

Markov Decision Process II

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science

Learning

Robot
Decision
Making

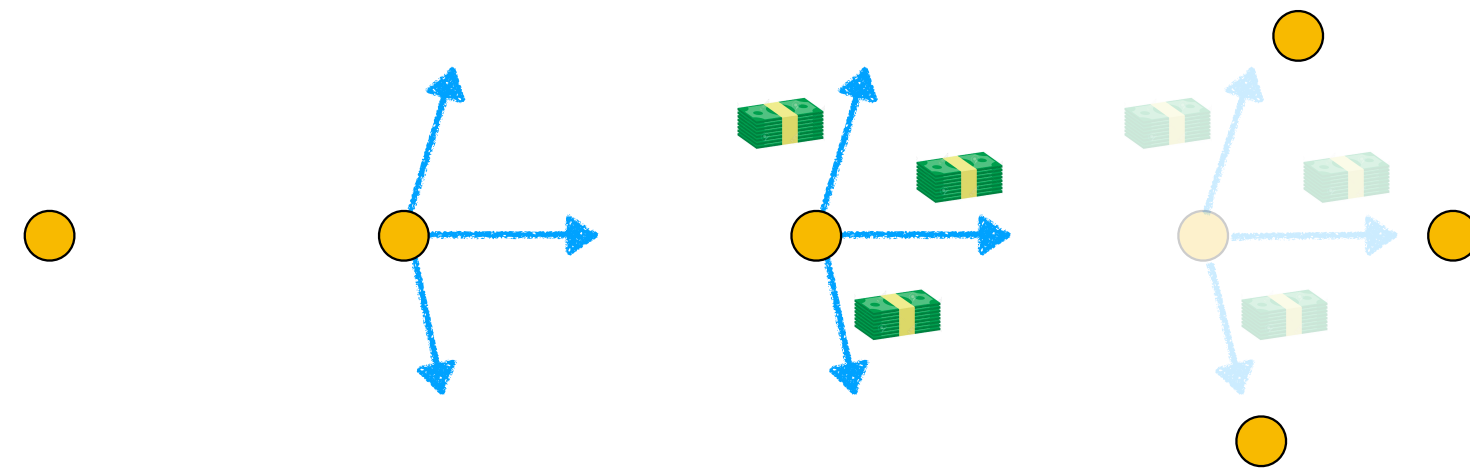
Today!

Recap

Markov Decision Process

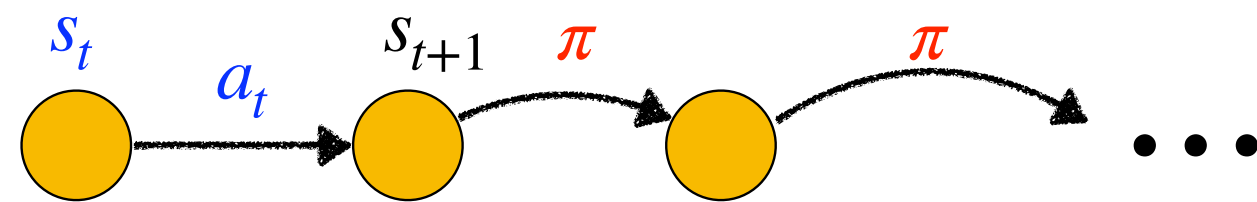
A mathematical framework for modeling sequential decision making

$$\langle S, A, C, \mathcal{T} \rangle$$



x

Value of a state-action



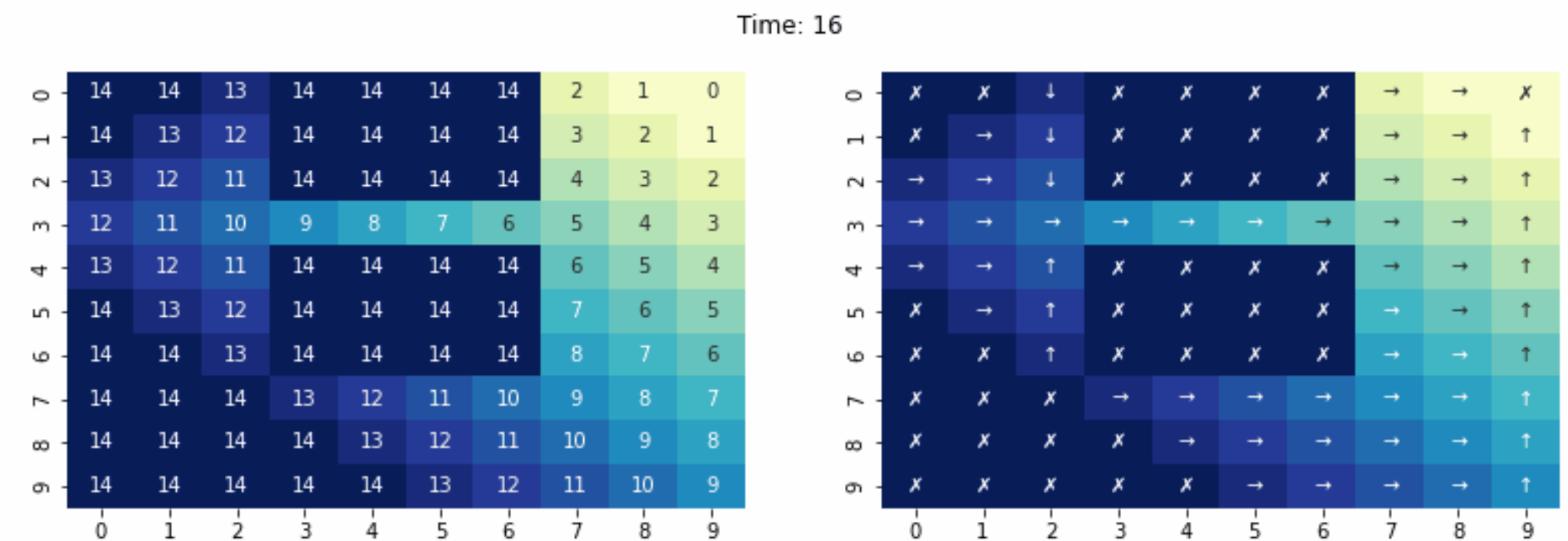
$$Q^\pi(s_t, a_t) = c_t + \gamma c_{t+1} + \gamma^2 c_{t+2} + \dots$$

Expected discounted sum of cost from starting at a state, executing action and following a policy from then on

$$Q^\pi(s_t, a_t) = c(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \mathcal{T}(s_t, a_t)} V^\pi(s_{t+1})$$

x

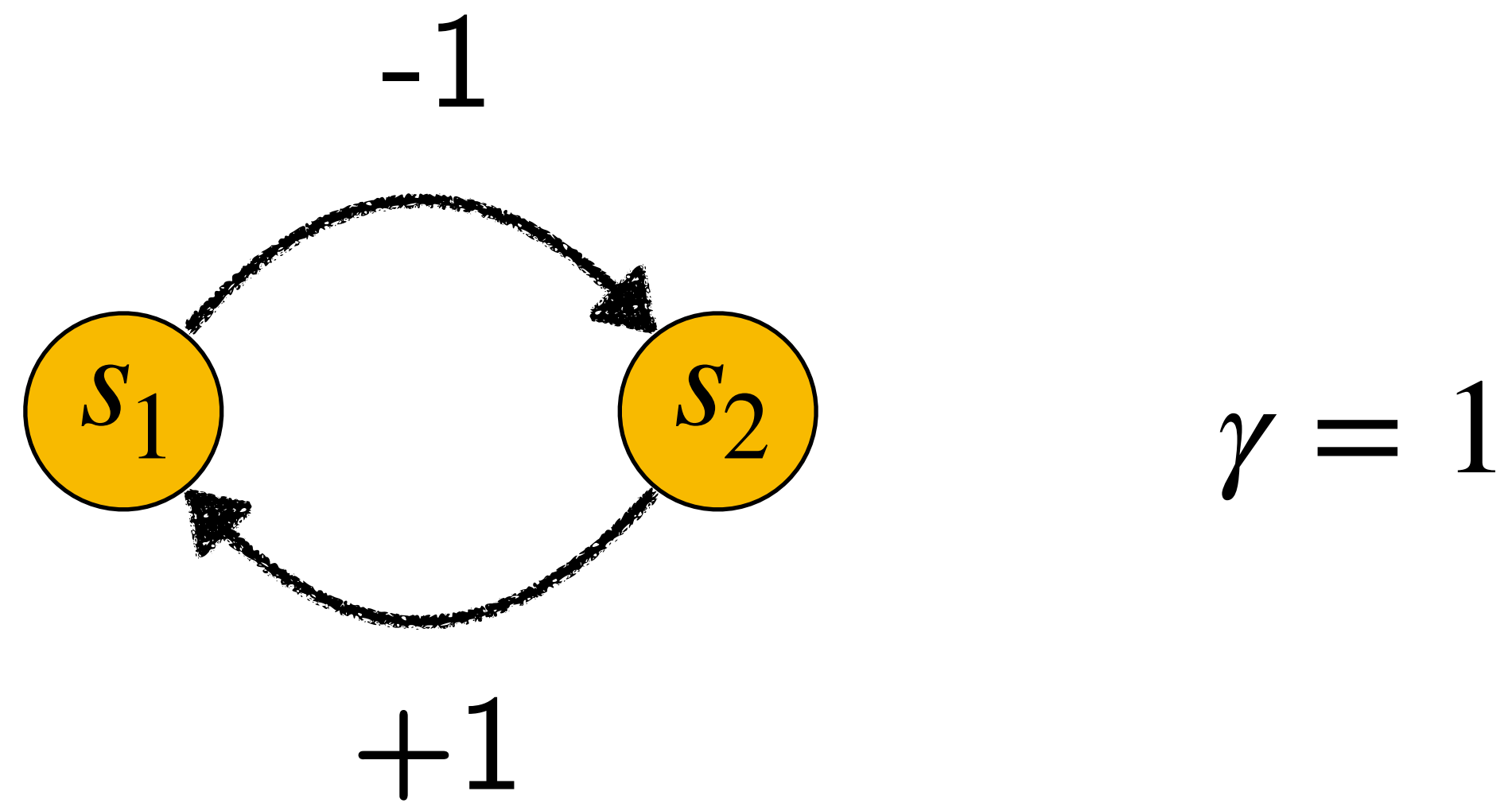
Dynamic Programming all the way!



