{\text{win}}$  and  $q_{\text{lose}}$ . The intuition is that non-distinguishing MCECs are not useful to refine  
 862 information in order to distinguish environments, so when a strategy visits such a MCEC,  
 863 one can assume that it will either stay inside forever (either if the MCEC is  $\varphi$ -winning, or if  
 864 there is no outgoing transition), or leave it as soon as possible (if the MCEC is  $\varphi$ -losing).

865 Observe that a distinguishing MCEC can contain a smaller non-distinguishing CEC. The  
 866 transformation described here only collapses MCECs that are non-distinguishing, and not  
 867 those smaller non-distinguishing CECs that are contained in MCECs.

868 Given an MEMDP  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$ , define the MEMDP  $\text{purge}(M) = \langle Q', A', (\delta'_e)_{e \in E} \rangle$   
 869 where  $Q'$  contains all states of  $Q$  except those that belong to non-distinguishing MCECs;  
 870 and for each non-distinguishing MCEC  $D$ , we add a fresh state  $s_D$  to  $Q'$ , and redirect all  
 871 transitions that enter a state of  $D$  in  $M$  to  $s_D$  in  $M'$ . We define the map  $f : Q \rightarrow Q'$  by



**Figure 5** An MEMDP  $M$  with two environments (left) and the construction  $\text{purge}(M)$  (right). Transition probabilities are uniform. Here  $D$  is the MCEC defined by the pairs  $\{(q_3, a), (q_4, a)\}$ , and  $D'$  is the MCEC defined by  $\{(q_5, a), (q_6, a)\}$ . The priority function is omitted, we assume that  $D$  is winning (e.g., by assigning priority 0 to  $q_3$  and  $q_4$ ) and that  $D'$  is losing (e.g., by assigning priority 1 to  $q_5$  and  $q_6$ ).

872 mapping all states of non-distinguishing MCECs  $D$  to  $s_D$ , and as the identity for other  
873 states.

874 We add a fresh action  $\text{stay}$  which from  $s_D$  goes to a winning absorbing state  $q_{\text{win}}$  if  
875  $D$  is  $\varphi$ -winning, and to a losing absorbing state  $q_{\text{lose}}$  otherwise. For each pair  $(q, a) \in D$   
876 such that  $\text{Supp}(\delta(q, a))$  is not included in  $D$ , we add a fresh action  $F_{(q,a)}$  from  $s_D$  with  
877  $\delta'_e(s_D, F_{(q,a)})(q') = \sum_{q'' \in f^{-1}(q')} \delta_e(q, a, q'')$  for all  $e \in E$ . (These state-action pairs can leave  
878  $D$  in some environments, so  $F$  stands for the *frontier* of  $D$ .)

879 Given the set of MCECs,  $\text{purge}(M)$  can be computed in polynomial time. However, the  
880  $\varepsilon$ -gap procedure we give does not actually compute  $\text{purge}(M)$ ; this construction is only used  
881 for proving the existence of a finite-memory strategy of bounded memory size.

882 **Example 16.** An example of this construction is given in Fig. 5 for MEMDP  $M$  with  
883 two environments  $e_1, e_2$ . Here  $\{(q_2, b)\}$  is an end-component in  $M[e_2]$  but not in  $M[e_1]$   
884 due to the edge to  $q_3$  so this is not a CEC, and is not collapsed in  $\text{purge}(M)$ . The  
885 MCEC  $D$  defined by  $\{(q_3, a), (q_4, a)\}$  has a single frontier action  $F_{(q_4,b)}$ . In  $M[e_1]$ , we  
886 have  $\delta'_{e_1}(s_D, F_{(q_4,b)}, s_{D'}) = 2/3$  since  $\delta_{e_1}(q_4, b, q_5) + \delta_{e_1}(q_4, b, q_6) = 2/3$  (since the probabil-  
887 ities are uniform), and  $\delta'_{e_1}(s_D, F_{(q_4,b)}, s_D) = 1/3$ . In  $M[e_2]$ , the latter edge is missing, so  
888  $\delta'_{e_2}(s_D, F_{(q_4,b)}, s_{D'}) = 1$ .

889 **Lemma 17.** For all MEMDPs  $M$ , the only non-distinguishing MCECs of  $\text{purge}(M)$  are  
890 the trivial  $q_{\text{win}}$  and  $q_{\text{lose}}$ .

891 **Proof.** Let  $D = (Q', A')$  be any non-distinguishing MCEC in  $M'$ .  $D$  must contain a state  
892 of the form  $s_{D'}$  since otherwise this is also a MCEC of  $M$ , and the construction would  
893 have collapsed it. We consider the component in  $M$  given by the inverse image of  $D$  by  $f$ .

## XX:24 The Value Problem for Multiple-Environment MDPs with Parity Objective

894 Formally, let  $Q'' = f^{-1}(Q') \subseteq Q$ , and for each  $q'' \in Q''$ , define  $A''(q'') = \{a \in A(q'') \mid \forall e \in$   
 895  $E : \text{Supp}(\delta_e(q, a)) \subseteq Q''\}$ .

896 Then for each state of the form  $q_{D'}$  in  $D$ ,  $(Q'', A'')$  contains all state-action pairs of  $D'$ .  
 897 But  $D$  is strongly connected in each  $M'[e]$ , and all non-distinguishing MCECs  $D'$  of  $M$   
 898 that were collapsed are also strongly connected in each  $M[e]$  by definition,  $(Q'', A'')$  is also  
 899 strongly connected in each  $M[e]$ , thus a MCEC in  $M$ .

900 Now,  $(Q'', A'')$  cannot be distinguishing, since the construction only collapses MCECs, so  
 901 no subset of  $(Q'', A'')$  can be collapsed in  $M'$ ; and  $(Q'', A'')$  would remain untouched and  
 902 be distinguishing in  $M'$  as well. So  $(Q'', A'')$  is non-distinguishing; but in this case, it is  
 903 collapsed into a trivial MCEC in  $M'$ , so  $D$  is trivial. ◀

904 To relate the histories of  $M$  to those of  $\text{purge}(M)$ , we introduce the function  $h \mapsto$   
 905  $\text{purge}(h)$  which, intuitively, maps the state of a non-distinguishing MCECs  $D$  to the state  
 906  $s_D$ , removes the state-actions pairs that stay in  $D$ , and replaces the state-action pairs  $(q, a)$   
 907 having a transition that leaves  $D$  by a new action  $F_{(q,a)}$ . Formally,  $\text{purge}(h)$  is obtained  
 908 from  $h = q_1 a_1 \dots q_n$  by applying the following transformation: for each non-distinguishing  
 909 MCEC  $D = (Q', A')$  of  $M$ ,

- 910 1. Replace the maximal suffix of  $h$  of the form  $q_i a_i \dots q_n$  such that for all  $i \leq k \leq n$ ,  $q_k \in Q'$   
 911 and  $a_k \in A'(q_k)$ , if such a suffix exists, by  $s_D$ ;
- 912 2. Remove all maximal factors of  $h$  of the form  $q_i a_i \dots q_j a_j$  satisfying  $q_k \in Q'$  and  $a_k \in A'(q_k)$   
 913 for all  $i \leq k \leq j$ ;
- 914 3. Replace each pair  $q_i a_i$  with  $q_i \in Q'$  and  $a_i \notin A'(q_i)$  by  $s_D F_{(q_i, a_i)}$ ;

915 ► **Example 18.** In the MEMDP of Fig. 5, with  $D$  containing the pairs  $(q_3, a)$  and  $(q_4, a)$ ,  
 916 for  $h = q_1 a q_3 a q_4 a q_3 a q_4 b q_4 b q_5$ , we get  $\text{purge}(h) = q_1 a s_D F_{q_4, b} s_D F_{q_4, b} q_5$ . Here we first apply  
 917 rule 2 above to the factor  $q_3 a q_4 a q_3 a$ , and get  $q_1 a q_4 b q_4 b q_5$ ; then an application of rule 3  
 918 yields  $\text{purge}(h) = q_1 a s_D F_{q_4, b} s_D F_{q_4, b} q_5$ . For the history  $h' = q_1 a q_3 a q_4 a q_3 a q_4$ , we would get  
 919 by rule 1,  $\text{purge}(h') = q_1 a s_D$ .

920 We establish a relation between  $M$  and  $\text{purge}(M)$  in Lemmas 19 and 20. These will be  
 921 used to give a memory bound for strategies for the quantitative case in Lemma 24.

922 In the Lemma 19, we only establish an inequality. This is because a given strategy  $\sigma$  of  
 923  $M$  may not be optimal within a non-distinguishing MCEC, while the construction  $\text{purge}(M)$   
 924 is based on the assumption that optimal strategies are used within each MCECs.

925 ► **Lemma 19.** Consider an MEMDP  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$ , and objective  $\varphi = \text{Parity}(p)$ , and  
 926 the map  $f : Q \rightarrow Q'$  relating states of  $M$  and those of  $\text{purge}(M) = \langle Q', A', (\delta'_e)_{e \in E} \rangle$ . For all  
 927  $q \in Q$ , and strategy  $\sigma$  for  $M$ , there exists  $\sigma'$  with  $\mathbb{P}_q^\sigma(M, \varphi) \leq \mathbb{P}_{f(q)}^{\sigma'}(\text{purge}(M), \varphi)$ .

928 **Proof.** Let us write  $M' = \text{purge}(M)$ . Consider  $q \in Q$ , and a strategy  $\sigma$  for  $M$ . We define  $\sigma'$   
 929 for  $M'$  as follows. For all histories  $h$  of  $M'$ , and action  $a \in A'(\text{last}(h))$ , we define

$$930 \quad \sigma'(h)(a) = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(ha) \mid \text{purge}^{-1}(h)]$$

931 for some arbitrary  $e \in E$  for which  $\mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)] > 0$ , if such  $e \in E$  exists; and  
 932 otherwise define  $\sigma'(h)$  arbitrarily. This quantity does not depend on  $e$  since, assuming

$$933 \quad \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)] > 0,$$

$$934 \quad \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(ha) \mid \text{purge}^{-1}(h)]$$

$$935 \quad = \sum_{\rho \in \text{purge}^{-1}(h)} \mathbb{P}_q^\sigma [M[e], \rho a' \mid \rho] \mathbb{P}_q^\sigma [M[e], \rho \mid \text{purge}^{-1}(h)]$$

$$936 \quad = \sum_{\rho \in \text{purge}^{-1}(h)} \sigma(\rho)(a') \mathbb{P}_q^\sigma [M[e], \rho \mid \text{purge}^{-1}(h)],$$

937 where  $a' = b$  if  $a$  has the form  $F(\_, b)$ , and  $a' = a$  otherwise (in which case we have  
 938  $a \in A(\text{last}(\rho))$ ). Moreover,  $\mathbb{P}_q^\sigma [M[e], \rho \mid \text{purge}^{-1}(h)]$  does not depend on  $e$  here since  
 939  $\text{purge}^{-1}(h)$  determines the outcomes of all transitions whose probability distributions differ  
 940 among environments because these were not erased by  $\text{purge}(\cdot)$ , and these probability  
 941 distributions are identical in the remaining transitions since they belong to non-distinguishing  
 942 MCECs.

