

# Minimal Disclosure in Partially Observable Markov Decision Processes

**Nathalie Bertrand**<sup>1</sup> and Blaise Genest<sup>2</sup>

<sup>1</sup>INRIA Rennes, France

<sup>2</sup>CNRS UMI IPAL, Singapore

FSTTCS

December 14th 2011

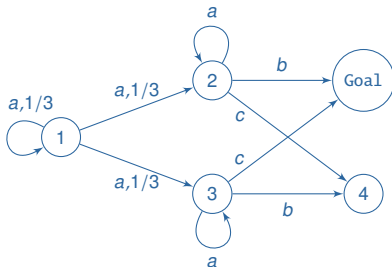
# Outline

---

- 1 Introduction
- 2 Worst-case cost
- 3 Average cost
- 4 Conclusion

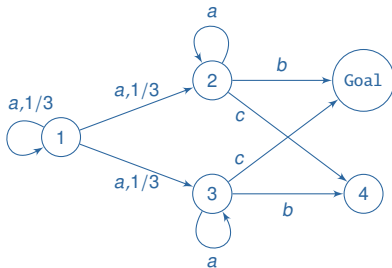
# Markov Decision Processes (MDP)

States:  $Q$ ; Actions:  $Act$ ; Probabilistic transition function:  $\Delta$



# Markov Decision Processes (MDP)

States:  $Q$ ; Actions:  $Act$ ; Probabilistic transition function:  $\Delta$

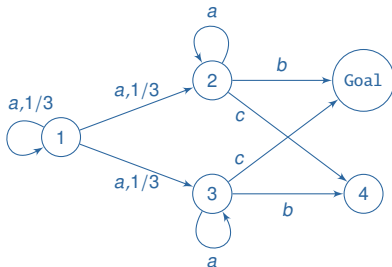


Strategy for the controller: based on actions and states

$$\sigma : Q \cdot (Act \cdot Q)^* \rightarrow Dist(Act)$$

# Markov Decision Processes (MDP)

States:  $Q$ ; Actions:  $Act$ ; Probabilistic transition function:  $\Delta$



Strategy for the controller: based on actions and states

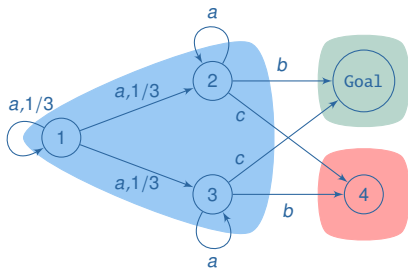
$$\sigma : Q \cdot (Act \cdot Q)^* \rightarrow Dist(Act)$$

Memoryless pure strategy to reach Goal almost-surely:

$$\sigma(1) = a, \sigma(2) = b, \sigma(3) = c$$

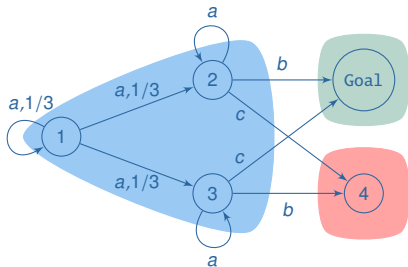
# Partially Observable MDP (POMDP)

Partial observation: induced by partition  $\mathcal{O}$



# Partially Observable MDP (POMDP)

Partial observation: induced by partition  $\mathcal{O}$

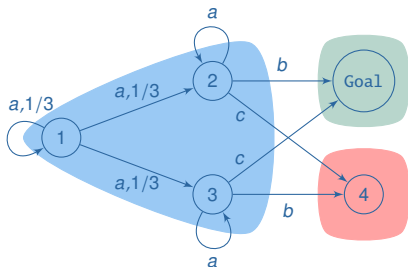


Strategy for the controller: based on actions and observations

$$\sigma : \mathcal{O} \cdot (\text{Act} \cdot \mathcal{O})^* \rightarrow \text{Dist}(\text{Act})$$

# Partially Observable MDP (POMDP)

Partial observation: induced by partition  $\mathcal{O}$



Strategy for the controller: based on actions and observations

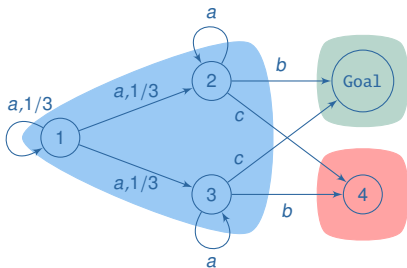
$$\sigma : \mathcal{O} \cdot (\text{Act} \cdot \mathcal{O})^* \rightarrow \text{Dist}(\text{Act})$$

No strategy to reach Goal almost-surely.



