

POMDP Value Iteration	$ S \sim 10-20$
SARSOP	$ S \sim 10,000$

POMDP Formulation Approximations

POMDP Computational Complexity

POMDP Computational Complexity

Sad facts ● ☹️

- Infinite horizon POMDPs are *undecidable*

POMDP Computational Complexity

Sad facts ●

- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*

POMDP Computational Complexity

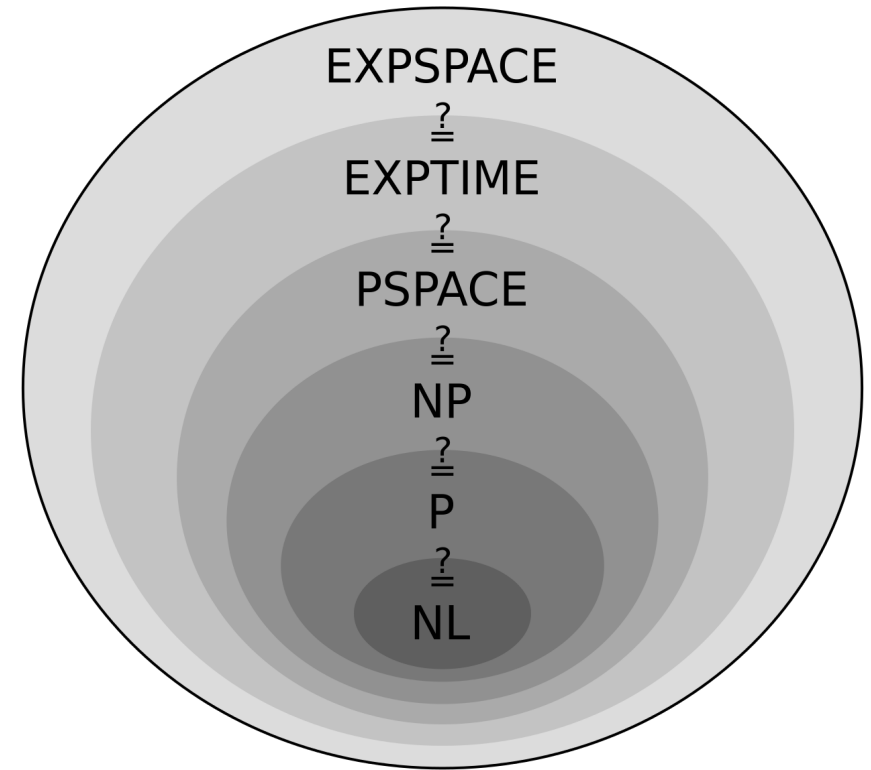
Sad facts ●

- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space

POMDP Computational Complexity

Sad facts ●

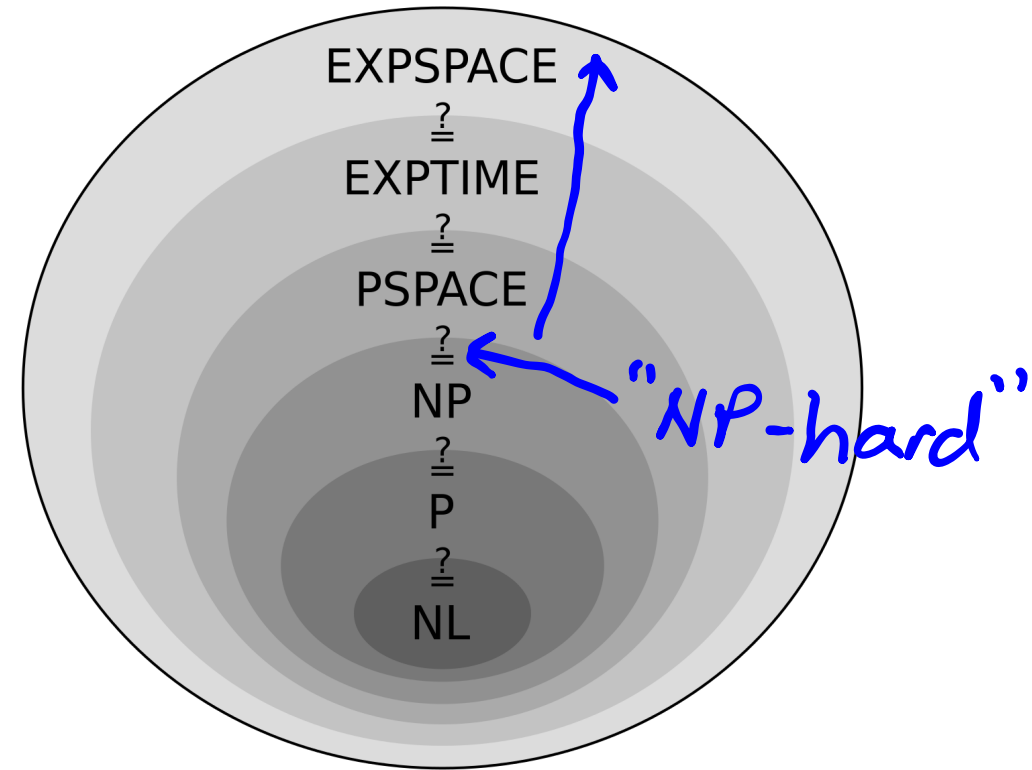
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

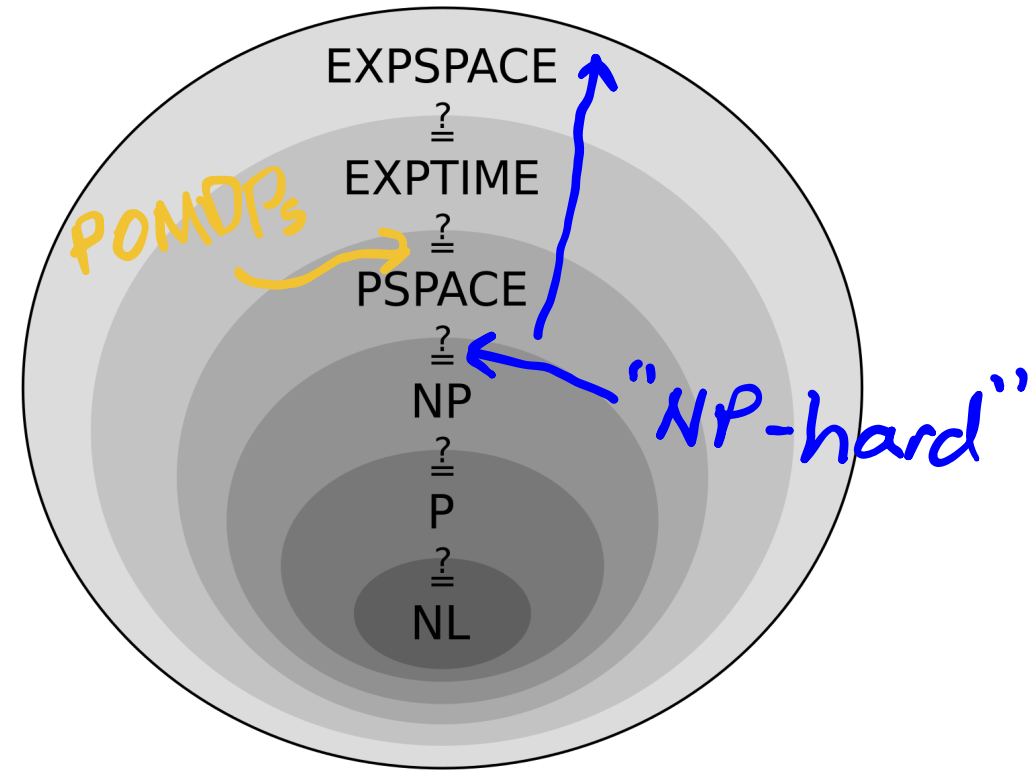
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

