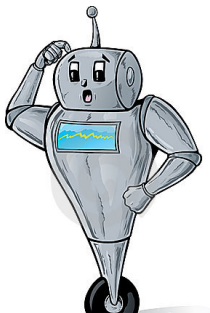


# QUANTUM POMDPs

Jenny Barry

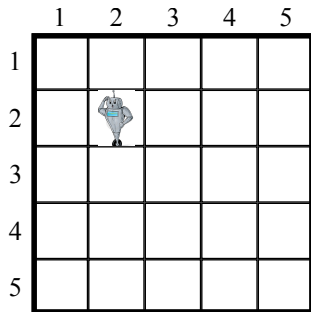
6.845 Final Project Presentation  
December 12, 2012



## ROBOTS...

- Don't know where they are.
- Don't know what they are doing.
- Don't understand what they are seeing.

# POMDPs

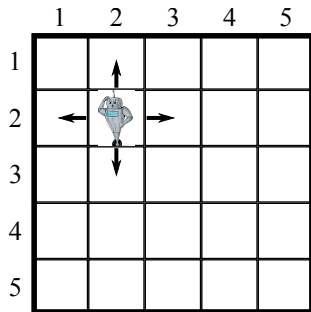


**States:**  $(i, j)$

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations

# POMDPs



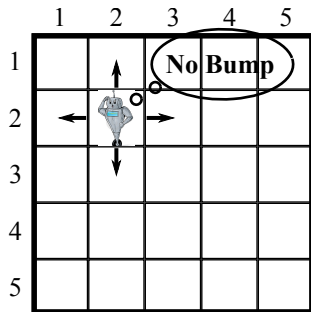
**States:**  $(i, j)$

**Actions:** (L, R, U, D, S)

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations

# POMDPs



**States:**  $(i, j)$

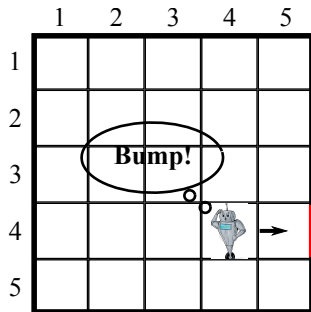
**Actions:** (L, R, U, D, S)

**Observations:** No Bump, Bump

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations

# POMDPs



**States:**  $(i, j)$



**Actions:** (L, R, U, D, S)

**Observations:** No Bump, Bump

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations


# POMDPs

	1	2	3	4	5
1	-1	-1	-1	-1	-1
2	-1		-1	-1	-1
3	-1	-1	-1	-1	-1
4	-1	-1	-1		-1
5	-1	-1	-1	-1	-1

**States:**  $(i, j)$

**Actions:** (L, R, U, D, S)

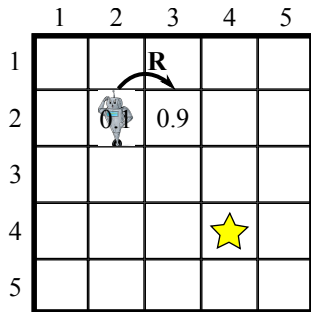
**Observations:** No Bump, Bump

**Rewards:** 0 at , -1 else

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations
- $R(s_i, a_j)$ : Reward for taking action  $a_j$  in state  $s_i$

# POMDPs



**States:**  $(i, j)$

**Actions:** (L, R, U, D, S)

**Observations:** No Bump, Bump

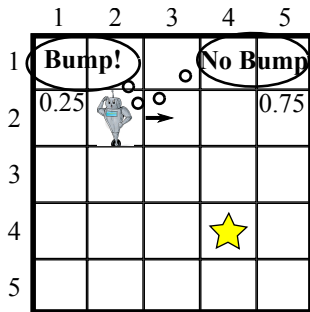
**Rewards:** 0 at ★, -1 else

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations
- $R(s_i, a_j)$ : Reward for taking action  $a_j$  in state  $s_i$
- $T(s_i|a_j, s_k)$ : Probability of transitioning to  $s_i$  starting in  $s_j$  taking action  $a_j$