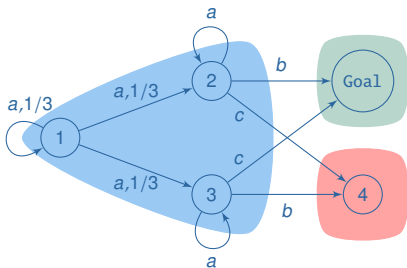
# POMDP with disclosure

Additional **request** action to reveal the precise state of system.  
Observations: partition + individual states



# POMDP with disclosure

Additional **request** action to reveal the precise state of system.  
 Observations: partition + individual states

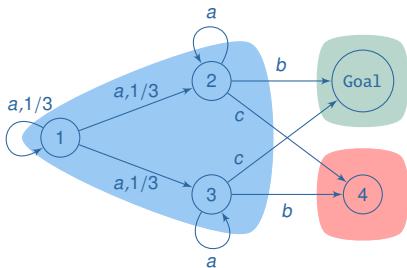


Strat. for the controller: based on **extended** actions and observations

$$\sigma : \mathcal{O}' \cdot (\text{Act}' \cdot \mathcal{O}')^* \rightarrow \text{Dist}(\text{Act}')$$

# POMDP with disclosure

Additional **request** action to reveal the precise state of system.  
 Observations: partition + individual states



Strat. for the controller: based on **extended** actions and observations

$$\sigma : \mathcal{O}' \cdot (\text{Act}' \cdot \mathcal{O}')^* \rightarrow \text{Dist}(\text{Act}')$$

Cheap strategy to reach Goal almost-surely?

# Problem statement

---

cost of a path = number of requests for disclosure

cost of a strategy  $\sigma =$

- ▶ worst-case cost along  $\sigma$ -paths (max number of requests)
- ▶ average cost along  $\sigma$ -paths (expected number of requests)

# Problem statement

---

cost of a path = number of requests for disclosure

cost of a strategy  $\sigma =$

- ▶ worst-case cost along  $\sigma$ -paths (max number of requests)
- ▶ average cost along  $\sigma$ -paths (expected number of requests)

## Problem statement

Finding almost-surely winning strategies that minimize:

- ▶ the worst-case cost, or
- ▶ the average cost

# Outline

---

- 1 Introduction
- 2 Worst-case cost**
- 3 Average cost
- 4 Conclusion

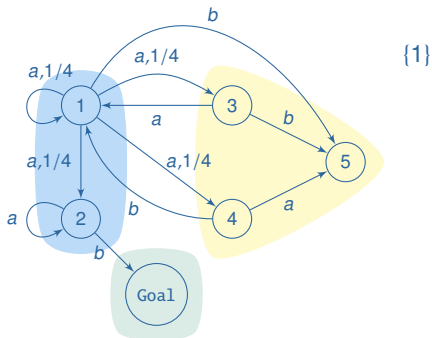
# Belief

---

Belief: (distribution over) states the system can be in

# Belief

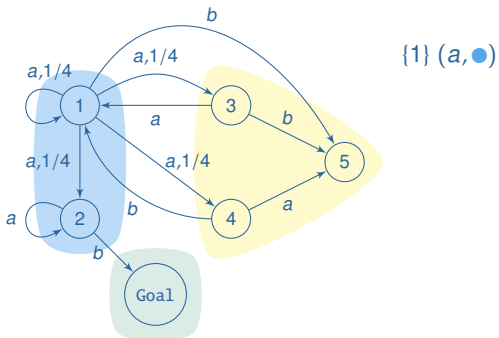
Belief: (distribution over) states the system can be in