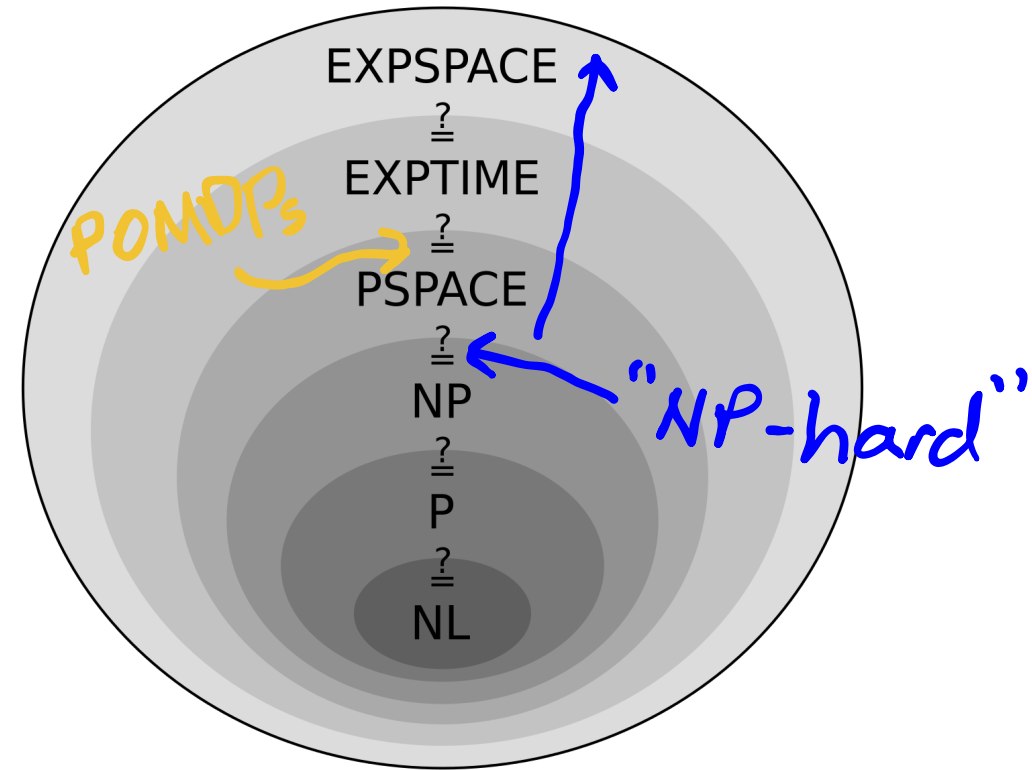
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

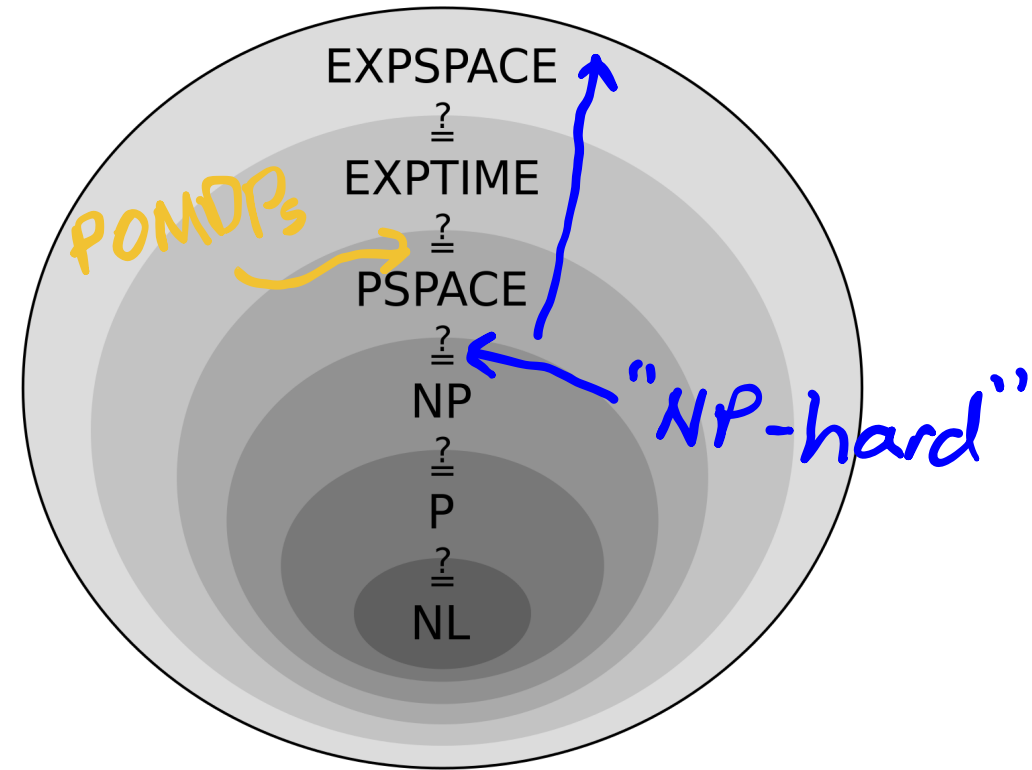
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space
 - Any algorithm that can solve a general POMDP will have exponential complexity



POMDP Computational Complexity

Sad facts ●

- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space
 - Any algorithm that can solve a general POMDP will have exponential complexity (we think)



Approximate POMDP Solutions

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Approximate POMDP Solutions

$$\bar{V} - \underline{V} = \varepsilon$$

Numerical Approximations

(approximately solve original problem)



Offline

~~Last week~~

Today

scalable to $|S| \sim 10,000$

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Last week



Online

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Last week



Online

Thursday

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Last week



Online

Thursday

Formulation Approximations

(solve a slightly different problem)

Approximate POMDP Solutions

Numerical Approximations

(approximately solve original problem)



Offline

Last week



Online

Thursday

Formulation Approximations

(solve a slightly different problem)

Today!

POMDP Objective

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$\underline{\underline{b' = \tau(b, a, o)}}$$

Certainty Equivalent

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

Certainty Equivalent

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

$$b' = \tau(b, a, o)$$

Certainty Equivalent

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

$$\pi_{\text{CE}}(b) = \pi_s(\mathbb{E}[s])$$

Handwritten annotations:
An arrow points from the text "or node or median" to the π_s term.
An arrow points from the text "or node or median" to the $\mathbb{E}[s]$ term.
A horizontal line is drawn under the π_s term.

$$b' = \tau(b, a, o)$$

Certainty Equivalent

MDP LQR

Optimal for LQG

LQG POMDP

$$T(\mathbf{s}' | \mathbf{s}, \mathbf{a}) = \mathcal{N}(\mathbf{s}' | \mathbf{T}_s \mathbf{s} + \mathbf{T}_a \mathbf{a}, \Sigma_s)$$

Linear Dynamics

$$O(\mathbf{o} | \mathbf{s}') = \mathcal{N}(\mathbf{o} | \mathbf{O}_s \mathbf{s}', \Sigma_o)$$