$$V^*(s_t) = \min_a [c(s_t, a) + V^*(s_{t+1})]$$

$$\pi^*(s_t) = \arg \min_a [c(s_t, a) + V^*(s_{t+1})]$$

Does value iteration converge?



What is $V^*(s_1)$? What is $V^*(s_2)$?

What is the effect of discount factor?

Gamma: 0.0

0	1	1	1	1	1	1	1	1	1	0
1	1	1	1	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1
3	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1
5	1	1	1	1	1	1	1	1	1	1
6	1	1	1	1	1	1	1	1	1	1
7	1	1	1	1	1	1	1	1	1	1
8	1	1	1	1	1	1	1	1	1	1
9	1	1	1	1	1	1	1	1	1	1

0	x	x	x	x	x	x	x	x	x	x
1	x	x	x	x	x	x	x	x	x	x
2	x	x	x	x	x	x	x	x	x	x
3	x	x	x	x	x	x	x	x	x	x
4	x	x	x	x	x	x	x	x	x	x
5	x	x	x	x	x	x	x	x	x	x
6	x	x	x	x	x	x	x	x	x	x
7	x	x	x	x	x	x	x	x	x	x
8	x	x	x	x	x	x	x	x	x	x
9	x	x	x	x	x	x	x	x	x	x

Activity!



Think-Pair-Share

Think (30 sec): What are some attributes of a **hard** MDP?

Pair: Find a partner

Share (45 sec): Partners exchange
ideas

Policy Iteration

How frequently does the best action change?

0	-	10	10	10	10	10	10	10	10	10	10
1	-	10	10	10	10	10	10	10	10	10	10
2	-	10	10	10	10	10	10	10	10	10	10
3	-	10	10	10	10	10	10	10	10	10	10
4	-	10	10	10	10	10	10	10	10	10	10
5	-	10	10	10	10	10	10	10	10	10	10
6	-	10	10	10	10	10	10	10	10	10	10
7	-	10	10	10	10	10	10	10	10	10	10
8	-	10	10	10	10	10	10	10	10	10	10
9	-	10	10	10	10	10	10	10	10	10	10
		0	1	2	3	4	5	6	7	8	9

Values

0	-	x	x	x	x	x	x	x	x	x	x
1	-	x	x	x	x	x	x	x	x	x	x
2	-	x	x	x	x	x	x	x	x	x	x
3	-	x	x	x	x	x	x	x	x	x	x
4	-	x	x	x	x	x	x	x	x	x	x
5	-	x	x	x	x	x	x	x	x	x	x
6	-	x	x	x	x	x	x	x	x	x	x
7	-	x	x	x	x	x	x	x	x	x	x
8	-	x	x	x	x	x	x	x	x	x	x
9	-	x	x	x	x	x	x	x	x	x	x
		0	1	2	3	4	5	6	7	8	9

Policy



Policy converges **faster**
than the value

Can we iterate over **policies**?

Policy Iteration

Init with some policy π

Repeat forever

Evaluate policy

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')$$

Improve policy

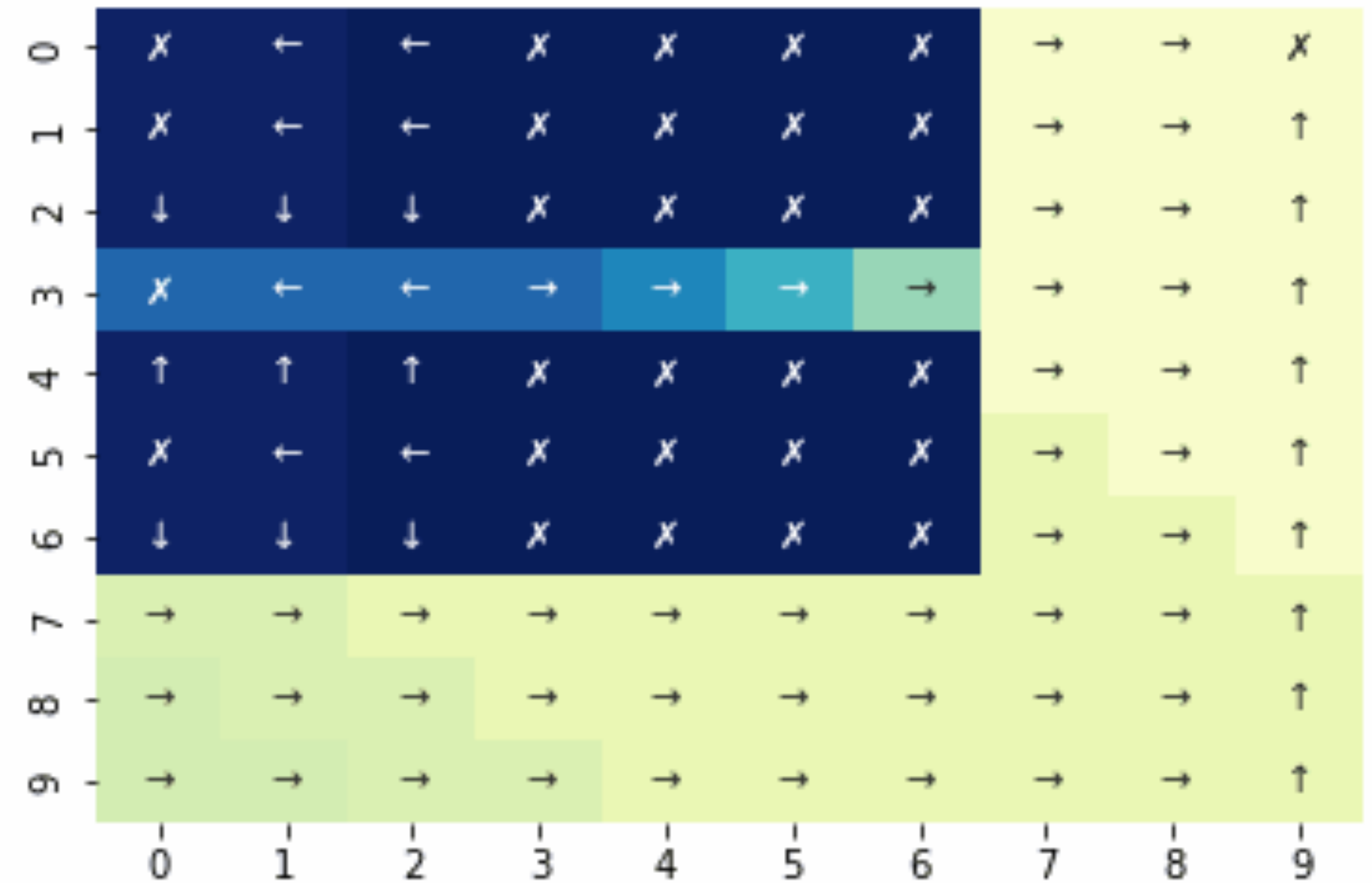
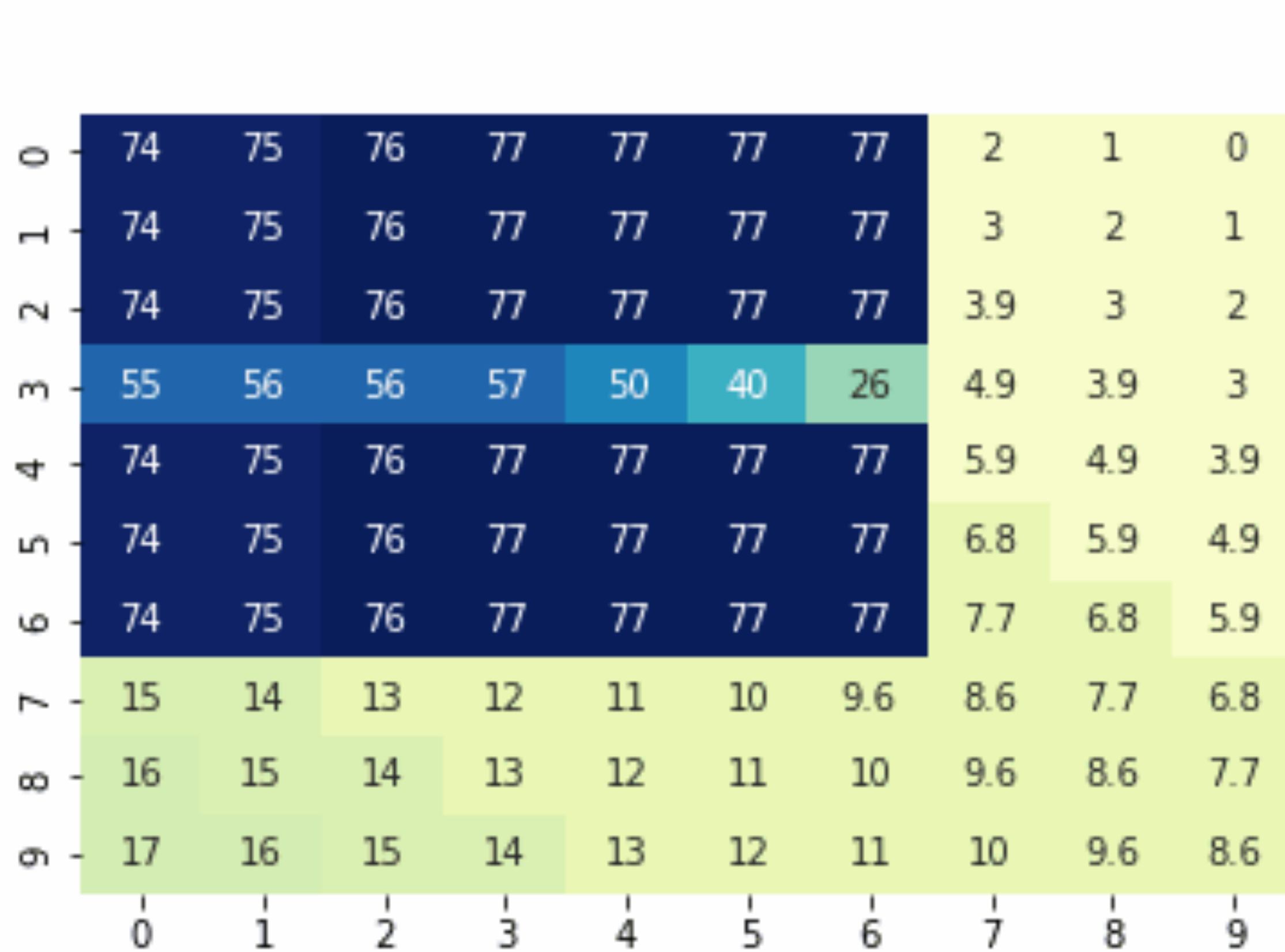
$$\pi^+(s) = \arg \min_a [c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')]$$