# POMDPs



**States:**  $(i, j)$

**Actions:** (L, R, U, D, S)

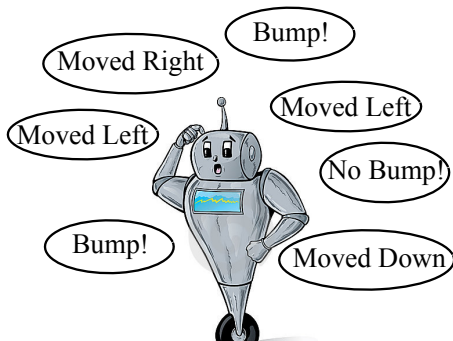
**Observations:** No Bump, Bump

**Rewards:** 0 at ★, -1 else

## PARTIALLY OBSERVABLE MARKOV DECISION PROCESS (POMDP)

- $S, A, \Omega$ : Possible states, actions, observations
- $R(s_i, a_j)$ : Reward for taking action  $a_j$  in state  $s_i$
- $T(s_i | a_j, s_k)$ : Probability of transitioning to  $s_i$  starting in  $s_j$  taking action  $a_j$
- $O(o_i | a_j, s_k)$ : Probability of observing  $o_i$  given that action  $a_j$  ended in  $s_k$

# BELIEF STATES



## DEFINITION: BELIEF STATE

POMDP  $P = \langle S, A, \Omega, R, T, O \rangle \Rightarrow$  Belief space  $B \subset \mathbb{R}^{|S|}$ :

- $\vec{b}_i = \Pr(s_i)$
- $\sum_i \vec{b}_i = |\vec{b}|_1 = 1$

# BELIEF STATES

	1	2	3	4	5
1	0.04	0.04	0.04	0.04	0.04
2	0.04	0.04	0.04	0.04	0.04
3	0.04	0.04	0.04	0.04	0.04
4	0.04	0.04	0.04	0.04	0.04
5	0.04	0.04	0.04	0.04	0.04


## DEFINITION: BELIEF STATE

POMDP  $P = \langle S, A, \Omega, R, T, O \rangle \Rightarrow$  Belief space  $B \subset \mathbb{R}^{|S|}$ :

- $\vec{b}_i = \Pr(s_i)$
- $\sum_i \vec{b}_i = |\vec{b}|_1 = 1$

# BELIEF STATES

	1	2	3	4	5
1	0.04	0.04	0.04	0.04	0.04
2	0.04	0.04	0.04	0.04	0.04
3	0.04	0.04	0.04	0.04	0.04
4	0.04	0.04	0.04	0.04	0.04
5	0.04	0.04	0.04	0.04	0.04

Move Right  
  
See No Bump

	1	2	3	4	5
1	0	0.07	0.07	0.07	0
2	0	0.07	0.07	0.07	0
3	0	0.07	0.07	0.07	0
4	0	0.07	0.07	0.07	0
5	0	0.07	0.07	0.07	0

## DEFINITION: BELIEF STATE

POMDP  $P = \langle S, A, \Omega, R, T, O \rangle \Rightarrow$  Belief space  $B \subset \mathbb{R}^{|S|}$ :

- $\vec{b}_i = \Pr(s_i)$
- $\sum_i \vec{b}_i = |\vec{b}|_1 = 1$

# BELIEF STATES

## DEFINITION: BELIEF STATE

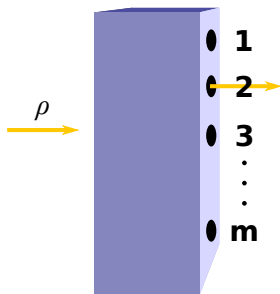
POMDP  $P = \langle S, A, \Omega, R, T, O \rangle \Rightarrow$  Belief space  $B \subset \mathbb{R}^{|S|}$ :