Gaussian Process Noise

Gaussian Observation Noise

$$b(\mathbf{s}) = \mathcal{N}(\mathbf{s} | \overset{\text{mean}}{\mu_b}, \overset{\text{covariance}}{\Sigma_b})$$

$$\mu_p \leftarrow \mathbf{T}_s \mu_b + \mathbf{T}_a \mathbf{a}$$

$$\Sigma_p \leftarrow \mathbf{T}_s \Sigma_b \mathbf{T}_s^\top + \Sigma_s$$

$$\mathbf{K} \leftarrow \Sigma_p \mathbf{O}_s^\top (\mathbf{O}_s \Sigma_p \mathbf{O}_s^\top + \Sigma_o)^{-1}$$

$$\mu_b \leftarrow \mu_p + \mathbf{K} (\mathbf{o} - \mathbf{O}_s \mu_p)$$

$$\Sigma_b \leftarrow (\mathbf{I} - \mathbf{K} \mathbf{O}_s) \Sigma_p$$

Bayesian Update
Kalman Filter

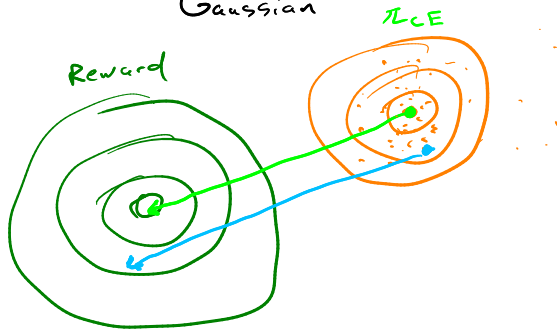
$$\pi^*(b) = -\mathbf{K}_{LQR} \mu_b$$

$$S = \mathbb{R}^2$$

Linear Dynamics + Observations
Gaussian

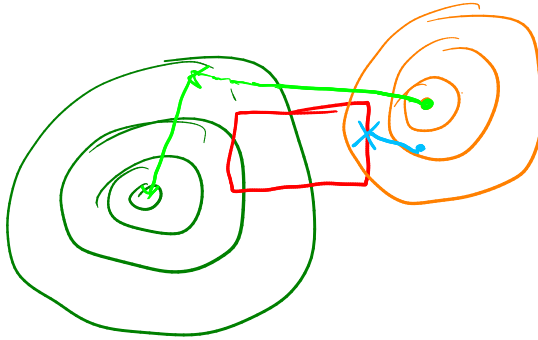
Reward Quadratic

Case 1:



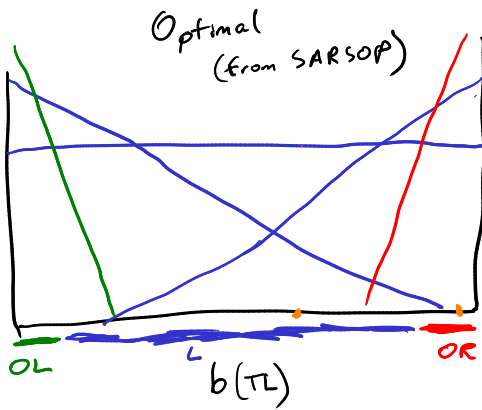
C.E. is optimal

Case 2:

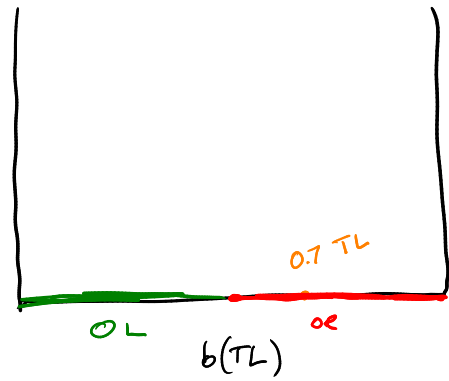


C.E. is bad

Tiger POMDP



C.E. using mode



QMDP

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

QMDP

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

$$b' = \tau(b, a, o)$$

QMDP


POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

$$\pi_{\text{QMDP}}(b) = \operatorname{argmax}_{a \in A} \mathbb{E}_{\underline{s \sim b}} [\underline{Q_{\text{MDP}}(s, a)}]$$

solve underlying
MDP
to get QMDP



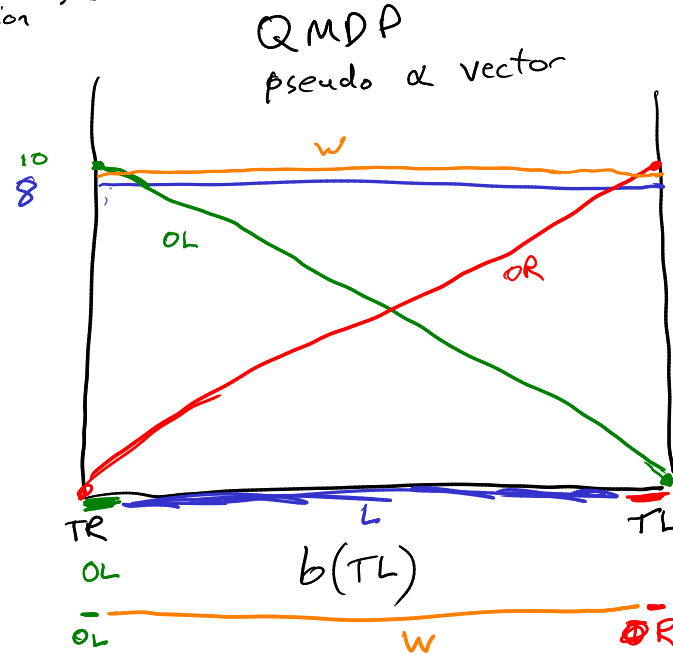
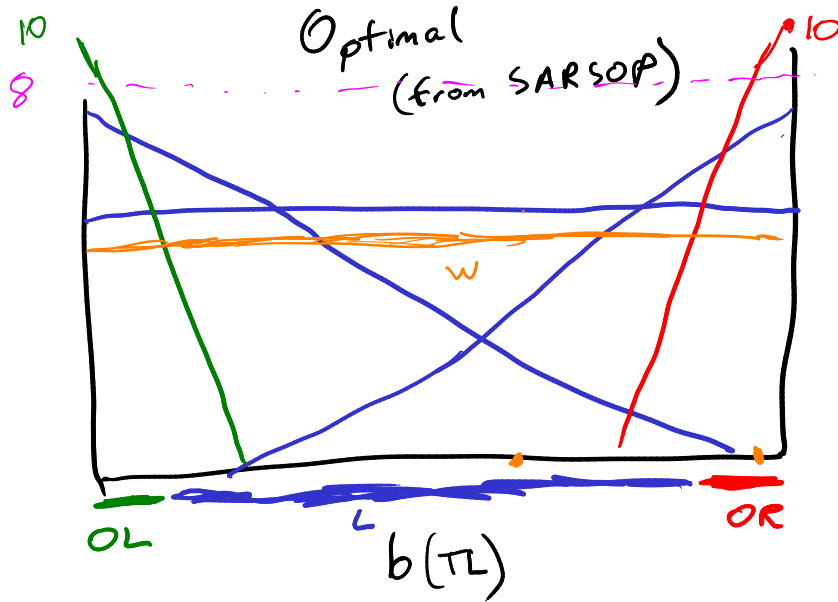
$$b' = \tau(b, a, o)$$