Init with some policy π

Iter: 0

0	-	→	→	→	→	→	→	→	→	→	↑
1	-	→	→	→	→	→	→	→	→	→	↑
2	-	→	→	→	→	→	→	→	→	→	↑
3	-	→	→	→	→	→	→	→	→	→	↑
4	-	→	→	→	→	→	→	→	→	→	↑
5	-	→	→	→	→	→	→	→	→	→	↑
6	-	→	→	→	→	→	→	→	→	→	↑
7	-	→	→	→	→	→	→	→	→	→	↑
8	-	→	→	→	→	→	→	→	→	→	↑
9	-	→	→	→	→	→	→	→	→	→	↑
		0	1	2	3	4	5	6	7	8	9

Iteration 1



$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')$$

$$\pi^+(s) = \arg \min_a c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')$$

Policy Iteration

Iter: 0

0	-	0	0	0	0	0	0	0	0	0	0
1	-	0	0	0	0	0	0	0	0	0	0
2	-	0	0	0	0	0	0	0	0	0	0
3	-	0	0	0	0	0	0	0	0	0	0
4	-	0	0	0	0	0	0	0	0	0	0
5	-	0	0	0	0	0	0	0	0	0	0
6	-	0	0	0	0	0	0	0	0	0	0
7	-	0	0	0	0	0	0	0	0	0	0
8	-	0	0	0	0	0	0	0	0	0	0
9	-	0	0	0	0	0	0	0	0	0	0
		0	1	2	3	4	5	6	7	8	9

0	-	→	→	→	→	→	→	→	→	→	↑
1	-	→	→	→	→	→	→	→	→	→	↑
2	-	→	→	→	→	→	→	→	→	→	↑
3	-	→	→	→	→	→	→	→	→	→	↑
4	-	→	→	→	→	→	→	→	→	→	↑
5	-	→	→	→	→	→	→	→	→	→	↑
6	-	→	→	→	→	→	→	→	→	→	↑
7	-	→	→	→	→	→	→	→	→	→	↑
8	-	→	→	→	→	→	→	→	→	→	↑
9	-	→	→	→	→	→	→	→	→	→	↑
		0	1	2	3	4	5	6	7	8	9

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')$$

$$\pi^+(s) = \arg \min_a c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')$$



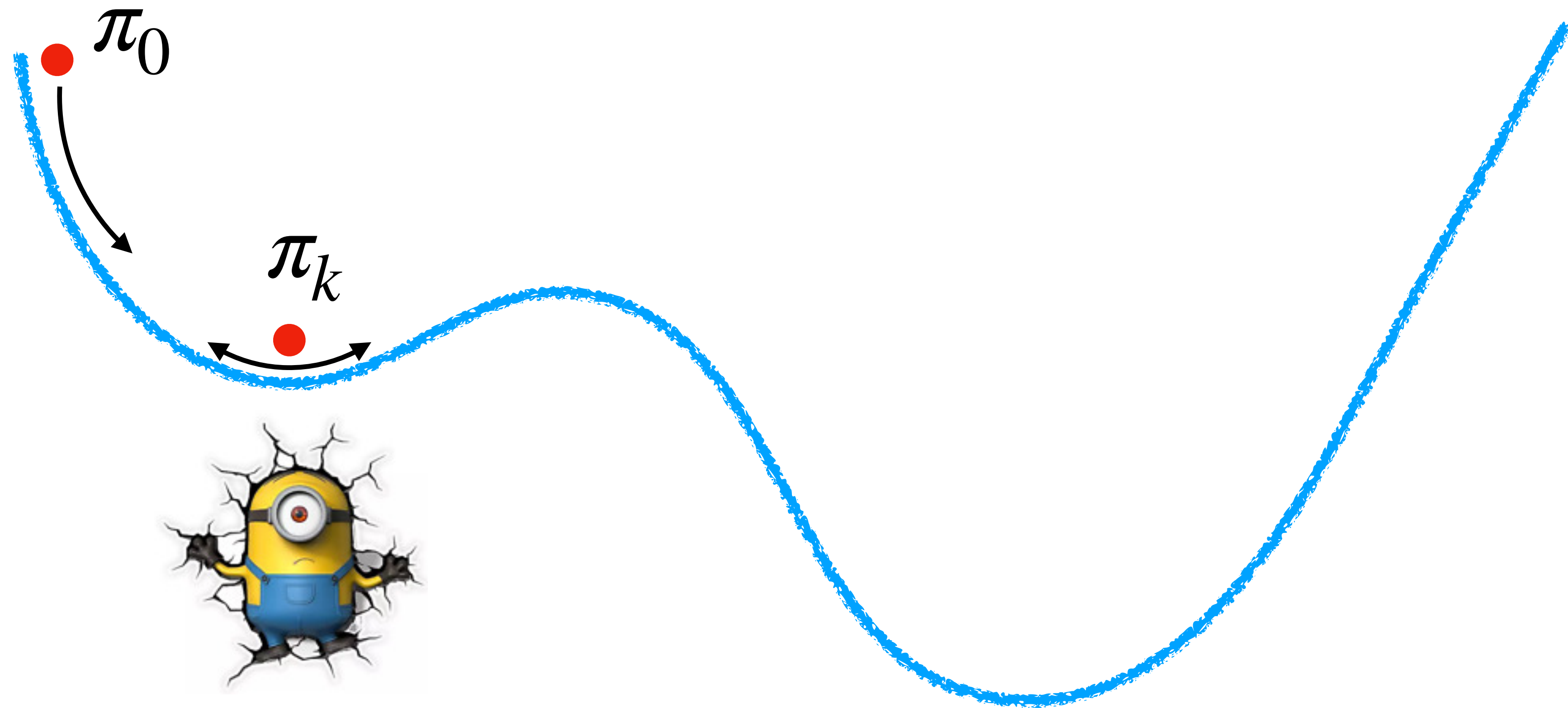
So many questions

Q1. Does policy iteration converge to the optimal policy?

