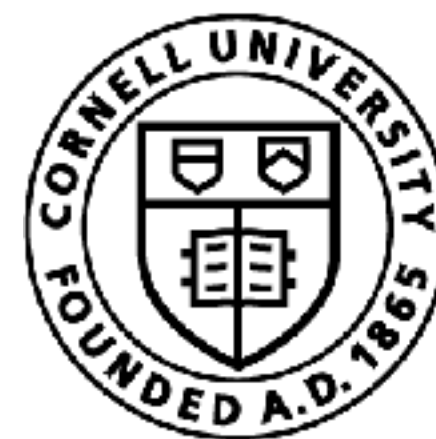


Temporal Difference & Q Learning

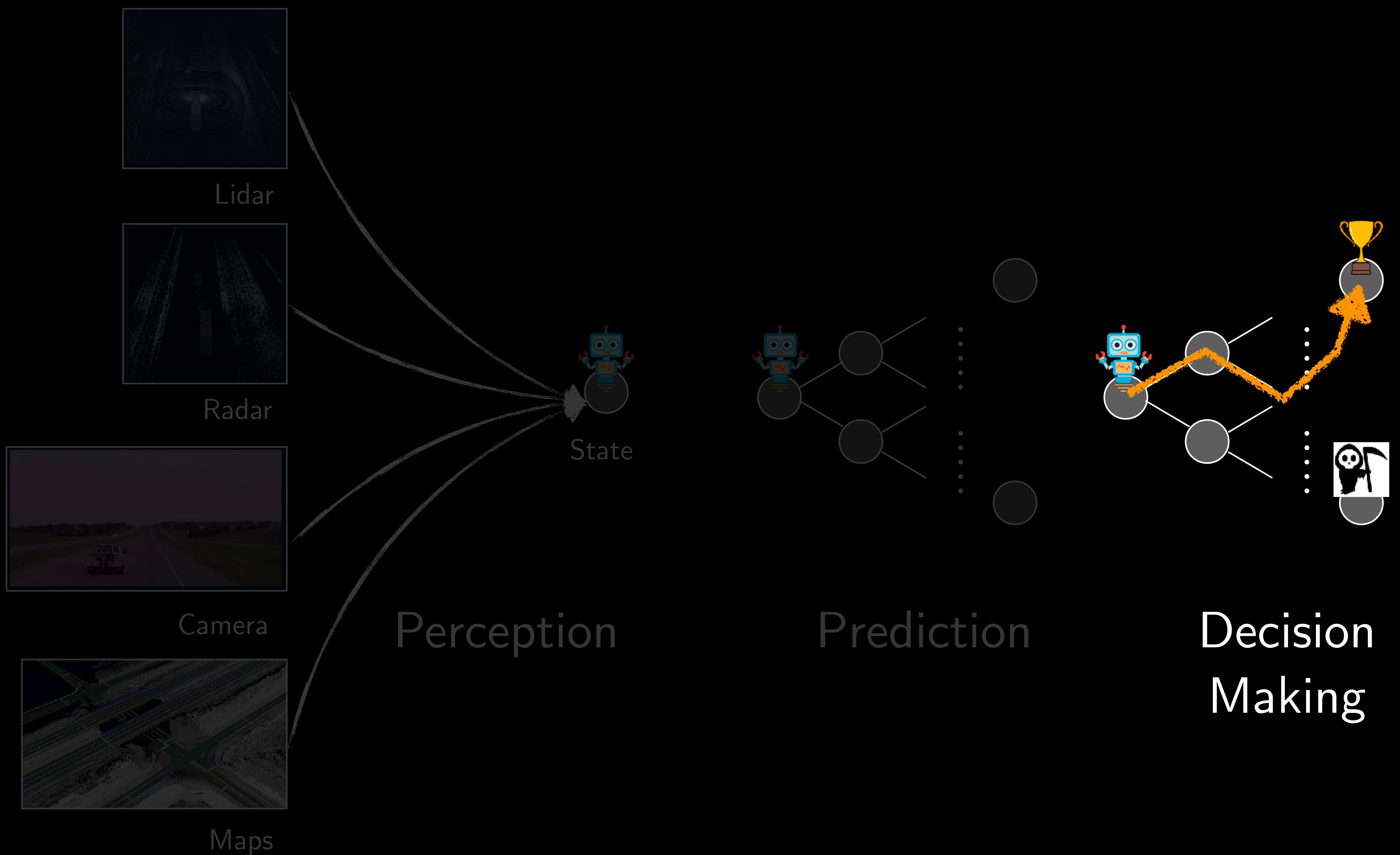
