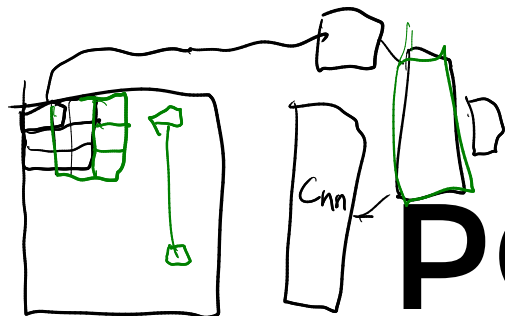


$$Q' \leftarrow Q$$

$$Q' \leftarrow (1-\tau)Q' + \tau Q$$

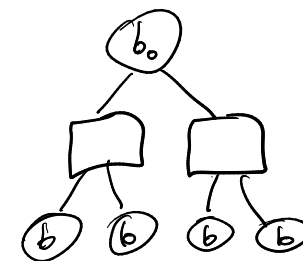
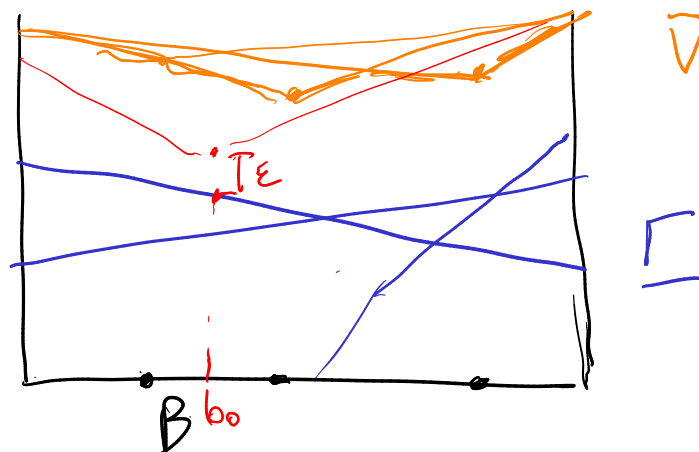


POMDP Formulation

Approximations

Convolutional
Neural
Network

SARSA
 $\sim 100,000$



POMDP Computational Complexity

POMDP Computational Complexity

Sad facts ● 🥲

POMDP Computational Complexity

Sad facts ●

- Infinite horizon POMDPs are *undecidable*

POMDP Computational Complexity

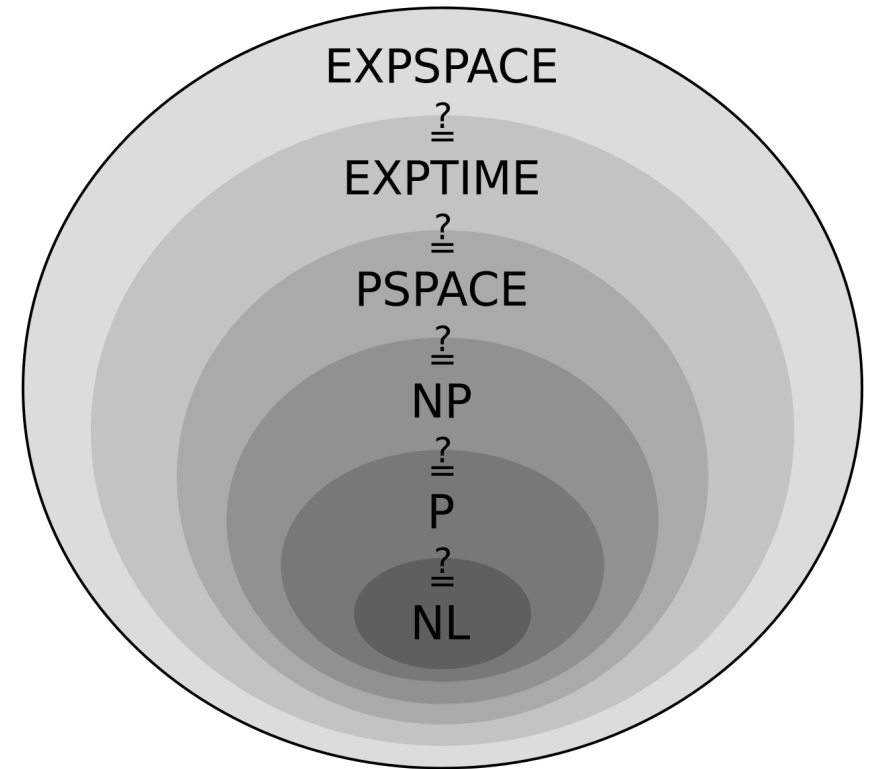
Sad facts ●

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

POMDP Computational Complexity

Sad facts ●

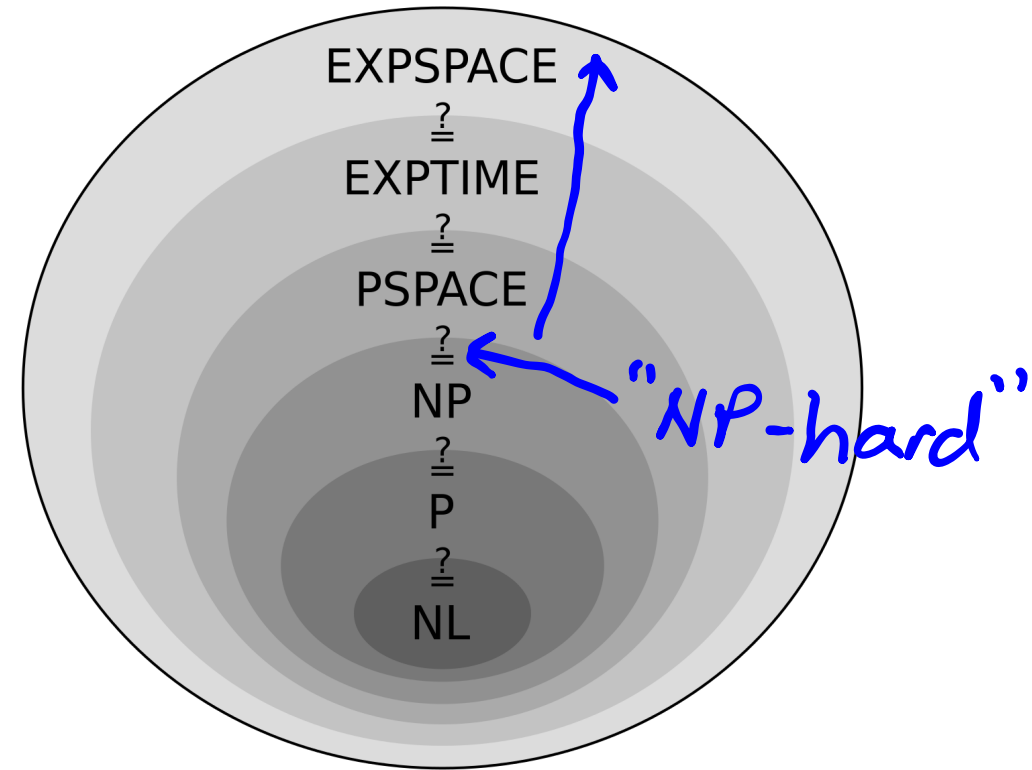
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