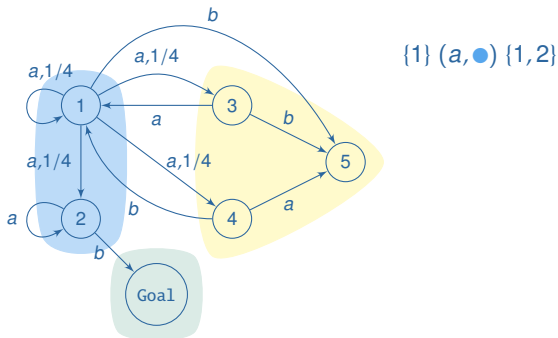
# Belief

Belief: (distribution over) states the system can be in



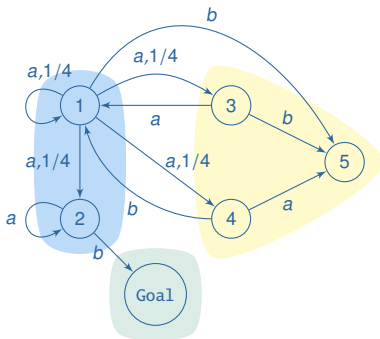
# Belief

Belief: (distribution over) states the system can be in



# Belief

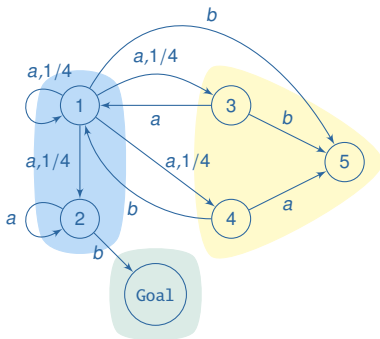
Belief: (distribution over) states the system can be in



$\{1\} (a, \bullet)$   $\{1, 2\} (a, \bullet)$

# Belief

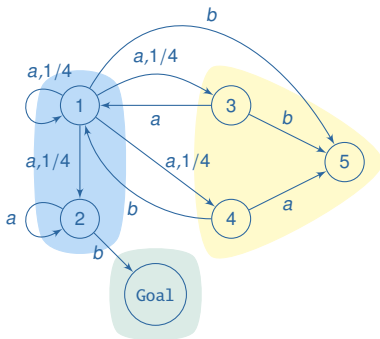
Belief: (distribution over) states the system can be in



$\{1\}$  (a, ●)  $\{1,2\}$  (a, ●)  $\{3,4\}$

# Belief

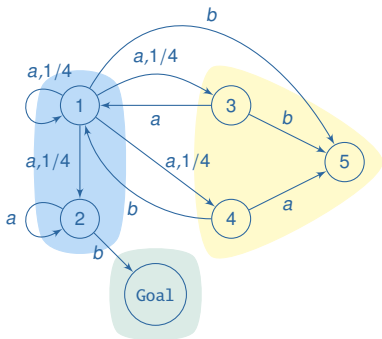
Belief: (distribution over) states the system can be in



$\{1\} (a, \bullet)$   $\{1,2\} (a, \bullet)$   $\{3,4\} (b, \bullet)$

# Belief

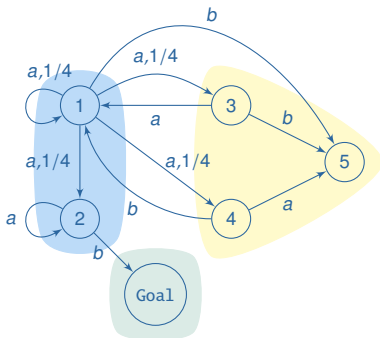
Belief: (distribution over) states the system can be in



$\{1\} (a, \bullet)$   $\{1,2\} (a, \bullet)$   $\{3,4\} (b, \bullet)$   $\{5\}$

# Belief

Belief: (distribution over) states the system can be in

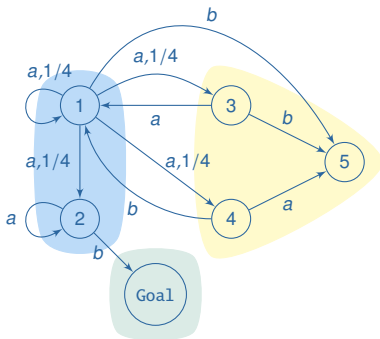


$\{1\} (a, \bullet) \{1,2\} (a, \bullet) \{3,4\} (b, \bullet) \{5\}$

$\{1\} (a, \bullet) \{3,4\} (req, \{4\}) \{4\} (b, \bullet) \{1\}$

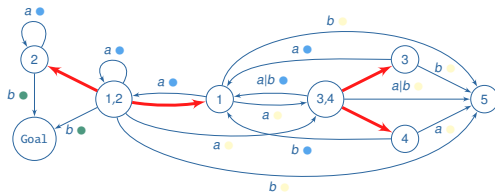
# Belief

Belief: (distribution over) states the system can be in



$\{1\} (a, \bullet) \{1,2\} (a, \bullet) \{3,4\} (b, \bullet) \{5\}$

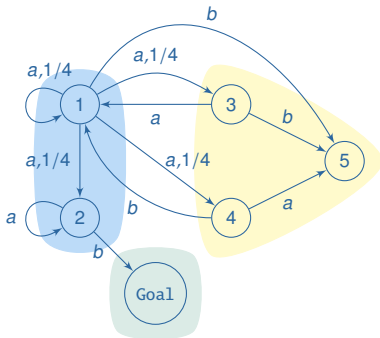
$\{1\} (a, \bullet) \{3,4\} (req, \{4\}) \{4\} (b, \bullet) \{1\}$