Q2. Does it converge faster than value iteration?

Q3. Can we bootstrap policy evaluation?

Q1. Does policy iteration converge to optimal policy?



Stuck at local minima??

π^*

Q1. Does policy iteration converge to optimal policy?

Proof has 2 step argument

Step 1: Policy iteration *monotonically* improves

Step 2: Once it reaches the optimal policy, it does not change

Performance Difference Lemma



Q1. Does policy iteration converge to optimal policy?

2 step argument

✓ Step 1: Policy iteration *monotonically* improves

All advantages ≤ 0 implies monotonic performance improvement

✓ Step 2: Once it reaches the optimal policy, it does not change

Advantage ≥ 0 for the optimal policy

PDL answers *all* ...

In Imitation Learning

In Model Free Reinforcement Learning

In Model Based Reinforcement Learning

How it partners with online learning



So many questions

✓ Q1. Does policy iteration converge to the optimal policy?

Q2. Does it converge faster than value iteration?

Q3. Can we bootstrap policy evaluation?



So many questions

✓ Q1. Does policy iteration converge to the optimal policy?

(ツ)/

Q2. Does it converge faster than value iteration?

Empirically yes, but no rigorous theory about when ..

Q3. Can we bootstrap policy evaluation?



So many questions

✓ Q1. Does policy iteration converge to the optimal policy?

(ツ)/

Q2. Does it converge faster than value iteration?

Empirically yes, but no rigorous theory about when ..



Q3. Can we bootstrap policy evaluation?

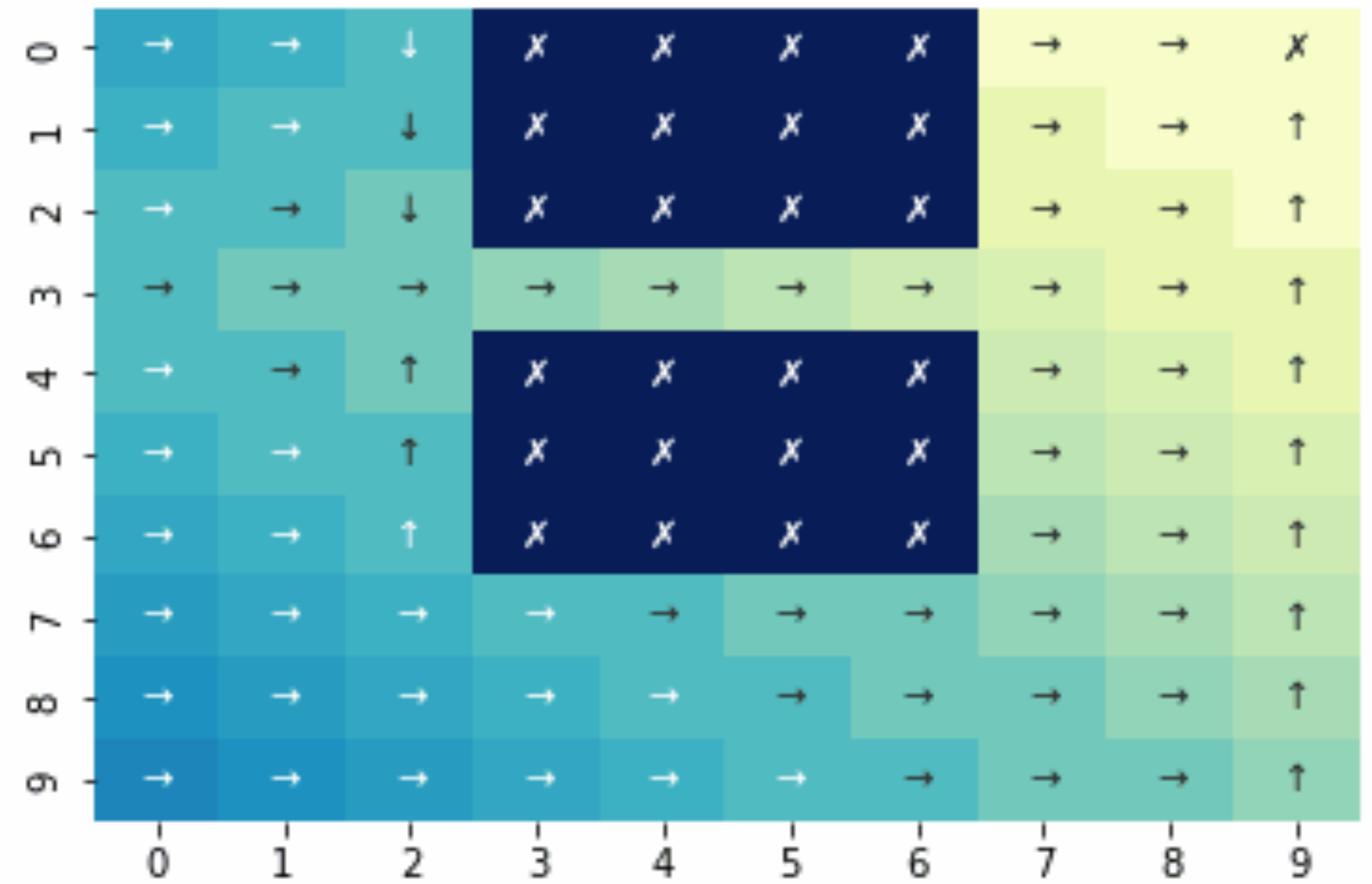
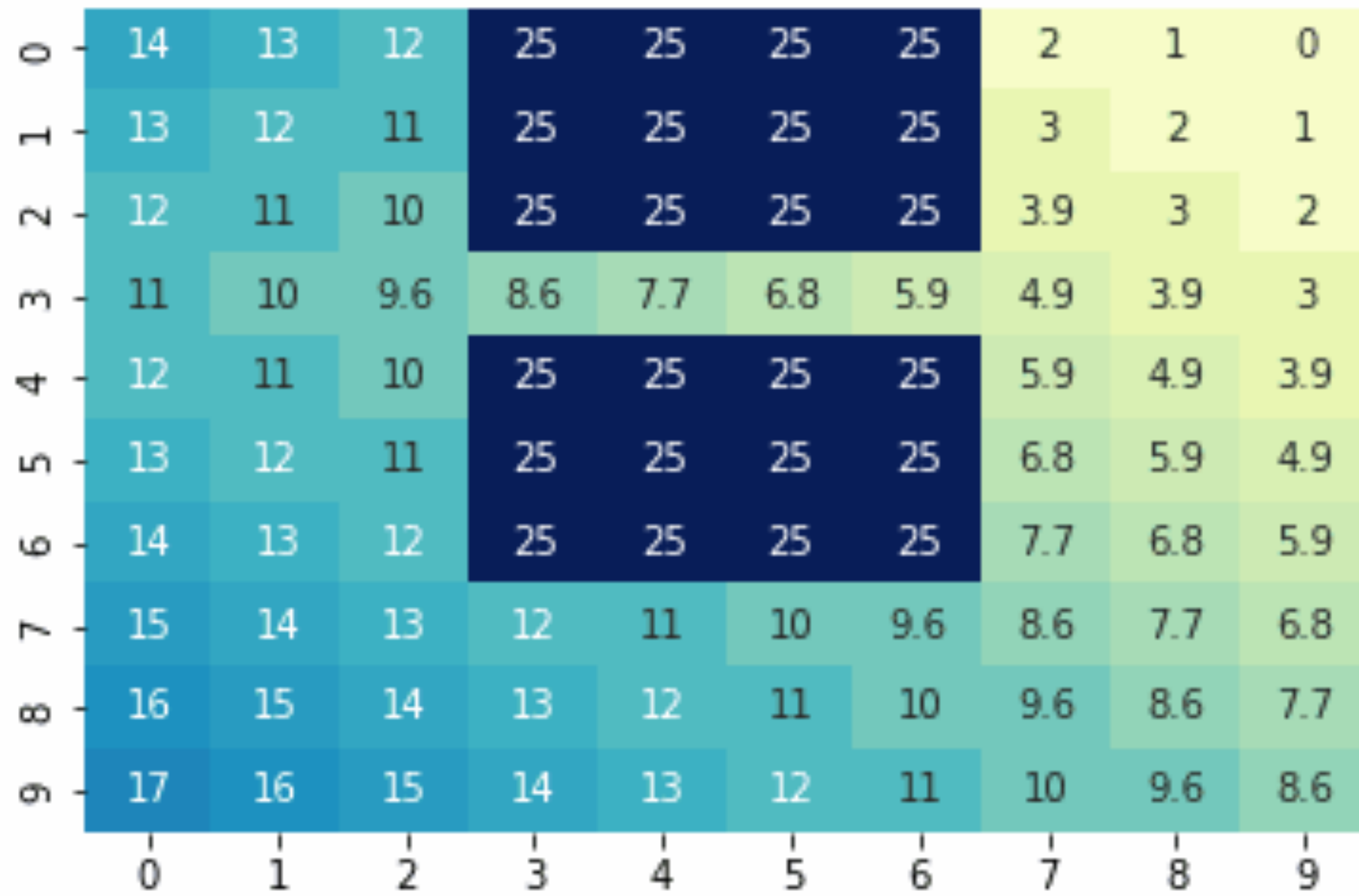
Yes \Rightarrow Modified policy iteration