943 For a history  $h$  of  $M'$  that ends in a state of the form  $s_D$ , and with  $\mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)] >$   
 944  $0$ , we let  $\sigma'$  take the action  $\text{stay}$  with probability  $\mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega \mid \text{purge}^{-1}(h)]$ , where  
 945  $D^\omega$  denotes the set of all runs that stay inside  $D$ . This probability is similarly independent  
 946 from the particular choice of  $e$ .

947 We prove that for all histories  $h$  of  $M'$  that do not contain  $\text{stay}$ ,  $a \in A'(\text{last}(h)) \setminus \{\text{stay}\}$ ,  
 948 and  $e \in E$ ,

$$949 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h] = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)], \quad (1)$$

$$950 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h \cdot a] = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h \cdot a)], \quad (2)$$

$$951 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h \cdot \text{stay}] = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega]. \quad (3)$$

952 We proceed by induction on the length of  $h$  to prove the above three properties.

953 Initially, if  $|h| = 1$ , then  $h = f(q)$  and  $\text{purge}^{-1}(h) = \{q\}$ . Then (1) follows trivially since  
 954 both sides are equal to 1. To see (2), note that, by definition of  $\sigma'$ ,

$$955 \quad \sigma'(h)(a) = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h \cdot a) \mid \text{purge}^{-1}(h)] = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h \cdot a)]$$

956 since  $h = f(q)$  here. Furthermore,

$$957 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h \cdot a] = \mathbb{P}_{f(q)}^{\sigma'} [M'[e], a \mid h] \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h]$$

$$958 \quad = \sigma'(h)(a)$$

959 since  $\mathbb{P}_{f(q)}^{\sigma'} [M'[e], h] = 1$ ; which yields (2).

960 Last, assume that  $\text{stay} \in A'(f(q))$ , that is  $f(q)$  has the form  $s_D$  for some non-distinguishing  
 961 MCEC  $D$ .

$$962 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h \cdot \text{stay}] = \sigma'(h)(\text{stay}) \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h]$$

$$963 \quad = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega \mid \text{purge}^{-1}(h)]$$

$$964 \quad = \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega],$$

965 which proves (3).

**XX:26 The Value Problem for Multiple-Environment MDPs with Parity Objective**

966 Assume now that  $|h| > 1$ , and let us write  $h = h'ar$  for a history  $h'$ ,  $a \in A'(\text{last}(h'))$ .

$$\begin{aligned}
 967 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'a] &= \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'a \mid h'] \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'] \\
 968 \quad &= \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'a \mid h'] \mathbb{P}_q^\sigma [M'[e], \text{purge}^{-1}(h')] \\
 969 \quad &= \sigma'(h')(a) \mathbb{P}_q^\sigma [M'[e], \text{purge}^{-1}(h')] \\
 970 \quad &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'a) \mid \text{purge}^{-1}(h')] \\
 971 \quad &\quad \cdot \mathbb{P}_q^\sigma [M'[e], \text{purge}^{-1}(h')], \\
 972 \quad &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'a)],
 \end{aligned}$$

973 where we used the induction hypothesis to apply (1) on the second line. This proves (2).

974 Consider now  $r \in Q'$ .

$$\begin{aligned}
 975 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'ar] &= \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'ar \mid h'a] \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'a] \\
 976 \quad &= \delta'_e(\text{last}(h'), a)(r) \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'a)].
 \end{aligned}$$

977 We distinguish two cases. If  $\text{last}(h)$  does not have the form of  $s_D$ , then it also belongs to  $Q$ ,  
 978  $a \in A(\text{last}(q))$ , with  $\delta_e(\text{last}(h'), a)(r) = \delta'_e(\text{last}(h'), a)(r)$ . In this case, the above is equal to  
 979  $\mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'ar)]$ . Assume now that  $\text{last}(h) = s_D$  for some non-distinguishing MCEC  
 980  $D$ , and that  $a = F_{(r', a')}$  for some pair  $(r', a')$ . Then  $\text{purge}^{-1}(h'a)$  only contains histories that  
 981 end at  $r'$ , followed by action  $a'$ . We have, moreover,  $\delta'_e(\text{last}(h'), a)(r) = \sum_{q \in f^{-1}(r)} \delta_e(r', a')(q)$ ,  
 982 by the definition of  $\text{purge}(M)$ , so

$$\begin{aligned}
 983 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h'ar] &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'a)f^{-1}(r)] \\
 984 \quad &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h'ar)].
 \end{aligned}$$

985 This proves (1).

986 Last, consider history  $h$  ending at some state  $s_D$ , and write

$$\begin{aligned}
 987 \quad \mathbb{P}_{f(q)}^{\sigma'} [M'[e], h \cdot \text{stay}] &= \mathbb{P}_{f(q)}^{\sigma'} [M[e], h \cdot \text{stay} \mid h] \mathbb{P}_{f(q)}^{\sigma'} [M[e], h] \\
 988 \quad &= \mathbb{P}_{f(q)}^{\sigma'} [M[e], h \cdot \text{stay} \mid h] \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)] \\
 989 \quad &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega \mid \text{purge}^{-1}(h)] \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)] \\
 990 \quad &= \mathbb{P}_q^\sigma [M[e], \text{purge}^{-1}(h)D^\omega],
 \end{aligned}$$

991 where we used the induction hypothesis on the second line, and the definition of  $\sigma'$  on the  
 992 third line. This proves (3).

993 We now show that  $\mathbb{P}_q^\sigma(M, \varphi) \leq \mathbb{P}_{f(q)}^{\sigma'}(M', \varphi)$  follows from these properties. In fact, for all  
 994  $e \in E$ , one can write

$$\begin{aligned}
 995 \quad \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \varphi) &= \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_{f(q)}^{\sigma'}(M'[e], (Q'A')^* \cdot s_D \cdot \text{stay} \cdot q_{\text{win}}) \\
 &\quad + \sum_{\substack{D \in \text{EC}(M'[e]), \varphi\text{-winning} \\ D \neq \{q_{\text{win}}, \dots\}}} \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \text{Inf} = D),
 \end{aligned} \tag{4}$$

996 by separating winning end-components of  $M'$  into two: the winning absorbing state  $q_{\text{win}}$   
 997 reached via some  $s_D$  for a non-distinguishing MCEC of  $M$ , and any other end-component of  
 998  $M'$ .

999 For the first term of (4), we have, from the above properties of  $\sigma'$ ,

$$\begin{aligned}
1000 & \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_{f(q)}^{\sigma'}(M'[e], (Q'A')^* \cdot s_D \cdot \text{stay} \cdot q_{\text{win}}) \\
1001 & = \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \sum_{\substack{D' \in \text{EC}(M[e]) \\ D' \subseteq D}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D'] \\
1002 & \geq \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \sum_{\substack{D' \in \text{EC}(M[e]), \varphi\text{-winning} \\ D' \subseteq D}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D'] \\
1003 & = \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D].
\end{aligned}$$

1004 For the second term of (4), let  $D \in \text{EC}(M') \setminus \{(q_{\text{win}}, \_)\}$ , and observe that  $\text{Inf} = D$  does  
1005 not contain the action **stay**. Notice how we only have an inequality because  $\sigma$  might actually  
1006 have a nonzero probability of realizing  $\text{Inf} = D'$  for some non-winning  $D'$  included in a  
1007 winning  $D$ .

1008 We established above that for all histories  $h$  of  $M'$  without the action **stay**,  $\mathbb{P}_{f(q)}^{\sigma'}[M'[e], h] =$   
1009  $\mathbb{P}_q^\sigma[M[e], \text{purge}^{-1}(h)]$ , that is, cylinders generated by  $h$  and  $\text{purge}^{-1}(h)$  have the same proba-  
1010 bilities in  $M'$  under  $\sigma'$ , and, respectively, in  $M$  under  $\sigma$ . It follows that

$$\begin{aligned}
1011 & \sum_{\substack{D \in \text{EC}(M'[e]), \varphi\text{-winning} \\ D \neq \{(q_{\text{win}}, \_)\}}} \mathbb{P}_{f(q)}^{\sigma'}[M'[e], \text{Inf} = D] \\
1012 & = \sum_{\substack{D \in \text{EC}(M'[e]), \varphi\text{-winning} \\ D \neq \{(q_{\text{win}}, \_)\}}} \mathbb{P}_q^\sigma[M[e], \text{purge}^{-1}(\text{Inf} = D)] \\
1013 & = \sum_{\substack{D \in \text{EC}(M'[e]), \varphi\text{-winning} \\ D \neq \{(q_{\text{win}}, \_)\}}} \sum_{D' \in \text{EC}(M), f(D')=D} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D'] \\
1014 & \geq \sum_{\substack{D \in \text{EC}(M'[e]), \varphi\text{-winning} \\ D \neq \{(q_{\text{win}}, \_)\}}} \sum_{\substack{D' \in \text{EC}(M), f(D')=D \\ D' \text{ is } \varphi\text{-winning}}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D'], \\
1015 & = \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{not a non-distinguishing MCEC}}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D].
\end{aligned}$$

1016 where we extend the definition of  $f$  to state-action pairs so that  $f(D')$  denotes an end-  
1017 component of  $M'$ .

1018 Combining these bounds on both terms of (4), we conclude

$$\begin{aligned}
1019 & \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \varphi) \geq \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D] \\
1020 & \quad + \sum_{\substack{D \in \text{EC}(M[e]), \varphi\text{-winning} \\ \text{not a non-distinguishing MCEC}}} \mathbb{P}_q^\sigma[M[e], \text{Inf} = D] \\
1021 & \geq \mathbb{P}_q^\sigma(M[e], \varphi).
\end{aligned}$$

1022



## XX:28 The Value Problem for Multiple-Environment MDPs with Parity Objective

1023 The following lemma is the dual, and shows that any strategy for  $\text{purge}(M)$  can be  
 1024 replicated in  $M$ , albeit with a bit more memory. The additional memory is required to  
 1025 implement behaviors inside non-distinguishing MCECs.