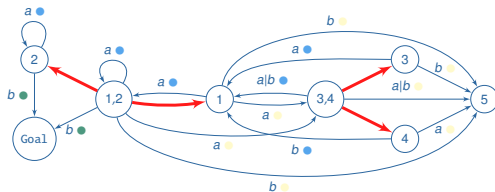
# Belief

Belief: (distribution over) states the system can be in



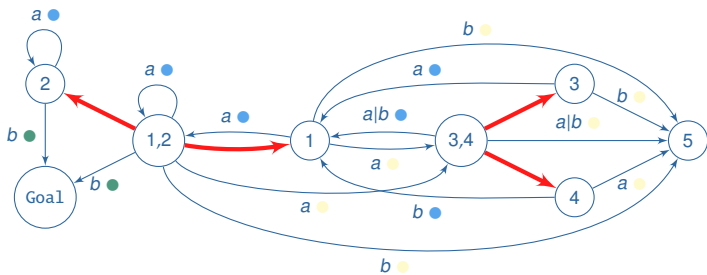
$\{1\} (a, \bullet) \{1,2\} (a, \bullet) \{3,4\} (b, \bullet) \{5\}$

$\{1\} (a, \bullet) \{3,4\} (req, \{4\}) \{4\} (b, \bullet) \{1\}$



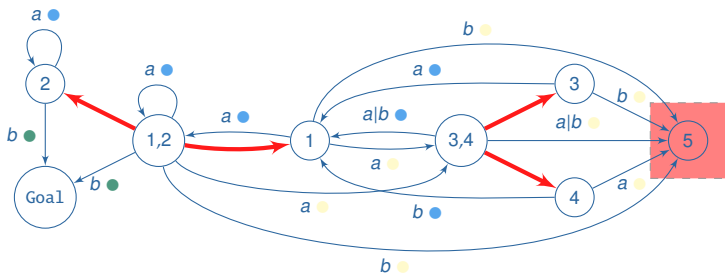
$up(S, a, O)$ : belief update from  $S$ , after action  $a$  and observation  $O$

# Reaching the goal almost-surely



# Reaching the goal almost-surely

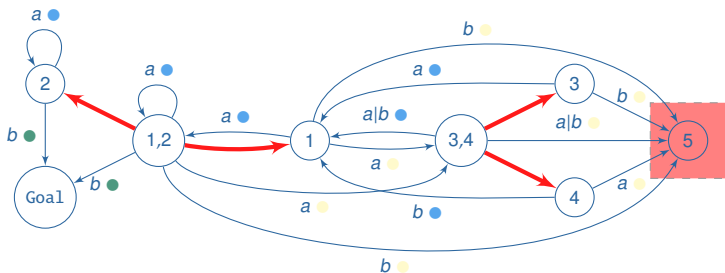
*Lose*: beliefs that contain a state losing in the (fully-observable) MDP



# Reaching the goal almost-surely

**Lose**: beliefs that contain a state losing in the (fully-observable) MDP

$$Win = \mathcal{B} \setminus Lose = W_{ok} \sqcup W_{req} \sqcup W_{safe}$$

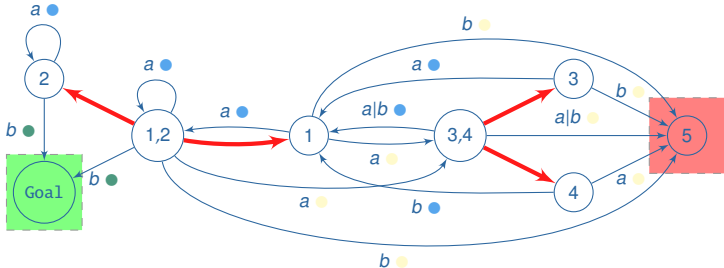


# Reaching the goal almost-surely

**Lose**: beliefs that contain a state losing in the (fully-observable) MDP

$$\text{Win} = \mathcal{B} \setminus \text{Lose} = W_{ok} \sqcup W_{req} \sqcup W_{safe}$$