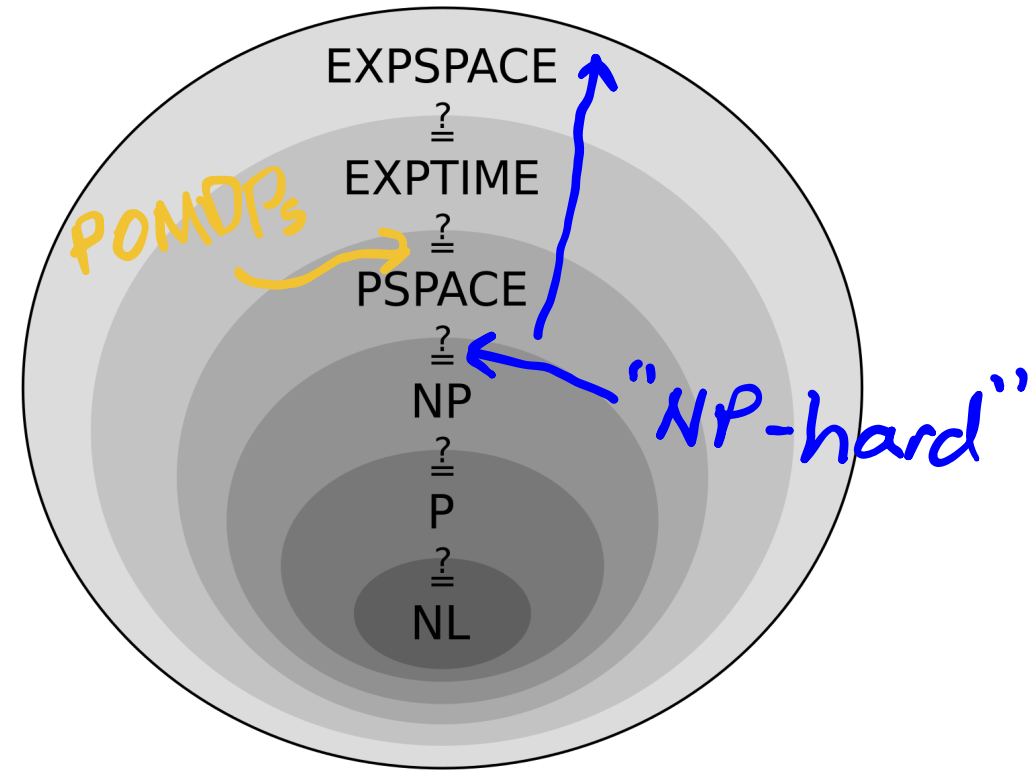
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

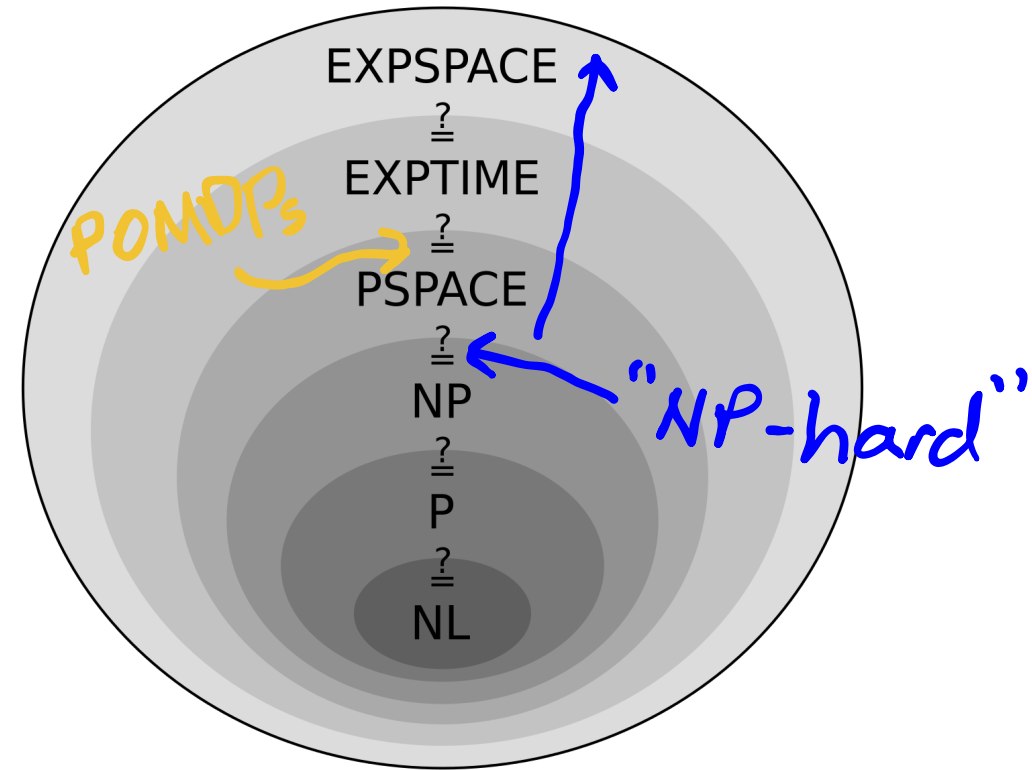
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space



POMDP Computational Complexity

Sad facts ●

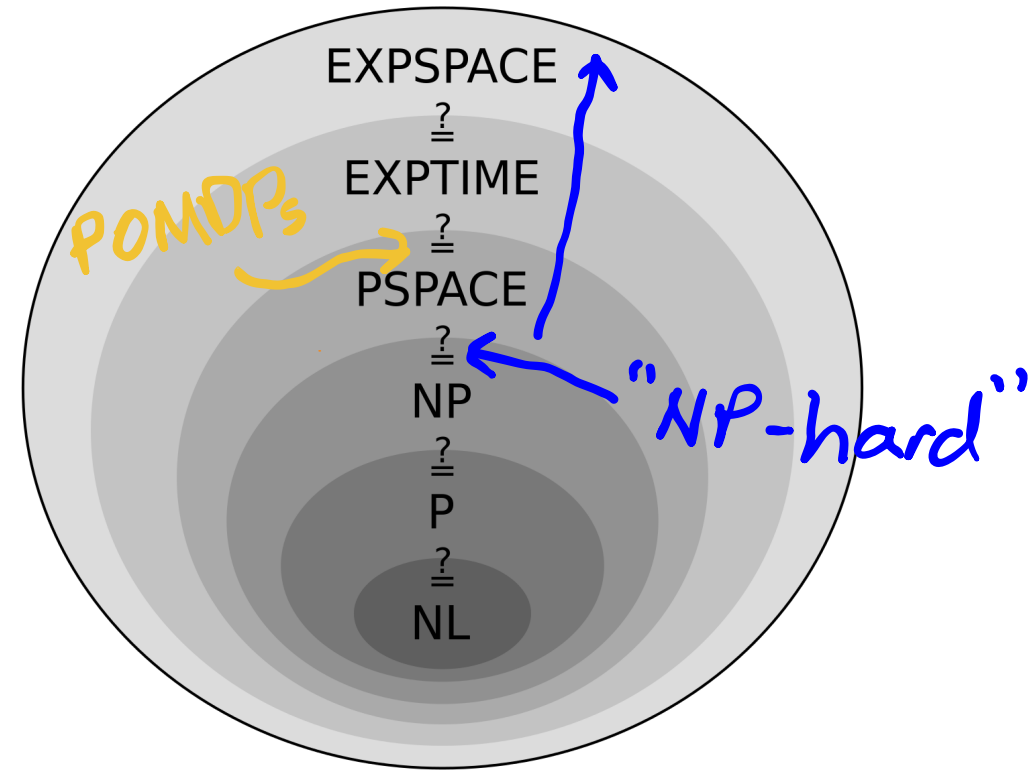
- Infinite horizon POMDPs are *undecidable*
- Finite horizon POMDPs are *PSPACE Complete*
 - Among the hardest problems that can be solved using a polynomial amount of space
 - Any algorithm that can solve a general POMDP will have exponential complexity