Example: Tiger POMDP with Waiting

Terminates when door is open

Always an upper bound on true POMDP Value function $\rightarrow E[Q_{MDP}(s,a)] \leq b(TL) Q_{MDP}(TL,a) + (1 - b(TL)) Q_{MDP}(TR,a)$

$\gamma = 0.9$



s	a	$Q_{MDP}(s,a)$
TL	OL	-100
TL	OR	+10
TR	OL	+10
TR	OR	-100
*	L	$-1 + \gamma 10 = 8$
*	W	$0 + \gamma 10 = 9$

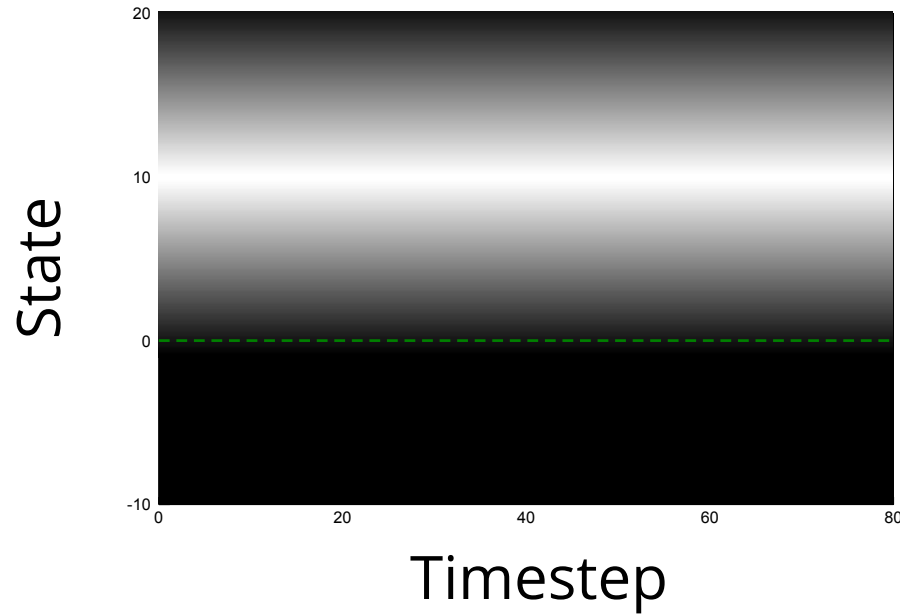
add new action Wait

50-50 observation

0 immediate reward

QMDP is bad at costly information gathering + long-lasting uncertainty
o.w. it's pretty good

