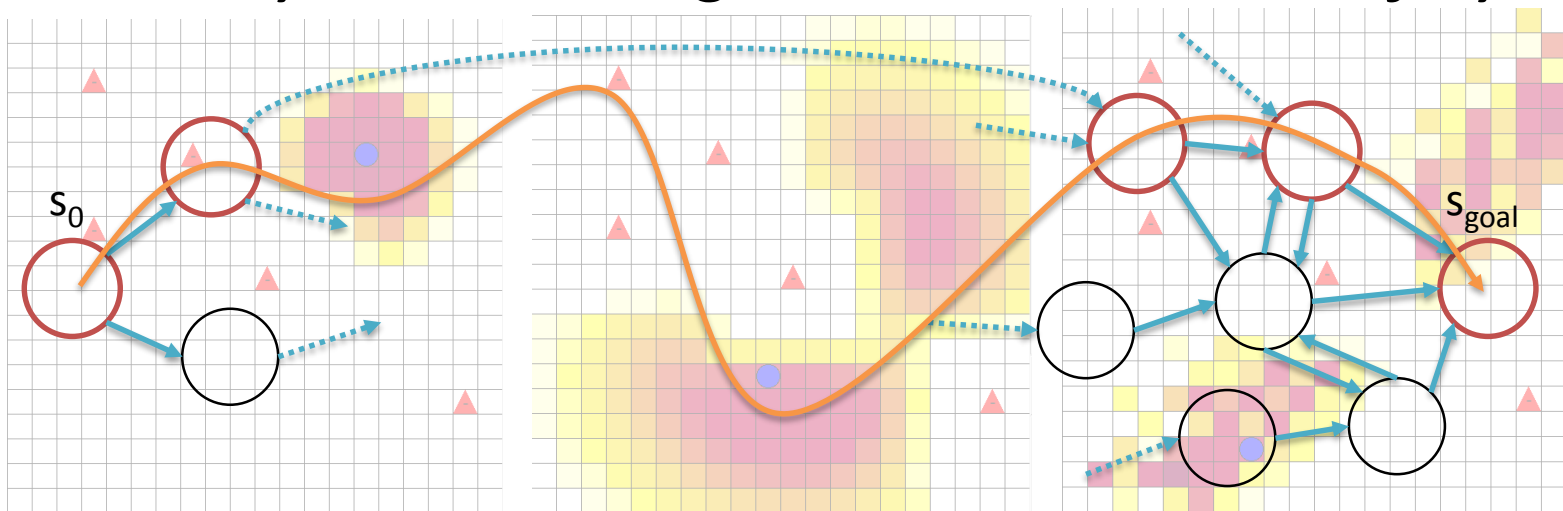


Risk-aware AO* (RAO*)

When you want to get there well and safely



16.412/6.834 Cognitive Robotics

Pedro Santana (psantana@mit.edu)

May the 4th be with you, 2016.





MIT
AEROASTRO

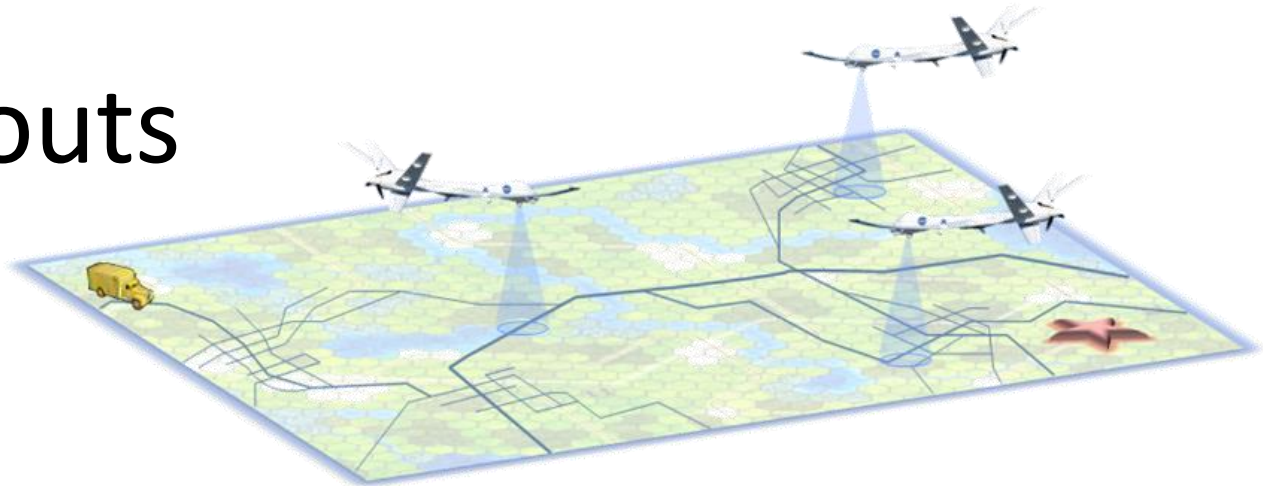
1. Motivation

Where can risk-aware planning be useful?



MIT
AEROASTRO

Science scouts



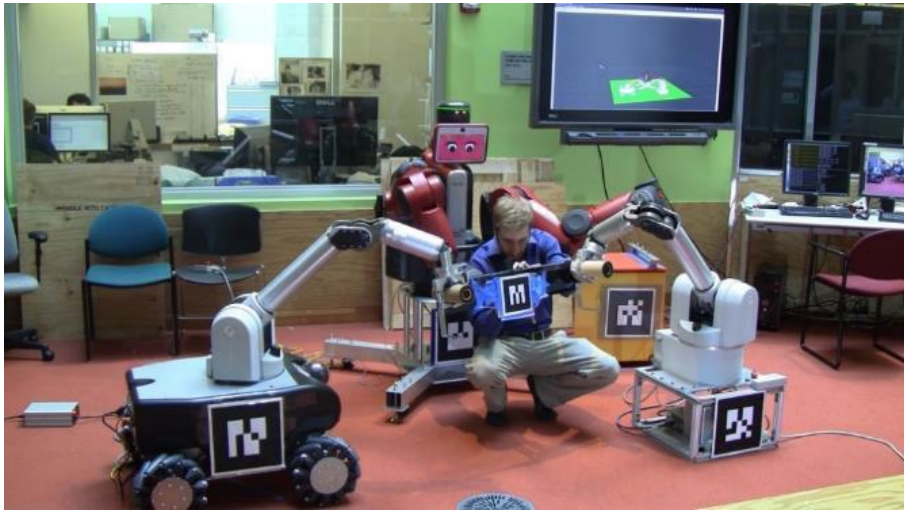
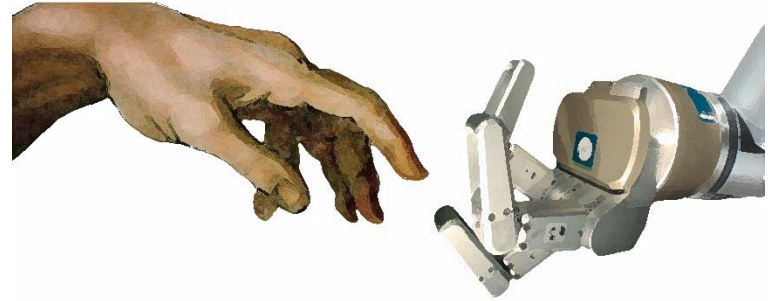
Courtesy of Andrew J. Wang





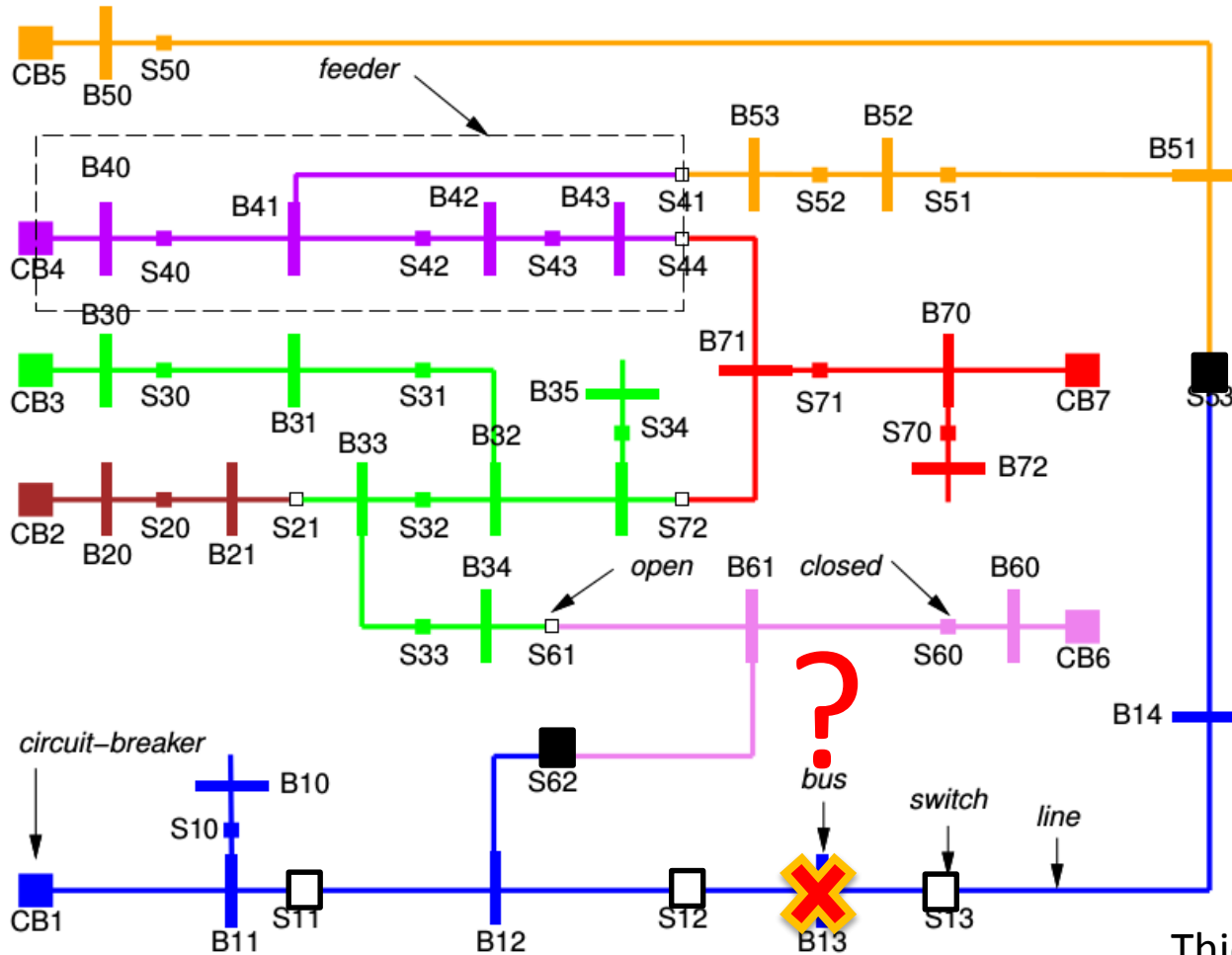
MIT
AEROASTRO

Collaborative manufacturing





Power supply restoration



Thiébaux & Cordier (2001)



Resilient Space Systems (RSS) demo

The image displays a ROS environment with a terminal window on the left and an RViz simulation on the right. The terminal window shows the execution of an activity plan, with a large text overlay that reads "Activity Plan Generated!". The RViz simulation shows a 3D environment with a rover and four locations (I1, I2, I3, I5) marked with yellow dots. A path is shown connecting these locations. The terminal window also shows the execution of the activity plan, with a list of tasks and their timestamps.

```
Activity Plan Generated!
```

```
0.0 : (move rover1 I1 I3)
270.0 : (turnon_mastcam rover1 I3)
290.0 : (take_pictures_mastcam rover1 I3 pic_req1)
310.0 : (move rover1 I3 I5)
650.0 : (turnon_mastcam rover1 I5)
670.0 : (take_pictures_mastcam rover1 I5 pic_req3)
690.0 : (survey_location rover1 I5)
740.0 : (collect_rock_sample rover1 I5 rock_req1)
790.0 : (move rover1 I5 I2)
1110.0 : (transmit_data rover1 I2 rock_req1)
1140.0 : (turnon_mastcam rover1 I2)
1160.0 : (take_pictures_mastcam rover1 I2 pic_req2)
1180.0 : (transmit_data rover1 I2 pic_req2)
1210.0 : (transmit_data rover1 I2 pic_req1)
1240.0 : (transmit_data rover1 I2 pic_req3)
1270.0 : (survey_location rover1 I2)
1320.0 : (collect_rock_sample rover1 I2 rock_req2)
1370.0 : (transmit_data rover1 I2 rock_req2)
```

Joint work between JPL, Caltech, WHOI, and MIT.

Our goal for today

*How can we generate **safe plans** that
optimize performance when
controlling a system with **stochastic
transitions** and **hidden state**?*