- ▶  $W_{ok} = \{S \mid S \subseteq \text{Goal}\}$

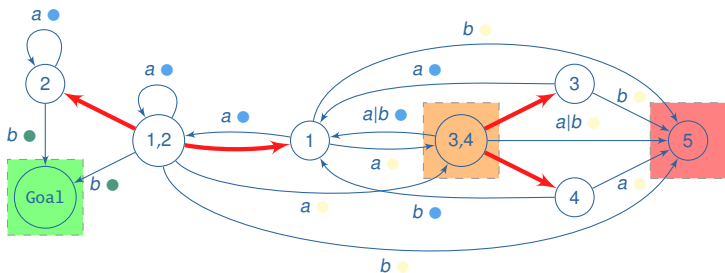


# Reaching the goal almost-surely

**Lose**: beliefs that contain a state losing in the (fully-observable) MDP

$$Win = \mathcal{B} \setminus Lose = W_{ok} \sqcup W_{req} \sqcup W_{safe}$$

- ▶  $W_{ok} = \{S \mid S \subseteq Goal\}$
- ▶  $W_{req} = \{S \mid \forall a \in Act \exists O \in \mathcal{O}, up(S, a, O) \in Lose\}$



# Reaching the goal almost-surely

**Lose**: beliefs that contain a state losing in the (fully-observable) MDP

$$Win = \mathcal{B} \setminus Lose = W_{ok} \sqcup W_{req} \sqcup W_{safe}$$

- ▶  $W_{ok} = \{S \mid S \subseteq Goal\}$
- ▶  $W_{req} = \{S \mid \forall a \in Act \exists O \in \mathcal{O}, up(S, a, O) \in Lose\}$

Canonical family of strategies  $(\sigma_n)_{n \in \mathbb{N}}$ :

- ▶ In  $W_{req}$ , play **req**, and
- ▶ in  $W_{safe}$ , play **req** with prob.  $1/n$  and unif. prob. on safe actions.

$a$  is safe from  $S$  if  $\forall O, up(S, a, O) \notin Lose$

## Lemma

$\sigma_n$  is almost-surely winning from  $Win$ .

# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$

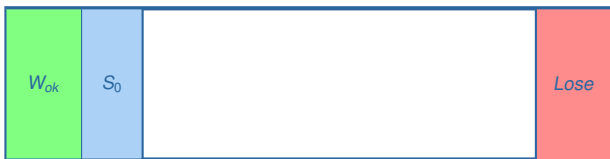


# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



Computation of  $S_0$ :  $S_0 = \text{reach}_{=1}(W_{ok})$

almost-sure reachability question for the belief-MDP without requests

Optimized strategy: no request from  $S \in S_0$

# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



Computation of  $S_1$

- ▶  $L_1 = \{S \mid \forall s \in S, \{s\} \in S_0\}$
- ▶  $S_1 = \text{reach}_{=1}(L_1 \cup S_0)$

Optimized Strategy:

request from  $S \in L_1 \setminus S_0$

uniform distribution on actions ensuring to stay in  $S_1$ , othw

# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



Stabilisation for  $N \leq |\mathcal{B}|$

# Optimizing the worst-case cost

---

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



Stabilisation for  $N \leq |\mathcal{B}|$

$$S_{\infty} = \text{Win} \setminus S_N$$

# Optimizing the worst-case cost

Iterative computation of  $S_k$ : set of beliefs where  $k$  req are sufficient.

$$S_0 \subseteq S_1 \subseteq S_2 \cdots \subseteq \text{Win}$$



## Proposition

The minimum worst-case cost can be computed in EXPTIME, together with a finite-memory strategy.

# Outline

---

- 1 Introduction
- 2 Worst-case cost
- 3 Average cost**
- 4 Conclusion

# Undecidability

---

Value: infimum of average cost over almost-surely winning strategies

$$\text{val}(G) = \inf\{\text{av\_cost}(\sigma) \mid \sigma \text{ almost-surely winning}\}$$



# Undecidability

---

Value: infimum of average cost over almost-surely winning strategies

$$\text{val}(G) = \inf\{\text{av\_cost}(\sigma) \mid \sigma \text{ almost-surely winning}\}$$

The value cannot be computed

For all  $K > 0$ , one cannot decide whether  $\text{val}(G) \leq K$ .

# Undecidability

---

Value: infimum of average cost over almost-surely winning strategies

$$val(G) = \inf\{av\_cost(\sigma) \mid \sigma \text{ almost-surely winning}\}$$

The value cannot be computed

For all  $K > 0$ , one cannot decide whether  $val(G) \leq K$ .