POMDP Example: Light-Dark



$$\mathcal{S} = \mathbb{Z}$$

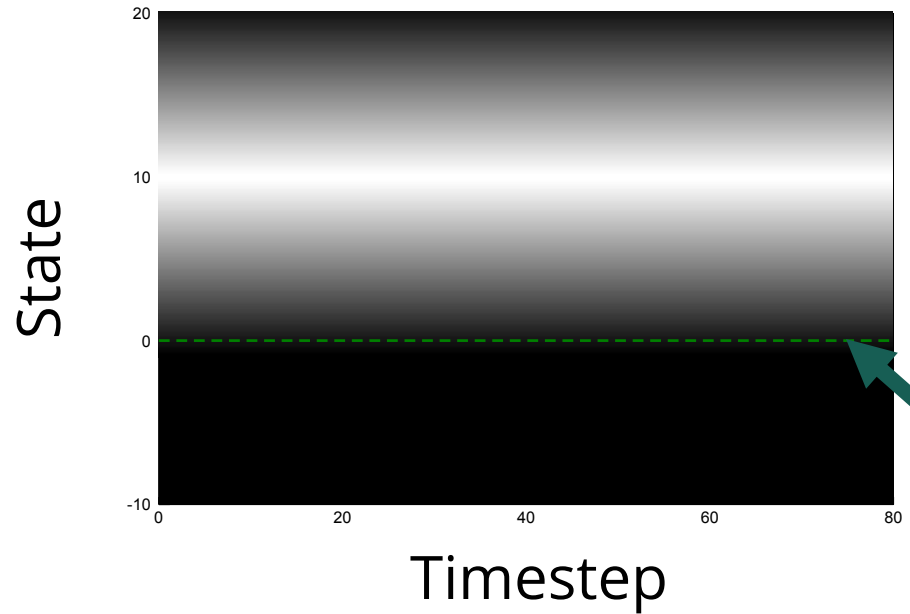
$$\mathcal{O} = \mathbb{R}$$

$$s' = s + a \quad o \sim \mathcal{N}(s, s - 10)$$

$$\mathcal{A} = \{-10, -1, 0, 1, 10\}$$

$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

POMDP Example: Light-Dark

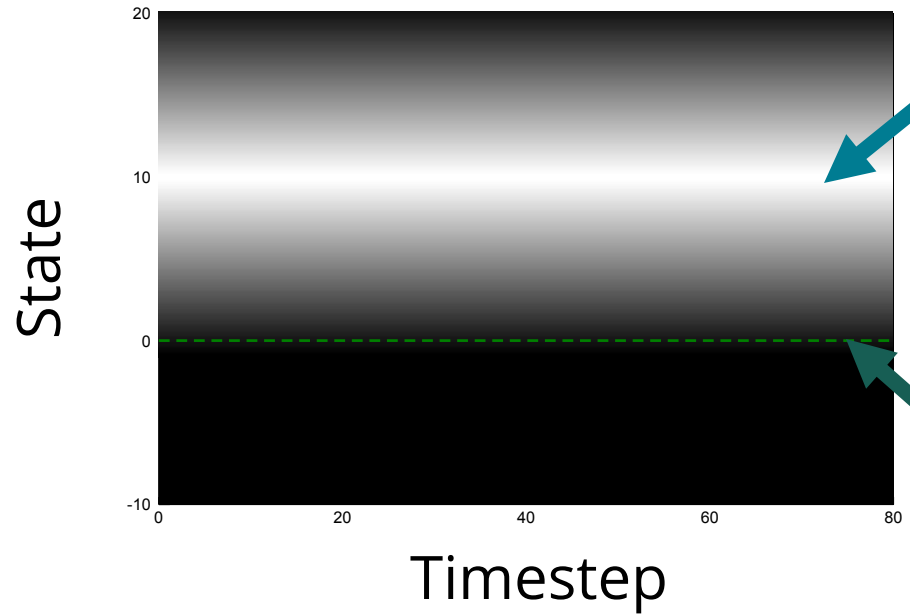


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations

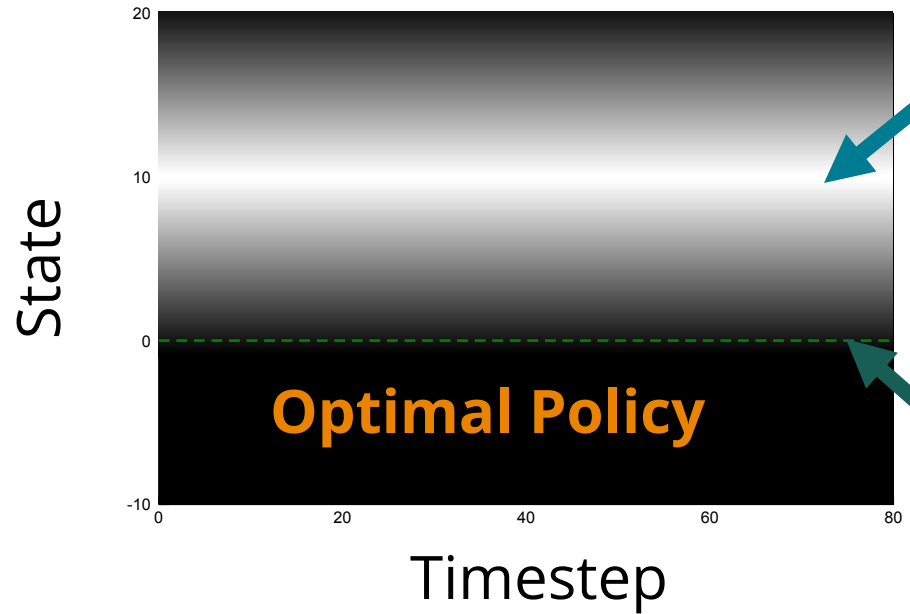


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations



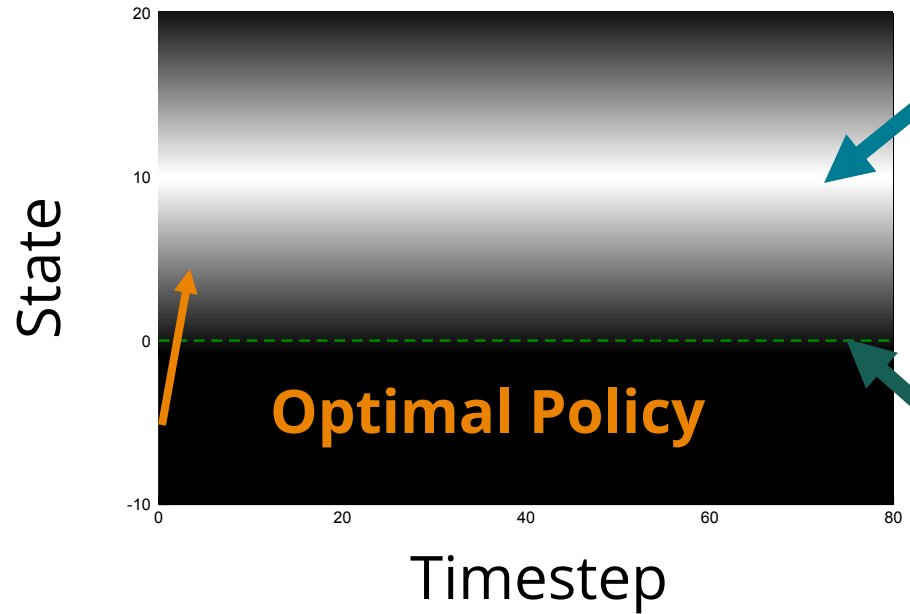
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations

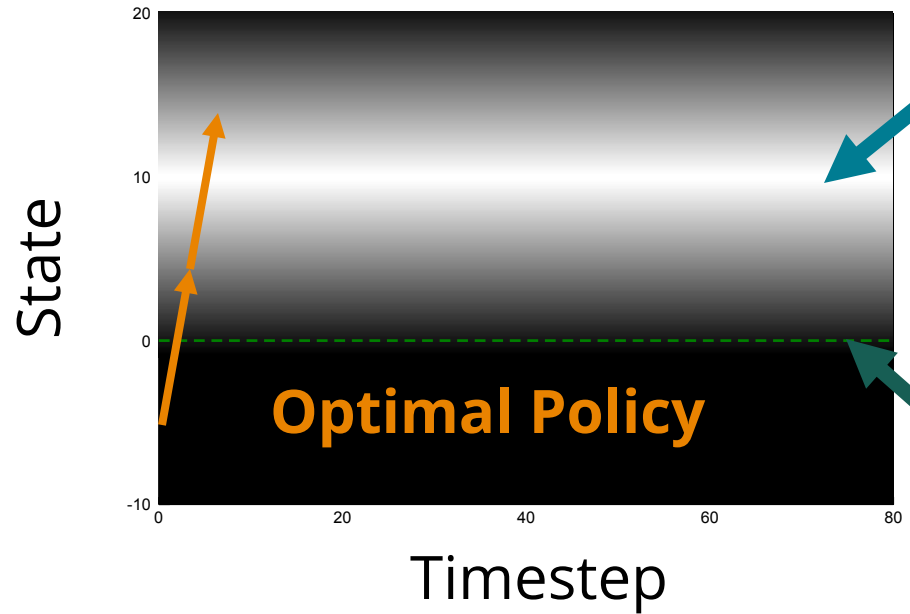


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations

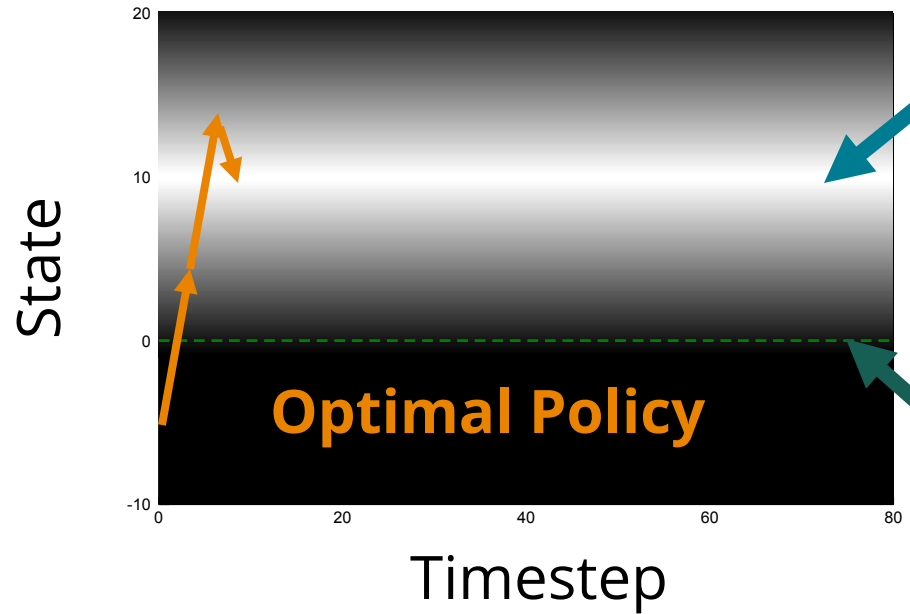


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations

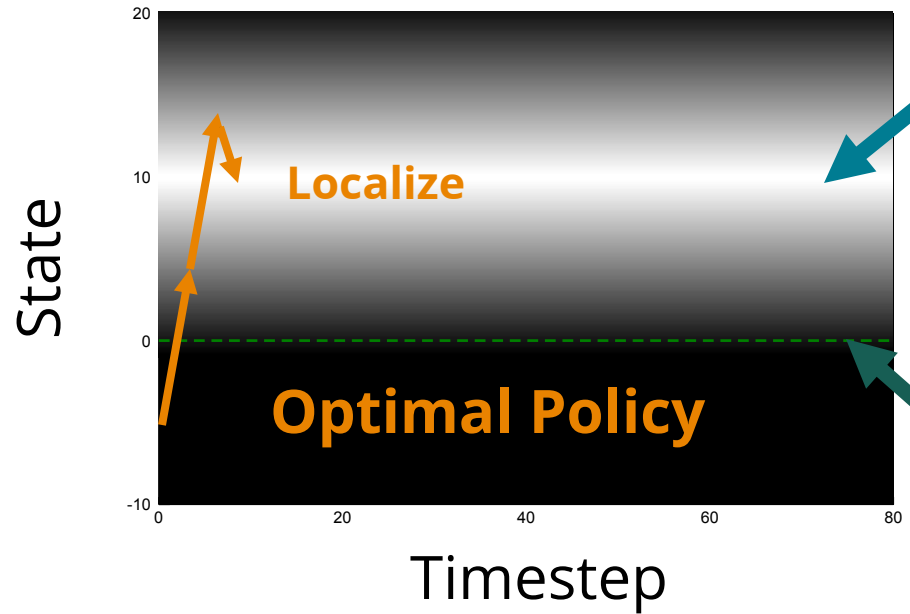


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations



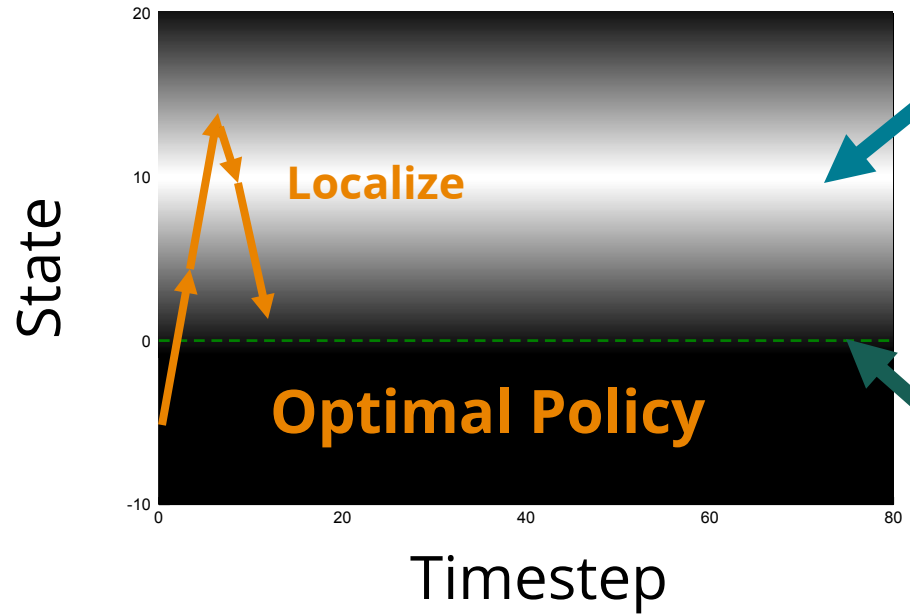
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations



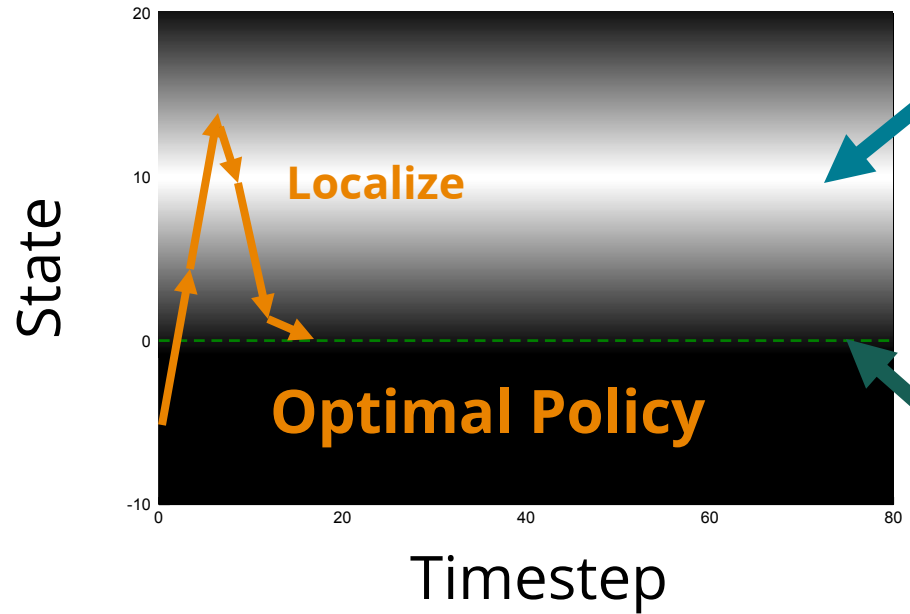
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

Accurate Observations



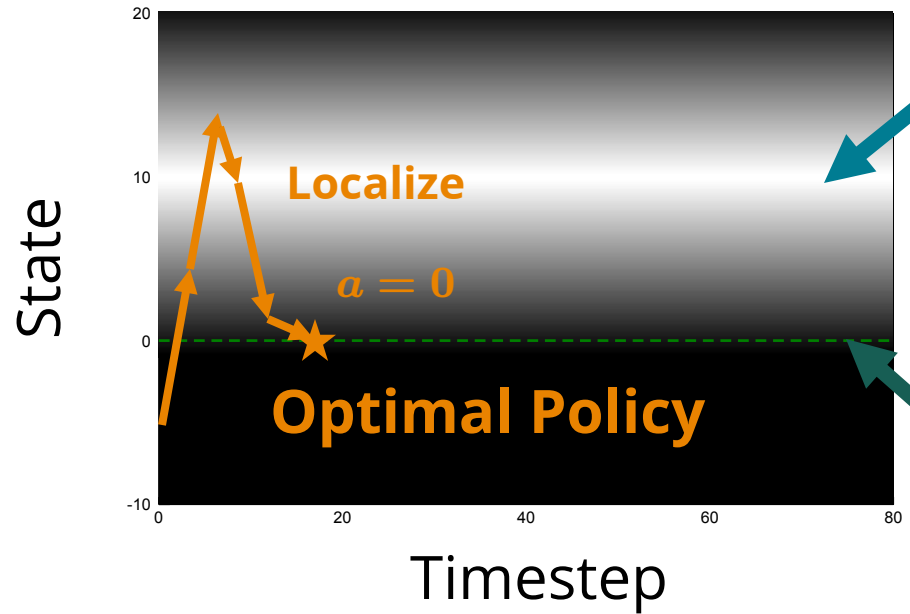
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