Today's topics

1. Motivation
2. Handling belief states
3. RAO*

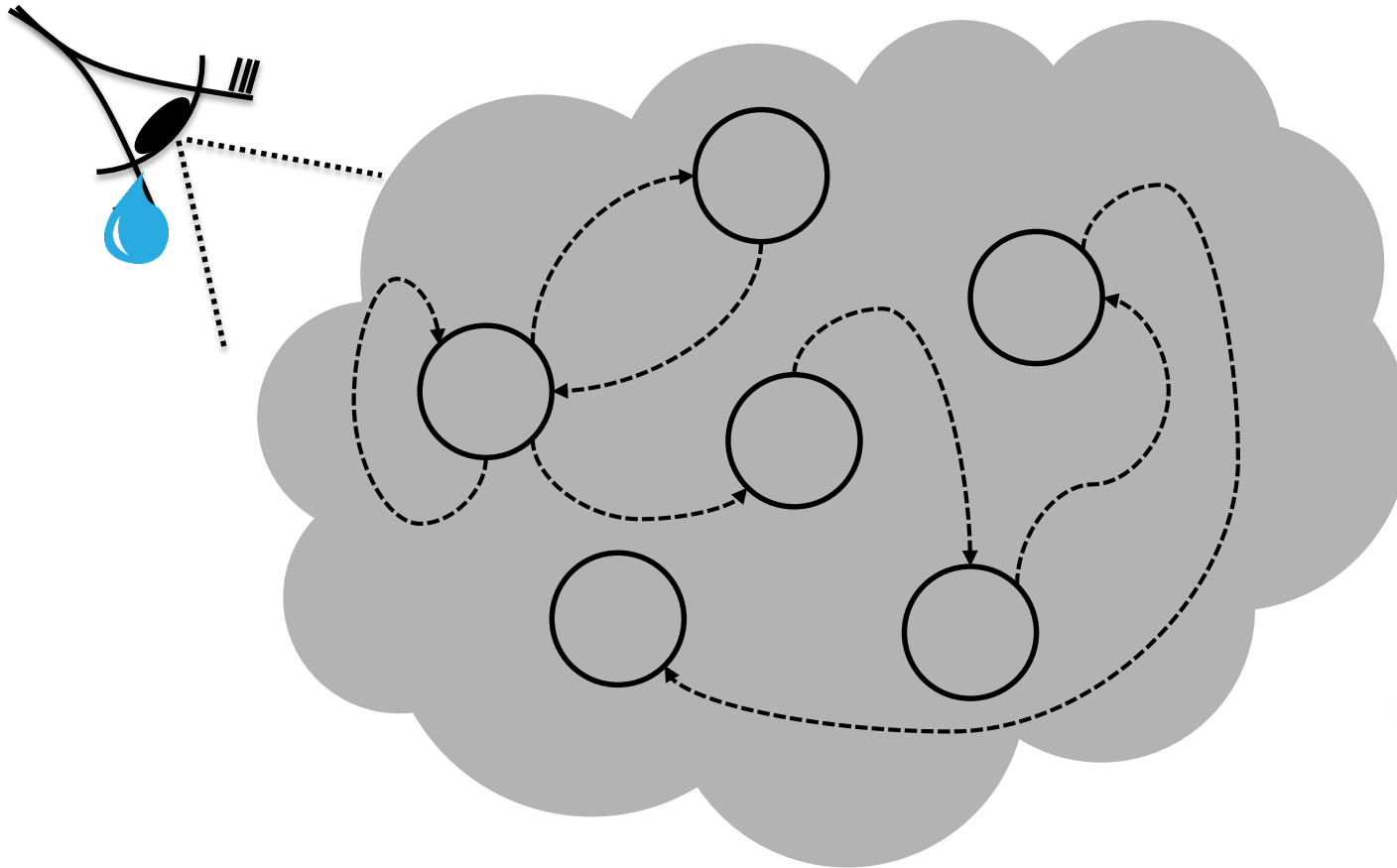
2. Handling belief states

“Probability is common sense reduced to calculation.”
— Pierre-Simon Laplace



MIT
AEROASTRO

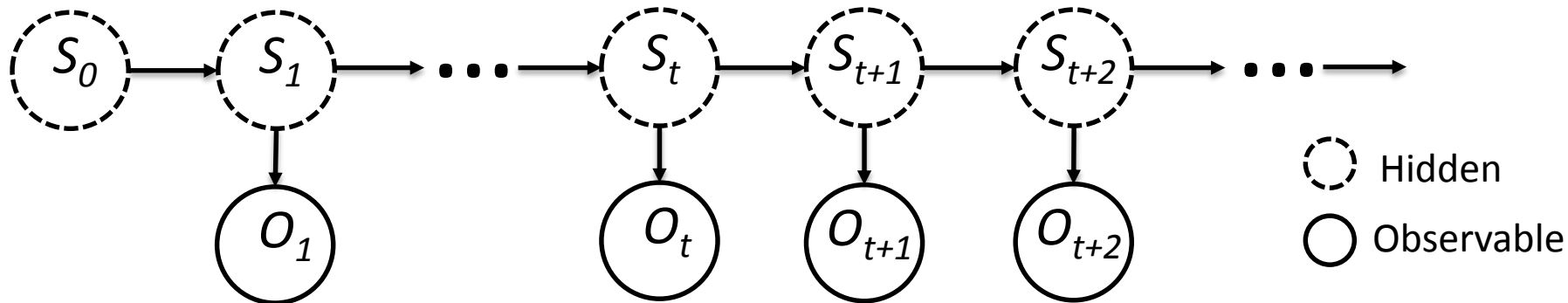
Hidden Markov models (HMMs)



Andrey Markov



Observing hidden Markov chains



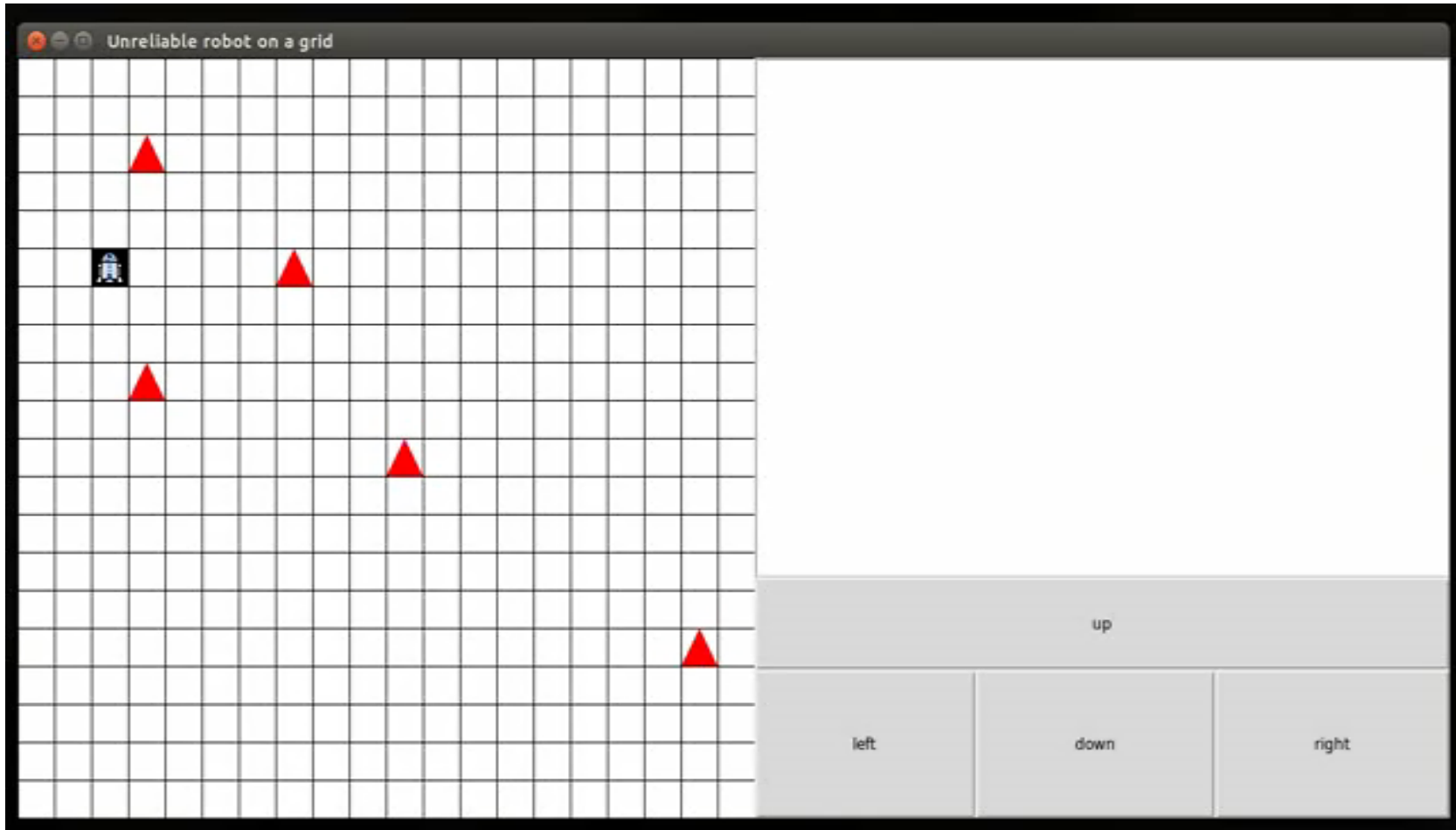
Definition: Hidden Markov Model (HMM)

A sequence of random variables $O_1, O_2, \dots, O_t, \dots$, is an HMM if the distribution of O_t is completely defined by the current (hidden) state S_t according to

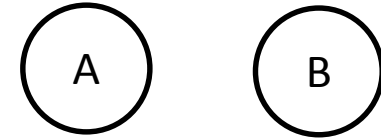
$$\Pr(O_t | S_t),$$

where S_t is part of an underlying Markov chain.

Robot navigation



Bayes' rule



A, B : random variables

$$\underbrace{\Pr(A, B)}_{\text{Joint}} = \underbrace{\Pr(A|B)}_{\text{Conditional}} \underbrace{\Pr(B)}_{\text{Marginal}}$$

$$\Pr(A, B) = \Pr(B|A)\Pr(A)$$

$$\Pr(A|B)\Pr(B) = \Pr(B|A)\Pr(A)$$

$$\Pr(A|B) = \frac{\Pr(B|A)\Pr(A)}{\Pr(B)} \propto \Pr(B|A)\Pr(A)$$

Bayes' rule!

Notation

Random variable!

$\Pr(S_t | \cdot)$ \longrightarrow Probability distribution of S_t

Vector of d
probability values.

$\Pr(S_t = s | \cdot) = \Pr(s_t | \cdot)$ \longrightarrow Probability of observing
 $S_t = s$ according to $\Pr(S_t | \cdot)$

Probability $\in [0,1]$

Filtering (*forward*)

“Given the available history of observations, what’s the belief about the current hidden state?”

$$\Pr(S_t | o_{1:t}) = \hat{p}_t$$

$$\begin{aligned} \Pr(S_t | o_{1:t}) &= \Pr(S_t | o_t, o_{1:t-1}) \\ &\propto \Pr(o_t | S_t, o_{1:t-1}) \Pr(S_t | o_{1:t-1}) && \text{Bayes} \\ &= \Pr(o_t | S_t) \Pr(S_t | o_{1:t-1}) && \text{Obs. model} \end{aligned}$$

$$\begin{aligned} \Pr(S_t | o_{1:t-1}) &= \sum_{i=1}^d \Pr(S_t | S_{t-1} = i, o_{1:t-1}) \Pr(S_{t-1} = i | o_{1:t-1}) && \text{Marg.} \\ &= \sum_{i=1}^d \Pr(S_t | S_{t-1} = i) \underbrace{\Pr(S_{t-1} = i | o_{1:t-1})}_{\text{Recursion!}} && \text{Trans. model} \end{aligned}$$

Filtering

“Given the available history of observations, what’s the belief about the current hidden state?”

$$\Pr(S_t | o_{1:t}) = \hat{p}_t$$

1. One-step prediction:

$$\Pr(S_t | o_{1:t-1}) = \bar{p}_t = \sum_{i=1}^d \Pr(S_t | S_{t-1} = i) \Pr(S_{t-1} = i | o_{1:t-1}) = T \hat{p}_{t-1}$$

2. Measurement update:

$$\hat{p}_t[i] = \eta \Pr(o_t | S_t = i) \bar{p}_t[i]$$

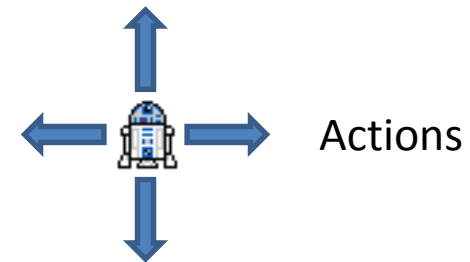
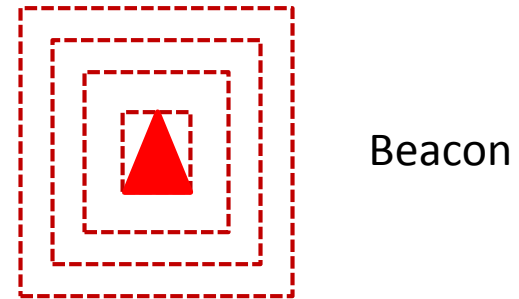
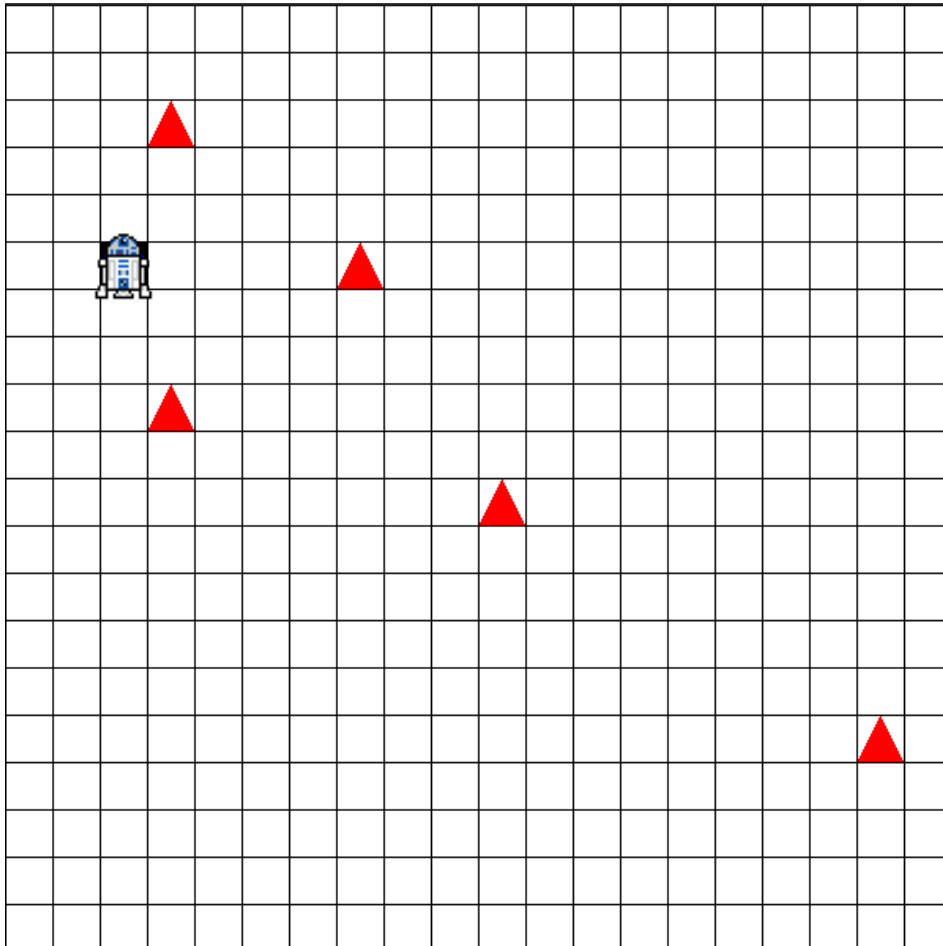
3. Normalize belief (to get rid of η):

$$\hat{p}_t[i] \leftarrow \frac{\hat{p}_t[i]}{\eta}, \eta = \sum_{j=1}^d \hat{p}_t[j]$$



MIT
AEROASTRO

Grid World



Prediction example

0	0	0	0
0	0.6	0.2 ^s	0
0	0.1	0.1	0
0	0	0	0

Robot model

“If told to perform an action, the robot will execute it with probability 90% or do nothing with probability 10%.”

Action = “Move right!” 

What is the probability of being in the red square?

$$\Pr(S_{t+1} = s | a_t = \text{right})?$$

$$\begin{aligned} \Pr(S_{t+1} = s | a_t = \text{right}) &= \Pr(S_{t+1} = s | S_t = s', a_t = \text{right}) \times \Pr(S_t = s') \\ &+ \Pr(S_{t+1} = s | S_t = s'', a_t = \text{right}) \times \Pr(S_t = s'') \\ &+ \Pr(S_{t+1} = s | S_t = s', a_t = \text{right}) \times \Pr(S_t = s') \\ &+ \Pr(S_{t+1} = s | S_t = s'', a_t = \text{right}) \times \Pr(S_t = s'') \end{aligned}$$

Prediction example

0	0	0	0
0	0.6	0.2 ^s	0
0	0.1	0.1	0
0	0	0	0

Robot model

“If told to perform an action, the robot will execute it with probability 90% or do nothing with probability 10%.”

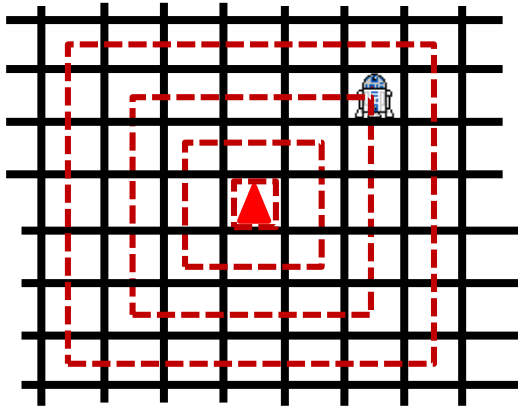
Action = “Move right!” 