Not too surprising: optimizing cost functions for POMDP is undecidable

▶ Skip proof

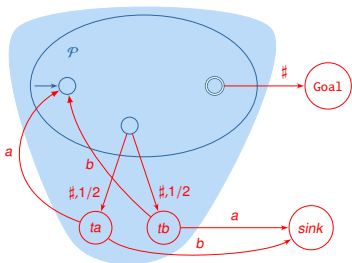
# Proof idea

---

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]

# Proof idea

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]



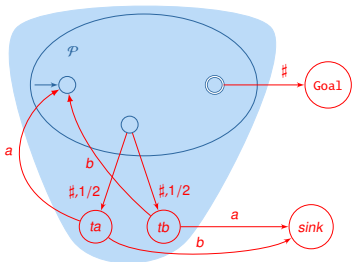
$\mathcal{P}$  accepts a word with probability greater than  $1 - \varepsilon$

iff

$$\text{val}(G) < \frac{\varepsilon}{1-\varepsilon}$$

# Proof idea

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]



$\mathcal{P}$  accepts a word with probability greater than  $1 - \varepsilon$

iff

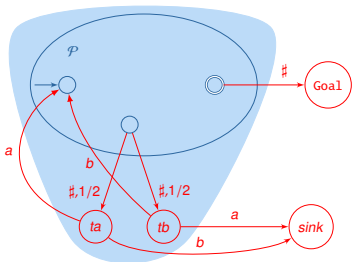
$$\text{val}(G) < \frac{\varepsilon}{1-\varepsilon}$$

$(\Rightarrow) \sigma$  plays  $(w \# \text{req } a|b)^*$  for  $w$  with  $\mathbb{P}(w) > 1 - \varepsilon$

$$\text{val}(\sigma) < 0 \times (1 - \varepsilon) + 1 \times \varepsilon(1 - \varepsilon) + 2 \times \varepsilon^2(1 - \varepsilon) \dots = \varepsilon/(1 - \varepsilon)$$

# Proof idea

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]



$\mathcal{P}$  accepts a word with probability greater than  $1 - \varepsilon$

iff

$$\text{val}(G) < \frac{\varepsilon}{1-\varepsilon}$$

$(\Rightarrow) \sigma$  plays  $(w \# \text{req } a|b)^*$  for  $w$  with  $\mathbb{P}(w) > 1 - \varepsilon$

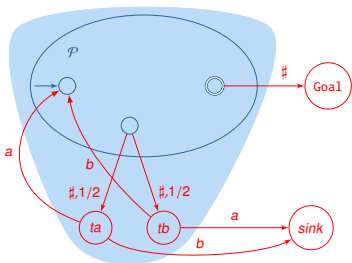
$$\text{val}(\sigma) < 0 \times (1 - \varepsilon) + 1 \times \varepsilon(1 - \varepsilon) + 2 \times \varepsilon^2(1 - \varepsilon) \cdots = \varepsilon/(1 - \varepsilon)$$

$(\Leftarrow) p$  probability in  $\sigma$  to have  $\#$  before  $\text{req}$

$$\text{val}(\sigma) > (1 - p) + p(1 - \varepsilon) \geq 1 - \varepsilon$$

# Proof idea

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]



$\mathcal{P}$  accepts a word with probability greater than  $1 - \varepsilon$

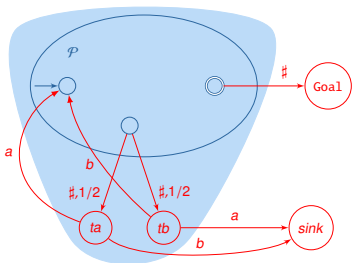
iff

$$val(G) < \frac{\varepsilon}{1-\varepsilon}$$

best approximation:  $|v - val(G)| = (\varepsilon/(1 - \varepsilon) + (1 - \varepsilon))/2$

# Proof idea

$\mathcal{P}$  PFA s.t. either all words have probability  $\leq \varepsilon$ , or some word has probability  $> 1 - \varepsilon$ . Which holds is undecidable. [Madani Hanks Condon 03]



$\mathcal{P}$  accepts a word with probability greater than  $1 - \varepsilon$

iff

$$val(G) < \frac{\varepsilon}{1-\varepsilon}$$

best approximation:  $|v - val(G)| = (\varepsilon/(1 - \varepsilon) + (1 - \varepsilon))/2$

approximation factor:  $\frac{|v - val(G)|}{val(G)} = \frac{(1-\varepsilon)(1/(1-\varepsilon)-\varepsilon)}{2\varepsilon} \xrightarrow{\varepsilon \rightarrow 0} \infty$



# Non-approximability

---

Corollary: For every  $\delta$  it is undecidable to approximate  $val(G)$  within  $\delta$ .