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

Numerical Approximations

(approximately solve original problem)



Offline

Last week

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

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


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


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

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

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

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

Handwritten notes in orange:

- or $\operatorname{argmax}_{s \in S} b(s)$ (with an arrow pointing to $\mathbb{E}[s]$)
- mode (with an arrow pointing to $\mathbb{E}[s]$)
- mean (with an arrow pointing to $s \sim b$)

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

Certainty Equivalent

Certainty Equivalent

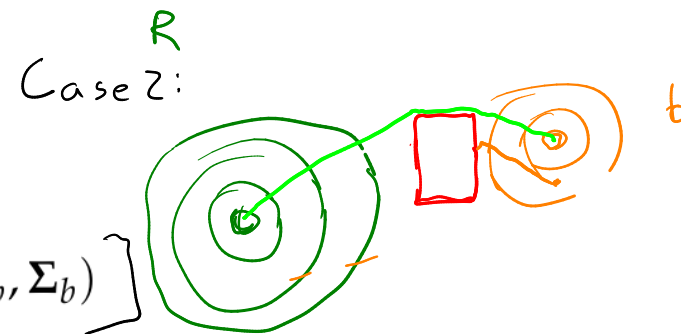
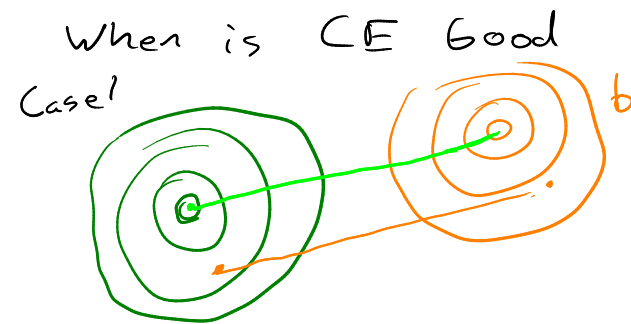
MDP analogy
LQR

$$s \in \mathbb{R}^n$$

$$a \in \mathbb{R}^m$$

Optimal for LQG

Linear Quadratic Gaussian



$$b(s) = \mathcal{N}(s \mid \mu_b, \Sigma_b)$$

$$T(s' \mid s, a) = \mathcal{N}(s' \mid T_s s + T_a a, \Sigma_s)$$

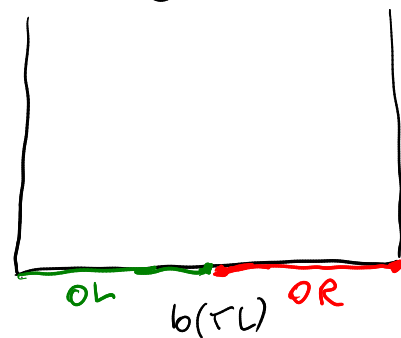
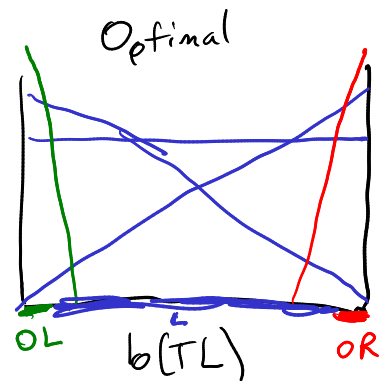
$$O(o \mid s') = \mathcal{N}(o \mid O_s s', \Sigma_o)$$

$$\underline{\mu}_p \leftarrow T_s \underline{\mu}_b + T_a a$$

$$\underline{\Sigma}_p \leftarrow T_s \underline{\Sigma}_b T_s^T + \Sigma_s$$

$$R(s, a) = -s^T R_s s - a^T R_a a$$

CE



$$\underline{K} \leftarrow \underline{\Sigma}_p \underline{O}_s^T (\underline{O}_s \underline{\Sigma}_p \underline{O}_s^T + \underline{\Sigma}_o)^{-1}$$

$$\underline{\mu}_b \leftarrow \underline{\mu}_p + \underline{K} (o - \underline{O}_s \underline{\mu}_p)$$

$$\underline{\Sigma}_b \leftarrow (I - \underline{K} \underline{O}_s) \underline{\Sigma}_p$$

Kalman Filter
(belief update)

$$\pi_{LQG}^*(b) = -K_{LQR} \mu_b$$

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

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

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

Given $(S, A, O, T, Z, R, \gamma)$
 (S, A, T, R, γ)

Q^*



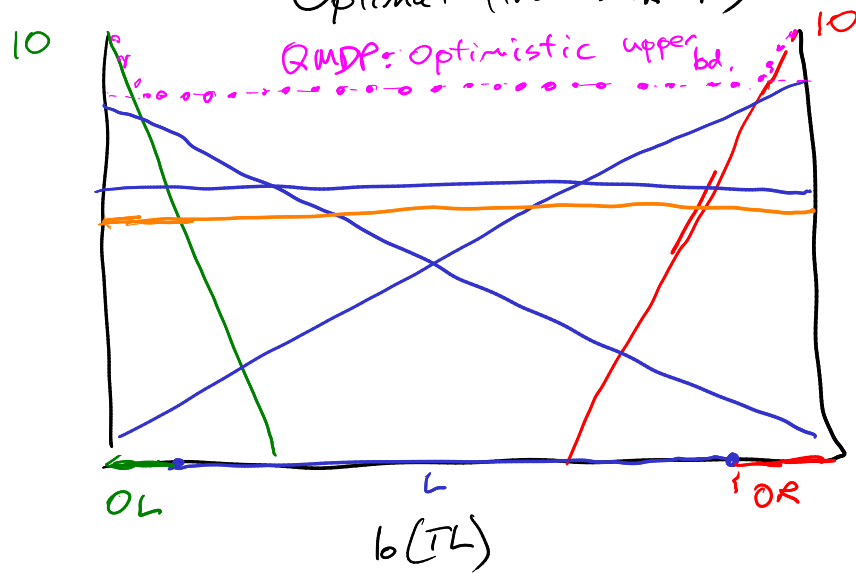
$$\pi_{\text{QMDP}}(b) = \operatorname{argmax}_{a \in A} \mathbb{E}_{s \sim b} [\underline{Q_{\text{MDP}}(s, a)}]$$

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

Example: Tiger POMDP with Waiting

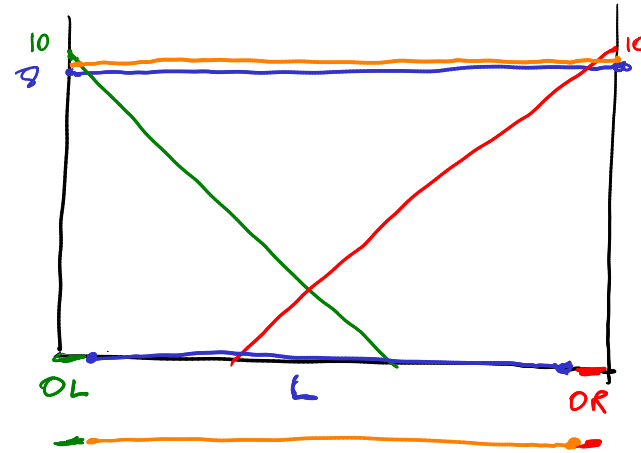
Terminates when door is open, $\gamma=0.9$

Optimal (from SARSOP)



$$\arg \max_a \mathbb{E}_{s \sim b} [Q_{MDP}(s, a)]$$

$$\mathbb{E}_{s \sim b} [Q_{MDP}(s, a)] = b(TL) Q_{MDP}(TL, a) + (1 - b(TL)) Q_{MDP}(TR, a)$$



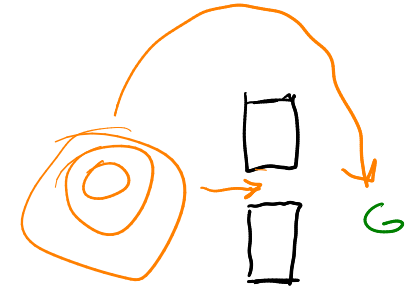
s	a	Q_{MDP}
TL	OL	-100
TL	OR	+10
TR	OL	+10
TR	OR	-100
*	L	$-1 + \gamma 10 = 8$
*	W	$0 + \gamma 10 = 9$

Is QMDP good for tiger?