What is the probability of being in the red square?

$$\Pr(S_{t+1} = s | a_t = \text{right})?$$

$$\begin{aligned} \Pr(S_{t+1} = s | a_t = \text{right}) &= 0.9 \times 0.6 &= 0.56 \\ &+ 0.0 \times 0.1 \\ &+ 0.1 \times 0.2 \\ &+ 0.0 \times 0.1 \end{aligned}$$

Sensor update



“What is my confidence about being in state s at time $t+1$, given that I took action a at time t and received observation o from beacon b at time $t+1$?”

$$\Pr(S_{t+1} = s | a_t = a, o_{t+1} = o_b)?$$


$$\Pr(S_{t+1} = s | a_t = a, o_{t+1} = o_b) \propto \Pr(o_{t+1} = o_b | S_{t+1} = s) \Pr(S_{t+1} = s | a_t = a)$$

→ Compute unnormalized probabilities.


$$\Pr(S_{t+1} = s | a_t = a, o_{t+1} = o_b) \leftarrow \frac{\Pr(S_{t+1} = s | a_t = a, o_{t+1} = o_b)}{\sum_{s'} \Pr(S_{t+1} = s' | a_t = a, o_{t+1} = o_b)}$$

→ Normalize their sum to 1.

Sensor update example

0	0	0	0
0	0.06	0.56 ^s	0.18
0	0.01	0.1	0.09
0	0	0	

Prediction example. Go check it!

Action = "Move right!" 


Observation = "You're 2 squares away."

Beacon Model

- If at the beacon's location, returns 0 with probability 95%;
- If 1 square away, returns '1' with probability 70%;
- If 2 squares away, returns '2' with probability 80%;
- If 3 squares away, returns '3' with probability 60%;
- In all other cases, the beacon returns no reading ('-').

$$\Pr(S_{t+1} = s | a_t = \text{right}, o'_{t+1} = 2) \propto \Pr(o_{t+1} = 2 | s_{t+1} = s) \Pr(S_{t+1} = s | a_t = \text{right})$$

Sensor update example

0	0	0	0
0	0.06	0.56 ^s	0.18
0	0.01	0.1	0.09
0	0	0	

Prediction example. Go check it!

Action = "Move right!" 


Observation = "You're 2 squares away."

Beacon Model

- If at the beacon's location, returns 0 with probability 95%;
- If 1 square away, returns '1' with probability 70%;
- If 2 squares away, returns '2' with probability 80%;
- If 3 squares away, returns '3' with probability 60%;
- In all other cases, the beacon returns no reading ('-').

$$\begin{aligned}
 \Pr(S_{t+1} = s | a_t = \text{right}, o'_{t+1} = 2) &\propto \Pr(o_{t+1} = 2 | S_{t+1} = s) \Pr(S_{t+1} = s | a_t = \text{right}) \\
 &= 0.8 \times 0.56 \\
 &= 0.448
 \end{aligned}$$

Sensor update example

0	0	0	0
0	0.048	0.448 ^S	0.144
0	0.008	0.0	0.0
0	0	0	

Action = "Move right!" 


Observation = "You're 2 squares away."

Beacon Model

- If at the beacon's location, returns 0 with probability 95%;
- If 1 square away, returns '1' with probability 70%;
- If 2 squares away, returns '2' with probability 80%;
- If 3 squares away, returns '3' with probability 60%;
- In all other cases, the beacon returns no reading ('-').

Normalize the probabilities on the grid!

Sensor update example

0	0	0	0
0	0.074	0.691 ^S	0.222
0	0.013	0.0	0.0
0	0	0	

Action = "Move right!" 

Observation = "You're 2 squares away."

Beacon Model

- If at the beacon's location, returns 0 with probability 95%;
- If 1 square away, returns '1' with probability 70%;
- If 2 squares away, returns '2' with probability 80%;
- If 3 squares away, returns '3' with probability 60%;
- In all other cases, the beacon returns no reading ('-').

Done with belief state updates!



MIT
AEROASTRO

3. RAO*

“It’s a trap!”

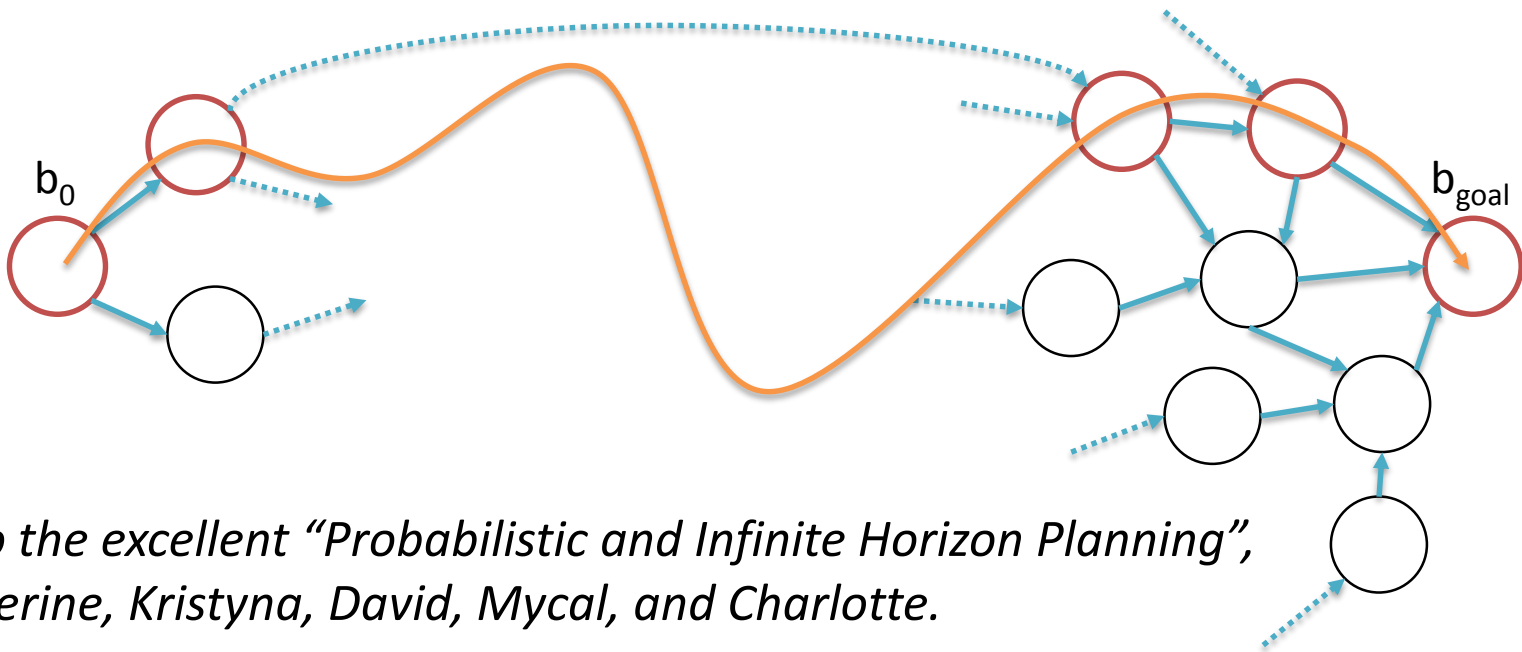
— RAO*, after determining that a policy was too risky.

RAO* = AO* + Belief states + Execution risk

“Probability of violating constraints c during execution.” $\leq \Delta \quad \equiv \quad er(b_0, c|\pi) \leq \Delta$

Key 1: (admissible) value heuristic guiding search towards “promising” policies;

Key 2: (admissible) *execution risk* heuristic allowing risk bounds to be propagated forward.



See also the excellent “Probabilistic and Infinite Horizon Planning”, by Katherine, Kristyna, David, Mycal, and Charlotte.

Elements of a CC-POMDP

S : discrete set of states.

A : discrete set of actions.

\mathcal{O} : discrete set of observations.

$T: S \times A \times S \rightarrow [0,1]$, transition function

$O: S \times \mathcal{O} \rightarrow [0,1]$, observation function

$R: S \times A \rightarrow \mathbb{R}$, reward function.

\mathcal{C} : set of state constraints.

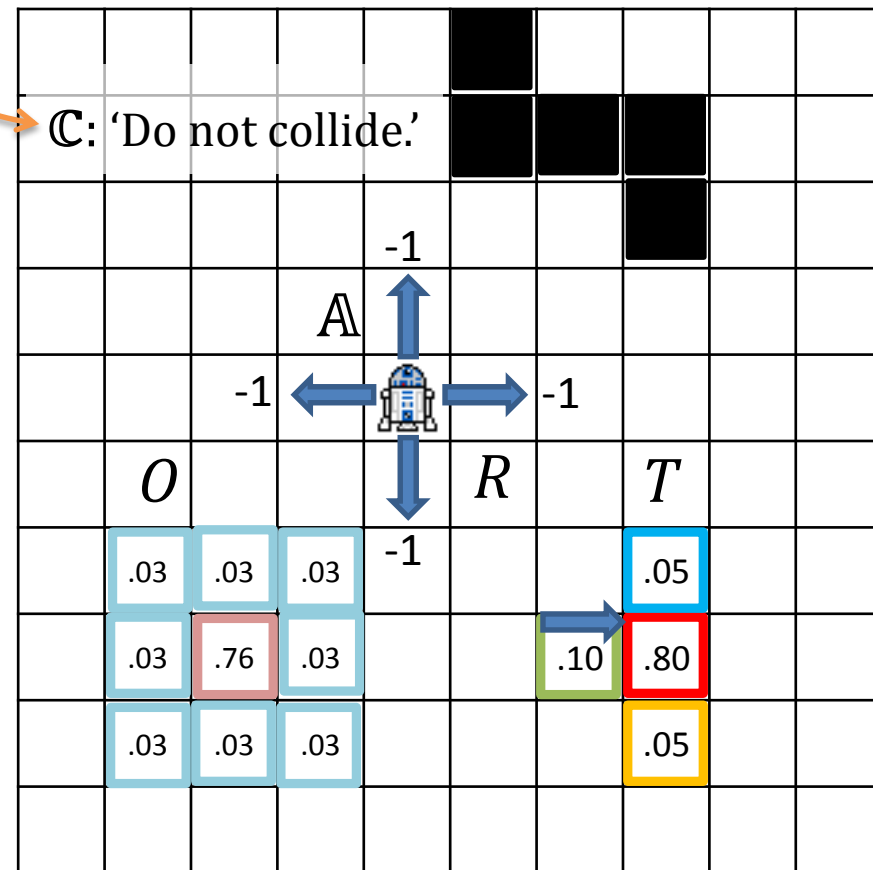
Δ : risk bound.

$$T(s_k, a_k, s_{k+1}) = \Pr(s_{k+1} | s_k, a_k)$$

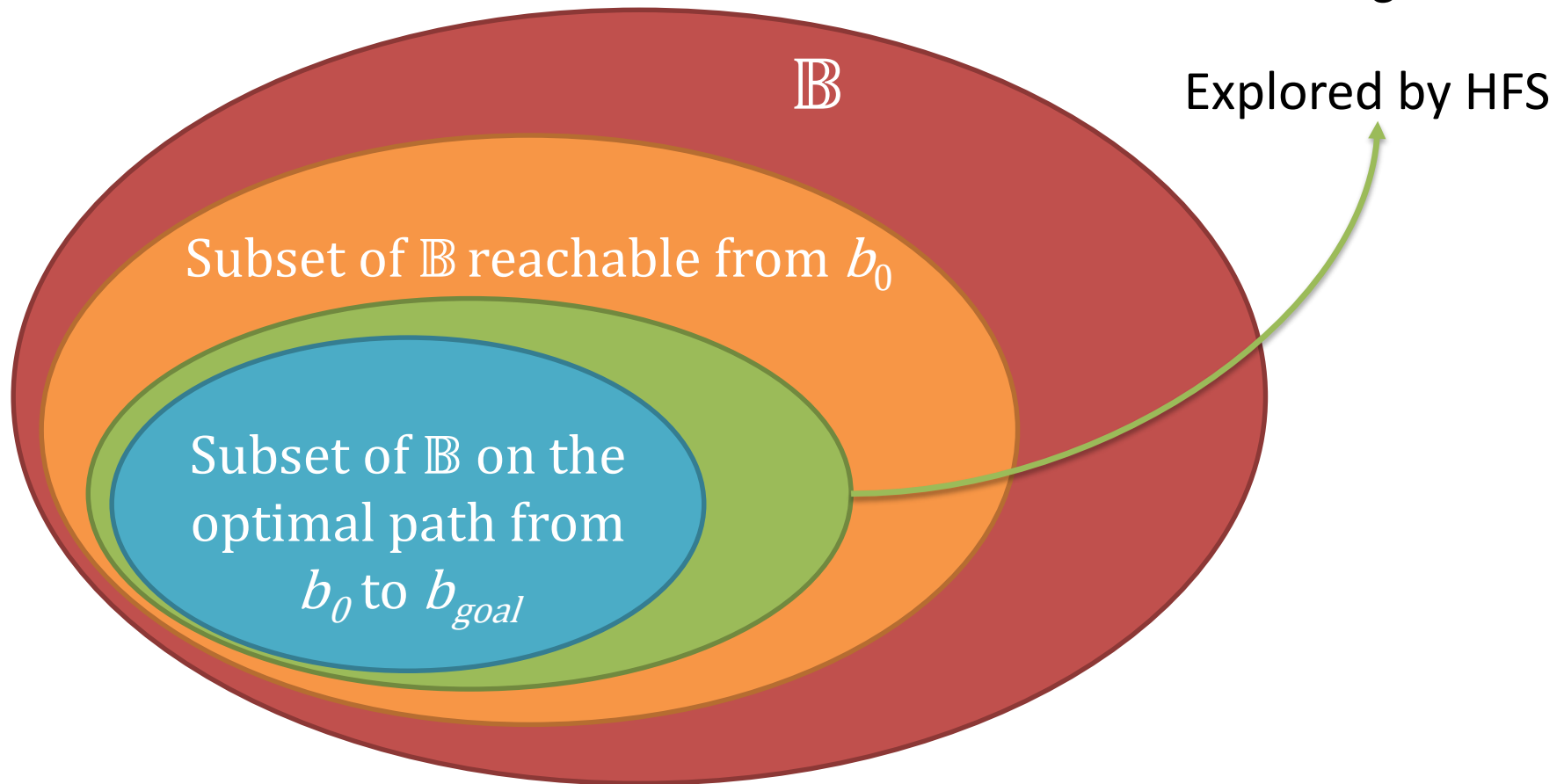
$$O(s_k, o_k) = \Pr(o_k | s_k)$$

$\Delta=0.01 \rightarrow$ 'Collision probability must be less than 1%.'