Messing with MDPs

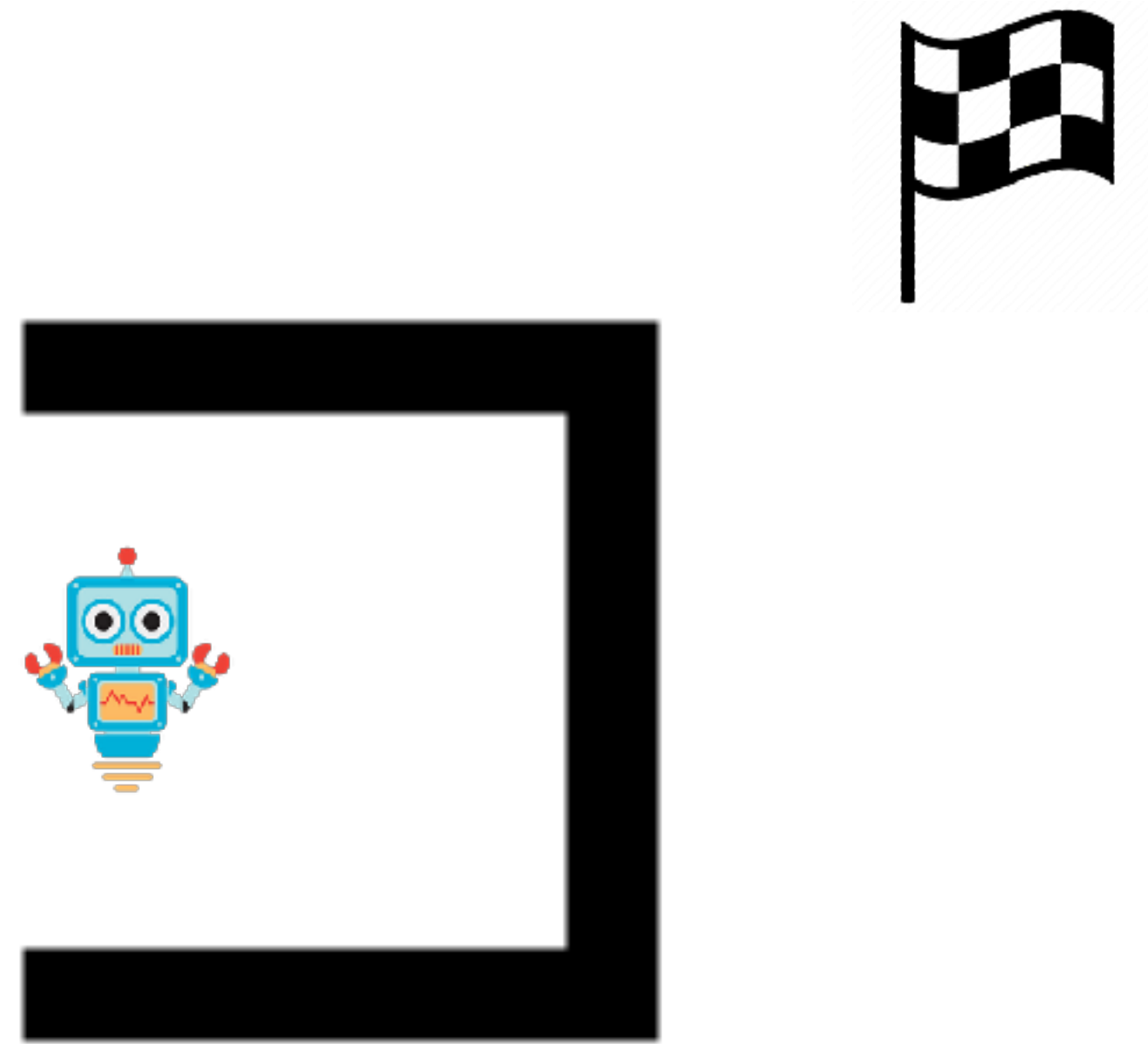


What happens as you increase the slipperiness of bridge?

p_slip: 0.0



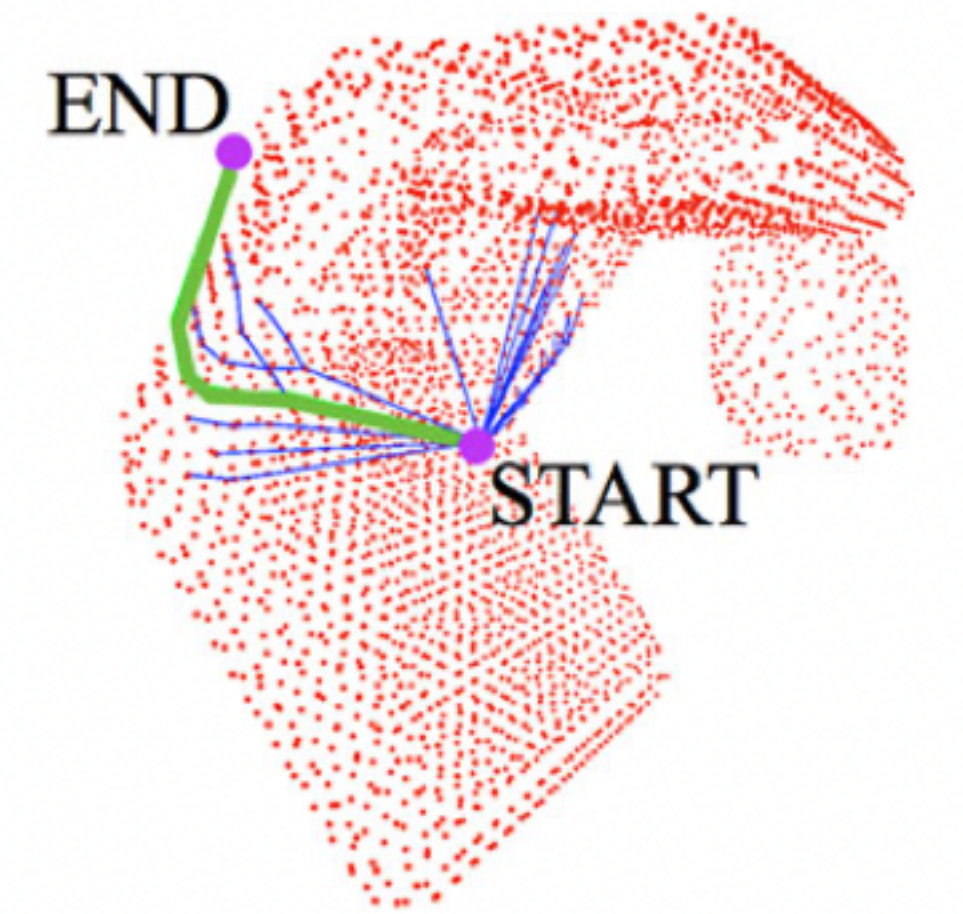
What happens if you change the cost function?