1026 **► Lemma 20.** *Consider an MEMDP  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$ , and objective  $\varphi = \text{Parity}(p)$ , and  
 1027 the map  $f : Q \rightarrow Q'$  relating states of  $M$  and that of  $\text{purge}(M) = \langle Q', A', (\delta'_e)_{e \in E} \rangle$ . For all  
 1028 states  $q' \in Q'$  and strategies  $\sigma'$  for  $\text{purge}(M)$ , and all  $q \in f^{-1}(q')$ , there exists a strategy  $\sigma$   
 1029 with  $\mathbb{P}_q^\sigma(M, \varphi) = \mathbb{P}_{q'}^{\sigma'}(\text{purge}(M), \varphi)$ . Furthermore, if  $\sigma'$  is a  $m$ -memory strategy,  $\sigma$  can be  
 1030 chosen to be a  $(m + |Q||A|)$ -memory strategy.*

1031 **Proof.** Consider  $q' \in Q'$ , an  $m$ -memory strategy  $\sigma'$  for  $M'$ , and  $q \in f^{-1}(q')$ , where  $m$  can  
 1032 be finite or infinite. We show that there exists an  $(m + |Q||A|)$ -memory strategy  $\sigma$  with  
 1033  $\mathbb{P}_q^\sigma(M, \varphi) = \mathbb{P}_{q'}^{\sigma'}(M', \varphi)$ . We define  $\sigma$  as follows. Consider a history  $h$  of  $M$ .

- 1034 **■** If  $f(\text{last}(h)) \in Q$ , we let  $\sigma(h) = \sigma'(\text{purge}(h))$ .
- 1035 **■** Assume  $f(\text{last}(h)) = s_D$  for some non-distinguishing MCEC  $D$ . With probability  
 1036  $\sigma'(\text{purge}(h))(\text{stay})$ , we let  $\sigma$  switch to a pure memoryless strategy that maximizes the  
 1037 probability of  $\varphi$  inside  $D$  (this strategy is independent from the environment). For each  
 1038  $F_{(q,a)}$ , with probability  $\sigma'(\text{purge}(h))(F_{(q,a)})$  we let  $\sigma$  run a pure memoryless strategy until  
 1039 state  $q$  is reached (which happens probability 1), and from  $q$  take  $a$ .

1040 The memory bound for  $\sigma$  is  $m + |Q||A|$  where  $m$  is the memory size of  $\sigma'$ , because inside  
 1041 each collapsed MCEC, and for each pair  $F_{(q,a)}$ , a pure memoryless strategy is executed until  
 1042 reaching  $q$  and taking action  $a$ .

1043 By construction, for all histories  $h$  that start at  $q$  in  $M$  and end outside of non-  
 1044 distinguishing end-components, we have:

$$1045 \quad \mathbb{P}_q^\sigma(M[e], h) = \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \text{purge}(h)) \text{ for all environments } e \in E. \quad (5)$$

1046 So if  $R$  denotes a measurable set of infinite runs of  $M$  such that for all  $\rho \in R$ ,  $\text{purge}(\rho)$  is  
 1047 infinite (in other terms,  $\rho$  does not stay inside a non-distinguishing MCEC), then

$$1048 \quad \mathbb{P}_q^\sigma(M[e], R) = \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \text{purge}(R)), \quad (6)$$

1049 writing  $\text{purge}(R) = \{\text{purge}(\rho) \mid \rho \in R\}$ .

1050 Furthermore, for those histories  $h = h'as$  where  $q \in D$  is a non-distinguishing MCEC  
 1051 and  $\text{last}(h') \notin D$ , we have:

$$1052 \quad \mathbb{P}_q^\sigma(M[e], h) = \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \text{purge}(h)) \text{ for all environments } e \in E.$$

1053 Then, by definition of  $\sigma$ , for a non-distinguishing MCEC  $D$ ,

$$1054 \quad \sum_{D' \in \text{EC}(M), D' \subseteq D} \mathbb{P}_q^\sigma(M[e], \text{Inf} = D') = \mathbb{P}_{f(q)}^{\sigma'}(M'[e], (Q'A')^* \cdot s_D \cdot \text{stay}). \quad (7)$$

1055 Observe that any end-component  $D$  of  $M$  that is not a non-distinguishing MCEC maps

1056 to an end-component of  $M'$ . We have for all  $e \in E$ ,

$$\begin{aligned}
1057 \quad \mathbb{P}_q^\sigma(M[e], \varphi) &= \sum_{\substack{D \in \text{EC}(M), \varphi\text{-winning} \\ \text{not a non-distinguishing MCEC}}} \mathbb{P}_q^\sigma(M[e], \text{Inf} = D) \\
1058 \quad &+ \sum_{\substack{D \in \text{EC}(M), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_q^\sigma(M[e], \text{Inf} = D) \\
1059 \quad &= \sum_{\substack{D \in \text{EC}(M'), \varphi\text{-winning} \\ \text{not a non-distinguishing MCEC}}} \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \text{Inf} = D) \\
1060 \quad &+ \sum_{\substack{D \in \text{EC}(M'), \varphi\text{-winning} \\ \text{non-distinguishing MCEC}}} \mathbb{P}_{f(q)}^{\sigma'}(M'[e], (Q' A')^* \cdot s_D \cdot \text{stay}) \\
1061 \quad &= \mathbb{P}_{f(q)}^{\sigma'}(M'[e], \varphi),
\end{aligned}$$

1062 using (6) and (7). ◀

## 1063 5.2 Learning While Playing

1064 In this section, we show that after collapsing non-distinguishing MCECs, over  $n$  steps (for  $n$   
1065 large enough), with high probability, we either reach a MCEC (which is either distinguishing  
1066 or trivial) or collect a large number of samples of distinguishing transitions whose empirical  
1067 average is close to their mean. Intuitively, this means that either the knowledge can be  
1068 improved after  $n$  steps using the collected samples while bounding the probability of error,  
1069 or a MCEC is reached.

1070 If the MCEC is distinguishing, the strategy can improve the knowledge as in Lemma 9,  
1071 and if not, then the MCEC is trivial and there is a unique way to play. These results will  
1072 be used in the next section to build a finite-memory strategy with approximately the same  
1073 probability of winning, given any arbitrary strategy.

1074 For a history  $h$ , let  $|h|_{q,a}$  denote the number of occurrences of the state-action pair  $(q, a)$ ,  
1075 and  $|h|_{q,a,q'}$  the number of times these are followed by  $q'$ , where  $q' \in \text{Supp}(\delta(q, a))$ . For a  
1076 distinguishing transition  $t = (q, a, q')$ , we say that a history  $h$  is a *bad*  $(t, \eta)$ -*classification* in  
1077 MDP  $M[e]$  if  $\left| \frac{|h|_{q,a,q'}}{|h|_{q,a}} - \delta_e(q, a)(q') \right| \geq \eta/2$ , that is the measured and theoretical frequency of  
1078  $t$  are too far apart. It is a *good*  $(t, \eta)$ -*classification* otherwise. Intuitively, over long histories,  
1079 good classifications have high probability.

1080 We first prove the following technical lemma, bounding the difference between the  
1081 empirical average and the mean when sampling among a finite number of transitions, when  
1082 the transitions to sample are chosen at each step by an *adversary*. This adversary corresponds  
1083 to strategies in an MDP, is arbitrary, and can depend on the history and use randomization.

1084 We state the following lemma for (single-environment) MDPs, and apply it to each  
1085 environment in an MEMDP.

1086 ► **Lemma 21.** *Consider MDP  $M$ , state  $q_0$ , and  $T = \{t_i = (q_i, a_i, q'_i)\}_{1 \leq i \leq k}$  a subset of*  
1087 *transitions such that  $(q_i, a_i) = (q_j, a_j)$  implies  $q'_i = q'_j$  for all  $i, j$ . For all  $\eta, \varepsilon > 0$ , all  $n_0 > \frac{k^3}{\varepsilon \eta^2}$ ,*  
1088 *and any strategy  $\sigma$  with  $\mathbb{P}_{q_0}^\sigma \left[ \{h : \sum_{(q,a,q') \in T} |h_{q,a}| \geq n_0 \} \right] = 1$ , the set of histories  $h$  that*  
1089 *satisfy the following conditions has probability at most  $\varepsilon$ :*

- 1090 ■  $\sum_{(q,a,q') \in T} |h_{q,a}| \geq n_0$
- 1091 ■ there exists  $1 \leq i \leq k$  such that  $|h|_{q_i, a_i} = \max_{i'} |h|_{q_i', a_i'}$  and  $h$  is a *bad*  $(t_i, \eta)$ -*classification*.

## XX:30 The Value Problem for Multiple-Environment MDPs with Parity Objective

1092 Here the assumption on  $T$  simplifies the proofs since it means that for each state-action  
 1093 pair, we will be observing the frequency of a unique successor state. The lemma also requires  
 1094 that at least  $n_0$  occurrences of  $T$  is visited with probability 1. This hypothesis ensures that  
 1095 we have enough samples to obtain a good approximation (that is, a good  $(t_i, \eta)$ -classification)  
 1096 with high probability (at least  $1 - \varepsilon$ ). In fact, if a strategy  $\sigma$  avoids visiting transitions  
 1097 from  $T$ , say, with probability  $1/2$ , then it cannot ensure a good approximation with high  
 1098 probability because half the cases, there are just not enough samples of  $T$ .

1099 The lemma is easy for  $k = 1$ . In fact, all trials are identical and independent, so one  
 1100 can use e.g. Hoeffding's inequality to derive a bound. When  $k > 1$ , trials are no longer  
 1101 independent since  $\sigma$  might react to the success or failure of a given transition to make  
 1102 its decisions in the future. In fact, the lemma is not trivial to prove due to the possible  
 1103 dependency between the trials.

1104 Here is such a situation of dependency. Consider a state  $q$  from which action  $a$  leads to  
 1105 either to  $q_1$  or  $q_2$ , each with probability 0.5, from which a deterministic transition comes back  
 1106 to  $q$ . Another action  $b$  from  $q$  deterministically loops back at  $q$ . Consider  $\sigma$  that picks  $(q, a)$   
 1107 first. As long as we reach  $q_1$ ,  $\sigma$  continues to pick  $(q, a)$ . Whenever  $q_2$  is reached,  $\sigma$  switches  
 1108 definitively to  $(q, b)$ . Now the probability of observing  $(q, a, q_1)$  at step  $n > 1$  depends on  
 1109 the result of the first  $n - 1$  trials. For example, conditioned on observing  $(q, a, q_1)$  on the  
 1110 first  $n - 1$  trials, the probability of observing  $(q, a, q_1)$  again is 0.5. But conditioned on not  
 1111 observing  $(q, a, q_1)$  on the  $(n - 1)$ -th trial, this probability is 0. This shows that given such  $\sigma$ ,  
 1112 the successive trials are not independent, and theorems such as Hoeffding's inequality cannot  
 1113 be applied.

1114 It turns out that although the trials can be dependent, their covariance is 0. We exploit  
 1115 this observation to derive a good bound using Chebyshev's inequality:

1116 ► **Theorem 22** (Chebyshev's Inequality). *Let  $X$  be a random variable with mean  $\mu$ , and*  
 1117 *standard deviation  $q$ . Then, for all  $a > 0$ , we have  $\mathbb{P}[|X - \mu| \geq sa] \leq \frac{1}{a^2}$ .*

1118 This inequality clearly also applies if  $q$  is an upper bound on the standard deviation of  $X$ .

1119 **Proof of Lemma 21.** We consider a slightly more abstract setting where there are  $k$  inde-  
 1120 pendent arms, each with a probability of success of  $p_i$ . In MDPs, each arm corresponds to a  
 1121 state-action pair  $(q_i, a_i)$  and it succeeds when reaching  $q'_i$ , with probability  $p_i = \delta(q_i, a_i)(q'_i)$ .