Wait

└ 50-50 observation

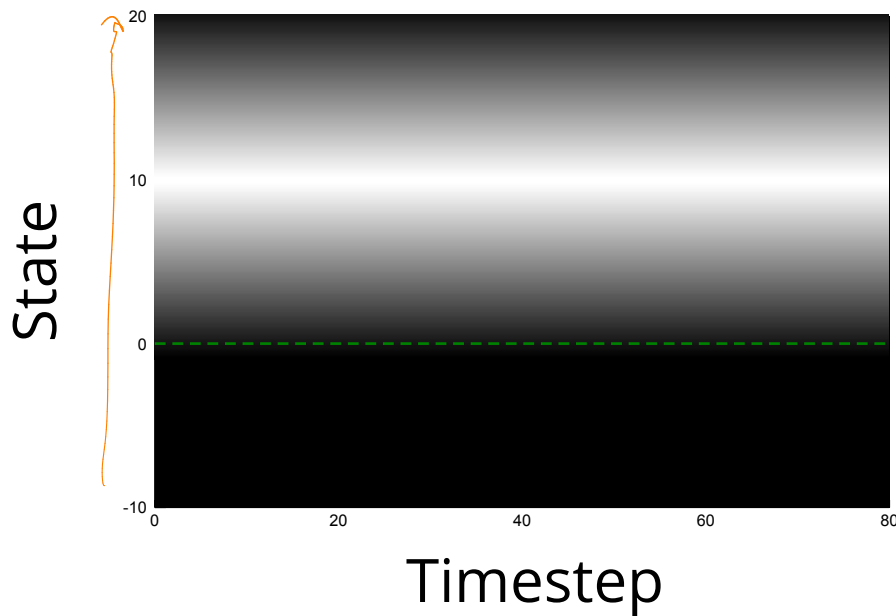
└ 0 immediate reward



QMDP is bad at costly information gathering + long-lasting uncertainty

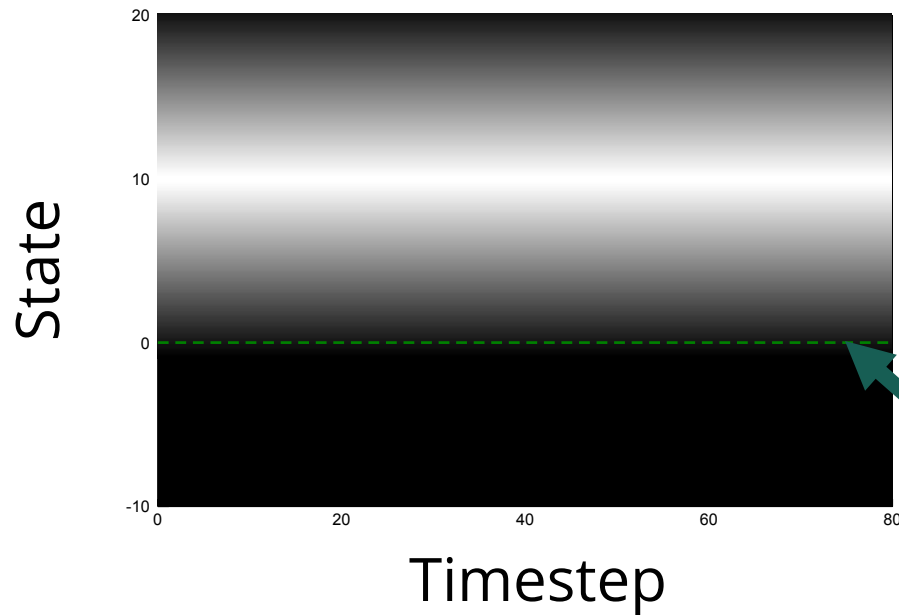
O.W. QMDP is pretty good + Much easier to solve