S, \mathcal{O}



Searching from an initial belief b_0

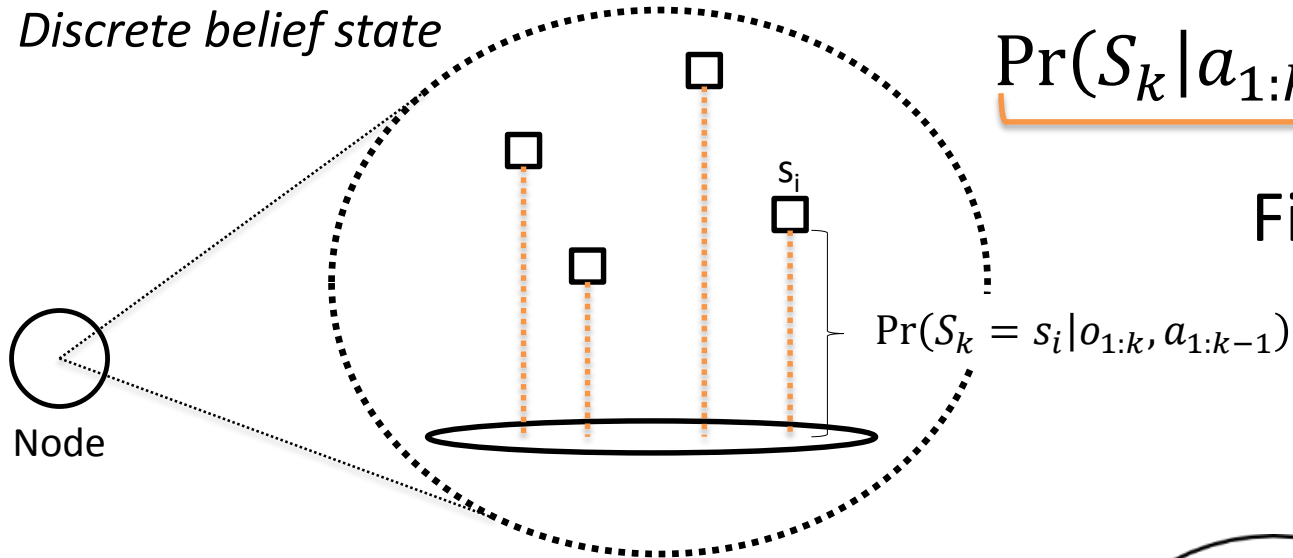


Good heuristic:  \approx 

Bad heuristic:  \approx 

RAO* nodes are belief states

Discrete belief state



$$\Pr(S_k | a_{1:k-1}, o_{1:k}) = \hat{p}_k$$

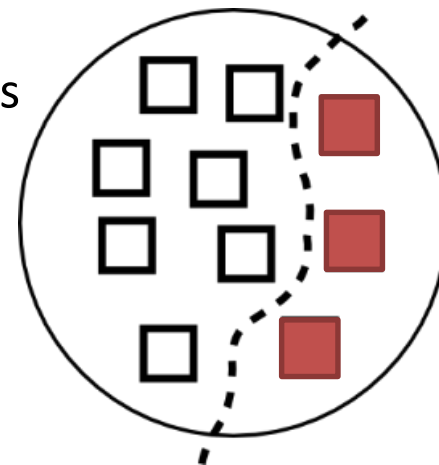
Filtering

Safe states

State is "safe"

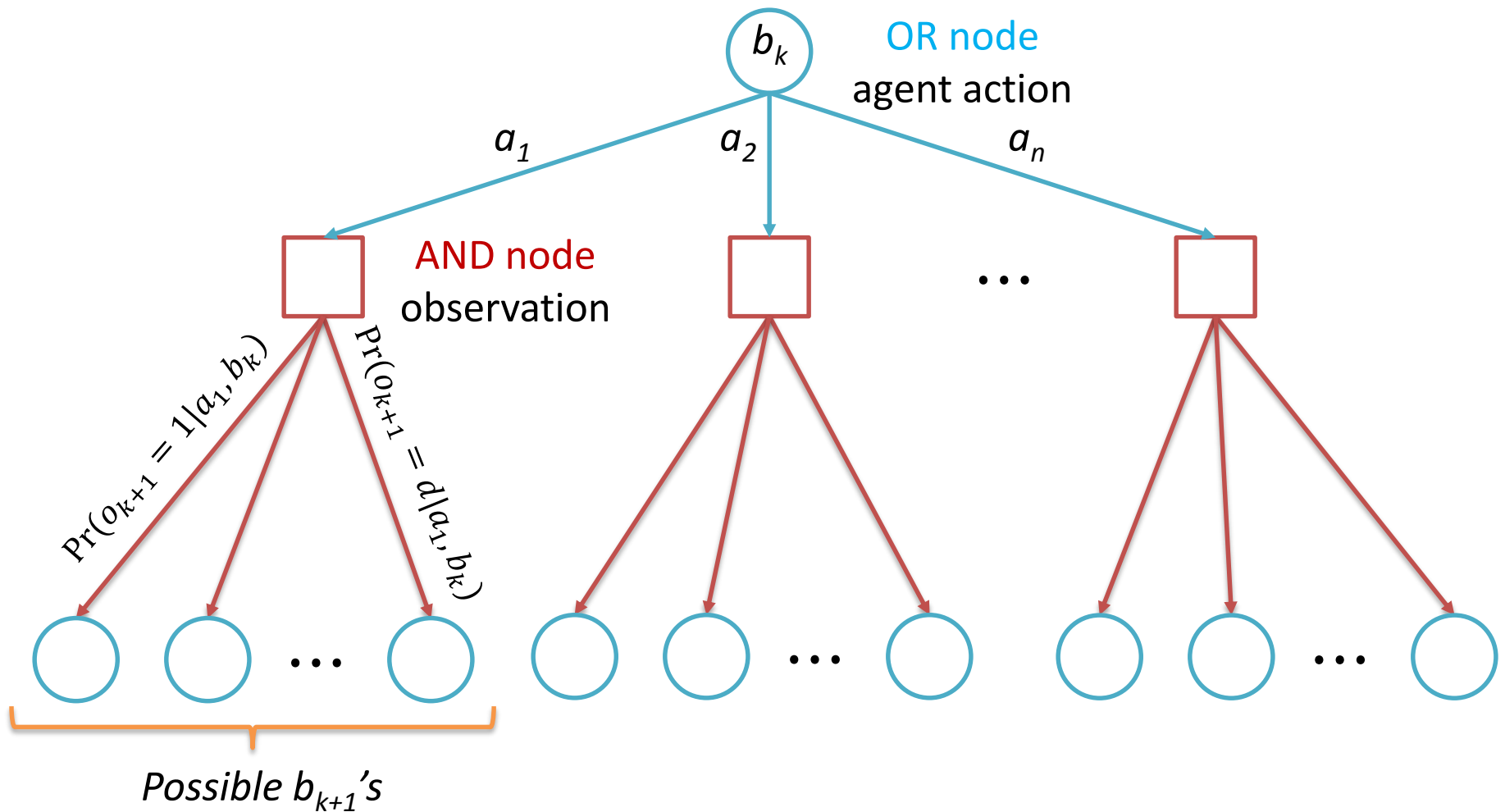


itself and its ancestors do not violate constraints



Unsafe states

Partially observable AND-OR search

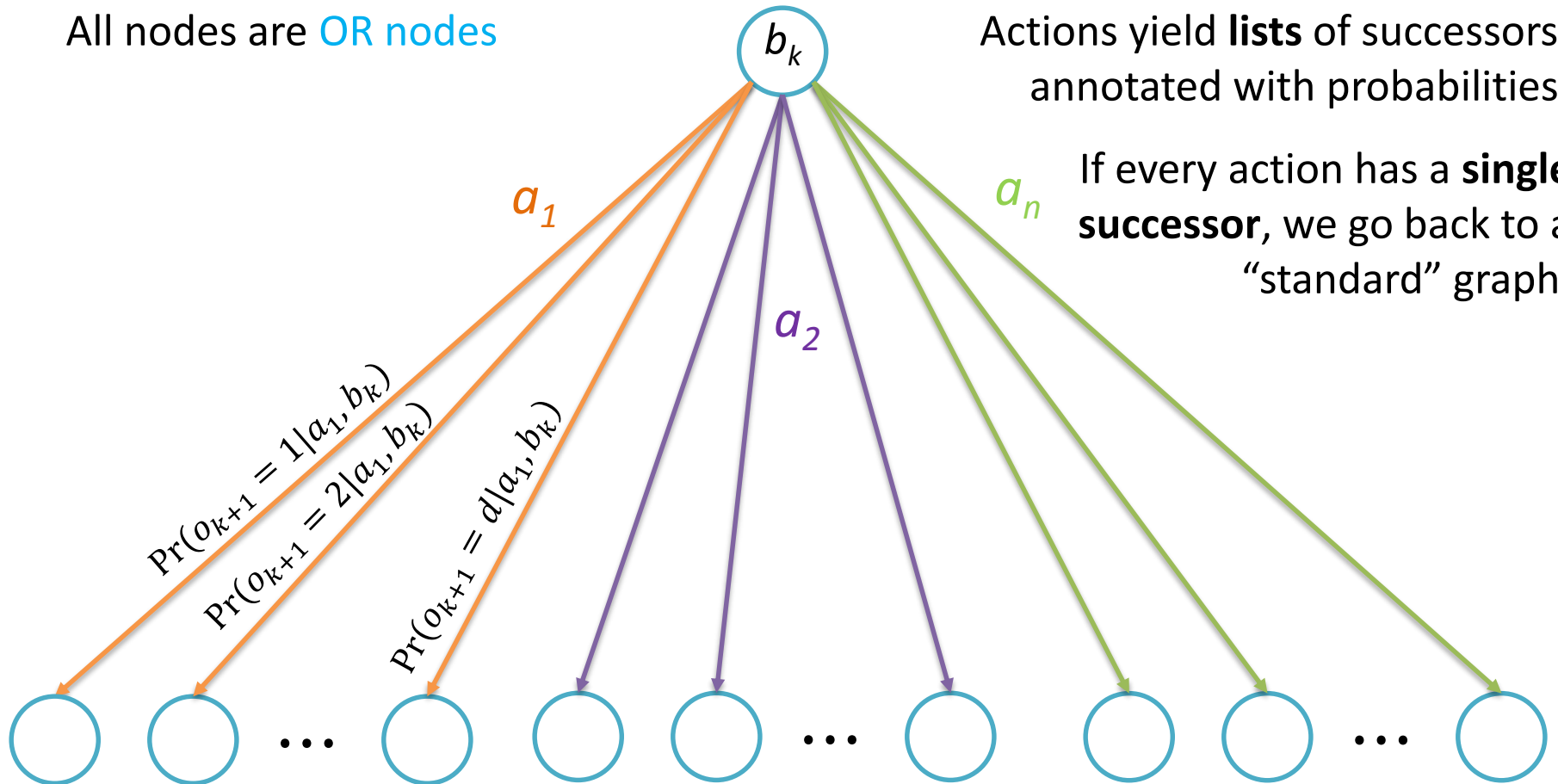


Hypergraph representation

All nodes are **OR nodes**

Actions yield **lists** of successors annotated with probabilities

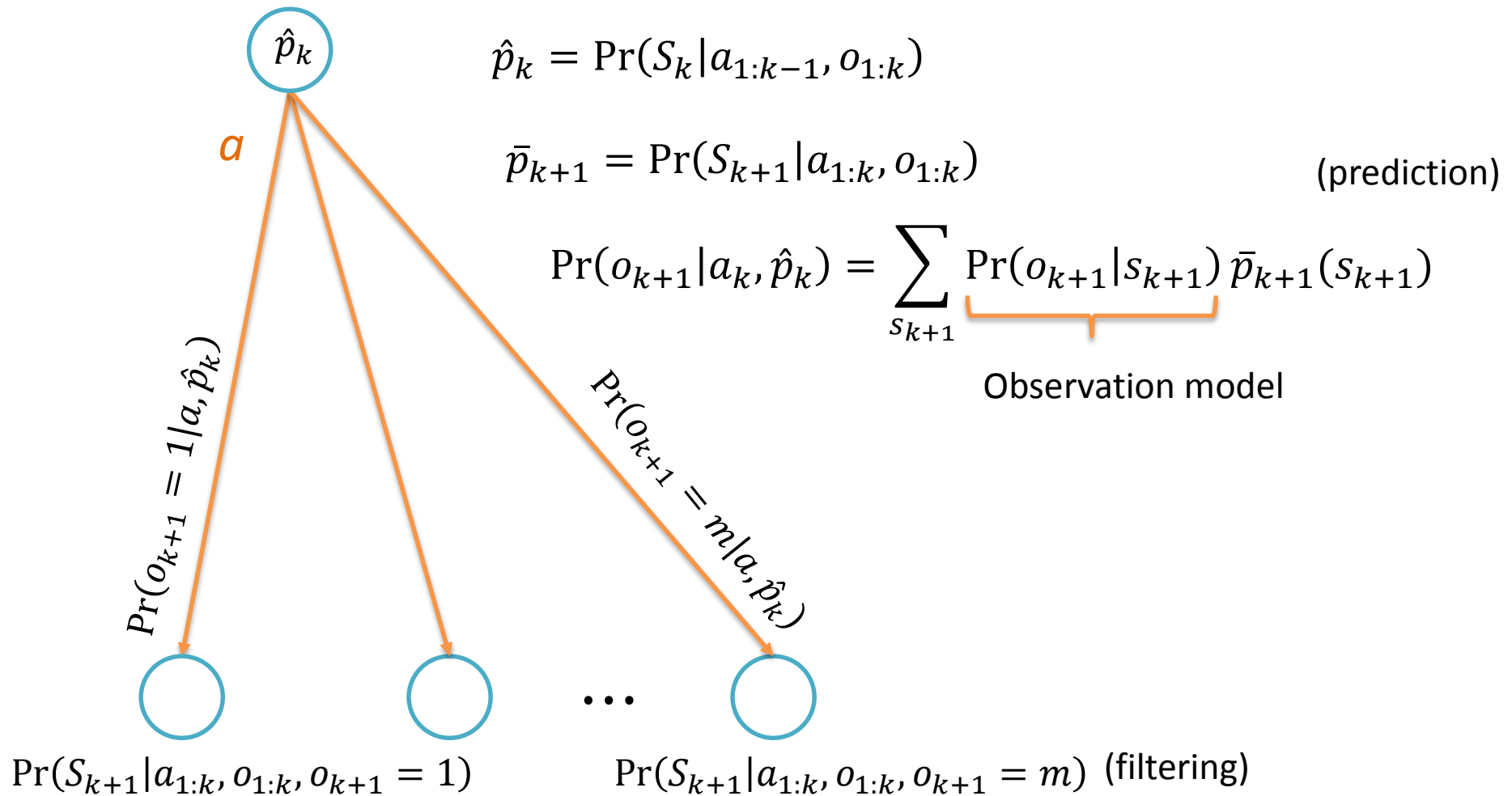
If every action has a **single successor**, we go back to a “standard” graph.