1122 Consider a strategy  $\sigma$  that chooses, at each step,  $i \in \{1, \dots, k\}$ , an arm to pull based on  
 1123 the full history and randomization. Consider  $\varepsilon, \eta > 0$ .

1124 We model the problem as follows. For each  $i \in \{1, \dots, k\}$ , define a sequence  $X_1^{(i)}, X_2^{(i)}, \dots$   
 1125 of identical and independent Bernoulli variables with probability  $p_i$ . Let  $\text{Choice}_j$  denote the  
 1126 arm selected by  $\sigma$  at step  $j$ . At each step  $j$ ,  $\text{Choice}_j$  selects an arm, and all types of arms  
 1127 are pulled. While  $\text{Choice}_j$  can depend on the history,  $X_j^{(i)}$  does not depend on the history,  
 1128 and in particular on  $\text{Choice}_j$ .

1129 Define the *weight* of arm  $i$  at step  $j$  as the following random variable.

$$1130 \quad \text{wgt}_j^{(i)} = \begin{cases} X_j^{(i)} - p_i & \text{if } \text{Choice}_j = i, \\ 0 & \text{otherwise.} \end{cases}$$

1131 Define  $\text{wgt}_{\leq n}^{(i)} = \sum_{j=1}^n \text{wgt}_j^{(i)}$ , for any  $n \geq 1$ . Let us also define  $\text{occ}_j^{(i)} = 1$  iff  $\text{Choice}_j = i$ , and  
 1132  $\text{occ}_{\leq n}^{(i)} = \sum_{j=1}^n \text{occ}_j^{(i)}$ . Observe that

$$1133 \quad \text{wgt}_{\leq n}^{(i)} = \sum_{1 \leq j \leq n, \text{Choice}_j = i} X_j^{(i)} - \text{occ}_{\leq n}^{(i)} p_i,$$

1134 that is, this is the difference between the empirical sum and the mean of the sum of the  
1135 subsequence of  $X_j^{(i)}$  where  $\text{Choice}_j = i$ .

1136 Then  $\frac{\text{wgt}_{\leq n}^{(i)}}{\text{occ}_{\leq n}^{(i)}}$  is the difference between the empirical average of the  $X_j^{(i)}$  and  $p_i$ , assuming

1137 that  $\text{occ}_{\leq n}^{(i)} > 0$ .

1138 We have, by the definition of variance,

$$1139 \quad \mathbb{E}^\sigma[\text{wgt}_j^{(i)}] = \mathbb{P}^\sigma[\text{Choice}_j = i](p_i(1 - p_i) + (1 - p_i)(-p_i)) = 0,$$

1140 so  $\mathbb{E}^\sigma[\text{wgt}_{\leq n}^{(i)}] = 0$  as well.

1141 We are going to apply Theorem 22 on the variable  $\text{wgt}_{\leq n}^{(i)}$ ; so we need a bound on the  
1142 variance of  $\text{wgt}_{\leq n}^{(i)}$ . We show that  $\mathbb{V}^\sigma[\text{wgt}_{\leq n}^{(i)}] \leq np_i(1 - p_i)$ . We have

$$1143 \quad \mathbb{V}^\sigma[\text{wgt}_{\leq n}^{(i)}] = \sum_{j=1}^n \mathbb{V}^\sigma[\text{wgt}_j^{(i)}] + 2 \sum_{1 \leq j < j' \leq n} \text{Cov}(\text{wgt}_j^{(i)}, \text{wgt}_{j'}^{(i)})$$

1144 For each  $j$ , because  $\mathbb{E}^\sigma[\text{wgt}_j^{(i)}] = 0$ , we have  $\mathbb{V}^\sigma[\text{wgt}_j^{(i)}] = \mathbb{E}^\sigma[(\text{wgt}_j^{(i)})^2]$ , which can be  
1145 calculated as

$$\begin{aligned} 1146 \quad & \mathbb{P}^\sigma[\text{Choice}_j = i](p_i(1 - p_i)^2 + (1 - p_i)(-p_i)^2) \\ 1147 \quad & = \mathbb{P}^\sigma[\text{Choice}_j = i]p_i(1 - p_i)((1 - p_i) + p_i) \\ 1148 \quad & \leq p_i(1 - p_i), \end{aligned}$$

1149 so that the first term of the variance is at most  $np_i(1 - p_i)$ .

1150 Now, as noted above,  $\text{wgt}_j^{(i)}$  and  $\text{wgt}_{j'}^{(i)}$  are not independent variables since  $\sigma$  can choose  
1151 the arm at step  $j'$  depending on the result of  $\text{wgt}_j^{(i)}$ ; we nevertheless show that the covariance  
1152 is equal to 0. We have  $\text{Cov}(\text{wgt}_j^{(i)}, \text{wgt}_{j'}^{(i)}) = \mathbb{E}^\sigma[\text{wgt}_j^{(i)} \cdot \text{wgt}_{j'}^{(i)}] - \mathbb{E}^\sigma[\text{wgt}_j^{(i)}]\mathbb{E}^\sigma[\text{wgt}_{j'}^{(i)}]$  by definition  
1153 of covariance; so this is equal to  $\mathbb{E}^\sigma[\text{wgt}_j^{(i)} \cdot \text{wgt}_{j'}^{(i)}]$  which can be calculated as follows.

$$\begin{aligned} 1154 \quad & \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1 \wedge X_{j'}^{(i)} = 1](1 - p_i)^2 \\ 1155 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1 \wedge X_{j'}^{(i)} = 0](1 - p_i)(-p_i) \\ 1156 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 0 \wedge X_{j'}^{(i)} = 1](-p_i)(1 - p_i) \\ 1157 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 0 \wedge X_{j'}^{(i)} = 0](-p_i)^2. \end{aligned}$$

1158 Now  $X_{j'}^{(i)}$  and the variables  $\text{Choice}_j, \text{Choice}_{j'}, X_j^{(i)}$  are independent; in fact, the values of  
1159  $\text{Choice}_j, \text{Choice}_{j'}$  cannot depend on  $X_{j'}^{(i)}$  since the latter is revealed after  $\text{Choice}_j, \text{Choice}_{j'}$ .  
1160 In contrast,  $X_j^{(i)}$  and  $\text{Choice}_{j'}$  can be dependent since the latter can depend on the value of  
1161  $X_j^{(i)}$ .

1162 We can rewrite  $\mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1 \wedge X_{j'}^{(i)} = 1]$  as follows.

$$\begin{aligned} 1163 \quad & \mathbb{P}^\sigma[X_{j'}^{(i)} = 1 \mid \text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1] \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1] \\ 1164 \quad & = \mathbb{P}^\sigma[X_{j'}^{(i)} = 1] \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1] \\ 1165 \quad & = p_i \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1] \end{aligned}$$

1166 by independence.

**XX:32 The Value Problem for Multiple-Environment MDPs with Parity Objective**

1167 Applying this to all four terms,  $\mathbb{E}^\sigma[\text{wgt}_j^{(i)} \cdot \text{wgt}_{j'}^{(i)}]$  can be written as

$$\begin{aligned}
 1168 \quad & \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1]p_i(1-p_i)^2 \\
 1169 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1](1-p_i)(1-p_i)(-p_i) \\
 1170 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 0]p_i(-p_i)(1-p_i) \\
 1171 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 0](1-p_i)(-p_i)^2, \\
 1172 \quad & = \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 1](p_i(1-p_i)^2 + (1-p_i)(1-p_i)(-p_i)) \\
 1173 \quad & + \mathbb{P}^\sigma[\text{Choice}_j = i \wedge \text{Choice}_{j'} = i \wedge X_j^{(i)} = 0](p_i(-p_i)(1-p_i) + (1-p_i)(-p_i)^2) \\
 1174 \quad & = 0.
 \end{aligned}$$

1175 So all covariance terms are 0, and we have  $\mathbb{V}^\sigma[\text{wgt}_{\leq n}^{(i)}] \leq np_i(1-p_i)$ .

1176 We now apply Theorem 22: For all  $1 \leq i \leq k$ , and for all  $a > 0$ ,

$$1177 \quad \mathbb{P}^\sigma \left[ \left| \text{wgt}_{\leq n}^{(i)} \right| \geq a\sqrt{np_i(1-p_i)} \right] \leq \frac{1}{a^2}.$$

1178 Using  $\mathbb{P}[X \cup Y] = \mathbb{P}[X] + \mathbb{P}[Y] - \mathbb{P}[X \cdot Y]$ , it follows that

$$1179 \quad \mathbb{P}^\sigma \left[ \exists i, \left| \text{wgt}_{\leq n}^{(i)} \right| \geq a\sqrt{np_i(1-p_i)} \right] \leq \frac{k}{a^2}.$$

1180 We have,

$$1181 \quad \mathbb{P}^\sigma \left[ \exists i, \text{occ}_{\leq n}^{(i)} = \max_{i'} \text{occ}_{\leq n}^{(i')} \wedge \left| \frac{\text{wgt}_{\leq n}^{(i)}}{\text{occ}_{\leq n}^{(i)}} \right| \geq \frac{a\sqrt{np_i(1-p_i)}}{\text{occ}_{\leq n}^{(i)}} \right] \leq \frac{k}{a^2}.$$

1182 Here we divided the inequality by  $\text{occ}_{\leq n}^{(i)}$  (since for  $n > 0$ ,  $\max_{i'} \text{occ}_{\leq n}^{(i')} > 0$ ); moreover,  
1183 the probability bound holds since each event has become smaller.

1184 Notice that for all  $n > 0$ ,  $\sum_{i=1}^k \text{occ}_{\leq n}^{(i)} = n$  since  $\sigma$  picks one of the arms at each step; so

1185  $\max_{i'} \text{occ}_{\leq n}^{(i')} \geq n/k$  with probability 1. We get

$$1186 \quad \mathbb{P}^\sigma \left[ \exists i, \text{occ}_{\leq n}^{(i)} = \max_{i'} \text{occ}_{\leq n}^{(i')} \wedge \left| \frac{\text{wgt}_{\leq n}^{(i)}}{\text{occ}_{\leq n}^{(i)}} \right| \geq \frac{ak\sqrt{p_i(1-p_i)}}{\sqrt{n}} \right] \leq \frac{k}{a^2}.$$

1187 Now, given  $\varepsilon, \eta > 0$ , we pick  $a = \sqrt{k/\varepsilon}$  so that  $k/a^2 \leq \varepsilon$ ; and then  $n$  large enough so that

$$1188 \quad \frac{ak\sqrt{p_i(1-p_i)}}{\sqrt{n}} \leq \eta/2; \text{ this means it suffices to pick } n \text{ such that } \max_{1 \leq i \leq k} \left( \frac{ak\sqrt{p_i(1-p_i)}}{\eta/2} \right)^2 \leq n,$$

1189 so  $(\frac{ak}{\eta/2})^2 = \frac{4k^3}{\varepsilon\eta^2} \leq n$  suffices.  $\blacktriangleleft$

1190 We use Lemma 21 to prove that in MEMDPs without non-trivial and non-distinguishing

1191 MCECs (for example, obtained by `purge(\cdot)`), after  $n$  steps, we either reach a MCEC or collect

1192 a large number of samples of distinguishing transitions whose empirical average is close to

1193 their mean. Given MEMDP  $M$ , let  $T_M$  denote a set obtained by selecting one distinguishing

1194 transition  $(q, a, q')$  for each state-action pair  $(q, a)$  whose probability distribution differs in a

1195 pair of different environments. We select one representative distinguishing transition  $(q, a, q')$

1196 for each pair  $(q, a)$  because this simplifies the calculations. We let  $|h|_{T_M} = \sum_{(q,a,q') \in T_M} |h|_{q,a}$ .