POMDP Example: Light-Dark



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

POMDP Example: Light-Dark



$$\mathcal{S} = \mathbb{Z} \quad \mathcal{O} = \mathbb{R}$$

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

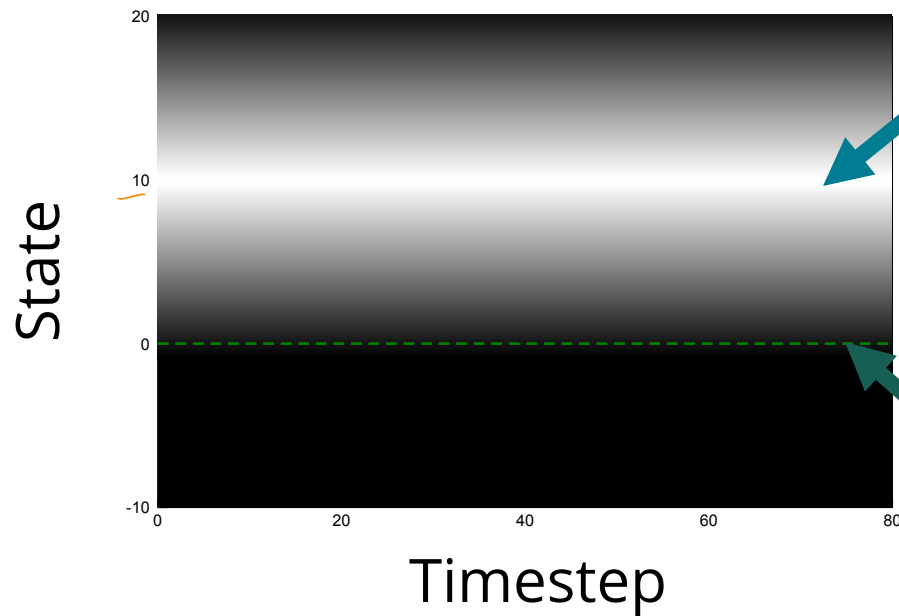
$$\mathcal{A} = \{-10, -1, 0, 1, 10\}$$

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

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

POMDP Example: Light-Dark

Accurate Observations



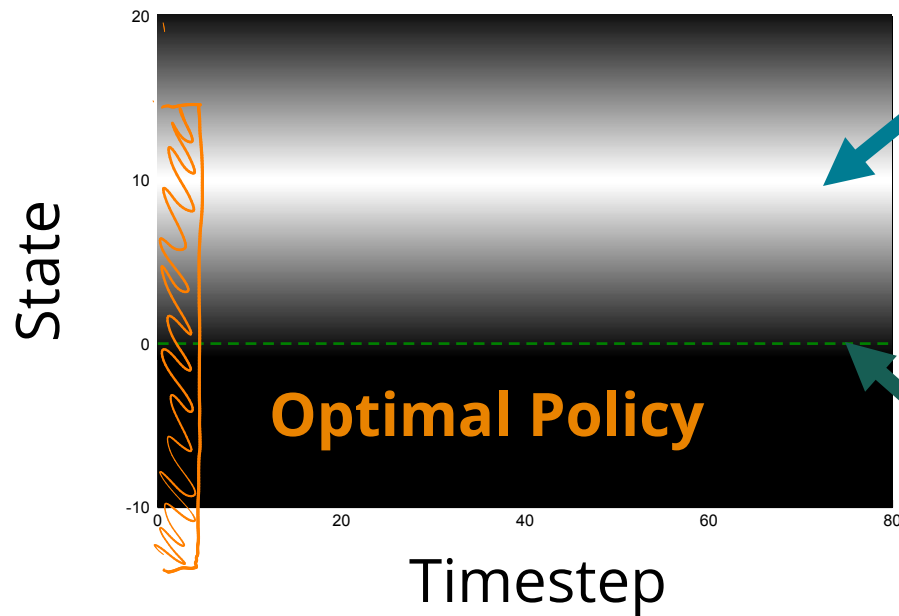
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\}\end{aligned}$$