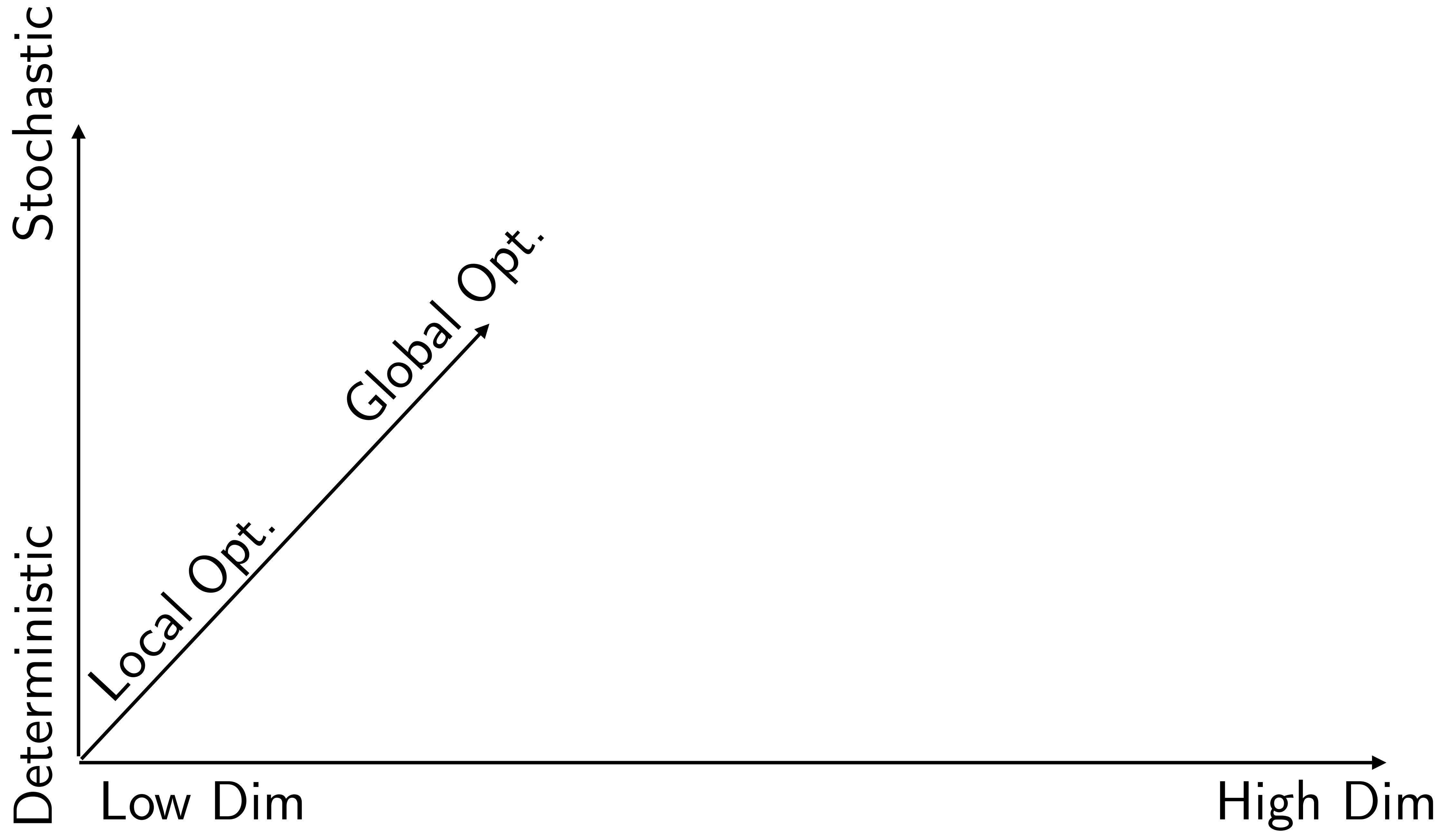


$$C(s) = ||s - s_{goal}||$$

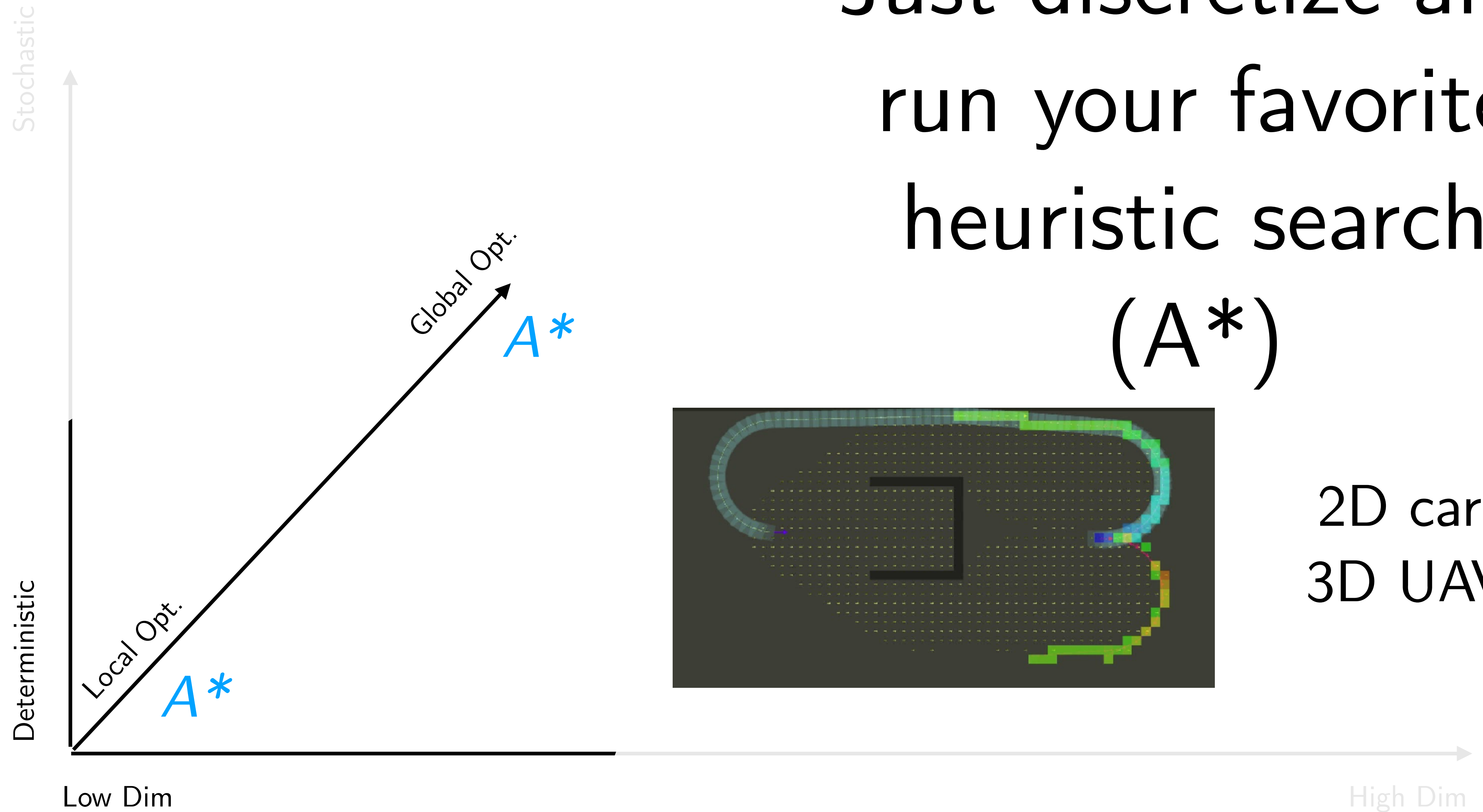
How do you convert a HARD MDP into an EASY one?