1197 Let us fix  $\eta$  as follows

$$1198 \quad \eta < \frac{1}{2} \min (\{ |\delta_e(q, a)(q') - \delta_f(q, a)(q')| \mid e, f \in E, q, q' \in Q, a \in A \} \setminus \{0\}).$$

1199 Let us define the set of *good histories with  $n_0$  samples*, denoted  $\text{Good}_{n_0}$ , as the set of  
1200 histories  $h$  satisfying

1201 ■  $|h|_{T_M} \geq n_0$ ,

1202 ■ for all  $t = (q, a, \_ ) \in T_M$  satisfying  $|h|_{q,a} = \max_{(q',a',\_) \in T_M} |h|_{q',a'}$ ,  $h$  is a good  $(t, \eta)$ -  
1203 classification.

1204 ► **Lemma 23.** *Consider an MEMDP  $M$  whose only non-distinguishing MCECs are trivial,*  
1205 *and fix  $\varepsilon > 0$ . Let  $n_0 = \lceil \frac{2(|Q||A|)^3}{\varepsilon\eta^2} \rceil$ , and  $n \geq 2p^{-2|Q|} \max(\log(\frac{4}{\varepsilon}), n_0)$  where  $p$  is the smallest*  
1206 *nonzero probability that appears in  $M$ . Then, from any starting state, and under any strategy,*  
1207 *with probability at least  $1 - \varepsilon$ , within  $n$  steps, the history either visits a MCEC (distinguishing*  
1208 *or trivial), or belongs to  $\text{Good}_{n_0}$ .*

1209 **Proof.** We show that in all  $M[e]$ , under any strategy  $\sigma$ , from every state  $q_0$ , there is a path  
1210 of size at most  $|Q|$  compatible with the strategy that reaches a MCEC or a distinguishing  
1211 transition. Consider first the case of a pure strategy  $\sigma$ . To prove this, towards a contradiction,  
1212 assume that MCECs and distinguishing transitions are not visited within  $|Q|$  steps under  
1213  $\sigma$ . Consider the execution tree that starts at  $q_0$  in  $M$  under  $\sigma$ : this is a tree labeled by  
1214  $Q$ , in which the children of a given node at history  $h$  are labeled by all possible successors  
1215  $\text{Supp}(\delta_e(\text{last}(h), \sigma(h)))$  for some  $e \in E$ . Since all transitions are non-distinguishing, the  
1216 choice of  $e$  is irrelevant here. We build this tree and cut each branch whenever a MCEC or a  
1217 distinguishing transition is seen, or a state is repeated. Since we assumed that MCECs and  
1218 distinguishing transitions are not reachable under  $\sigma$ , all branches of this tree are cut only  
1219 when a state is repeated. It follows that the set of states in this tree, together with the actions  
1220 prescribed by  $\sigma$  from these histories form a closed set of states. But then a strongly-connected  
1221 subset must exist, which is a non-distinguishing CEC. This is thus included in a MCEC,  
1222 contradicting our assumption.

1223 If  $\sigma$  is pure, then in all  $M[e]$ , from every history, there is a probability of at least  $p^{|Q|}$  of  
1224 either taking a distinguishing transition, or visiting a MCEC within  $|Q|$  steps, independently  
1225 of the current state. If  $\sigma$  is randomized, then the probability of such a single run can be  
1226 smaller since  $\sigma$  might assign small probabilities to its actions. In this case, since we are  
1227 only interested in the behaviors in the first  $n$  steps, we can see  $\sigma$  as a *mixed* strategy which  
1228 consists in randomly choosing among a set of pure strategies that stop after  $n$  steps. Since  
1229 the above argument can be applied to each pure strategy in the support of  $\sigma$  (when  $\sigma$  is seen  
1230 as a mixed strategy), it follows that under  $\sigma$ , there is a probability of at least  $p^{|Q|}$  of taking  
1231 a distinguishing transition or visiting a MCEC within the next  $|Q|$  steps, as well.

1232 Viewing runs as the concatenation of finite segments of size  $|Q|$ , we call each such segment  
1233 a trial. Consider the random Bernoulli variables  $X_1, X_2, \dots$  such that the value of  $X_i$  is 1 iff  
1234 a MCEC or a distinguishing transition is visited at the  $i$ -th trial.

1235 So by Hoeffding's inequality, for all states  $q_0$ ,  $n > 0$  and  $t > 0$ ,<sup>2</sup>

1236 
$$\mathbb{P}_{q_0}^{\sigma} \left[ \sum_{i=1}^n X_i \leq \sum_{i=1}^n \mathbb{E}^{\sigma}[X_i] - t \right] \leq 2e^{-2\frac{t^2}{n}}.$$

<sup>2</sup> Note that Hoeffding's inequality requires an independent sequence of random variables which is not the case of the  $X_i$ 's. We can nevertheless still apply this inequality here using a coupling argument: Define  $U_i$  as a sequence of independent and continuous variables uniformly distributed over  $[0, 1]$ . Define the Bernoulli variable  $Y_i = 1$  iff  $U_i \leq p^{|Q|}$ . Furthermore, define the sequence of Bernoulli variables  $\tilde{X}_1, \tilde{X}_2, \dots$  inductively, by  $\tilde{X}_i = 1$  iff  $U_i \leq \mathbb{P}(X_i = 1 \mid X_1 = \tilde{X}_1, \dots, X_{i-1} = \tilde{X}_{i-1})$ . Because  $\mathbb{P}(X_i = 1 \mid h) \geq p^{|Q|}$  regardless of the history  $h$ , we have  $Y_i \leq \tilde{X}_i$ . Furthermore,  $\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \mathbb{P}(\tilde{X}_1 = x_1, \dots, \tilde{X}_n = x_n)$  for all  $x_1, \dots, x_n \in \{0, 1\}$ . It follows that Hoeffding's inequality can be applied on the i.i.d. sequence  $Y_i$ , and we get for all  $A > 0$ ,  $\mathbb{P}[\sum_i X_i \leq A] \leq \mathbb{P}[\sum_i Y_i \leq A]$ .

## XX:34 The Value Problem for Multiple-Environment MDPs with Parity Objective

1237 Given  $n > 0$ , we choose here  $t = np^{|Q|}/2$ . This yields,

$$1238 \quad \mathbb{P}_{q_0}^\sigma \left[ \sum_{i=1}^n X_i \leq \sum_{i=1}^n \mathbb{E}^\sigma[X_i] - np^{|Q|}/2 \right] \leq 2e^{-2\frac{n^2 p^{2|Q|}}{4n}} \leq \varepsilon/2,$$

1239 which is the case since, by taking the log of both sides,

$$1240 \quad -\frac{np^{2|Q|}}{2} \leq \log(\varepsilon/4)$$

$$1241 \quad \Leftrightarrow n \geq 2 \log(4/\varepsilon) p^{-2|Q|}.$$

1242 Because  $\mathbb{E}^\sigma(X_i) \geq p^{|Q|}$ ,  $\sum_{i=1}^n \mathbb{E}^\sigma[X_i] \geq np^{|Q|}$ . This means that with probability at least  
 1243  $1 - \varepsilon/2$ ,  $\sum_{i=1}^n X_i \geq np^{|Q|}/2$ . As  $n \geq 2 \lceil \frac{2(|Q||A|^3)}{\varepsilon\eta^2} \rceil p^{-2|Q|}$ , we have  $\sum_{i=1}^n X_i \geq \lceil \frac{2(|Q||A|^3)}{\varepsilon\eta^2} \rceil$ ,  
 1244 that is, with probability at least  $1 - \varepsilon/2$ , either a MCEC or  $\lceil \frac{2(|Q||A|^3)}{\varepsilon\eta^2} \rceil$  occurrences of  
 1245 distinguishing transitions are seen (which can be good or bad classifications).

1246 Let us write  $n_0 = \lceil \frac{2(|Q||A|^3)}{\varepsilon\eta^2} \rceil$ . It remains to bound the probability of visiting either  
 1247 a MCEC or  $\text{Good}_{n_0}$ . Let us define a tree-shaped MDP  $M_n$  from  $M$  as follows. First, we  
 1248 unfold  $M$  by stopping each branch either when a MCEC is reached, or after  $n$  steps. Then,  
 1249 each leaf that belongs to a MCEC is extended with fresh states and transitions so that the  
 1250 branch contains  $n_0$  instances of distinguishing transitions. More precisely, we pick some  
 1251 distinguishing transition  $(q, a, q')$  of  $M$ , and extend a given leaf  $l_0$  of  $M_n$  as follows. The only  
 1252 enabled action at  $l_0$  is  $a$ , and it goes to  $l'_0$  with probability  $\delta_e(q, a, q')$  to  $l'_0$  in  $M_n[e]$ , and to  
 1253  $l''_0$  with probability  $1 - \delta_e(q, a, q')$ ; and both  $l'_0, l''_0$  deterministically go to  $l_1$ . We repeat this  
 1254 until  $n_0$  occurrences of distinguishing transitions are obtained. Last, all leafs are made into  
 1255 absorbing states.

1256 Let  $\diamond\text{CEC}$  denote the set of histories that reach a MCEC. For all  $e \in E$ ,

$$1257 \quad \mathbb{P}_{q_0}^\sigma [M[e], \diamond\text{CEC} \vee \text{Good}_{n_0}] \geq \mathbb{P}_{q_0}^\sigma [M_n[e], \text{Good}_{n_0}]$$

1258 where  $q'_0$  is the root of  $M_n[e]$ , since MCECs are replaced with a gadget that might not satisfy  
 1259  $\text{Good}_{n_0}$  with probability 1.