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

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

POMDP Example: Light-Dark

Accurate Observations

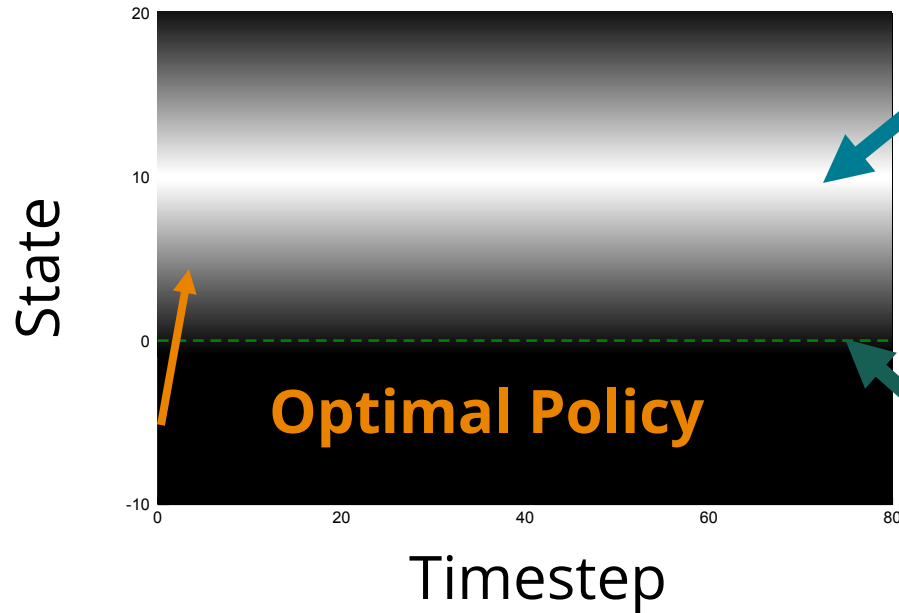


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

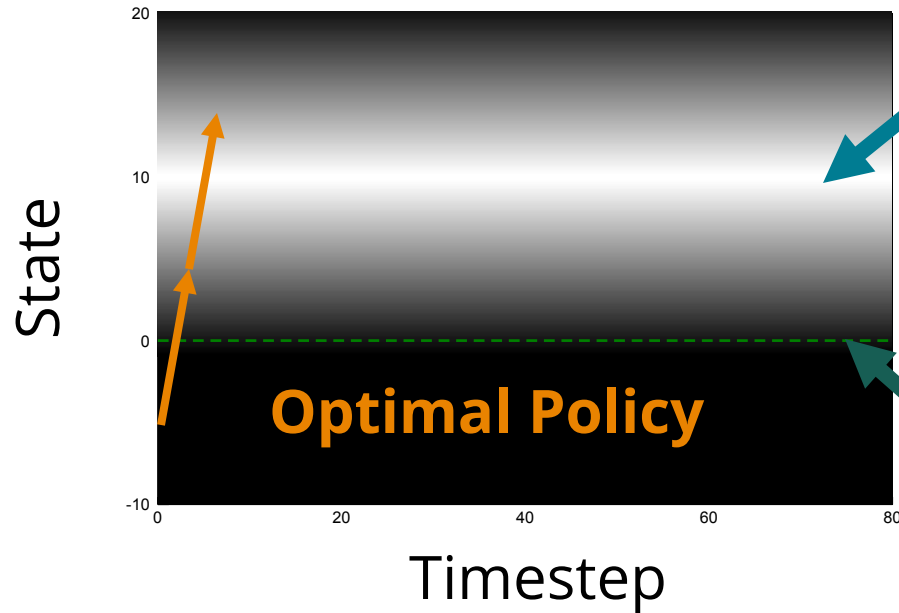


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

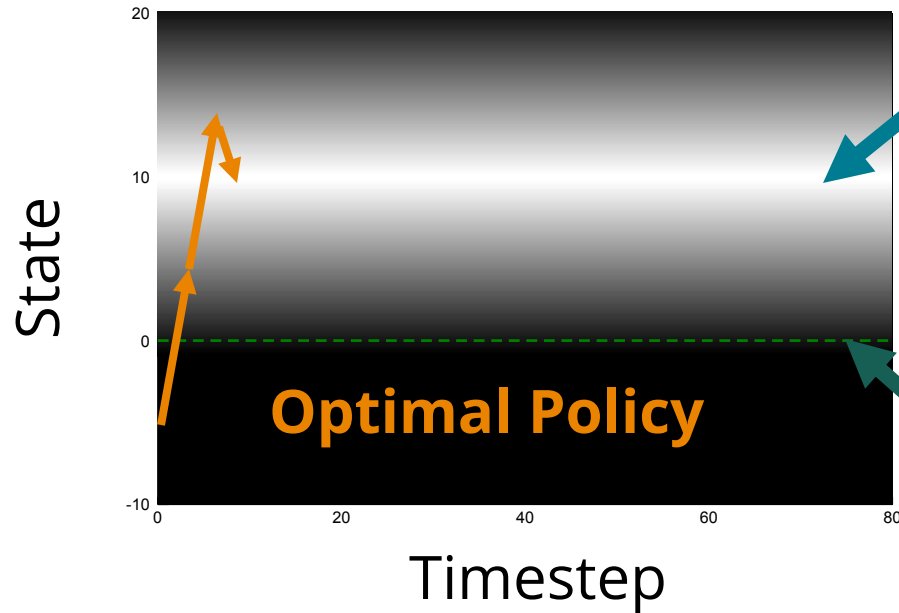


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

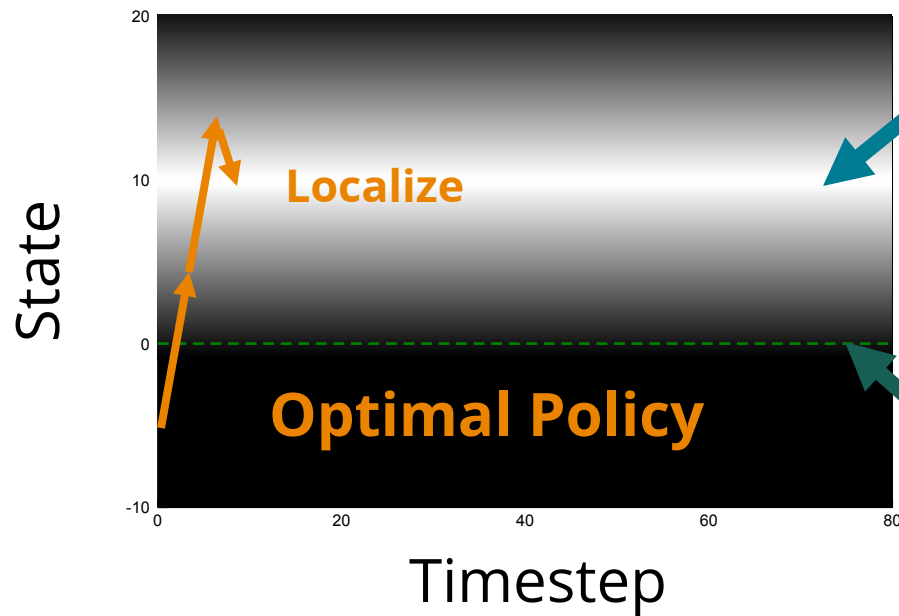


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

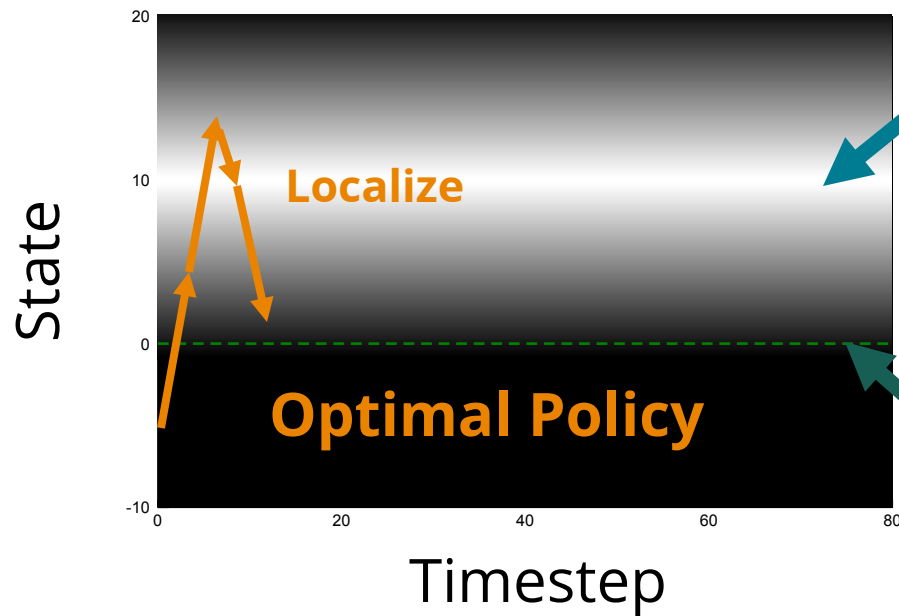


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