Key concepts

1. How to compute transition probabilities
2. How to compute (admissible) utility estimates
3. How to compute (admissible) execution risk estimates
4. How to put them all together

3.1 Hyperedge probabilities



Computing hyperedge probabilities

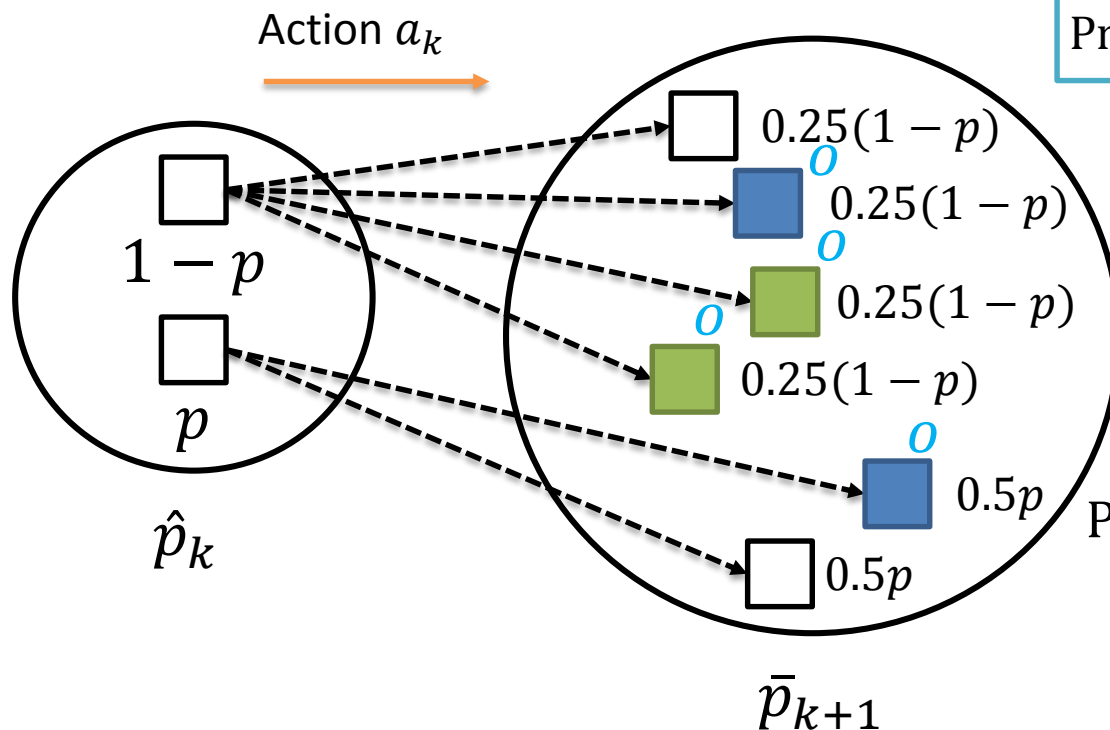
$$\bar{p}_{k+1} = \Pr(S_{k+1}|a_{1:k}, o_{1:k}) = T(a_k)\hat{p}_k \quad \Pr(o_{k+1}|a_k, \hat{p}_k) = \sum_{s_{k+1}} \Pr(o_{k+1}|s_{k+1}) \bar{p}_{k+1}(s_{k+1})$$

$$\Pr(o|\blacksquare) = 0.8$$

$$\Pr(-|\blacksquare) = 0.2$$

$$\Pr(o|\blacksquare) = 0.6$$

$$\Pr(\times|\blacksquare) = 0.4$$



$$\Pr(o|a_k, \hat{p}_k)?$$

$$\Pr(o|a_k, \hat{p}_k) = 0.25(1-p) \times 0.8$$

$$+ 0.25(1-p) \times 0.6$$

$$+ 0.25(1-p) \times 0.6$$

$$+ 0.5p \times 0.8$$

3.2 Admissible utility estimates

$$\underbrace{V(s_k, a_k)}_{\substack{\text{Expected reward} \\ \text{of executing } a_k \\ \text{at } s_k}} = \underbrace{R(s_k, a_k)}_{\substack{\text{Immediate} \\ \text{reward}}} + \sum_{s_{k+1}} \underbrace{V^*(s_{k+1})}_{\substack{\text{Optimal reward} \\ \text{at } s_{k+1}}} T(s_k, a_k, s_{k+1})$$

Unknown \rightarrow

Expected optimal future reward

$h(s_{k+1}) \geq V^*(s_{k+1})$

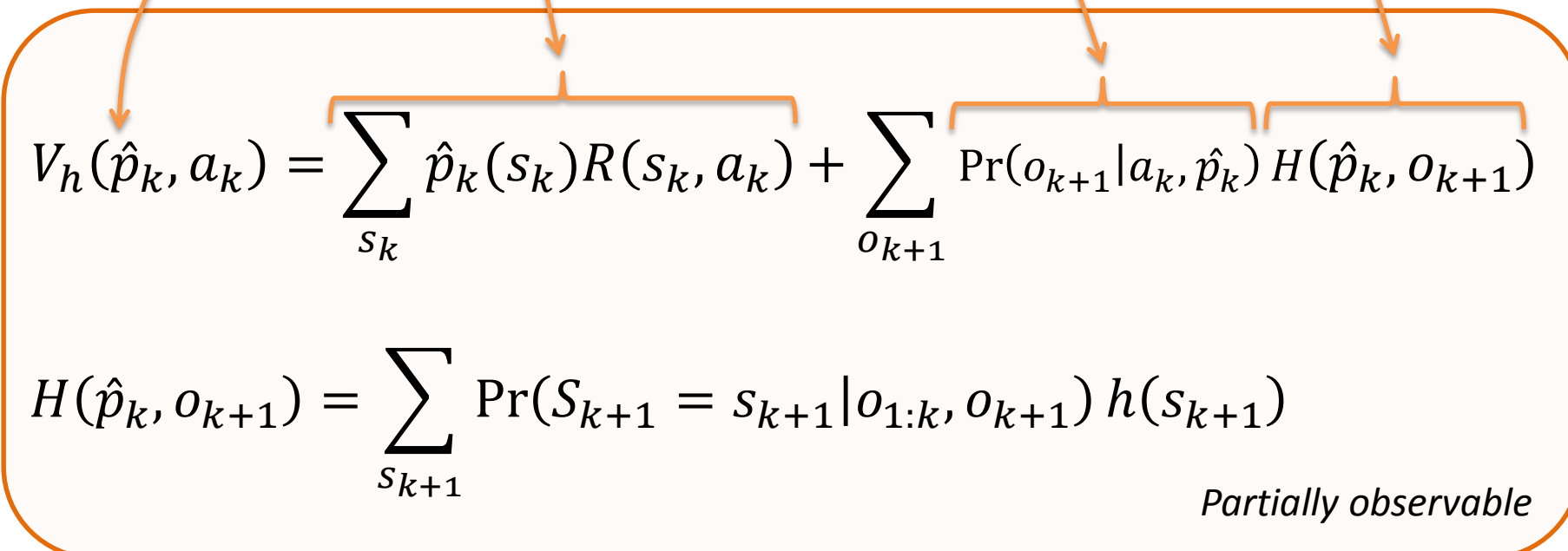
- \rightarrow Admissible (“optimistic”) estimate of future reward.
- \rightarrow Should be “easy” to compute.

$$V_h(s_k, a_k) = R(s_k, a_k) + \sum_{s_{k+1}} h(s_{k+1}) T(s_k, a_k, s_{k+1}) \geq V(s_k, a_k)$$

Estimating belief state utility

Fully observable

$$V_h(s_k, a_k) = R(s_k, a_k) + \sum_{s_{k+1}} T(s_k, a_k, s_{k+1}) h(s_{k+1})$$



$$V_h(\hat{p}_k, a_k) = \sum_{s_k} \hat{p}_k(s_k) R(s_k, a_k) + \sum_{o_{k+1}} \Pr(o_{k+1} | a_k, \hat{p}_k) H(\hat{p}_k, o_{k+1})$$

$$H(\hat{p}_k, o_{k+1}) = \sum_{s_{k+1}} \Pr(S_{k+1} = s_{k+1} | o_{1:k}, o_{k+1}) h(s_{k+1})$$

Partially observable

3.3 Execution risk

$Sa_i = 1$: agent hasn't violated \mathbb{c} until i -th step

$$er(b_k, \mathbb{c}|\pi) = 1 - \Pr \left(\bigwedge_{i=k}^T Sa_i | b_k, \pi \right) \quad (\text{Ono et al., 2012})$$

Probability of violating constraints from k onwards.

Probability of remaining safe from k onwards.

$$er(b_k|\pi) = r_b(b_k) + (1 - r_b(b_k)) \sum_{o_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) er(b_{k+1}|\pi)$$

Immediate risk at the current belief state.

Observations originated from safe states.

Execution risk

$$er(b_k|\pi) = \underbrace{r_b(b_k)}_{\text{Immediate risk at the current belief state.}} + (1 - r_b(b_k)) \sum_{o_{k+1}} \underbrace{Pr^{sa}(o_{k+1}|\pi(b_k), b_k)}_{\text{Observations originated from safe states.}} er(b_{k+1}|\pi)$$

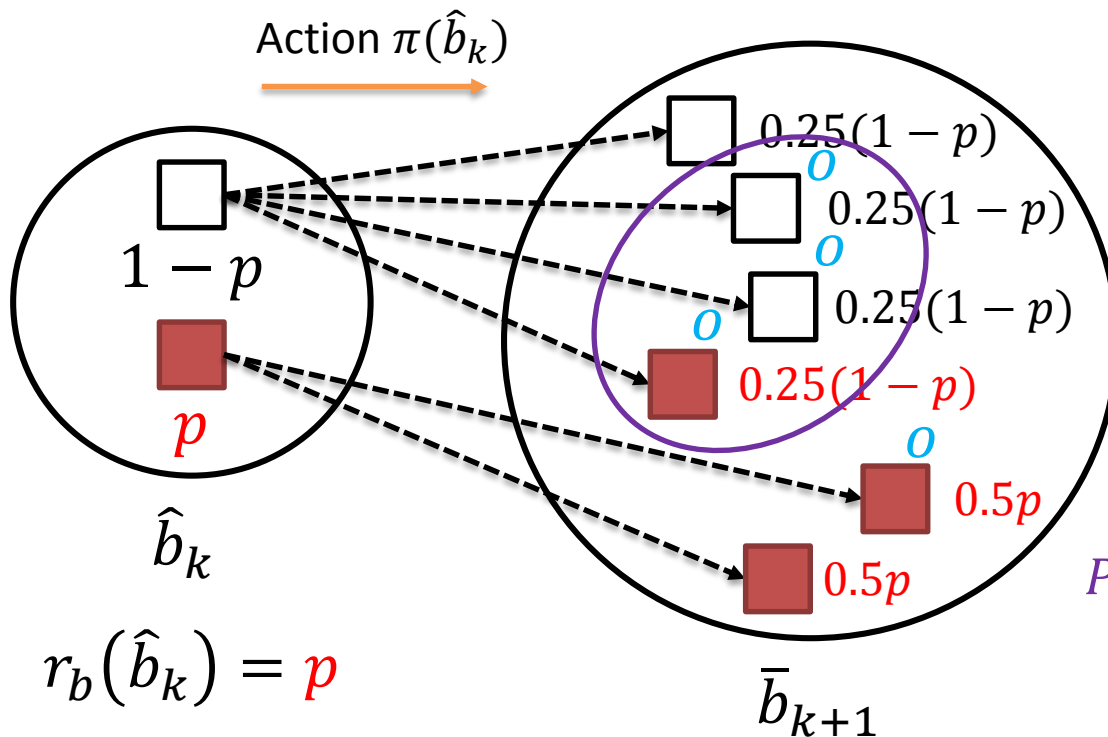
$$r_b(b_k) = \sum_{s_k \in \mathcal{S}} b(s_k) \underbrace{c_v(p(s_k), \mathbb{C})}_{\substack{\text{1 if the sequence of states ("path")} \\ \text{leading to } s_k \text{ violates } \mathbb{C}}}$$

$$\bar{b}^{sa}(s_{k+1}|a_k) = \Pr(s_{k+1}|Sa_k, a_k, b_k) = \frac{\sum_{s_k: c_v(p(s_k), \mathbb{C})=0} T(s_k, a_k, s_{k+1}) b(s_k)}{1 - r_b(b_k)}$$

$$\Pr^{sa}(o_{k+1}|a_k, b_k) = \Pr(o_{k+1}|Sa_k, a_k, b_k) = \sum_{s_{k+1}} O(s_{k+1}, o_{k+1}) \bar{b}^{sa}(s_{k+1}|a_k)$$

Execution risk in pictures

$$er(b_k|\pi) = r_b(b_k) + (1 - r_b(b_k)) \sum_{o_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) er(b_{k+1}|\pi)$$



What about this term?

$$Pr^{sa}(o|\cdot) = \frac{3 \times 0.25(1-p)}{1-p} = 0.75$$

Estimating execution risk

$$er(b_k|\pi) = r_b(b_k) + (1 - r_b(b_k)) \sum_{o_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) er(b_{k+1}|\pi)$$

 Admissible (“optimistic”) estimate of future execution risk

$$h_{er}(b_{k+1}|\pi) \leq er(b_{k+1}|\pi)$$

 Should be “easy” to compute.

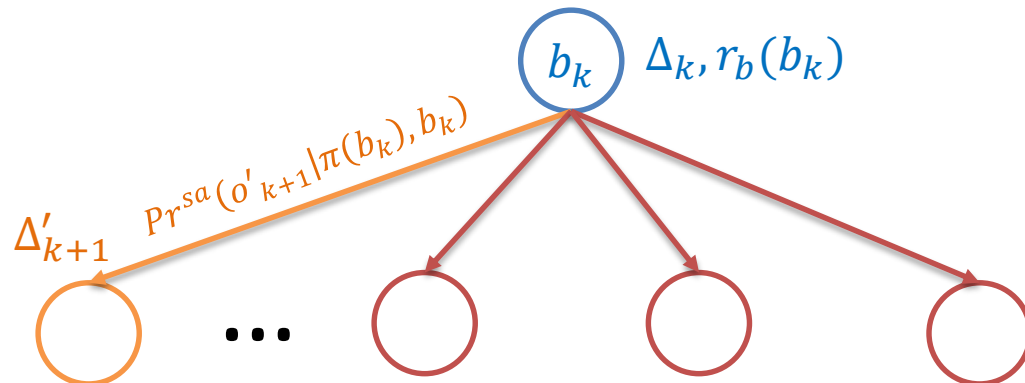
$$er_h(b_k|\pi) = r_b(b_k) + (1 - r_b(b_k)) \sum_{o_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) h_{er}(b_{k+1}|\pi)$$

$h_{er}(b_{k+1}|\pi) = r_b(b_{k+1})$ is always admissible

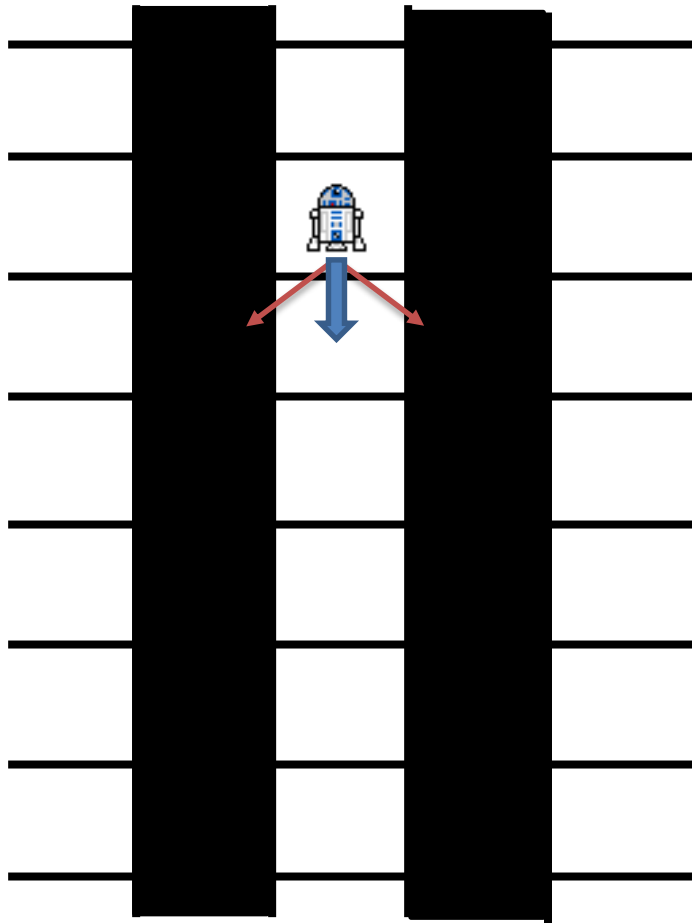
Propagating execution risk estimates

$$er_h(b_k|\pi) = r_b(b_k) + (1 - r_b(b_k)) \sum_{o_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) h_{er}(b_{k+1}|\pi)$$

$$\Delta'_{k+1} = \frac{1}{Pr^{sa}(o'_{k+1}|\pi(b_k), b_k)} \left(\frac{\Delta_k - r_b(b_k)}{1 - r_b(b_k)} - \sum_{o_{k+1} \neq o'_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) h_{er}(b_{k+1}|\pi) \right)$$