$$R(s, a) = \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}$$

Goal: $a = 0$ at $s = 0$

POMDP Example: Light-Dark

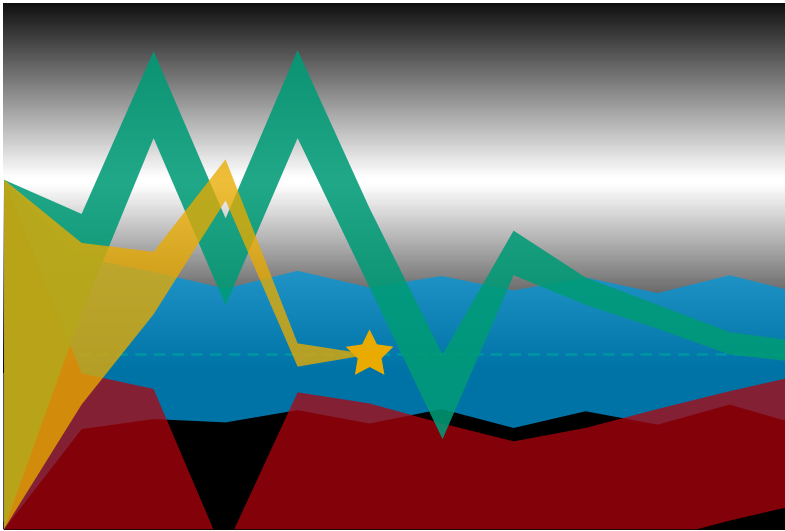
Accurate Observations



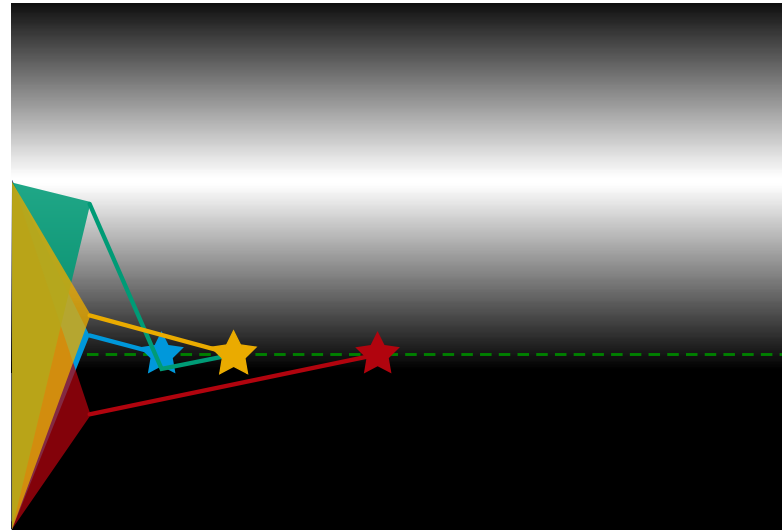
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