1260 As an additional step, we obtain  $M'_n$  by modifying  $M_n$  as follows: we extend each leaf  
 1261 whose branch does not contain  $n_0$  occurrences of distinguishing transitions (nor visit a  
 1262 MCEC), by adding fresh states and transitions as described above so that a total of  $n_0$   
 1263 distinguishing transitions is obtained at each branch. We get for all  $e \in E$ .

$$1264 \quad \mathbb{P}_{q'_0}^\sigma [M_n[e], \text{Good}_{n_0}] \geq \mathbb{P}_{q'_0}^\sigma [M'_n[e], \text{Good}_{n_0}] - \varepsilon/2$$

1265 since the probability of the modified branches was shown to be at most  $\varepsilon$  above.

1266 Now, by construction, for all strategies  $\sigma$  and  $e \in E$ ,  $n_0$  occurrences of distinguishing  
 1267 transitions are seen in  $M'_n[e]$  with probability 1. By Lemma 21 with  $k = |Q||A|$ , applied for  
 1268  $\varepsilon/2$ , we get

$$1269 \quad \mathbb{P}_{q'_0}^\sigma [M'_n[e], \text{Good}_{n_0}] \geq 1 - \varepsilon/2.$$

1270 It follows that  $\mathbb{P}_q^\sigma [M[e], \diamond\text{CEC} \vee \text{Good}_{n_0}] \geq \mathbb{P}_{q'_0}^\sigma [M'_n[e], \text{Good}_{n_0}] \geq 1 - \varepsilon$  for all  $e \in E$ , as  
 1271 required.  $\blacktriangleleft$

### 1272 5.3 Constructing Approximate Finite-Memory Strategies

1273 We are now ready to construct a finite-memory strategy that approximates an arbitrary  
 1274 strategy  $\sigma$ . We construct a finite-memory strategy for  $\text{purge}(M)$  and then transfer it to  $M$

1275 using Lemmas 19-20. The finite-memory strategy we construct consists in imitating the  
 1276 strategy  $\sigma$  for  $n$  steps, where  $n$  is defined in Lemma 23. Because all nontrivial MCECs of  
 1277  $\text{purge}(M)$  are distinguishing, when we play for  $n$  steps, with high probability, we either visit a  
 1278 trivial MCEC (which is either winning for all environments or losing for all environments), or  
 1279 reach a distinguishing MCEC, or observe enough samples of distinguishing transitions. The  
 1280 strategy is extended arbitrarily in trivial MCECs. Inside distinguishing MCECs, it gathers  
 1281 samples of distinguishing transitions as in Lemma 9, which improves the knowledge (with an  
 1282 arbitrarily small probability of error). The knowledge is also correctly improved with high  
 1283 probability if enough samples are gathered outside of MCECs. In both cases, the strategy  
 1284 switches to a finite-memory strategy for the improved knowledge constructed recursively for  
 1285 smaller sets of environments.

1286 Lemma 24 formalizes this reasoning and gives a bound  $N$  on the memory of the resulting  
 1287 strategy. In the memory bound, the term  $\lceil 8 \frac{\log(8/\varepsilon)}{\eta^2} \rceil$  comes from the application of Lemma 9  
 1288 for distinguishing MCECs for each subset of  $E$ ; and the term  $(2|Q|)^{n(|E|+1)}$  corresponds to  
 1289 the recursive analysis, since the strategy is defined inductively for each subset of  $E$ .

1290 **► Lemma 24.** *Consider an MEMDP  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$ , state  $q \in Q$ , parity objec-*  
 1291 *tive  $\varphi$ . For all strategies  $\sigma$ , and  $\varepsilon > 0$ , there exists a strategy  $\sigma'$  using at most  $N =$*   
 1292  *$(2|Q|)^{n(|E|+1)} |A| \lceil 8 \frac{\log(8/\varepsilon)}{\eta^2} \rceil$  memory where  $n = \lceil 2p^{-2|Q|} \max(\frac{8(|Q||A|)^3}{\varepsilon\eta^2}, \log(16/\varepsilon)) \rceil$ , with  $p$*   
 1293 *the smallest nonzero probability in  $M$ , and that satisfies  $\mathbb{P}_q^{\sigma'}(M, \varphi) \geq \mathbb{P}_q^\sigma(M, \varphi) - \varepsilon$ .*

1294 **Proof.** Given  $\varepsilon > 0$ , let

$$1295 \quad n_0 = \left\lceil \frac{8(|Q||A|)^3}{\varepsilon\eta^2} \right\rceil,$$

$$1296 \quad n = \left\lceil 2p^{-2|Q|} \max(n_0, \log(16/\varepsilon)) \right\rceil,$$

1297 Notice that the bounds on  $n$  and  $n_0$  come from Lemma 23 applied for  $\varepsilon/4$ . Define the  
 1298 sequence  $(g_i)_{i \geq 1}$  by  $g_1 = 1$ , and  $g_i = \alpha(g_{i-1} + \beta) + \gamma$  where  $\alpha = 2|Q|^n$ ,  $\beta = \lceil 8 \frac{\log(8/\varepsilon)}{\eta^2} \rceil$ ,  
 1299 and  $\gamma = |Q||A|$ . Note that we have, for  $i > 1$ ,  $g_i = \alpha^{i-1} + (\gamma + \alpha\beta)(\frac{\alpha^{i-1}-1}{\alpha-1})$ . Observe that  
 1300  $g_i \leq \alpha^{i-1}(1 + \gamma + \alpha\beta)$ .

1301 We prove, by induction on  $|E|$ , that for all states  $q$ , strategies  $\sigma$ , there exists a  $g_{|E|}$ -memory  
 1302 strategy  $\sigma'$  such that  $\mathbb{P}_q^{\sigma'}(M, \varphi) \geq \mathbb{P}_q^\sigma(M, \varphi) - \varepsilon$ .

1303 We have  $g_{|E|} \leq \alpha^{|E|-1}(1 + \gamma + \alpha\beta) \leq \alpha^{|E|}(1 + |Q||A| + 2|Q|^n\beta) \leq \alpha^{|E|}(3|Q|^n|A|\beta) \leq$   
 1304  $\alpha^{|E|}(2\alpha|A|\beta)$ , which is at most  $(2|Q|)^{n(|E|+1)} |A| \lceil 8 \frac{\log(8/\varepsilon)}{\eta^2} \rceil$ , and proves the lemma.

1305 The base case  $|E| = 1$  is obvious since there exists optimal memoryless strategies for  
 1306 parity objectives in MDPs. Assume  $|E| \geq 2$ .

1307 Let  $M' = \text{purge}(M)$  and  $\sigma'$  be given by Lemma 19 such that  $\mathbb{P}_q^\sigma(M, \varphi) \leq \mathbb{P}_{q'}^{\sigma'}(M', \varphi)$   
 1308 where  $q' = f(q)$ . We prove the property for  $M'$  and transfer the result back to  $M$  using  
 1309 Lemma 20. More precisely, we show below that there exists a  $(g_{|E|} - \gamma)$ -memory strategy  $\sigma''$   
 1310 with  $\mathbb{P}_{q'}^{\sigma''}(M', \varphi) \geq \mathbb{P}_{q'}^{\sigma'}(M', \varphi) - \varepsilon$ . It follows, by Lemma 19, that there exists a  $g_{|E|}$ -memory  
 1311 strategy  $\sigma'''$  for  $M$  such that

$$1312 \quad \mathbb{P}_q^{\sigma'''}(M, \varphi) = \mathbb{P}_{q'}^{\sigma''}(M', \varphi) \geq \mathbb{P}_{q'}^{\sigma'}(M', \varphi) - \varepsilon \geq \mathbb{P}_q^\sigma(M, \varphi) - \varepsilon$$

1313 which proves the result.

1314 We construct  $\sigma''$  by imitating  $\sigma'$  for  $n$  steps, and stopping if a MCEC is reached (thus,  
 1315 either trivial or distinguishing, by Lemma 17). More precisely, consider history  $h$  in  $M'$  that  
 1316 starts at  $q'$ . We define  $\sigma''(h) = \sigma'(h)$ , except in the following cases where  $\sigma''$  switches to a  
 1317 strategy as described below:

**XX:36 The Value Problem for Multiple-Environment MDPs with Parity Objective**

- 1318 1. If  $\text{last}(h)$  belongs to a trivial MCEC, then  $\sigma''$  is memoryless from that history (as there  
 1319 is only one possible action to choose). Notice that  $\mathbb{P}_{q'}^{\sigma''}(M', \varphi \mid h) = \mathbb{P}_{q'}^{\sigma'}(M', \varphi \mid h)$  since  
 1320 this MCEC is either winning or losing with probability 1, in each  $e \in E$ .  
 1321 2. Assume  $\text{last}(h)$  belongs to a distinguishing MCEC  $D$  with partition  $(K_1, K_2)$ . Let  
 1322  $\vec{\beta} = \mathbb{P}^{\sigma'}(M', \varphi \mid h)$ , the probability values achieved from history  $h$  under strategy  $\sigma'$   
 1323 starting with history  $h$ . One can define a strategy  $\sigma'_h$  such that  $\vec{\beta} = \mathbb{P}_{\text{last}(h)}^{\sigma'_h}(M', \varphi)$ , by  
 1324  $\sigma'_h : h' \mapsto \sigma'(h \cdot h')$ . By induction applied to  $M'$ ,  $\text{last}(h)$ ,  $\sigma'_h$ , environment set  $K_i$ , and  
 1325  $\varepsilon/8$ , there exist  $g_{|K_i|}$ -memory strategies  $\sigma_i$ , with  $\mathbb{P}_{\text{last}(h)}^{\sigma_i}(M'[K_i], \varphi) \geq \vec{\beta}|_{K_i} - \varepsilon/8$ . We  
 1326 apply Lemma 9 to build strategy  $\sigma''_h$  satisfying the following:

1327 
$$\mathbb{P}_{\text{last}(h)}^{\sigma''_h}(M'[e], \varphi) \geq \mathbb{P}_{\text{last}(h)}^{\sigma_i}(M'[e], \varphi) - \varepsilon/4 \text{ for all environments } e \in K_i. \quad (8)$$

- 1328 At  $h$ , we let  $\sigma''$  switch to  $\sigma''_h$ . It follows that  $\mathbb{P}_{q'}^{\sigma''}(M', \varphi \mid h) \geq \mathbb{P}_{q'}^{\sigma'}(M', \varphi \mid h) - \varepsilon/4$ .  
 1329 3. Assume that  $h$  contains  $n_0$  occurrences of distinguishing state-action pairs, that is,  
 1330  $|h|_{T_{M'}} = n_0$ . Let  $(q, a, q') \in T_M$  be a distinguishing transition with the largest number of  
 1331 occurrences in  $h$ ; and let  $(K_1, K_2)$  be the partition of  $E$  induced by this transition. For  
 1332 each  $i = 1, 2$ , let  $\sigma_i$  be the  $g_{|K_i|}$ -memory strategy given by induction hypothesis applied  
 1333 to  $M'$ , state  $\text{last}(h)$ , environment set  $K_i$ , bound  $\varepsilon/4$ , and strategy  $\sigma'_h : h' \mapsto \sigma'(h \cdot h')$   
 1334 that achieves  $\mathbb{P}_{\text{last}(h)}^{\sigma_i}(M'[K_i], \varphi) \geq \mathbb{P}_{\text{last}(h)}^{\sigma'_h}(M'[K_i], \varphi) - \varepsilon/4$  for each  $i = 1, 2$ . We let  $\sigma''$   
 1335 switch to:

- 1336 =  $\sigma_1$  if  $\left| \frac{|h|_{q,a,q'}}{|h|_{q,a}} - \delta_e(q, a)(q') \right| < \eta/2$  for some  $e \in K_1$ ,  
 1337 =  $\sigma_2$  otherwise.