Sanjiban Choudhury

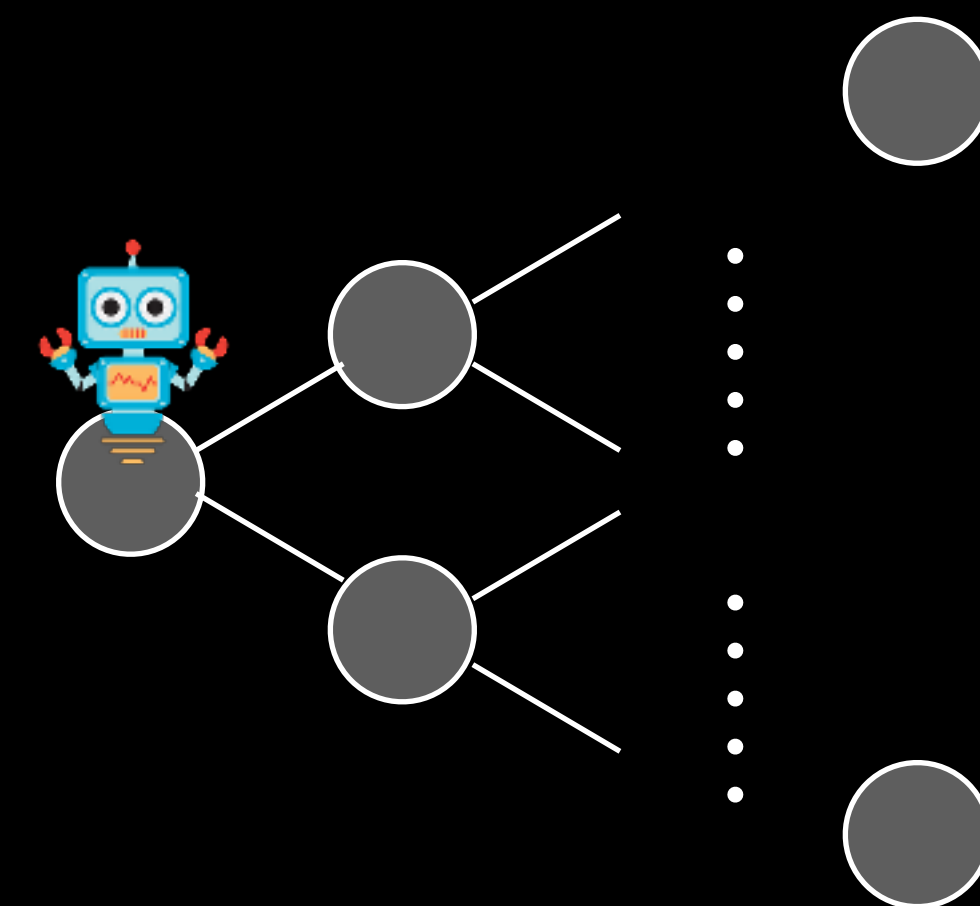
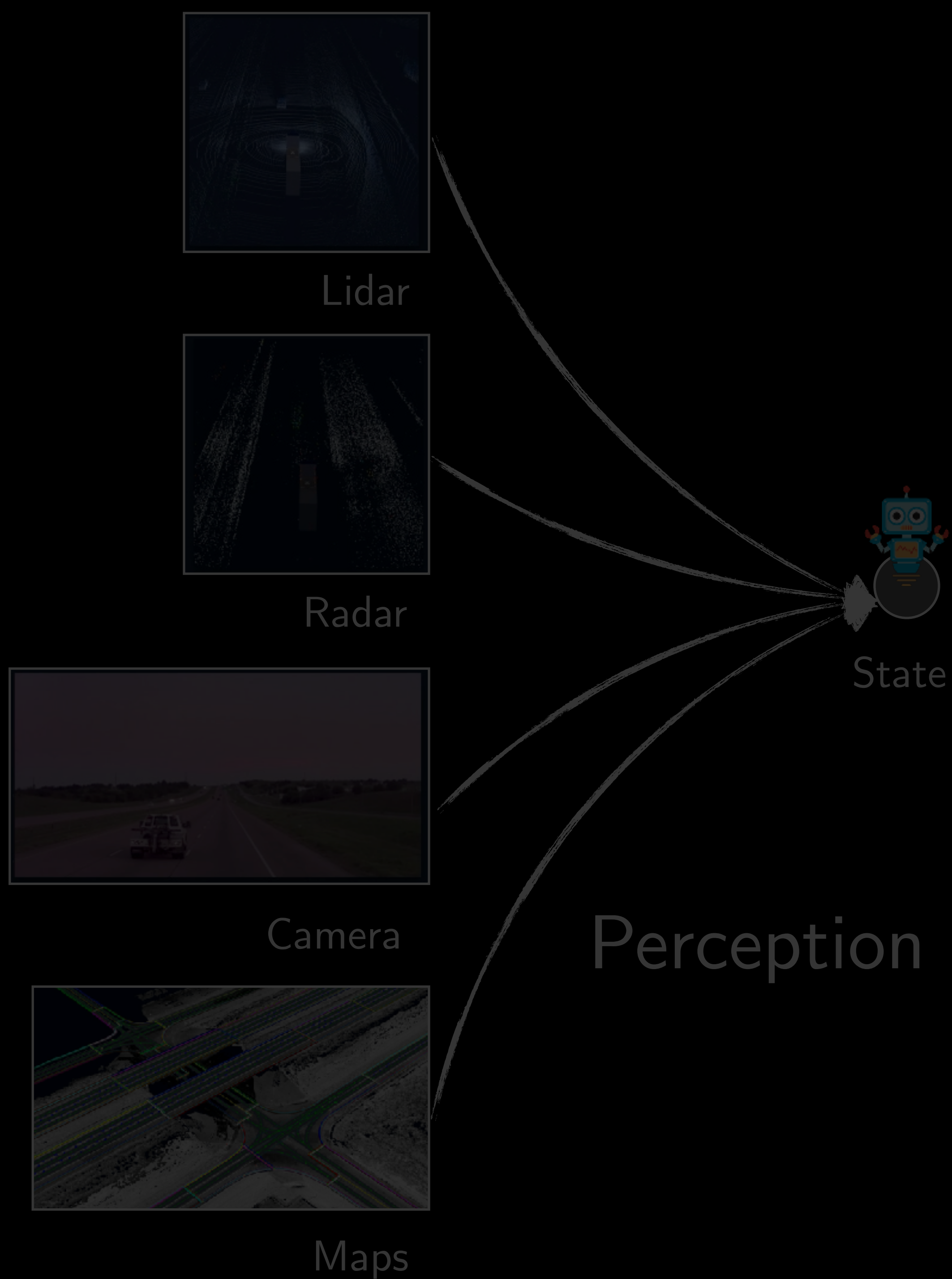