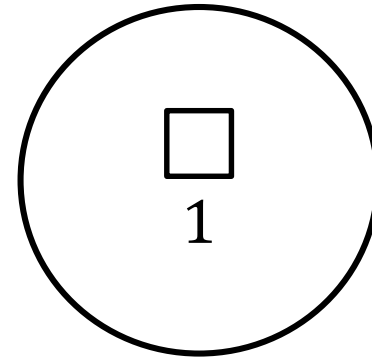


Risk example



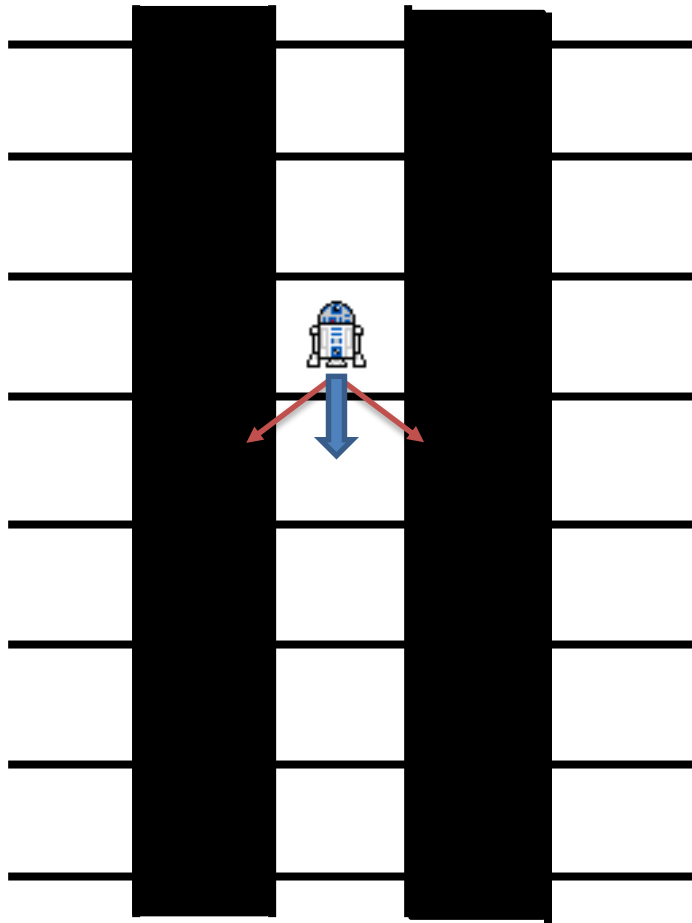
Robot model

“If told to move, R2D2 achieves the desired cell with probability 90%, or slips to either side with probability 5%.”



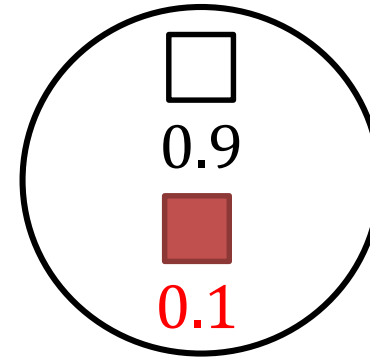


Risk example



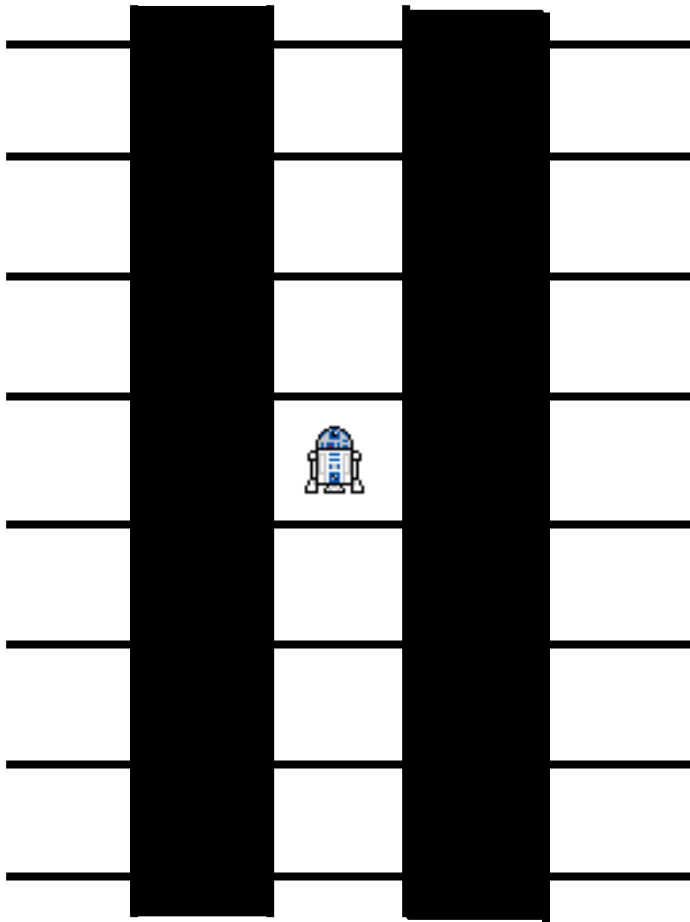
Robot model

“If told to move, R2D2 achieves the desired cell with probability 90%, or slips to either side with probability 5%.”



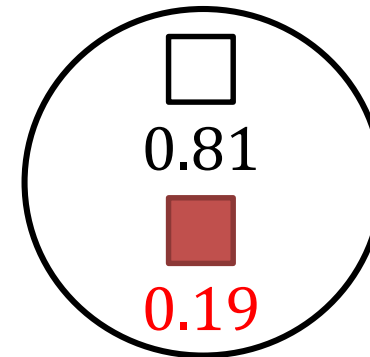
Both particles share the same position

Risk example



Robot model

“If told to move, R2D2 achieves the desired cell with probability 90%, or slips to either side with probability 5%.”



Both particles share the same position

3.4 RAO* in a nutshell

Additions to AO shown in red*

- *Input*: implicit **partially observable** AND-OR search problem
 $\langle S, A, \mathbb{O}, T, O, R, c, \Delta \rangle, b_0$
- *Output*: **optimal** policy in the form of an **acyclic hypergraph** mapping **belief states** to actions.
- *Strategy*: incrementally build solutions forward from b_0 , using h to estimate future utilities (just like A*!) and h_{er} to estimate policy risk. The **set of explored solutions** form the **explicit** hypergraph G , and the subset of G corresponding to the **current estimate** of the **best policy** is called the **greedy** hypergraph g .

RAO*'s pseudocode

Additions to AO shown in red*

Input: $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T, O, R, c, \Delta \rangle, b_0$

Output: Policy $\pi: \mathcal{B} \rightarrow \mathcal{A}$

Explicit graph $G \leftarrow b_0, g \leftarrow$ Best partial policy of G

while best partial policy graph g has nonterminal leafs

$m \leftarrow$ Expand any nonterminal leaf from g

Add to G the children in m which do not violate risk bound

$Z \leftarrow$ set containing m and all of its predecessors that are part of g

while Z is not empty

$n \leftarrow$ Remove from Z a node with no descendants in Z

Update utility and **execution risk** for n

$\pi \leftarrow$ Choose best action at n **not violating risk bound**

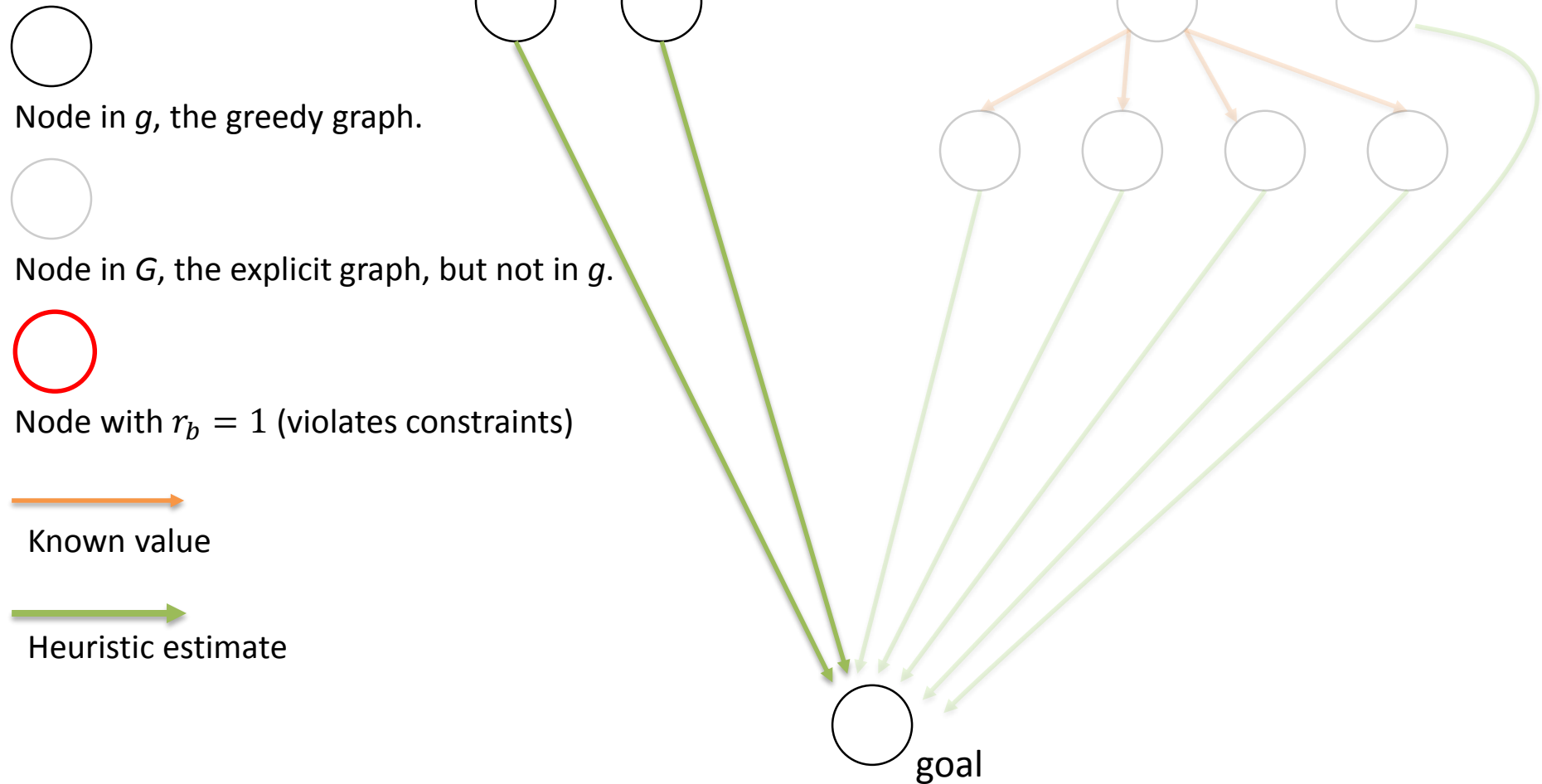
Update g with the new π

Heuristics for utility and **execution risk**

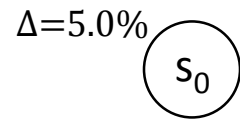
g is the graph obtained by following π from b_0

Bellman backups

RAO* example



Start



Open nodes: $[s_0]$

g = $\{s_0: \text{None}\}$

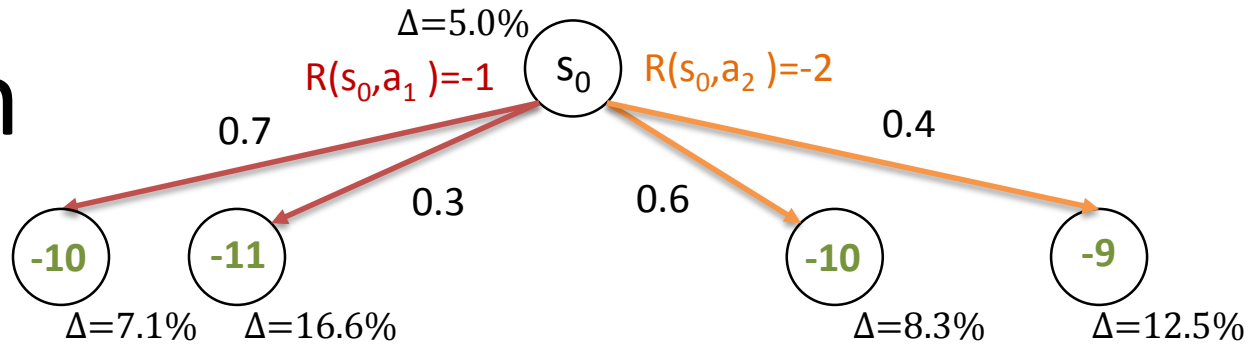
G starts just with just the initial state s_0



Expansion

Open nodes: $[s_0]$

$g = \{s_0: \text{None}\}$



1. Choose an open node to expand $\rightarrow s_0$

2. Estimate the **value** and **execution risk** of leaf nodes

$$h_{er} = r_b = 0 \text{ for all leaves}$$

3. **Propagate risk bounds**

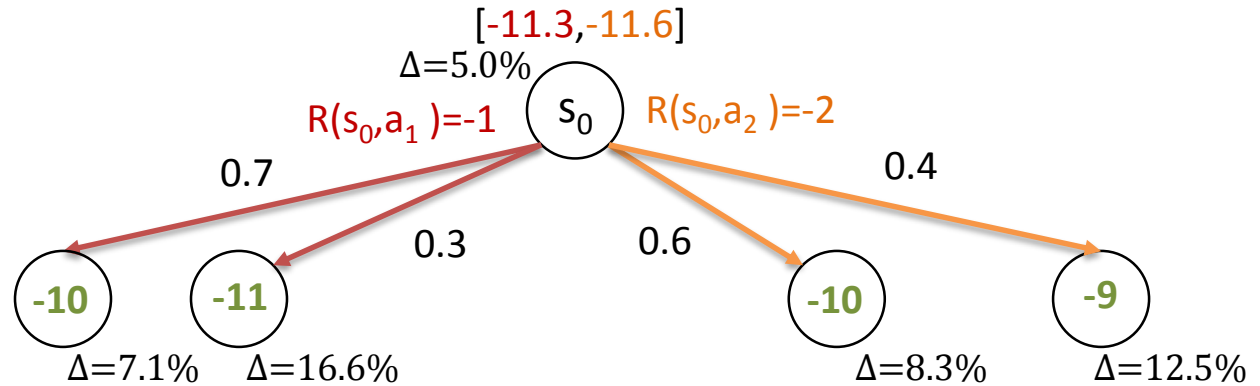
$$\Delta'_{k+1} = \frac{1}{Pr^{sa}(o'_{k+1} | \pi(b_k), b_k)} \left(\frac{\Delta_k - r_b(b_k)}{1 - r_b(b_k)} - \sum_{o_{k+1} \neq o'_{k+1}} Pr^{sa}(o_{k+1} | \pi(b_k), b_k) h_{er}(b_{k+1} | \pi) \right)$$

$$= \frac{1}{0.6} \left(\frac{0.05 - 0}{1 - 0} - 0.4 \times 0 \right) = 8.3\%$$

Backup

Open nodes: $[s_0]$

$g = \{s_0: \text{None}\}$



3. Backup value and *execution risk* for the currently expanded node (s_0) and all its ancestors that are part of g (no ancestors), recording the best value at each node.