- $\vec{b}_i = \Pr(s_i)$
- $\sum_i \vec{b}_i = |\vec{b}|_1 = 1$

## BELIEF MARKOV DECISION PROCESS

- $B$ : Belief space (continuous)
- $A$ : Robot's actions
- $\tau(\vec{b}'|a_i, \vec{b})$ : Probability of  $\vec{b}'$  after taking action  $a_i$  in state  $\vec{b}$ .
- $\rho(\vec{b}, a_i) = \sum_i \vec{b}_i R(s_i, a_i)$ : Reward for taking action  $a_i$  in state  $\vec{b}$
- $b_0$ : Starting belief state

*I know that I know nothing.* - Socrates

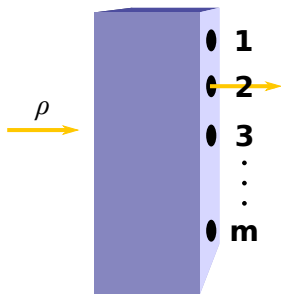


## DEFINITION: SUPEROPERATOR

$$\mathbf{S} = \{K_1, \dots, K_m\}$$

- $\sum_{i=1}^m K_i^\dagger K_i = \mathbb{I}$
- $\Pr[\text{Observation } i] = \text{Tr}(K_i \rho K_i^\dagger)$

$$\rho \rightarrow \frac{K_i \rho K_i^\dagger}{\text{Tr}(K_i \rho K_i^\dagger)}$$



## DEFINITION: SUPEROPERATOR

$$S = \{K_1, \dots, K_m\}$$

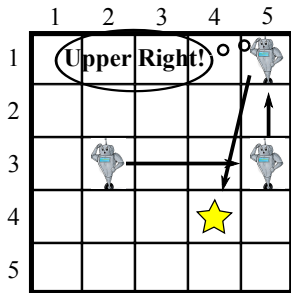
- $\sum_{i=1}^m K_i^\dagger K_i = \mathbb{I}$
- $\Pr[\text{Observation } i] = \text{Tr}(K_i \rho K_i^\dagger)$

$$\rho \rightarrow \frac{K_i \rho K_i^\dagger}{\text{Tr}(K_i \rho K_i^\dagger)}$$

## QUANTUM OBSERVABLE MARKOV DECISION PROCESS (QOMDP)

- $S$ : Hilbert space
- $\Omega$ : Set of observations
- $\mathcal{A}$ : Set of quantum superoperators
- $R$ : Reward function
- $\rho_0$ : Starting state

# POMDPs ARE HARD...

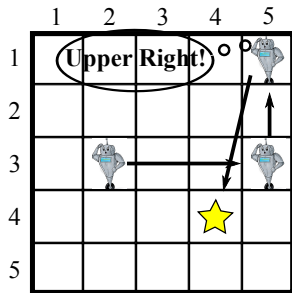


## Strategy:

- 1 Localize: go right until wall, then up
- 2 Go to goal



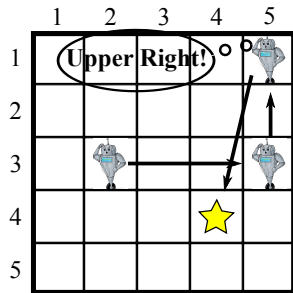
# POMDPs ARE HARD...



$$\pi(\vec{b}) = \begin{cases} R & \text{if } \sum_{s \in \text{Right Wall}} \vec{b}(s) < 1 \\ U & \text{if } \sum_{s \in \text{Right Wall}} \vec{b}(s) = 1 \\ & \text{and } \vec{b}(\text{Upper Right}) < 1 \\ \text{Go to goal} & \text{if } \|\vec{b}\|^2 = 1 \end{cases}$$