Cornell Bowers C-IS
Computer Science

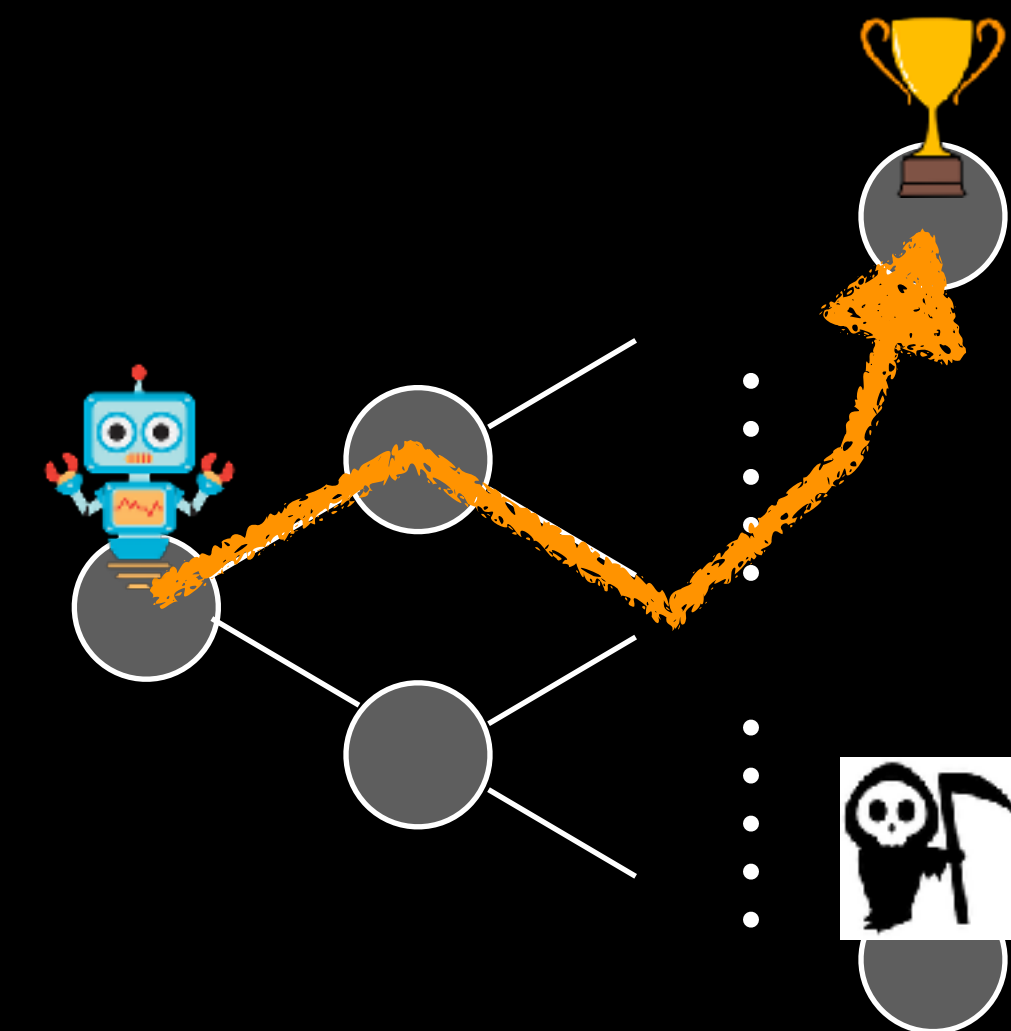
The story thus far ...