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 11 * 0.3) = -11.3$$

$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$

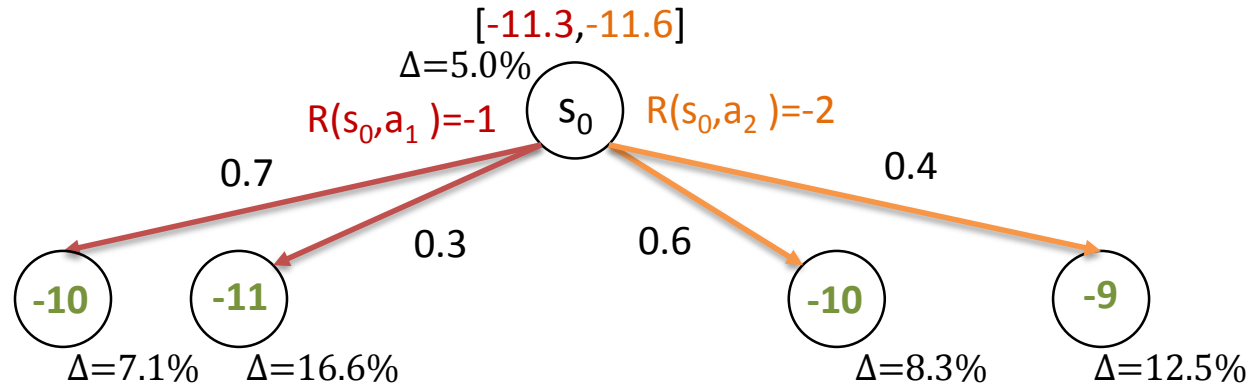
$$er(s_0, a_1) = er(s_0, a_2) = 0 < 5\%$$



Update g

Open nodes: $[s_0]$

$g = \{s_0: \text{None}\}$



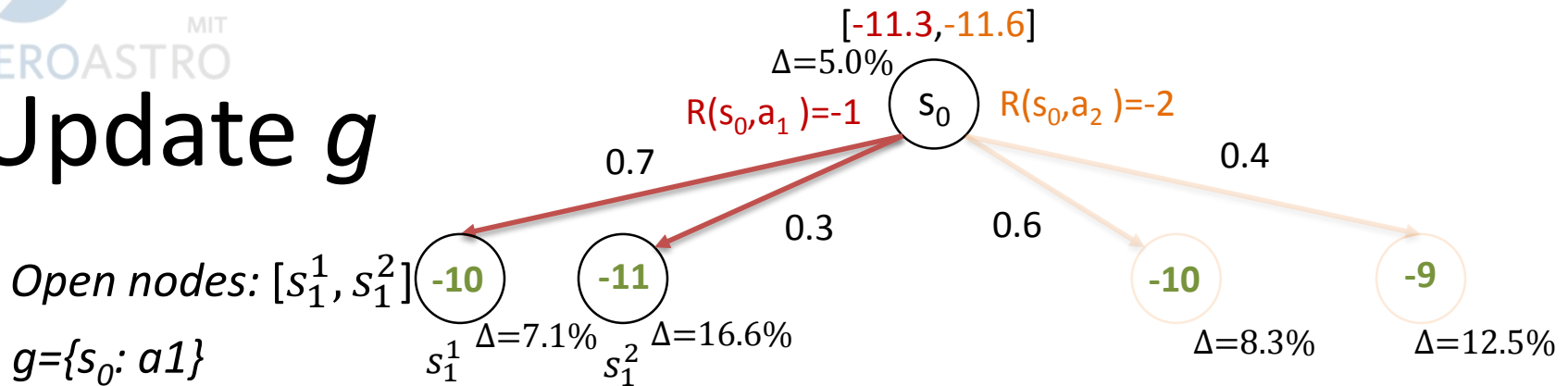
4. Update g and the list of open nodes (non-terminal) by selecting the best action at the nodes which got their values updated

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 11 * 0.3) = -11.3$$

$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$



Update g



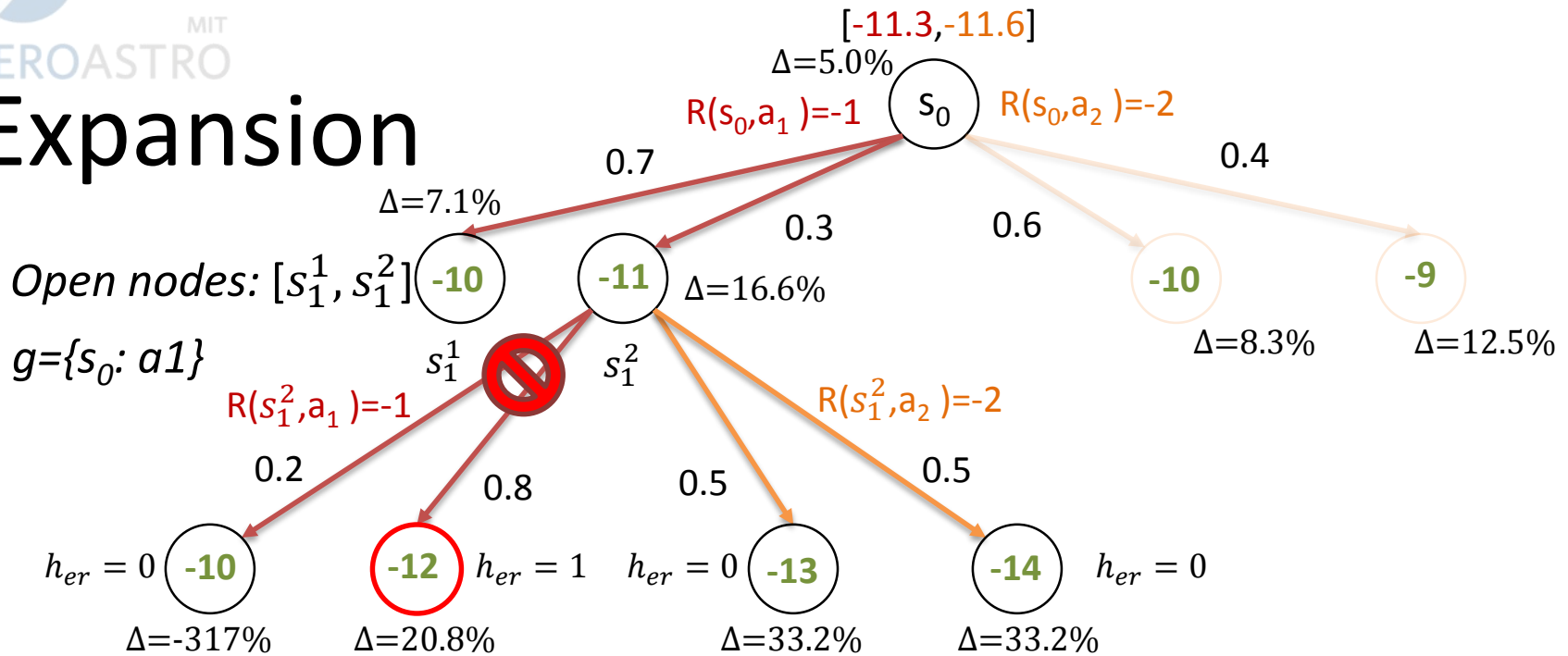
4. Update g and the list of open nodes (non-terminal) by selecting the best action at the nodes which got their values updated

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 11 * 0.3) = -11.3$$

$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$



Expansion



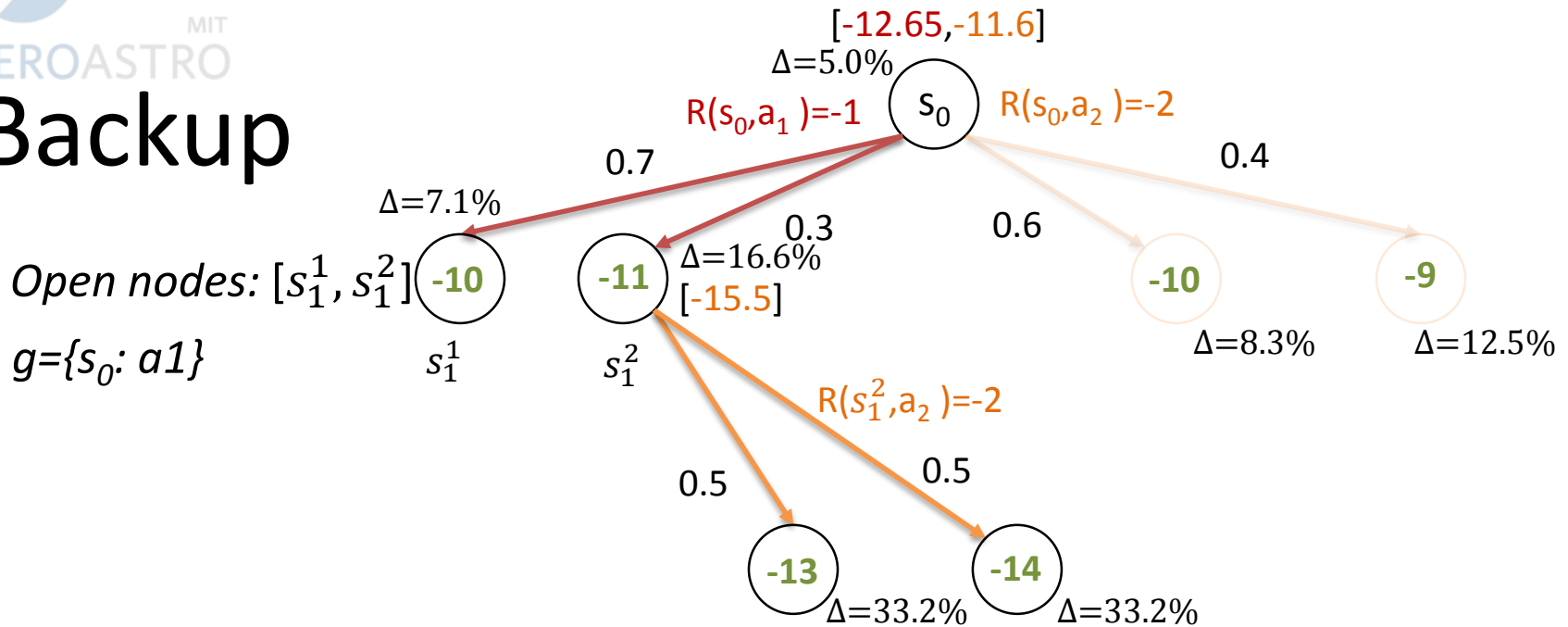
$$\frac{1}{0.2} \left(\frac{0.166 - 0}{1 - 0} - 0.8 \times 1 \right) = -317\%$$

$$\frac{1}{0.8} \left(\frac{0.166 - 0}{1 - 0} - 0.2 \times 0 \right) = 20.8\%$$

$$\Delta'_{k+1} = \frac{1}{Pr^{sa}(o'_{k+1} | \pi(b_k), b_k)} \left(\frac{\Delta_k - r_b(b_k)}{1 - r_b(b_k)} - \sum_{o_{k+1} \neq o'_{k+1}} Pr^{sa}(o_{k+1} | \pi(b_k), b_k) h_{er}(b_{k+1} | \pi) \right)$$

1. Choose any open node to expand $\rightarrow s_1^2$
2. Estimate the **value** and **execution risk** of leaf nodes
3. **Propagate risk bounds**

Backup



3. Backup values and *execution risk* for the currently expanded node (s_1^2) and all its ancestors that are part of g (s_0), recording the best value at each node.

$$V_h(s_1^2, a_2) = -2 + (-13 * 0.5 - 14 * 0.5) = -15.5$$

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 15.5 * 0.3) = -12.65$$

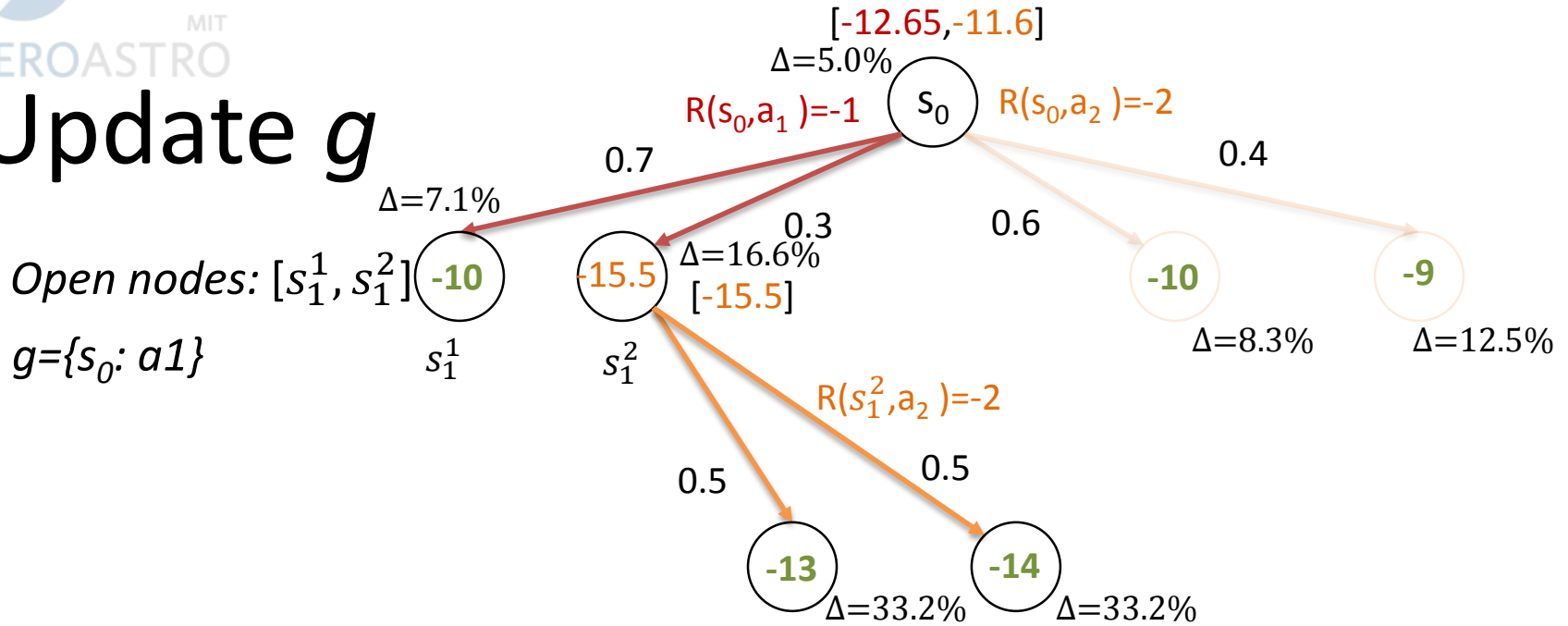
$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$

$$er(s_1^2, a_2) = 0 < 16.6\%, er(s_0, a_2) = 0 < 5\%$$

From leaves to the root



Update g



4. Update g and the list of open nodes (non-terminal) by selecting the best action at the nodes which got their values updated.