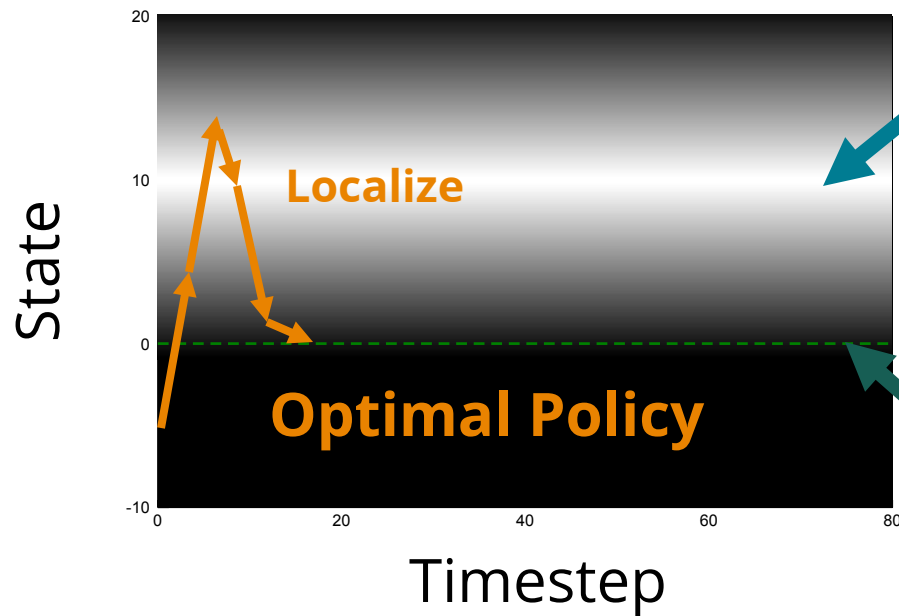


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

Accurate Observations

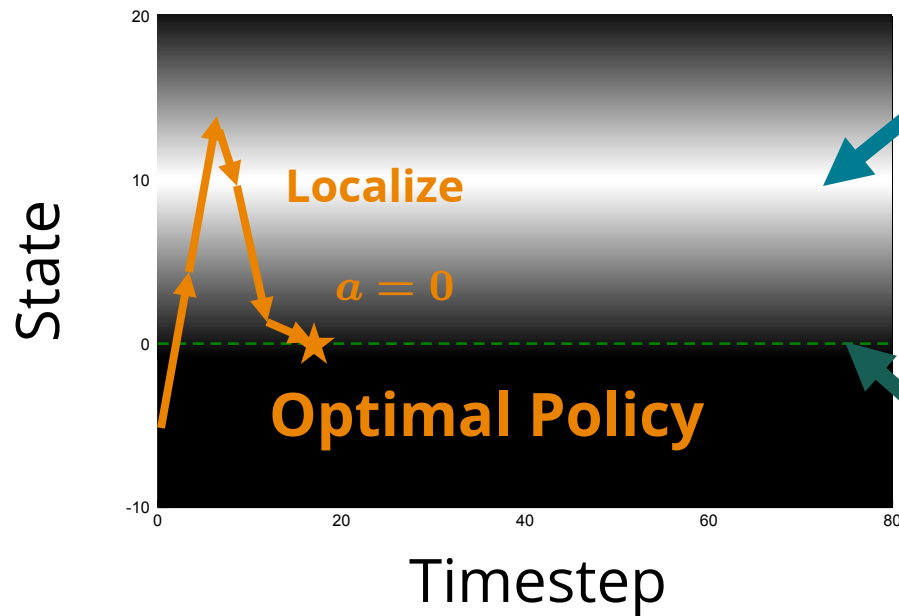


$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Example: Light-Dark

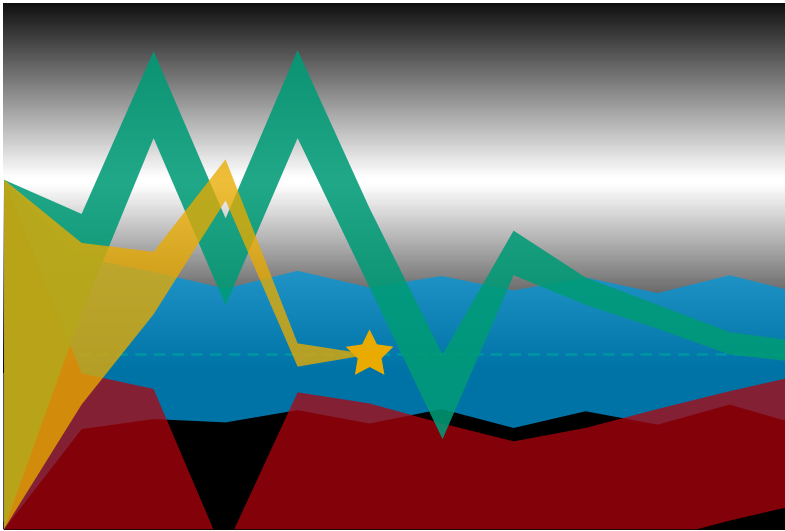
Accurate Observations



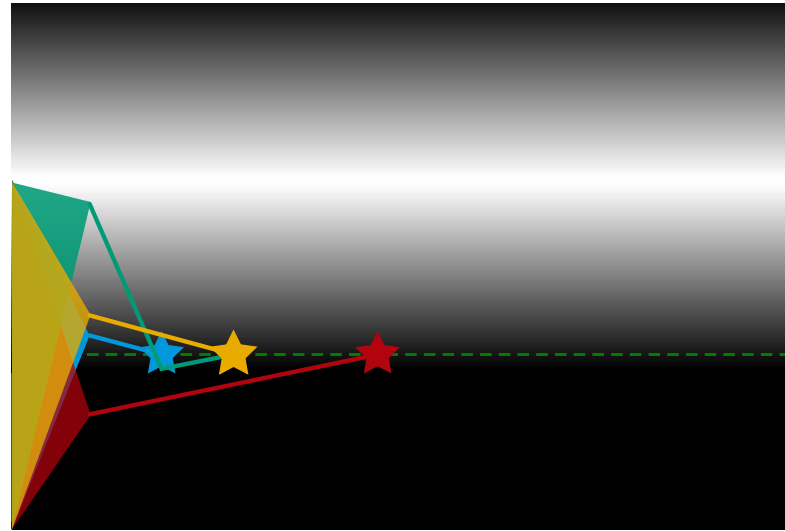
$$\begin{aligned}\mathcal{S} &= \mathbb{Z} & \mathcal{O} &= \mathbb{R} \\ s' &= s + a & o &\sim \mathcal{N}(s, s - 10) \\ \mathcal{A} &= \{-10, -1, 0, 1, 10\} \\ R(s, a) &= \begin{cases} 100 & \text{if } a = 0, s = 0 \\ -100 & \text{if } a = 0, s \neq 0 \\ -1 & \text{otherwise} \end{cases}\end{aligned}$$