1338 The above shows that if  $h$  is a good classification in  $e$ , then  $\mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi \mid h) \geq$   
 1339  $\mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi \mid h) - \varepsilon/4$ .

- 1340 4. If  $|h| = n$  and none of the above applies, then  $\sigma''$  switches to an arbitrary memoryless  
 1341 strategy. These histories that satisfy case 4 has probability at most  $\varepsilon/4$  by Lemma 23.

1342 Let us show that  $\mathbb{P}_{q'}^{\sigma''}(M', \varphi) \geq \mathbb{P}_{q'}^{\sigma'}(M', \varphi) - \varepsilon$ . To prove this, we distinguish histories  $h$   
 1343 according to the cases above, and relate  $\mathbb{P}_{q'}^{\sigma''}(M', \varphi \mid h)$  and  $\mathbb{P}_{q'}^{\sigma'}(M', \varphi \mid h)$ , and bound the  
 1344 probability of some histories  $h$ .

1345 Let us write

1346 
$$\mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi) = \sum_{h: \text{case 1}} \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi, h) + \sum_{h: \text{case 2}} \mathbb{P}_{q'}^{\sigma''}(M'[e], h) \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi \mid h)$$
  
 1347 
$$+ \sum_{\substack{h: \text{case 3} \\ \text{bad classification}}} \mathbb{P}_{q'}^{\sigma''}(M'[e], h) \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi \mid h)$$
  
 1348 
$$+ \sum_{\substack{h: \text{case 3} \\ \text{good classification}}} \mathbb{P}_{q'}^{\sigma''}(M'[e], h) \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi \mid h)$$
  
 1349 
$$+ \sum_{h: \text{case 4}} \mathbb{P}_{q'}^{\sigma''}(M'[e], h) \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi \mid h).$$

1350

1351 Since  $\mathbb{P}_{q'}^{\sigma''}(M', h) = \mathbb{P}_{q'}^{\sigma'}(M', h)$  for histories satisfying any of the cases (because  $\sigma''$  imitates  
 1352  $\sigma'$  until such a case occurs), and because the terms  $\mathbb{P}_{q'}^{\sigma''}(M', h)$  at the second and forth lines

1353 are each at most  $\varepsilon/4$ , using the cases above, we get

$$\begin{aligned}
1354 \quad \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi) &\geq \sum_{h: \text{case 1}} \mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi, h) + \sum_{h: \text{case 2}} \mathbb{P}_{q'}^{\sigma'}(M'[e], h) (\mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi | h) - \varepsilon/4) \\
1355 \quad &+ \left( \sum_{\substack{h: \text{case 3} \\ \text{bad classification}}} \mathbb{P}_{q'}^{\sigma'}(M'[e], h) \mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi | h) \right) - \varepsilon/4 \\
1356 \quad &+ \sum_{\substack{h: \text{case 3} \\ \text{good classification}}} \mathbb{P}_{q'}^{\sigma'}(M'[e], h) (\mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi | h) - \varepsilon/4) \\
1357 \quad &+ \left( \sum_{h: \text{case 4}} \mathbb{P}_{q'}^{\sigma''}(M'[e], h) \mathbb{P}_{q'}^{\sigma''}(M'[e], \varphi | h) \right) - \varepsilon/4. \\
1358 \quad &\geq \mathbb{P}_{q'}^{\sigma'}(M'[e], \varphi) - \varepsilon.
\end{aligned}$$

1359 Last, we argue that  $\sigma''$  uses memory of size  $g_{|E|}$ . Strategy  $\sigma''$  must store the histories until  
1360 one of the four cases occur: this happens in at most  $n$  steps, which means  $|Q|^n$  memory is  
1361 required for this phase. In addition, for each history of case 2,  $g_{|E|-1} + g_{|E|-1} + 8 \frac{\log(8/\varepsilon)}{\eta^2}$   
1362 memory states are needed by Lemma 9; where the terms  $g_{|E|-1}$  are upper bounds on the  
1363 memory requirement of the strategies to which we switch, given by induction. Case 3 does  
1364 not require additional memory since the decision is made depending on the current history,  
1365 which is already in the memory. In total, we thus need  $|Q|^n(2g_{|E|-1} + \beta)$  memory states,  
1366 which is at most  $\alpha(g_{|E|-1} + \beta) = g_{|E|} - \gamma$ .  $\blacktriangleleft$

## 1367 5.4 Approximation Algorithm

1368 We now provide a procedure solving the gap problem with threshold  $\alpha$  for parity objectives  
1369 in MEMDPs. Informally, given bound  $N$ , the procedure *guesses* an  $N$ -memory strategy by  
1370 solving a set of polynomial constraints over the reals, and checks that the strategy ensures  
1371 winning with probability at least  $\alpha - \varepsilon$  in all environments. We first give the construction  
1372 for reachability, then explain the extension to parity conditions.

1373 **Reachability in MDPs** Let us start by recalling the linear constraints that characterize  
1374 reachability probabilities in single-environment MDPs under memoryless strategies. Consider  
1375 an MDP  $M = \langle Q, A, \delta \rangle$  and objective  $\text{Reach}(T)$ . Let  $Q^{\text{no}} \subseteq Q$ , and  $Q^? = Q \setminus (Q^{\text{no}} \cup T)$ .  $Q^{\text{no}}$   
1376 will be the set of states from which the reachability probability is 0; it is necessary to make  
1377 sure all such states are in  $Q^{\text{no}}$  so that the equation given below has a unique solution. Define  
1378 the unknown  $x_q$  representing the probability of reaching  $T$  from  $q$  under the strategy that is  
1379 being guessed, and  $p_q(a)$  the probability of the strategy to pick action  $a$  from  $q$ , for  $a \in A_q$ .  
1380 Consider the following constraints:

$$\begin{aligned}
&x_q = 0 && \text{for all } q \in Q^{\text{no}}, \\
&x_q = 1 && \text{for all } q \in T, \\
1381 \quad &x_q = \sum_{a \in A_q} p_q(a) \cdot \sum_{q' \in Q} \delta(q, a, q') \cdot x_{q'} && \text{for all } q \in Q^?, \\
&0 \leq x_q \leq 1 \text{ and } 0 \leq p_q(a) \leq 1 && \text{for all } q \in Q, a \in A_q, \\
&\sum_{a \in A_q} p_q(a) = 1 && \text{for all } q \in Q.
\end{aligned} \tag{9}$$

1382 Any solution  $(\vec{x}, \vec{p})$  of (9) yields a strategy  $\sigma^{\vec{p}}$ , which is defined as picking action  $a$   
1383 from state  $q$  with probability  $p_q(a)$ . The following theorem shows that  $\vec{x}$  does capture  
1384 the reachability probabilities of  $\sigma^{\vec{p}}$ , provided that  $Q^{\text{no}}$  is the set of states from which the  
1385 reachability probability is 0.

1386 ► **Theorem 25** (Theorem 10.19, [3]). Consider any subset  $Q^{\text{no}} \subseteq Q$ , and a solution  
 1387  $(\vec{x}, \vec{p})$  of (9). If for all states  $q \in Q^{\text{no}}$ ,  $\mathbb{P}_q^{\sigma^{\vec{p}}}[M, \text{Reach}(T)] = 0$ , then for all  $q \in Q$ ,  
 1388  $x_q = \mathbb{P}_q^{\sigma^{\vec{p}}}[M, \text{Reach}(T)]$ . Conversely, for any memoryless strategy  $\tau$ , if  $Q^{\text{no}}$  denotes the  
 1389 set of states  $q$  with  $\mathbb{P}_q^\tau[M, \text{Reach}(T)] = 0$ , then (9) has a unique solution  $(\vec{x}, \vec{p})$  where  $\tau = \sigma^{\vec{p}}$ ,  
 1390 and  $x_q = \mathbb{P}_q^\tau[M, \text{Reach}(T)]$  for all  $q \in Q$ .

1391 **Finite-Memory Reachability in MEMDPs** We now show how to solve the gap problem  
 1392 for an instance of the quantitative reachability problem for MEMDPs. Consider MEMDP  
 1393  $M = \langle Q, A, (\delta_e)_{e \in E} \rangle$ , objective  $\text{Reach}(T)$ , a memory bound  $N$ , an initial state  $q_0$ , and bounds  
 1394  $\varepsilon, \alpha > 0$ . We want to check whether there exists a strategy  $\sigma$  such that, for all environments  
 1395  $e \in E$ , we have  $\mathbb{P}_{q_0}^\sigma[M[e], \text{Reach}(T)] \geq \alpha$ , or whether for all  $\sigma$ , there exists an environment  
 1396  $e \in E$  with  $\mathbb{P}_{q_0}^\sigma[M[e], \text{Reach}(T)] < \alpha - \varepsilon$ .

1397 We guess a memoryless randomized strategy on combined states  $(q, i)$  for  $q \in Q$  and  
 1398  $0 \leq i < N$ , which correspond to  $N$ -memory strategies on  $M$ . In the sequel, we write  
 1399  $[N] = \{0, 1, \dots, N - 1\}$ . We define the unknown variable  $x_{q,i}^e$  for each  $e \in E$ , and combined  
 1400 state  $(q, i)$  representing the probability of reaching  $T$  from state  $q$  and memory value  $i$  in  
 1401  $M[e]$ , under the strategy that is being guessed. Furthermore, define  $p_{q,i}(a, i')$  for each action  
 1402  $a \in A_q$  and  $i' \in [N]$ , as the unknown representing the probability of the strategy picking  
 1403 action  $a$  from  $(q, i)$  and updating the memory value to  $i'$ .