**POLICY:**  $\pi(\vec{b}, t) = a$  specifies action to take in belief  $\vec{b}$  at time  $t$

# POMDPs ARE HARD...



$$\pi(\vec{b}) = \begin{cases} R & \text{if } \sum_{s \in \text{Right Wall}} \vec{b}(s) < 1 \\ U & \text{if } \sum_{s \in \text{Right Wall}} \vec{b}(s) = 1 \\ & \text{and } \vec{b}(\text{Upper Right}) < 1 \\ \text{Go to goal} & \text{if } \|\vec{b}\|^2 = 1 \end{cases}$$

**POLICY:**  $\pi(\vec{b}, t) = a$  specifies action to take in belief  $\vec{b}$  at time  $t$

## POLICY EXISTENCE PROBLEM (PEP)

Given POMDP  $P = \langle S, A, \Omega, R, T, O \rangle$ , decide if there is some policy  $\pi$  that has expected future reward at least  $V$  over the next  $h$  timesteps.

- If  $h = \text{poly}(S)$ , PEP is in PSPACE and **PSPACE-COMPLETE**.
- If  $h = \infty$ , PEP is **UNDECIDABLE**.

# ...BUT QOMDPs ARE HARDER

**POMDPs  $\subseteq$  QOMDPs**

- PEP with  $h = \text{poly}(d)$  is at least PSPACE-Complete
- ✓ PEP with  $h = \infty$  is **UNDECIDABLE**

# ...BUT QOMDPs ARE HARDER

## POMDPs $\subseteq$ QOMDPs

- ✓ PEP with  $h = \text{poly}(d)$  is **PSPACE-COMPLETE**
- ✓ PEP with  $h = \infty$  is **UNDECIDABLE**

### THEOREM

PEP for QOMDPs with  $h = \text{poly}(d)$  is in PSPACE.

**PROOF SKETCH:** There are only  $O((|A||\Omega|)^h)$  policies. Try them all.

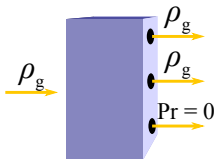
# ...BUT QOMDPs ARE HARDER

## POMDPs $\subseteq$ QOMDPs

- ✓ PEP with  $h = \text{poly}(d)$  is **PSPACE-COMPLETE**
- ✓ PEP with  $h = \infty$  is **UNDECIDABLE**

## GOAL-STATE REACHABILITY PROBLEM (GRP)

Assume the Q(P)OMDP has an absorbing goal state. Decide if there is a policy that reaches this goal state with probability 1.



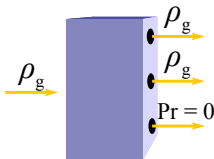
# ...BUT QOMDPs ARE HARDER

## POMDPs $\subseteq$ QOMDPs

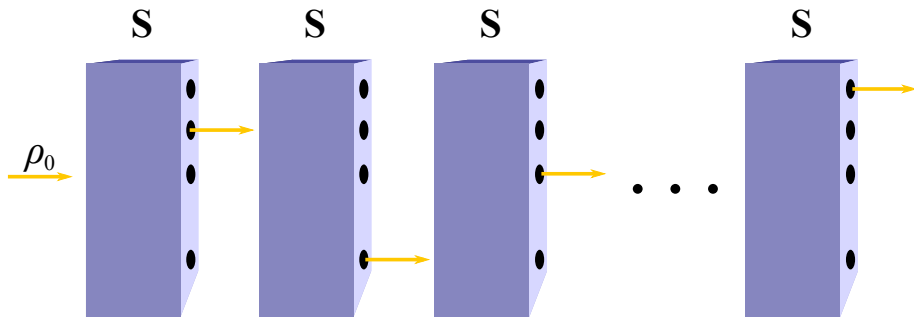
- ✓ PEP with  $h = \text{poly}(d)$  is **PSPACE-COMPLETE**
- ✓ PEP with  $h = \infty$  is **UNDECIDABLE**
- GRP is **DECIDABLE** for POMDPs
- GRP is **UNDECIDABLE** for QOMDPs

## GOAL-STATE REACHABILITY PROBLEM (GRP)

Assume the Q(P)OMDP has an absorbing goal state. Decide if there is a policy that reaches this goal state with probability 1.