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

POMDP Solution



QMDP

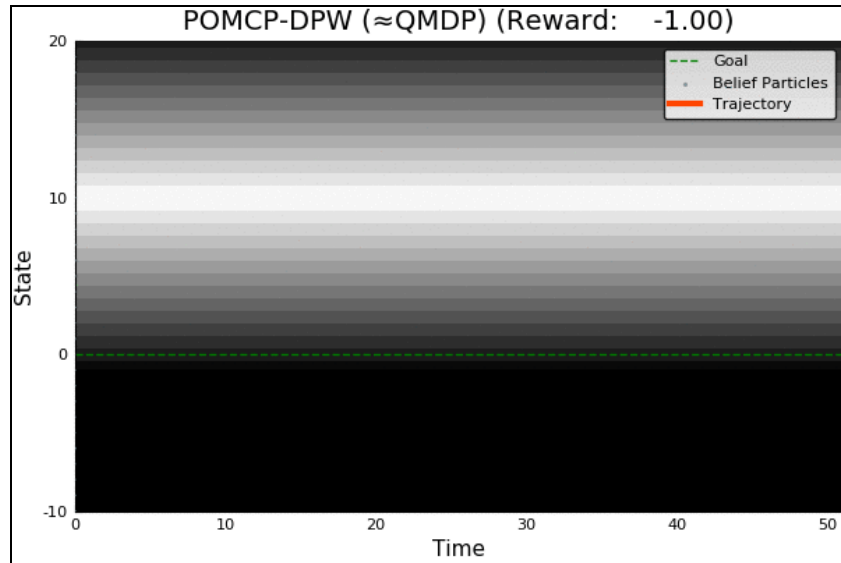


Same as **full observability**
on the next step

Information Gathering

QMDP

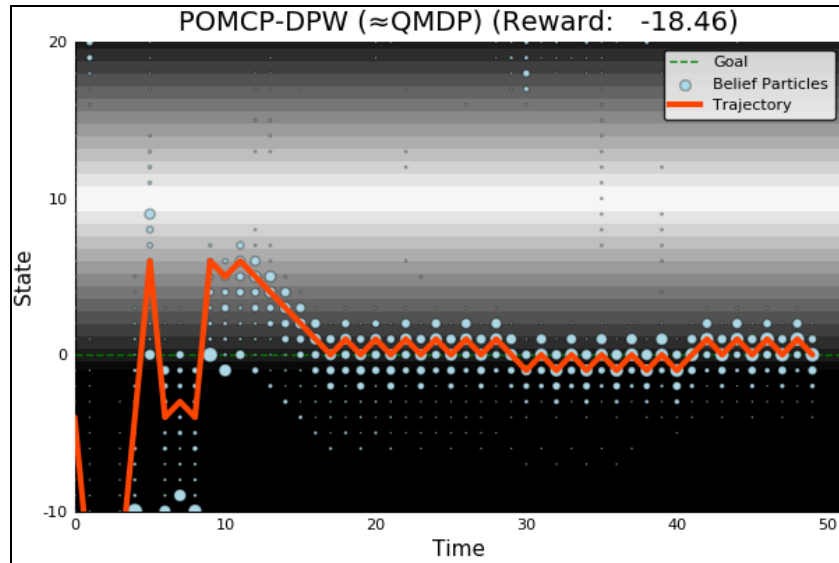
Full POMDP



Information Gathering

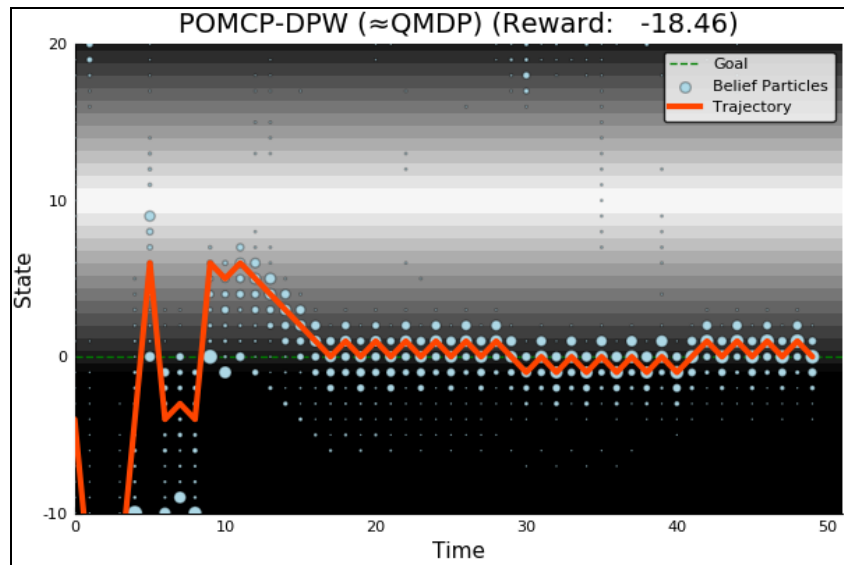
QMDP

Full POMDP

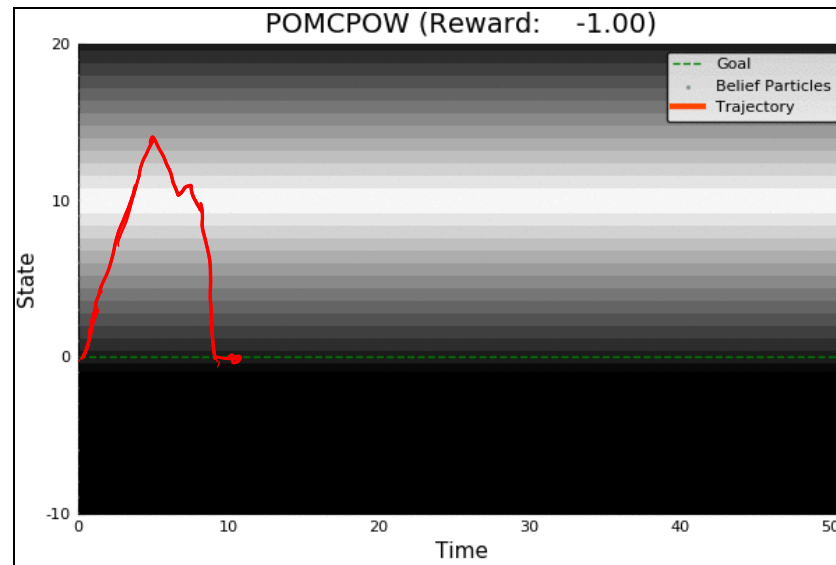


Information Gathering

QMDP

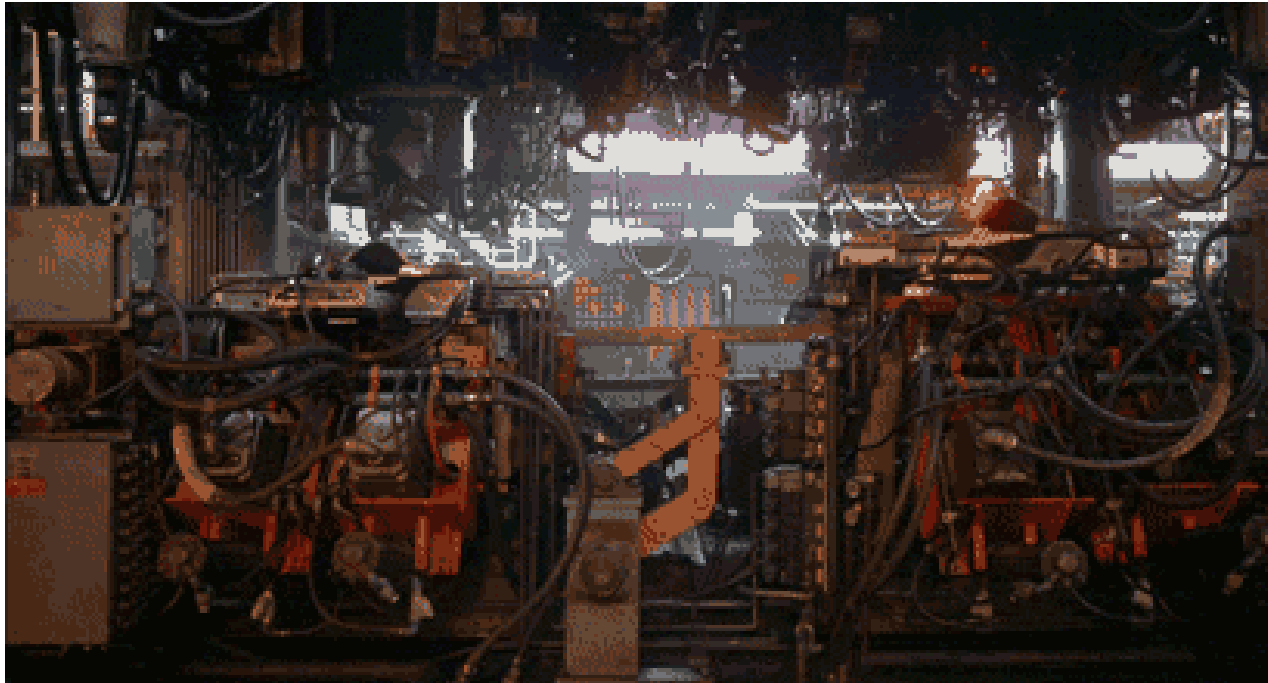


Full POMDP



QMDP

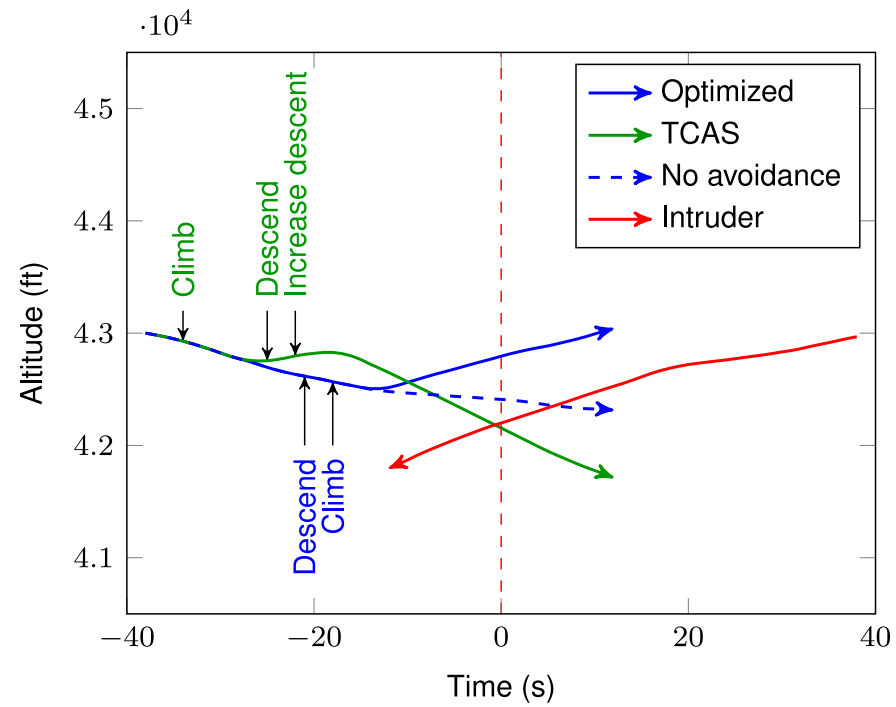
INDUSTRIAL GRADE



QMDP

ACAS X

[Kochenderfer, 2011]



Hindsight Optimization

POMDP Objective

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

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

K scenarios

pre-sample w_t^k

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

subject to $s_{t+1} = G(s_t, a_t, w_t^k)$

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

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

FIB

POMDP Objective

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

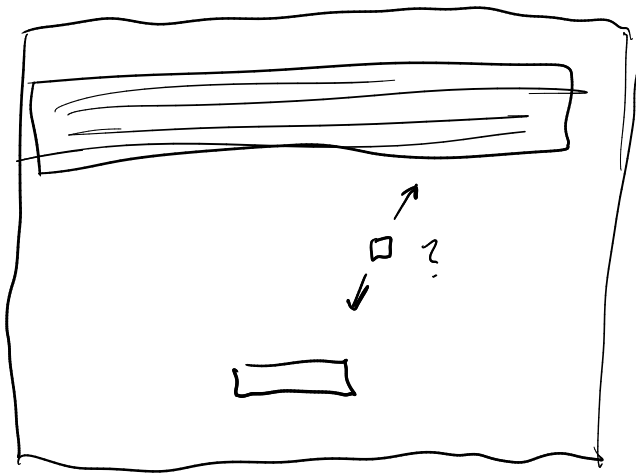
$$b' = \tau(b, a, o)$$

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
state is

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

Open Loop

POMDP Objective

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

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