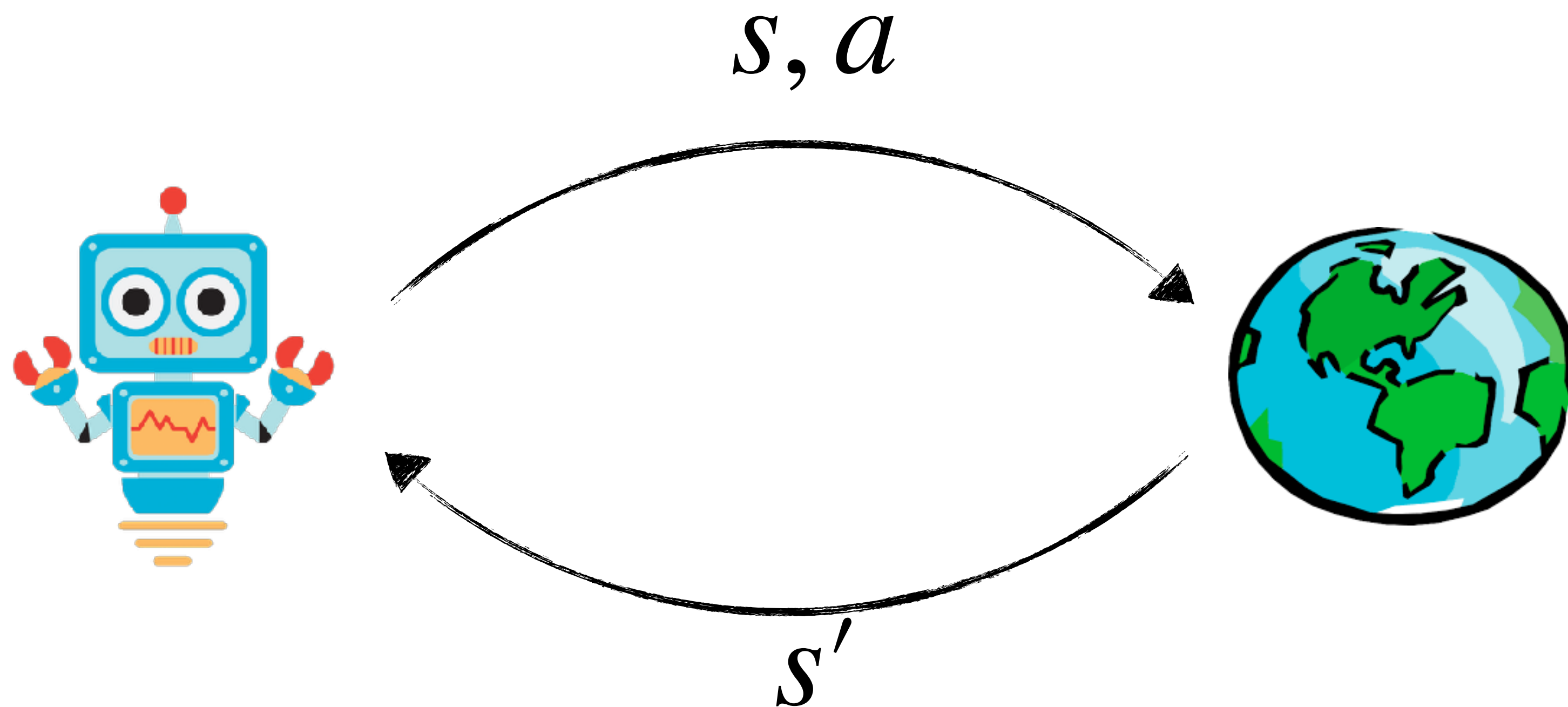
Prediction



Decision
Making

What if the transitions are unknown?

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





Activity!



Think-Pair-Share

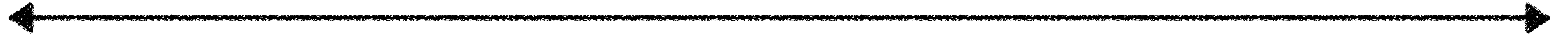