Solving continuous MDPs

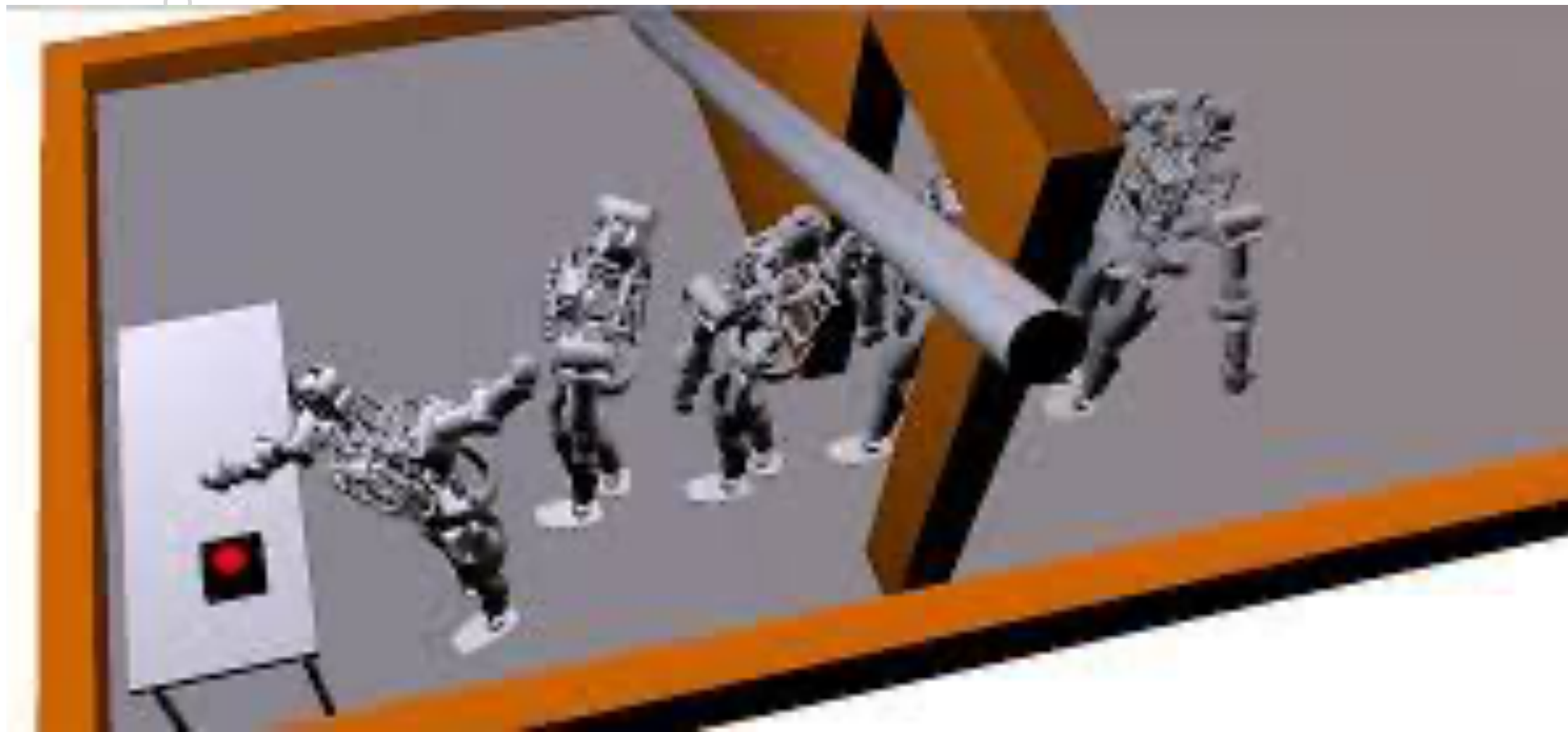




Just discretize and
run your favorite
heuristic search
(A^*)

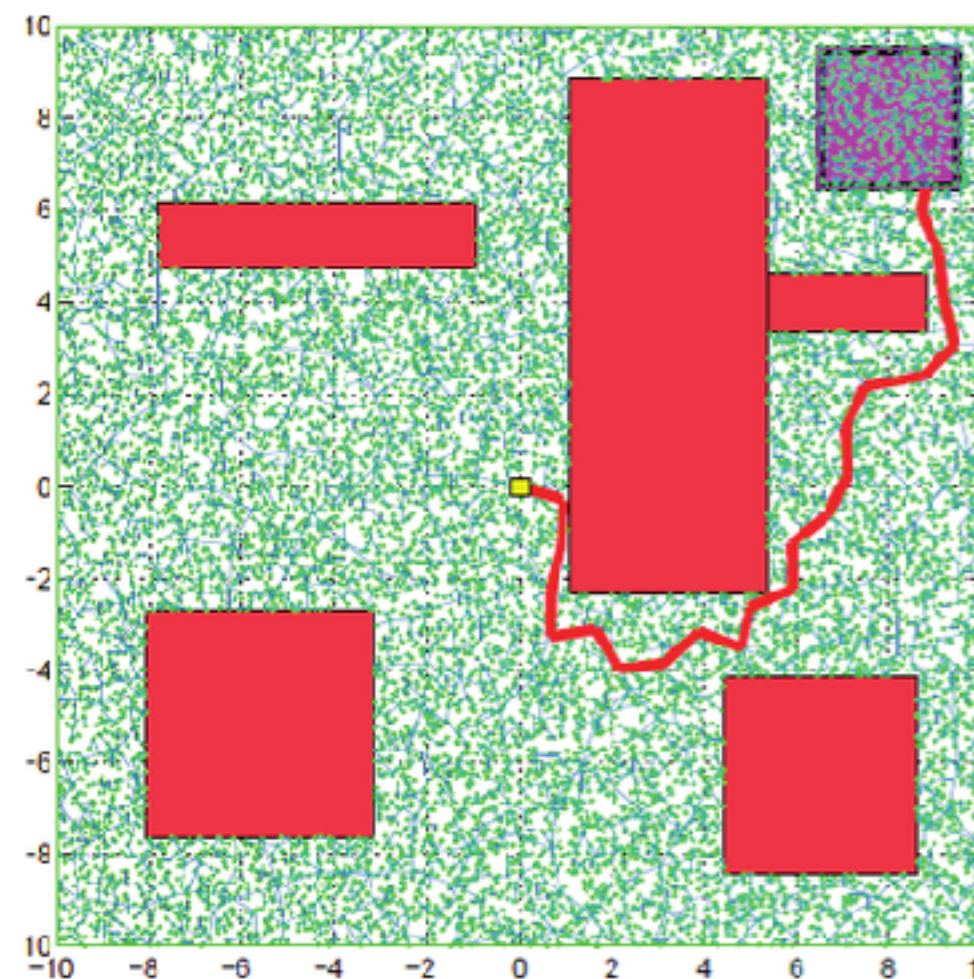
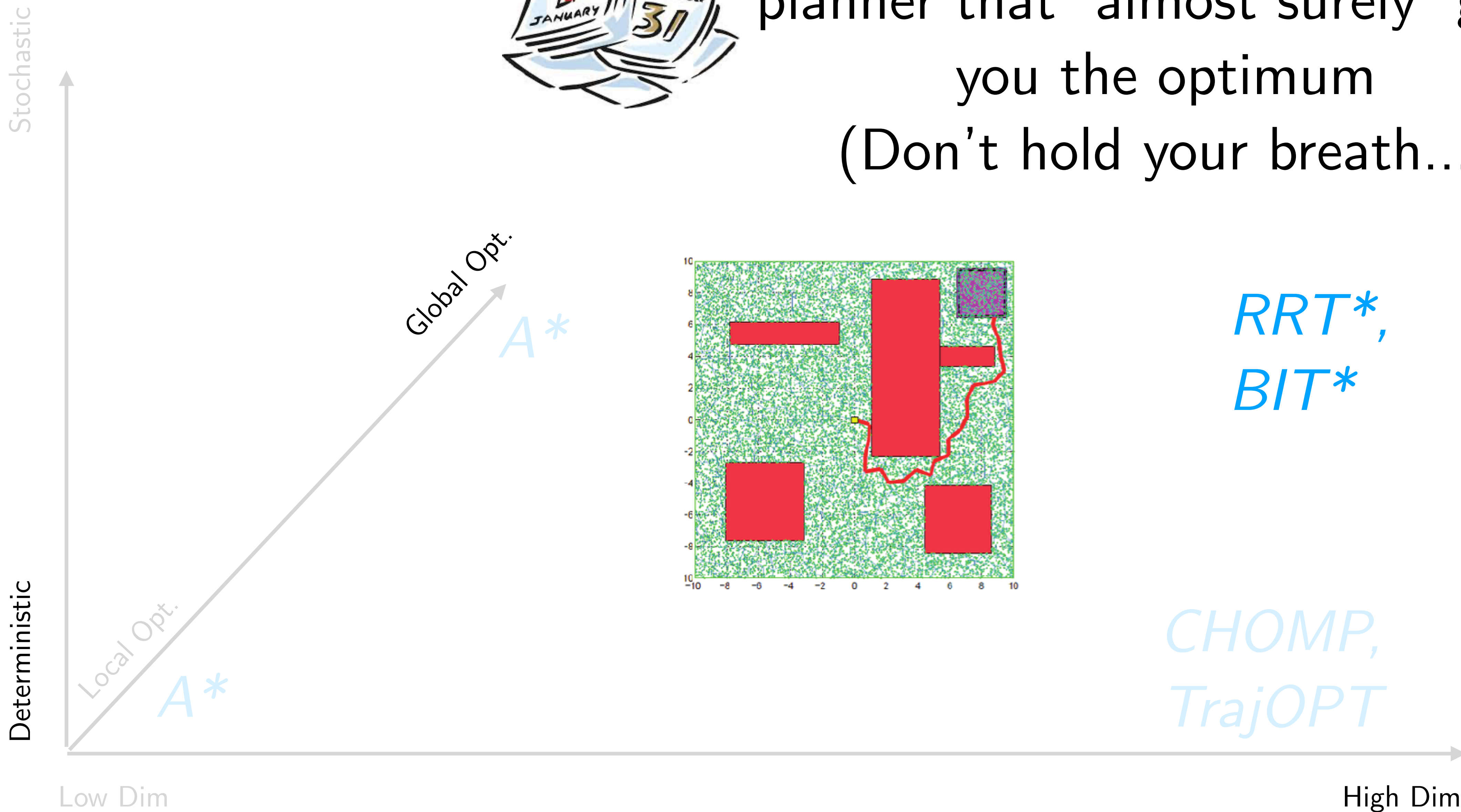


Run some
sequential convex
trajectory optim.