# GRP FOR QOMDPs: QMOP

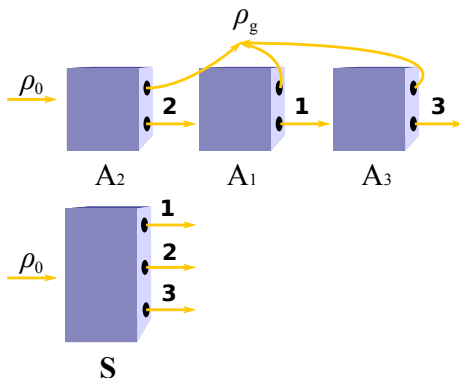


## QUANTUM MEASUREMENT OCCURRENCE PROBLEM (QMOP)

Given a superoperator  $S = \{K_1, \dots, K_m\}$  and starting state  $\rho_0$ , decide if there is some finite sequence of measurements that can never be observed if  $\rho_0$  is continually fed back into  $S$ .

QMOP is **UNDECIDABLE** [Eisert12]

# GRP FOR QOMDPs: REDUCTION FROM QMOP

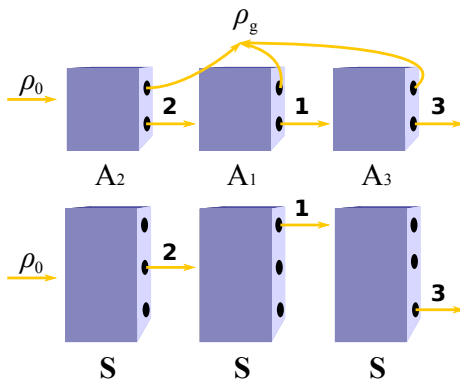


Given QMOP  $S = \{K_1, \dots, K_m\}$ :

- $m$  actions. Action  $i$  either:
  - Transitions according to  $K_i$
  - Transitions to goal state
- $m + 1$  observations:
  - ① At-Goal
  - ② Observation  $i$  from QMOP



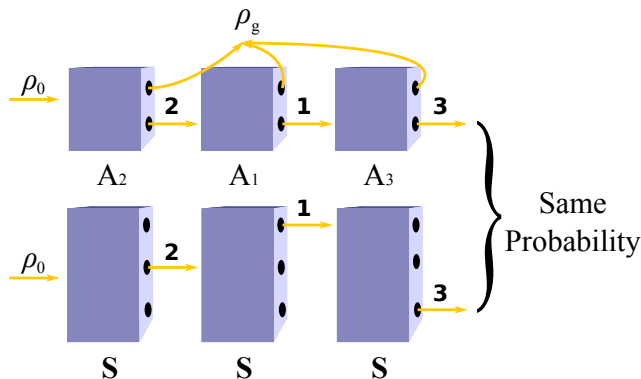
# GRP FOR QOMDPs: REDUCTION FROM QMOP



Given QMOP  $S = \{K_1, \dots, K_m\}$ :

- $m$  actions. Action  $i$  either:
  - Transitions according to  $K_i$
  - Transitions to goal state
- $m + 1$  observations:
  - ① At-Goal
  - ② Observation  $i$  from QMOP