Think (30 sec): What is the MDP $\langle S, A, C, T \rangle$ for this robot? Is the transition T known?

Pair: Find a partner

Share (45 sec): Partners exchange ideas

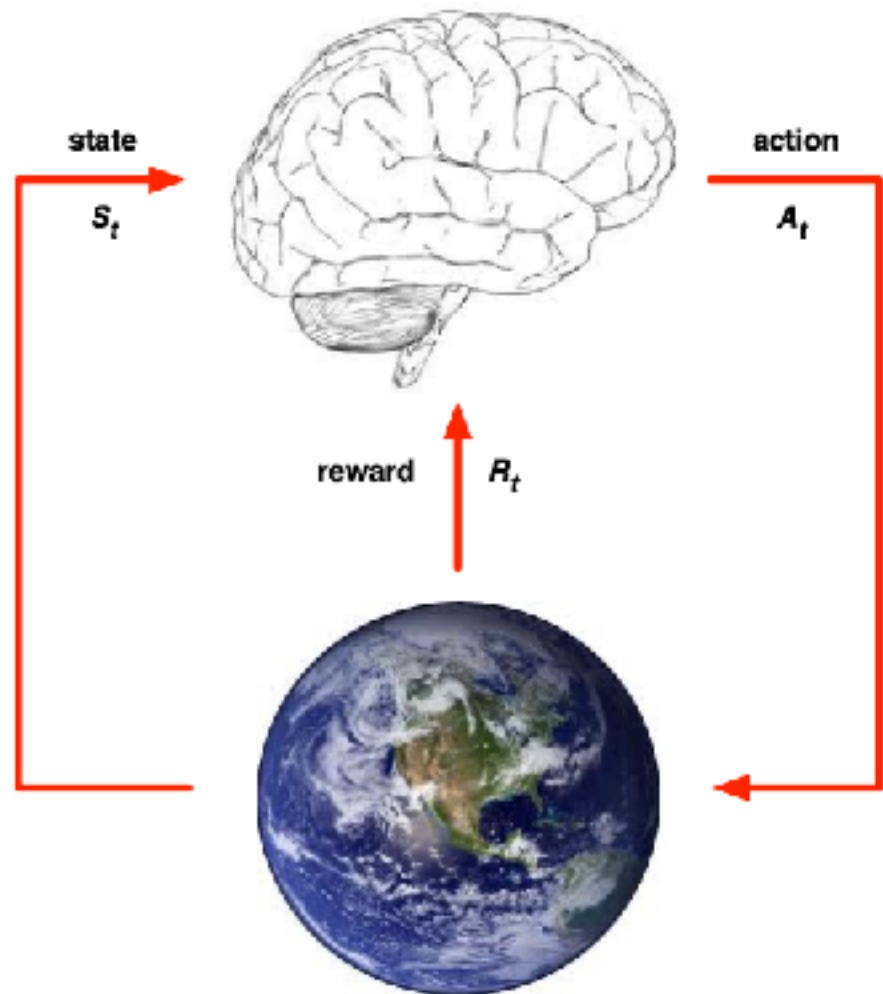


Model-Based OR Model Free?