NB: bigger  $\delta$  need bigger POMDP

# Non-approximability

---

Corollary: For every  $\delta$  it is undecidable to approximate  $val(G)$  within  $\delta$ .

NB: bigger  $\delta$  need bigger POMDP

## NP-hardness of good approximations

Assuming  $P \neq NP$  there is a POMDP  $G$  with

**few** reachable belief states (quadratic in  $n$ ) s.t.

any polynomial time algorithm  $\mathcal{A}$  returns for  $G$  a value  $v$  with

approximation factor:  $\frac{|v - val(G)|}{val(G)} \geq 2^{n-1}/n^2$ , and

absolute approximation error:  $|v - val(G)| \geq 2^{n-1}/n$ .

# Proof idea

---

$\varphi$  3-SAT instance with  $m$  clauses and  $k$  variables;  $n = mk$

$\varphi$  is satisfiable if for each clause  $C_i$ , one can choose a literal  $l_i$  and the choices do not conflict

POMDP behaviour:

- ▶ random choice of variable to monitor
- ▶ conflicts force a **request** not to lose

# Proof idea

---

$\varphi$  3-SAT instance with  $m$  clauses and  $k$  variables;  $n = mk$

$\varphi$  is satisfiable if for each clause  $C_i$ , one can choose a literal  $l_i$  and the choices do not conflict

POMDP behaviour:

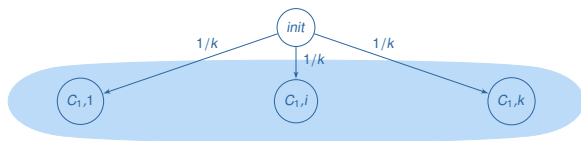
- ▶ random choice of variable to monitor
- ▶ conflicts force a **request** not to lose

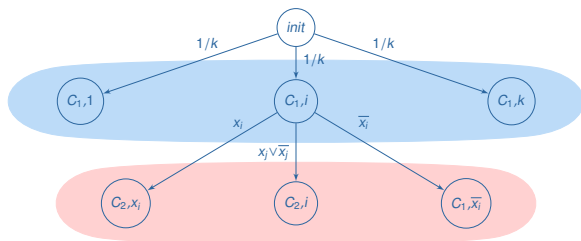
## Properties of the reduction

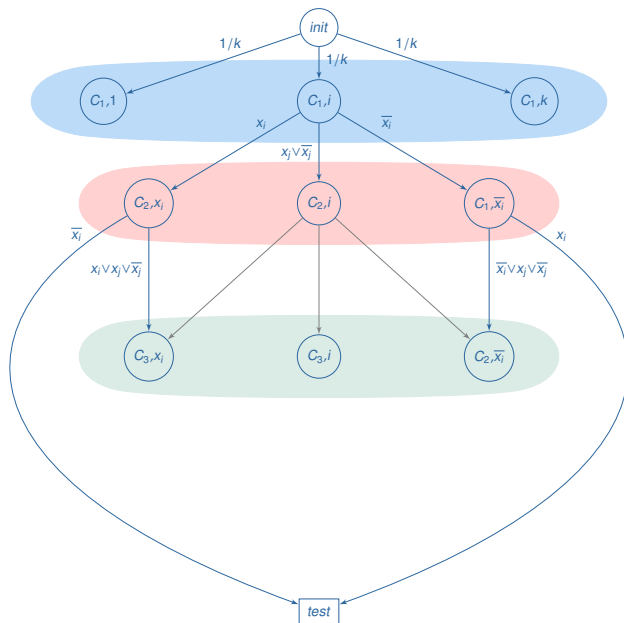
- ▶  $\varphi$  satisfiable  $\Rightarrow \text{val}(G) < n$
- ▶  $\varphi$  not satisfiable  $\Rightarrow \text{val}(G) > 2^n/n - 2$