1404 Consider subsets  $Q_e^{\text{no}} \subseteq Q$  for each  $e \in E$ , and let  $Q_e^? = Q \setminus (Q_e^{\text{no}} \cup T_e)$ . We write the  
 1405 following constraints in a slightly more general setting, where a possibly different target  
 1406 set  $T_e$  is considered for each environment  $e$  (this will be useful when generalizing to parity  
 1407 conditions below):

$$\begin{aligned}
 & x_{q,i}^e = 0 && \text{for all } e \in E, q \in Q_e^{\text{no}}, i \in [N], \\
 & x_{q,i}^e = 1 && \text{for all } e \in E, q \in T_e, i \in [N], \\
 & x_{q,i}^e = \sum_{a \in A_q} p_{q,i}(a, i') \cdot \sum_{q' \in Q} \delta_e(q, a, q') \cdot x_{q',i'}^e && \text{for all } e \in E, q \in Q_e^?, i \in [N], \\
 1408 & 0 \leq x_{q,i}^e \leq 1 && \text{for all } e \in E, q \in Q, i \in [N], \quad (10) \\
 & 0 \leq p_{q,i}(a, i') \leq 1 && \text{for all } q \in Q, a \in A_q, i, i' \in [N], \\
 & \sum_{a \in A_q} \sum_{i' \in [N]} p_{q,i}(a, i') = 1 && \text{for all } q \in Q, i \in [N], \\
 & x_{q_0}^e \geq \alpha - \varepsilon && \text{for all } e \in E.
 \end{aligned}$$

1409 Notice how the choice of the action and memory updates  $p_{q,i}$  does not depend on the  
 1410 environment. The constraints (10) simply combine  $|E|$  copies of (9) over a state space  
 1411 augmented with  $N$  memory values. In addition we added the constraints  $x_{q_0}^e \geq \alpha - \varepsilon$  for all  
 1412  $e \in E$ , which restrict the solution sets to those strategies that ensure the threshold  $\alpha - \varepsilon$ .

1413 **The Gap Problem for Reachability** The full procedure is as follows. We let  $T_e = T$  for  
 1414 all  $e \in E$ . Let  $N$  be the memory bound given in Lemma 24.

1415 We enumerate all possible choices for the sets  $Q_e^{\text{no}}$ . For each choice  $(Q_e^{\text{no}})_{e \in E}$ , we solve  
 1416 the corresponding constraints (10). If there is no solution, we continue with the next choice.  
 1417 Otherwise let  $\sigma^{\vec{p}}$  be the  $N$ -memory strategy given by the solution to this equation. If  
 1418  $\mathbb{P}_q^{\sigma^{\vec{p}}}[M[e], \text{Reach}(T_e)] = 0$  for each  $e \in E$  and  $q \in Q_e^{\text{no}}$ , then we return **Yes**; otherwise we  
 1419 continue with the next choice  $(Q_e^{\text{no}})_{e \in E}$ . We return **No** at the end of if no solution was found.

1420 Let us show that this procedure solves the gap problem. Assume that there exists a  
 1421 strategy  $\tau$  such that for all environments  $e \in E$ , we have  $\mathbb{P}_{q_0}^\tau[M[e], \text{Reach}(T)] \geq \alpha$ . By  
 1422 Lemma 24, there exists a  $N$ -memory strategy  $\tau'$  such that in all environments  $e \in E$ , we  
 1423 have  $\mathbb{P}_{q_0}^{\tau'}[M[e], \text{Reach}(T)] \geq \alpha - \varepsilon$ . Hence (10) must have a solution corresponding to this  
 1424 strategy for some choice of the sets  $(Q_e^{\text{no}})_{e \in E}$ , and the procedure returns **Yes**. Assume now  
 1425 that no strategy achieves the threshold  $\alpha - \varepsilon$ . In particular, no  $N$ -memory strategy achieves  
 1426 this threshold, and the procedure returns **No**.

1427 We now analyze the complexity of the procedure. The value  $N$  is double exponential in  
 1428 the size of the input, which means that the size of (10) is also double exponential. Polynomial  
 1429 equations can be solved in polynomial space in the size of the equations [5], so here we can  
 1430 solve (10) in double exponential space.

1431 **The Gap Problem for Parity** We now extend the previous procedure to solve the  
 1432 quantitative parity gap problem based on the following observations. In  $M[e]$ , any finite-  
 1433 memory strategy  $\sigma$  induces a Markov chain. Then the probability  $\mathbb{P}_{q_0}^\sigma[M[e], \varphi]$  of satisfying  
 1434 a parity condition  $\varphi$  is equal to the probability of reaching bottom strongly connected  
 1435 components (BSCC) that are winning<sup>3</sup> for  $\varphi$  in the induced Markov chain [3]. But the set  
 1436 of BSCCs only depends on the support of  $\sigma$ , that is, the set of state-action pairs that have  
 1437 positive probability. When considering an MDP under  $N$ -memory strategies, the support is  
 1438 the set of tuples  $(q, i, a, i')$  such that from combined state  $(q, i)$  the strategy has a nonzero  
 1439 probability of picking action  $a$  and updating memory to  $i'$ .

1440 We proceed as follows. Given MEMDP  $M$ , initial state  $q_0$ , parity condition  $\varphi$ , and  
 1441 bound  $N$ , we enumerate all supports  $S \subseteq Q \times [N] \times A \times [N]$ . For each support  $S$ , let  $T_e^S$   
 1442 be the set of  $\varphi$ -winning BSCCs in  $M[e]$  under a strategy with support  $S$ . We apply the  
 1443 reachability procedure described above based on (10) for the target sets  $(T_e^S)_e$  augmented  
 1444 with the following constraints: for all  $(q, i, a, i') \in S$ , we add the constraint  $p_{q,i}(a, i') > 0$ ,  
 1445 and for all others  $p_{q,i}(a, i') = 0$ . If the answer is Yes for some support  $S$ , then we return Yes;  
 1446 otherwise we return No.

1447 This solves the gap problem: if there is  $\tau$  such that for all environments  $e \in E$  we have  
 1448  $\mathbb{P}_{q_0}^\tau[M[e], \varphi] \geq \alpha$ , then by Lemma 24 there exists a  $N$ -memory strategy  $\tau'$  such that for all  
 1449 environments  $e \in E$  we have  $\mathbb{P}_{q_0}^{\tau'}[M[e], \varphi] \geq \alpha - \varepsilon$ . Let  $S$  denote the support of the strategy  
 1450  $\tau'$ , and let  $T_e^S$  be the set of winning BSCCs in  $M[e]$  under  $\tau$ . So (10), instantiated for  $S$  has  
 1451 a solution, for some choice of the sets  $(Q_e^{\text{no}})_{e \in E}$ , and the procedure returns Yes. Assume now  
 1452 that no strategy achieves the threshold  $\alpha - \varepsilon$ . In particular, there is no  $N$ -memory strategy  
 1453 with any support  $S$  that achieves this threshold, and the procedure returns No.

1454 There are an exponential number of possibilities for the choice of support. Moreover,  
 1455 given a support  $S$ , each set  $T_e^S$  can be determined in polynomial time. Overall, the procedure  
 1456 remains in double exponential space.

1457 **► Theorem 26.** *The gap problem can be solved in double exponential space for MEMDPs*  
 1458 *with parity objectives.*

## 1459 ——— References ———

- 1460 1 T. S. Badings, T. D. Simão, M. Suilen, and N. Jansen. Decision-making under uncertainty:  
 1461 beyond probabilities. *Int. J. Softw. Tools Technol. Transf.*, 25(3):375–391, 2023.
- 1462 2 C. Baier, M. Größer, and N. Bertrand. Probabilistic  $\omega$ -automata. *J. ACM*, 59(1):1, 2012.
- 1463 3 C. Baier and J.-P. Katoen. *Principles of Model Checking*. MIT Press, 2008.
- 1464 4 P. Buchholz and D. Scheftelowitsch. Computation of weighted sums of rewards for concurrent  
 1465 MDPs. *Math. Methods Oper. Res.*, 89(1):1–42, 2019.
- 1466 5 J. Canny. Some algebraic and geometric computations in PSPACE. In *Proc. of STOC:*  
 1467 *Symposium on Theory of Computing*, page 460–467. ACM, 1988.
- 1468 6 K. Chatterjee, M. Chmelík, D. Karkhanis, P. Novotný, and A. Royer. Multiple-environment  
 1469 Markov decision processes: Efficient analysis and applications. In *Proc. of ICAPS: Automated*  
 1470 *Planning and Scheduling*, pages 48–56. AAAI Press, 2020.

<sup>3</sup> Recall that a BSCC is winning for a parity condition if the smallest priority of its states is even.

- 1471 7 K. Chatterjee, L. Doyen, H. Gimbert, and T. A. Henzinger. Randomness for free. *Information and Computation*, 245:3–16, 2017.
- 1472 8 C. Courcoubetis and M. Yannakakis. The complexity of probabilistic verification. *J. ACM*, 42(4):857–907, July 1995.
- 1473 9 L. de Alfaro. *Formal verification of probabilistic systems*. Ph.d. thesis, Stanford University, 1997.
- 1474 10 S. Even, A. L. Selman, and Y. Yacobi. The complexity of promise problems with applications to public-key cryptography. *Information and Control*, 61(2):159 – 173, 1984.
- 1475 11 E. A. Feinberg and A. Schwartz, editors. *Handbook of Markov Decision Processes - Methods and Applications*. Kluwer, 2002.
- 1476 12 H. Gimbert and Y. Oualhadj. Probabilistic automata on finite words: Decidable and undecidable problems. In *Proc. of ICALP (2)*, LNCS 6199, pages 527–538. Springer, 2010.
- 1477 13 O. Goldreich. On promise problems (a survey in memory of Shimon Even [1935-2004]). *Manuscript*, 2005.
- 1478 14 W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- 1479 15 K. J. Åström. Optimal control of Markov processes with incomplete state information I. *Journal of Mathematical Analysis and Applications*, 10:174–205, 1965.
- 1480 16 O. Madani, S. Hanks, and A. Condon. On the undecidability of probabilistic planning and related stochastic optimization problems. *Artif. Intell.*, 147(1-2):5–34, 2003.
- 1481 17 D. A. Martin. The determinacy of Blackwell games. *The Journal of Symbolic Logic*, 63(4):1565–1581, 1998.
- 1482 18 C. H. Papadimitriou and J. N. Tsitsiklis. The complexity of Markov decision processes. *Math. Oper. Res.*, 12(3):441–450, 1987.
- 1483 19 M. L. Puterman. *Markov decision processes*. John Wiley and Sons, 1994.
- 1484 20 N. M. van Dijk R. J. Boucherie. *Markov decision processes in practice*. Springer, 2017.
- 1485 21 J.-F. Raskin and O. Sankur. Multiple-environment Markov decision processes. In *Proc. of FSTTCS: Foundation of Software Technology and Theoretical Computer Science*, volume 29 of *LIPICs*, pages 531–543. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2014.
- 1486 22 M. Suilen, M. van der Vegt, and S. Junges. A PSPACE algorithm for almost-sure Rabin objectives in multi-environment MDPs. In *Proc. of CONCUR: Concurrency Theory*, volume 311 of *LIPICs*, pages 40:1–40:17. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024.
- 1487 23 M. van der Vegt, N. Jansen, and S. Junges. Robust almost-sure reachability in multi-environment MDPs. In *Proc. of TACAS: Tools and Algorithms for the Construction and Analysis of Systems*, LNCS 13993, pages 508–526. Springer, 2023.
- 1488
- 1489
- 1490
- 1491
- 1492
- 1493
- 1494
- 1495
- 1496
- 1497
- 1498
- 1499
- 1500
- 1501
- 1502
- 1503
- 1504
- 1505