Goal: $a = 0$ at $s = 0$

POMDP Solution



QMDP

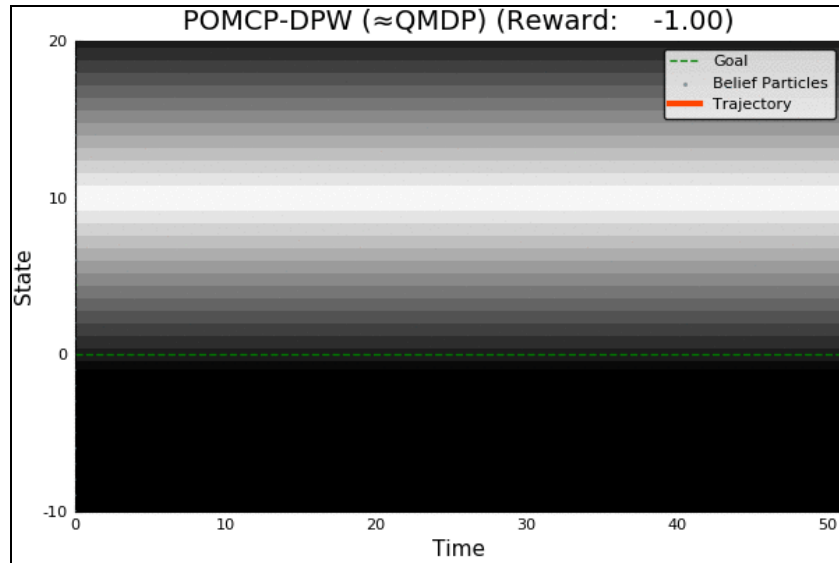


Same as **full observability**
on the next step

Information Gathering

QMDP

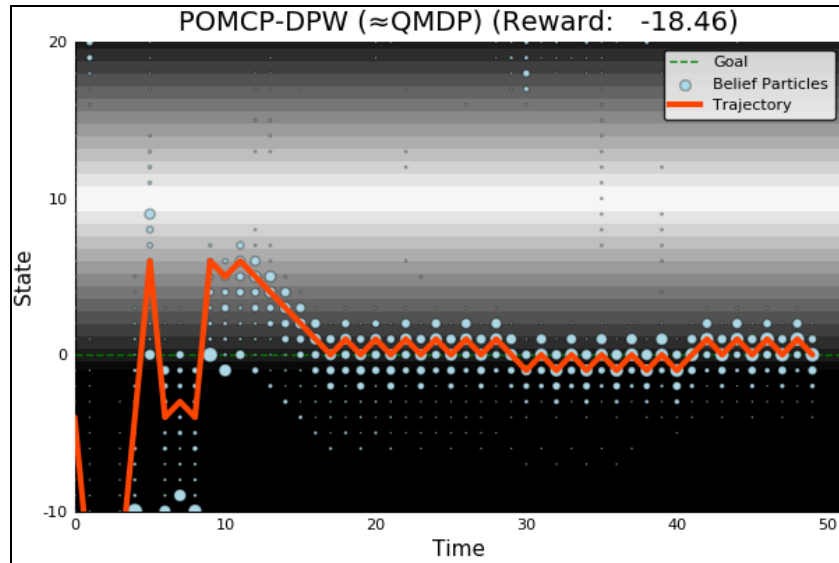
Full POMDP



Information Gathering

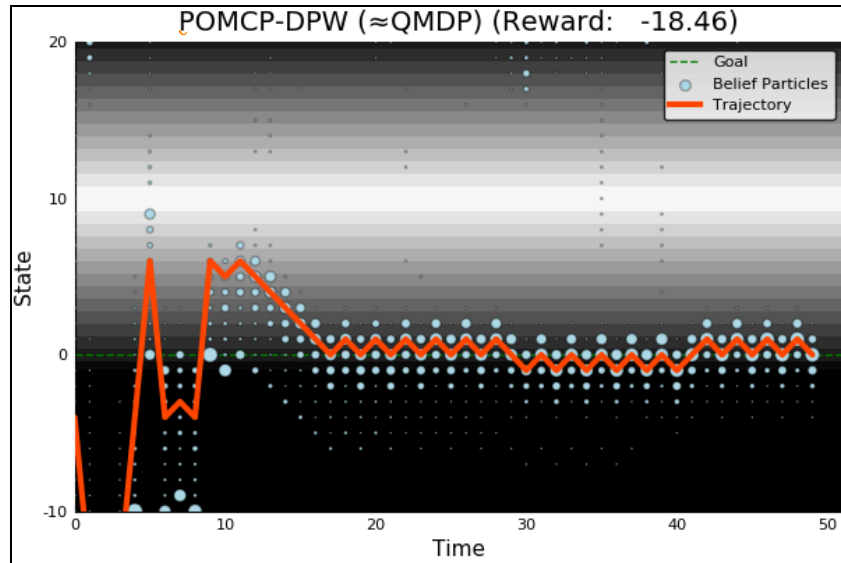
QMDP

Full POMDP

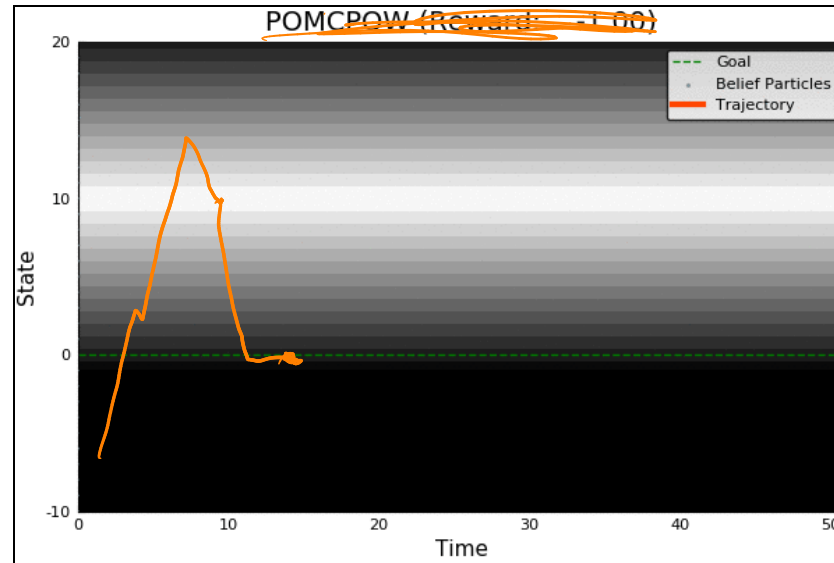


Information Gathering

QMDP

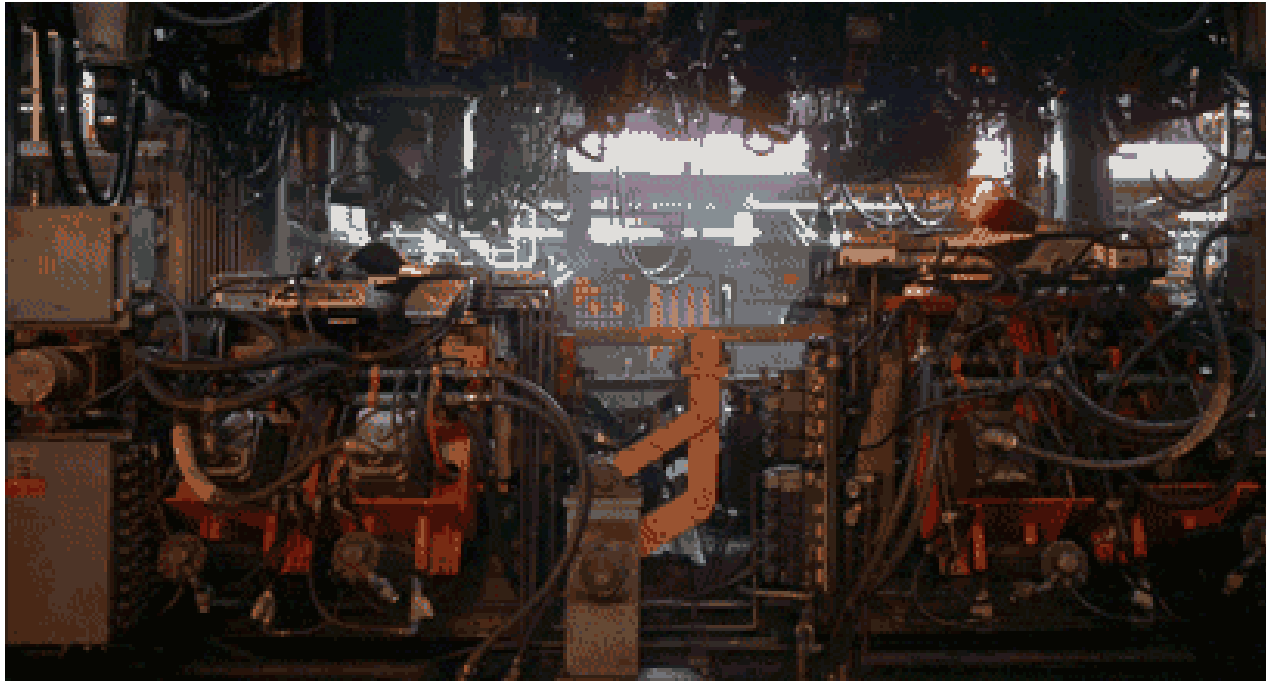


Full POMDP



QMDP

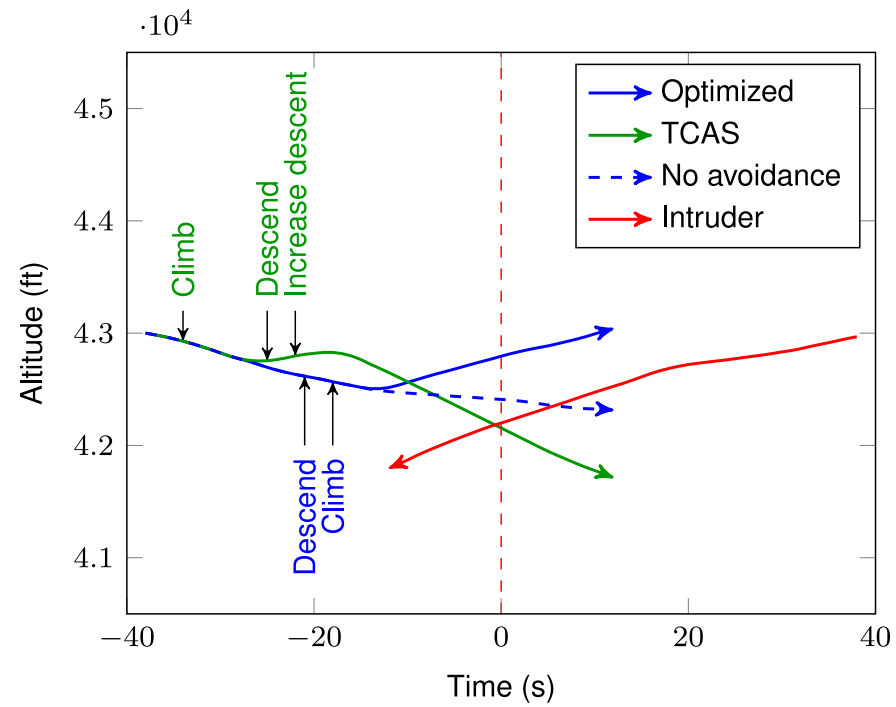
INDUSTRIAL GRADE



QMDP

ACAS X

[Kochenderfer, 2011]



Hindsight Optimization

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

pre-sample w_t^k

$$Q_{\text{HS}}(s, a) = \frac{1}{K} \sum_{k=1}^K \max_{a_{1:T}} \sum_{t=0}^T R(s_t, a_t)$$

s.t. $s_{t+1} = G(s_t, a_t, w_t^k)$
 $a_0 = a$

$$\pi_{\text{HS}} = \operatorname{argmax}_a \mathbb{E}_{s \sim b} [Q_{\text{HS}}(s, a)]$$

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$

FIB

FIB

loop

$$\alpha_a[s] \leftarrow R(s, a) + \gamma \sum_o \max_{a'} \sum_{s'} T(s' | s, a) Z(o | a, s') \alpha_{a'}[s']$$

↑

compared to QMDP

loop

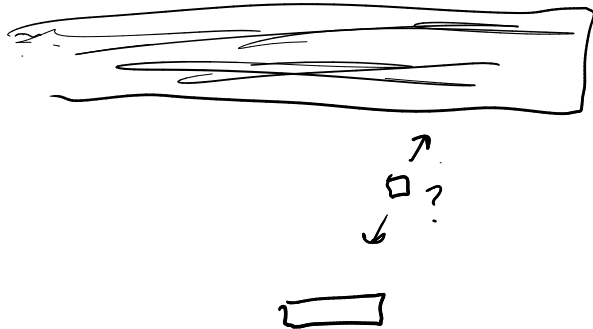
$$\alpha_a[s] \leftarrow R(s, a) + \gamma \sum_{s'} T(s' | s, a) \max_{a'} \alpha_{a'}[s']$$

k-Markov

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$



solve an MDP where

$$s_t = [o_t, o_{t-1}, \dots, o_{t-(k-1)}]$$

$k=4$ - markov approx

s_t = last 4 observations

Open Loop

POMDP Objective

$$\pi^* = \operatorname{argmax}_{\pi: B \rightarrow A} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(b_t)) \right]$$

$$b' = \tau(b, a, o)$$