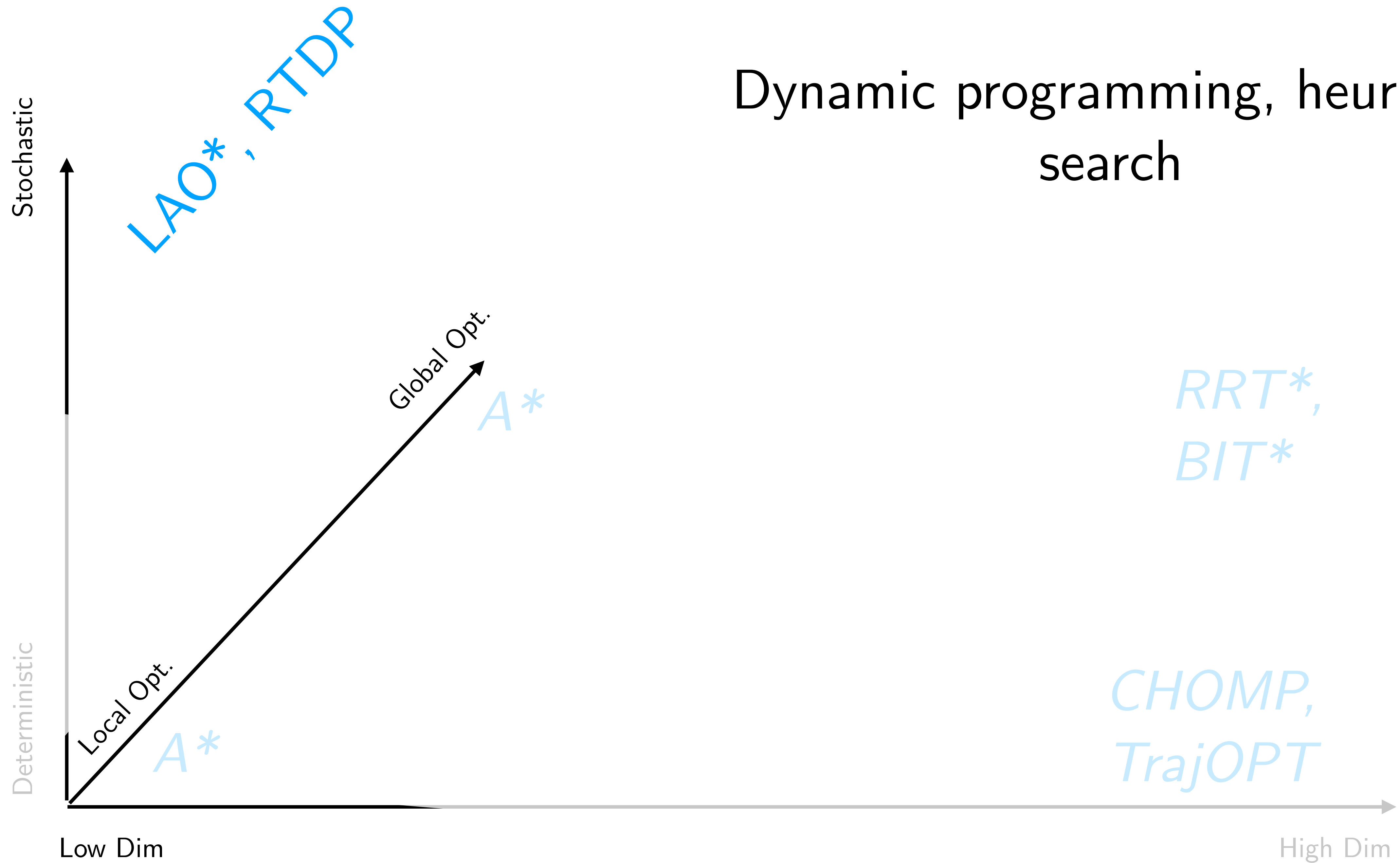


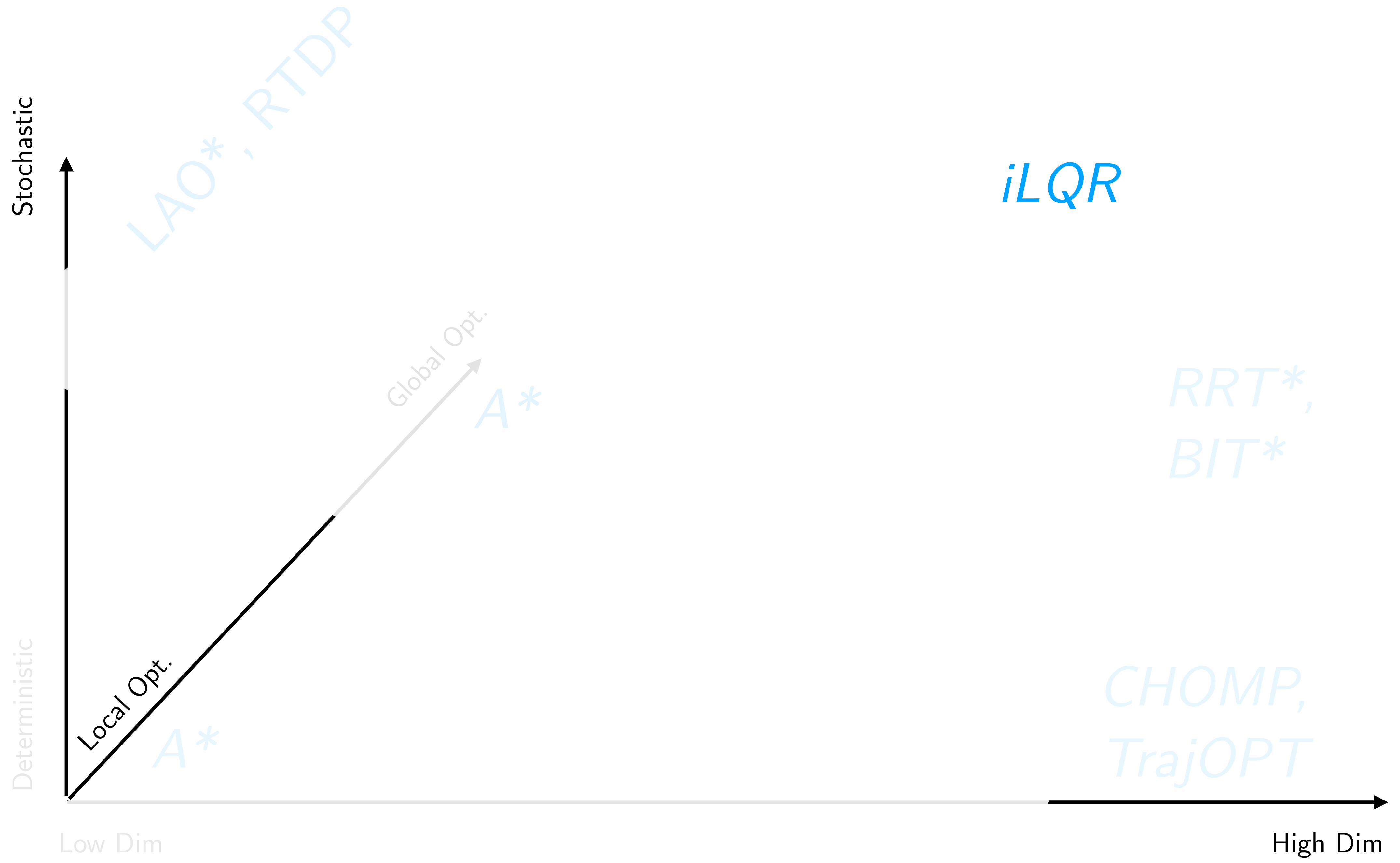


Call a (slow) probabilistic planner that “almost surely” gives you the optimum
(Don't hold your breath...)



Dynamic programming, heuristic search





GIVE UP!