$$V_h(s_1^2, a_2) = -2 + (-13 * 0.5 - 14 * 0.5) = -15.5$$

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 15.5 * 0.3) = -12.65$$

$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$

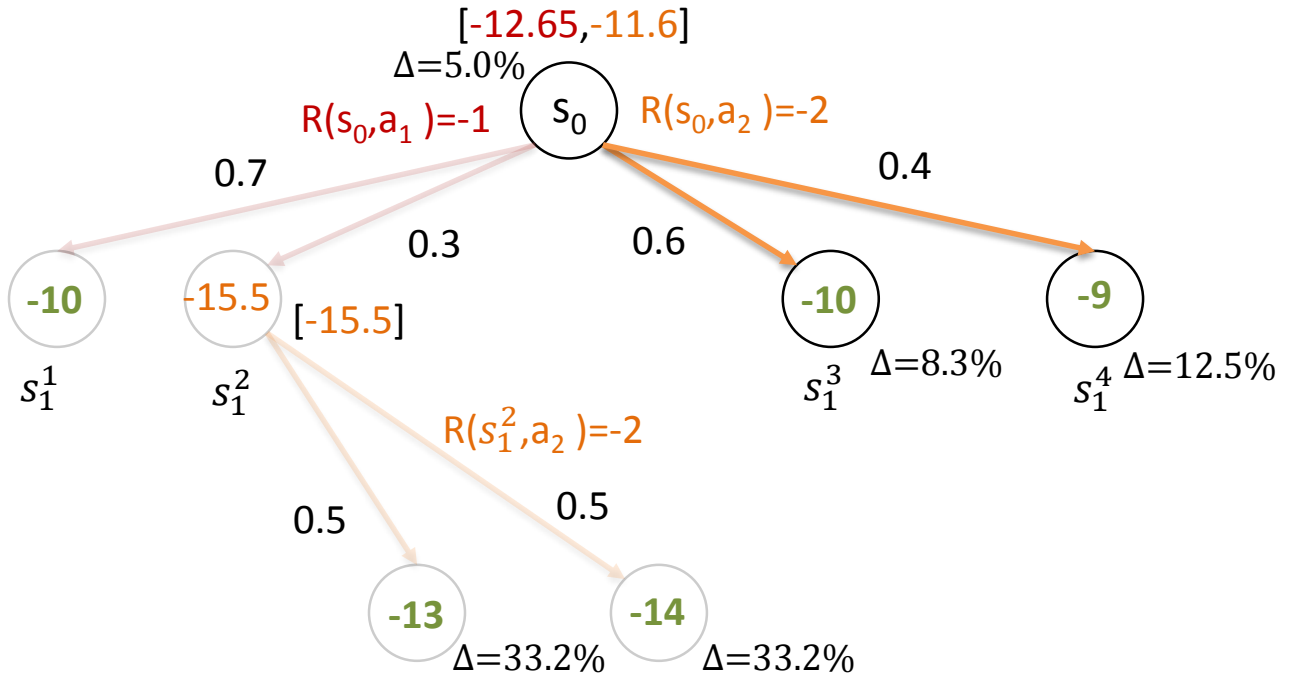
From leaves to the root



Update g

Open nodes: $[s_1^3, s_1^4]$

$g = \{s_0: a_2\}$



4. Update g and the list of open nodes (non-terminal) by selecting the best action at the nodes which got their values updated.

$$V_h(s_1^2, a_2) = -2 + (-13 * 0.5 - 14 * 0.5) = -15.5$$

$$V_h(s_0, a_1) = -1 + (-10 * 0.7 - 15.5 * 0.3) = -12.65$$

$$V_h(s_0, a_2) = -2 + (-10 * 0.6 - 9 * 0.4) = -11.6$$

From leaves to the root

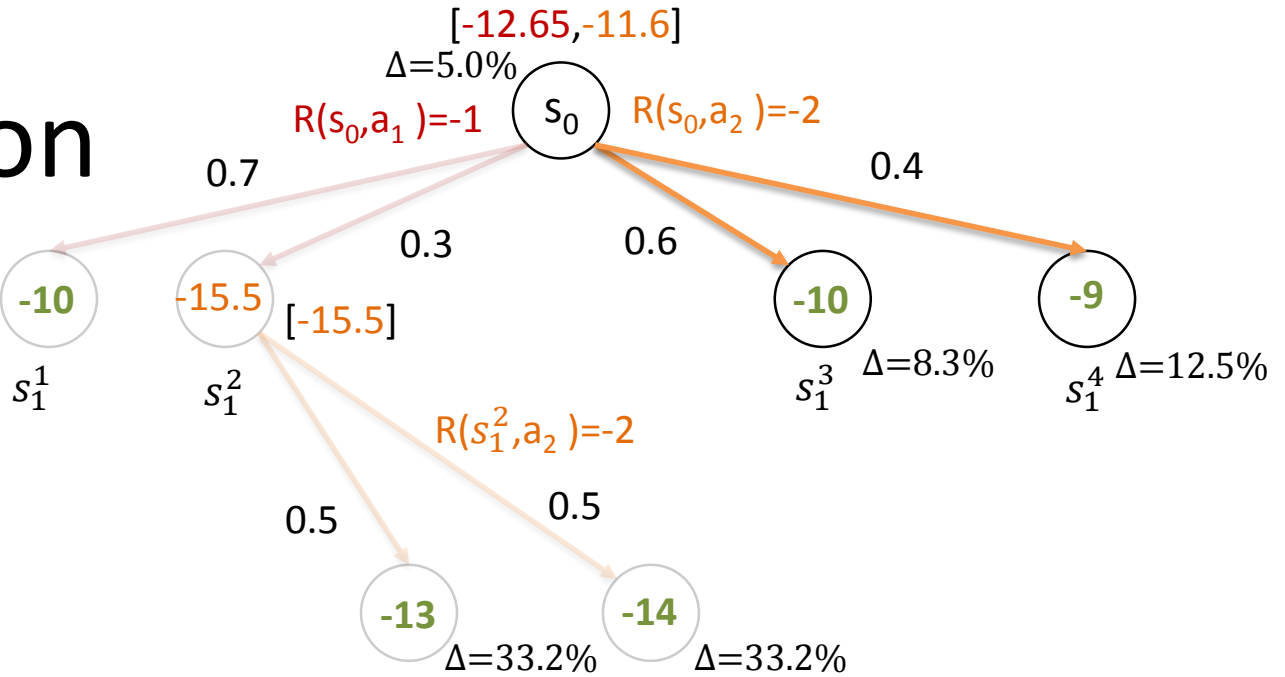




Termination

Open nodes: $[s_1^3, s_1^4]$

$g = \{s_0: a_2\}$

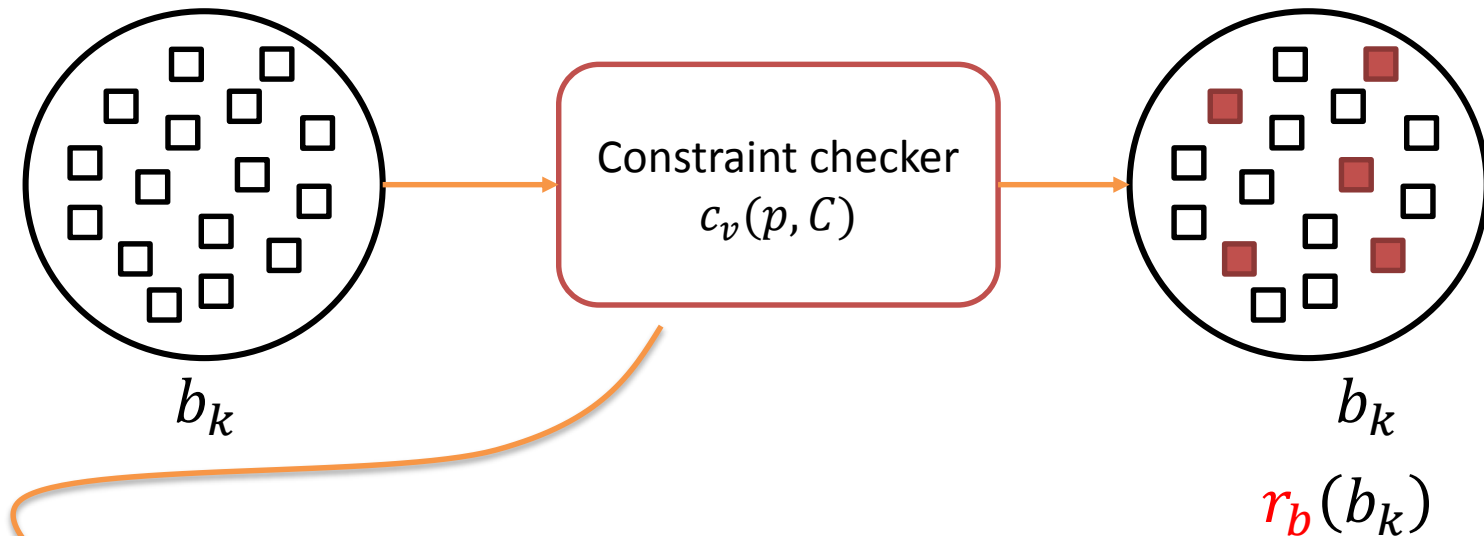


Search terminates when the list of open nodes is empty, i.e., all leaf nodes in g are terminal (goals).



At this point, return g as the **optimal** policy π .

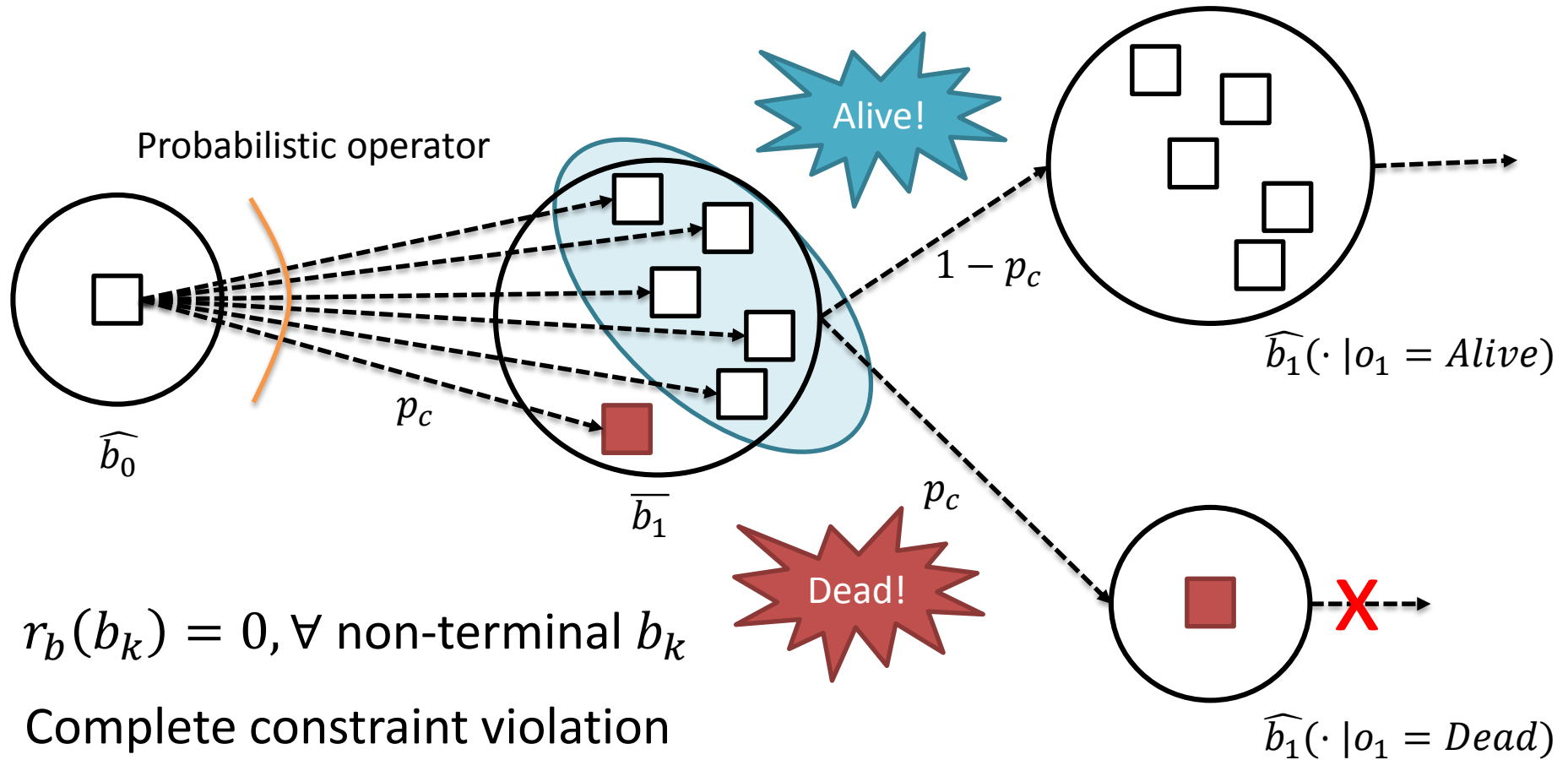
How hard is it to assess risk in RAO*?



Must run *really fast* for RAO* to be useful in practice!

e.g., see the probabilistic scheduling lecture by Andrew Wang and Cheng Fang.

Terminal violations = Observable violations

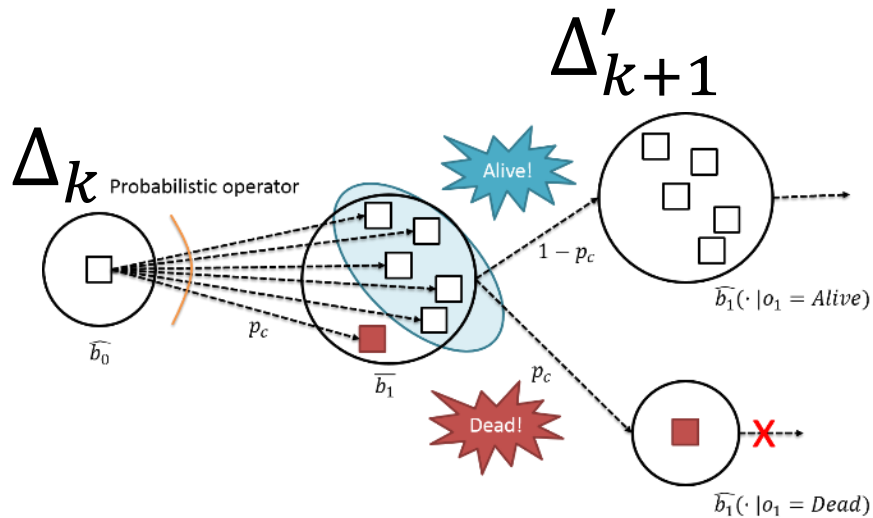


$$r_b(b_k) = 0, \forall \text{ non-terminal } b_k$$

Complete constraint violation awareness!

Dynamic vs Static risk allocation

$$\Delta'_{k+1} = \frac{1}{Pr^{sa}(o'_{k+1}|\pi(b_k), b_k)} \left(\frac{\Delta_k - r_b(b_k)}{1 - r_b(b_k)} - \sum_{o_{k+1} \neq o'_{k+1}} Pr^{sa}(o_{k+1}|\pi(b_k), b_k) h_{er}(b_{k+1}|\pi) \right)$$



Static allocation: $\Delta'_{k+1} = \Delta_k - p_c$

Dynamic allocation: $\Delta'_{k+1} = \frac{1}{1 - p_c} \left(\frac{\Delta_k - 0}{1 - 0} - p_c \cdot 1 \right) = \frac{\Delta_k - p_c}{1 - p_c}$

Some takeaways

1. *Execution risk* should be applicable to risk-aware planning in general and could be incorporated into other POMDP solvers to endow them with a keen sensitivity to risk;
2. Risk-bounded plan execution improves upon the conservatism of risk-minimal alternatives while offering strict safety guarantees;
3. Efficient risk-aware constraint solvers are necessary for risk-aware planning;