# GRP FOR QOMDPs: REDUCTION FROM QMOP

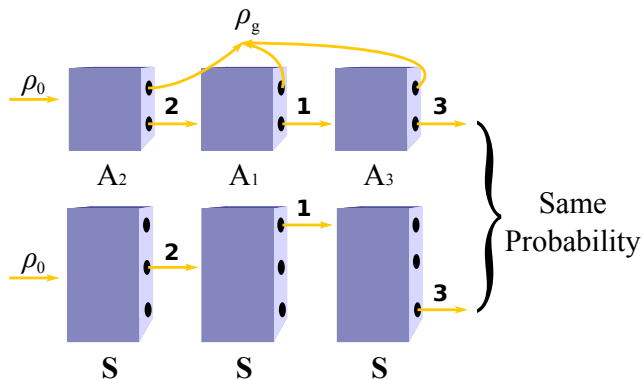


Given QMOP  $S = \{K_1, \dots, K_m\}$ :

- $m$  actions. Action  $i$  either:
  - Transitions according to  $K_i$
  - Transitions to goal state
- $m + 1$  observations:
  - ① At-Goal
  - ② Observation  $i$  from QMOP
- $\Pr(\rho_n \neq \text{goal} \mid \text{actions } j_1, \dots, j_n) = \Pr(\text{Observing sequence } j_1, \dots, j_n)$ .

$\Rightarrow$  Path to goal of probability 1 if and only some sequence unobservable.

# GRP FOR QOMDPs: REDUCTION FROM QMOP



## THEOREM

GRP for QOMDPs is undecidable.

# GOAL-STATE REACHABILITY FOR POMDPs

## CONVERSION TO PLUS/ZERO LAND

$$\begin{bmatrix} 0.2 & 0 & 0.8 \\ 0.3 & 0.1 & 0.6 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} + & 0 & + \\ + & + & + \\ 0 & 0 & + \end{bmatrix} \begin{bmatrix} + \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} + \\ + \\ 0 \end{bmatrix}$$

# GOAL-STATE REACHABILITY FOR POMDPs

## CONVERSION TO PLUS/ZERO LAND

$$\begin{bmatrix} 0.2 & 0 & 0.8 \\ 0.3 & 0.1 & 0.6 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} + & 0 & + \\ + & + & + \\ 0 & 0 & + \end{bmatrix} \begin{bmatrix} + \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} + \\ + \\ 0 \end{bmatrix}$$

- Convert POMDP probabilities to plus/zero
  - Finitely many  $(2^{|S|} - 1)$  states
  - Finitely many policies
- ⇒ We find the goal state or repeat a previously seen state in finite time.

# GOAL-STATE REACHABILITY FOR POMDPs

## CONVERSION TO PLUS/ZERO LAND

$$\begin{bmatrix} 0.2 & 0 & 0.8 \\ 0.3 & 0.1 & 0.6 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} + & 0 & + \\ + & + & + \\ 0 & 0 & + \end{bmatrix} \begin{bmatrix} + \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} + \\ + \\ 0 \end{bmatrix}$$

- Convert POMDP probabilities to plus/zero
- Finitely many  $(2^{|S|} - 1)$  states
- Finitely many policies

⇒ We find the goal state or repeat a previously seen state in finite time.

## THEOREM

GRP for POMDPs is decidable.

## COMPLEXITY PROBLEMS

- Complexity separations using non-negative properties of POMDPs
- Complexity separations using value function structure of POMDPs
- What if we don't know the starting state in a QOMDP?

## COMPLEXITY PROBLEMS

- Complexity separations using non-negative properties of POMDPs
- Complexity separations using value function structure of POMDPs
- What if we don't know the starting state in a QOMDP?

## ALGORITHMS

- Algorithms for solving QOMDPs
- Algorithms for approximating QOMDPs



# FUTURE WORK

## COMPLEXITY PROBLEMS

- Complexity separations using non-negative properties of POMDPs
- Complexity separations using value function structure of POMDPs
- What if we don't know the starting state in a QOMDP?

## ALGORITHMS

- Algorithms for solving QOMDPs
- Algorithms for approximating QOMDPs

## APPLICATIONS

- Reward structure for QOMDPs
- Practical applications of QOMDPs