▶ Skip details

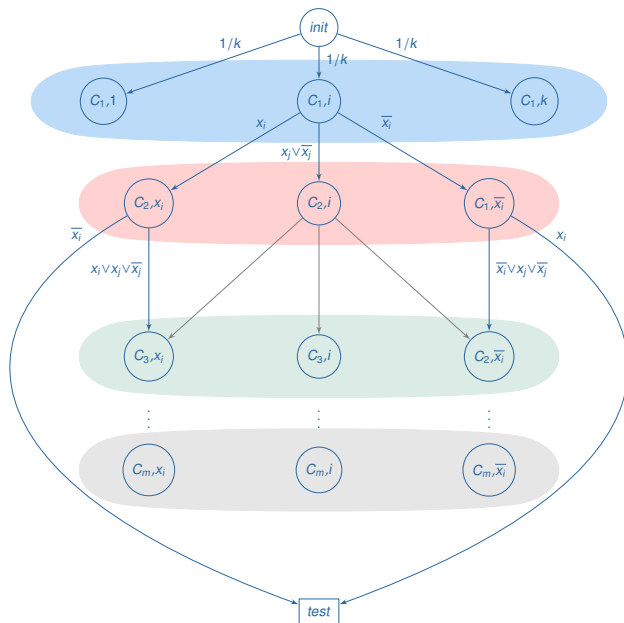


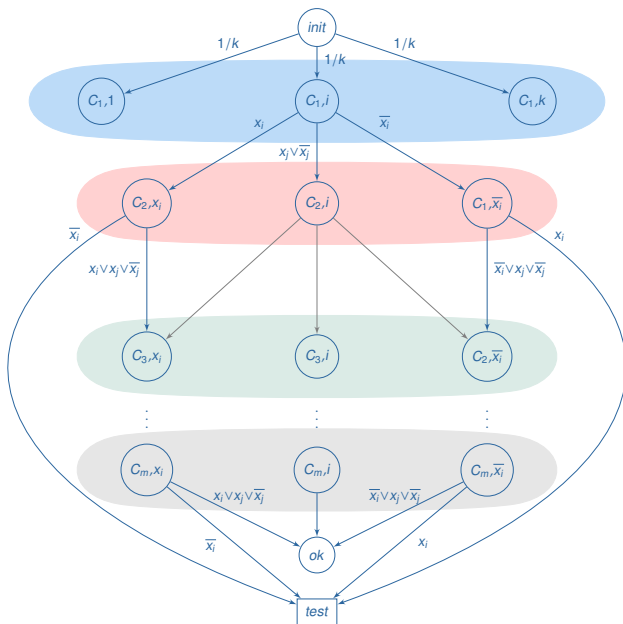


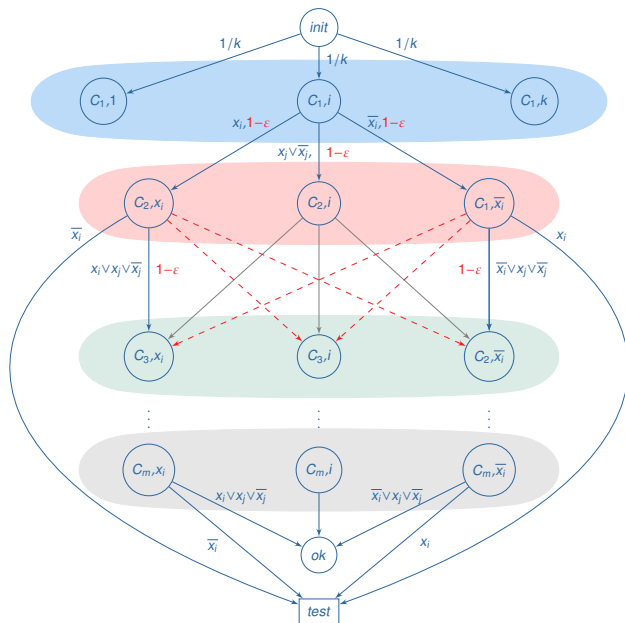












# Outline

---

- 1 Introduction
- 2 Worst-case cost
- 3 Average cost
- 4 Conclusion**

# Conclusion

---

## Contribution

Minimize requests for full-information in POMDP  
under an almost-sure reachability objective.

- ▶ Worst-case cost
  - ▶ computation in EXPTIME, together with an optimal strategy
- ▶ Average cost
  - ▶ computation undecidable
  - ▶ approximation unfeasible
  - ▶ large least approximation factors for polytime algorithms

# Conclusion

---

## Contribution

Minimize requests for full-information in POMDP  
under an almost-sure reachability objective.

- ▶ Worst-case cost
  - ▶ computation in EXPTIME, together with an optimal strategy
- ▶ Average cost
  - ▶ computation undecidable
  - ▶ approximation unfeasible
  - ▶ large least approximation factors for polytime algorithms

## Future work

- ▶ extend to several information levels
  - ▶ successive partition refinement
- ▶ tradeoff between objective (reachability probability) and cost