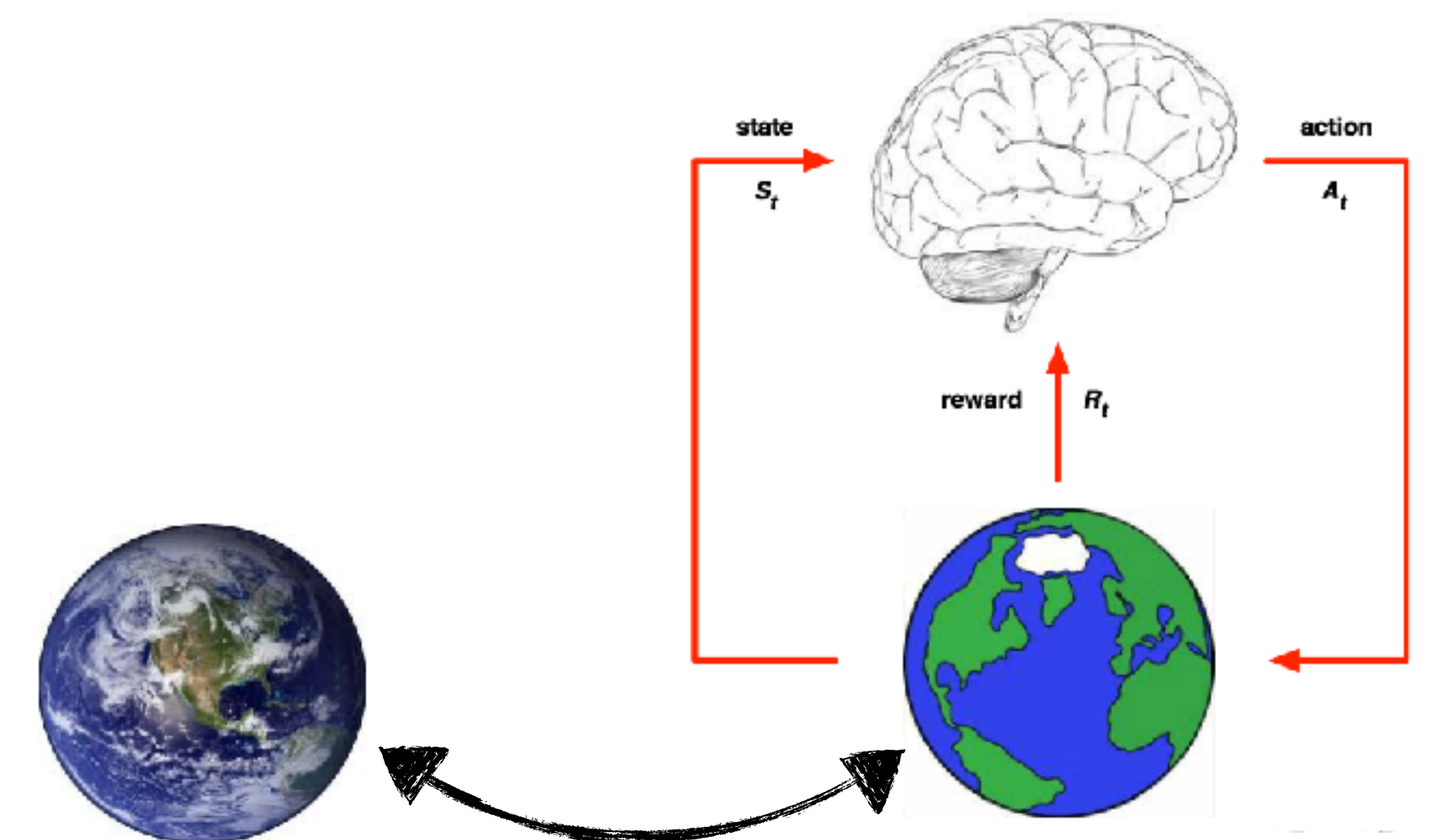
Model Free

Directly learn
 π or $Q(s, a)$



Model Based

Learn a model
 $T(s' | s, a)$, plan with
model to find π

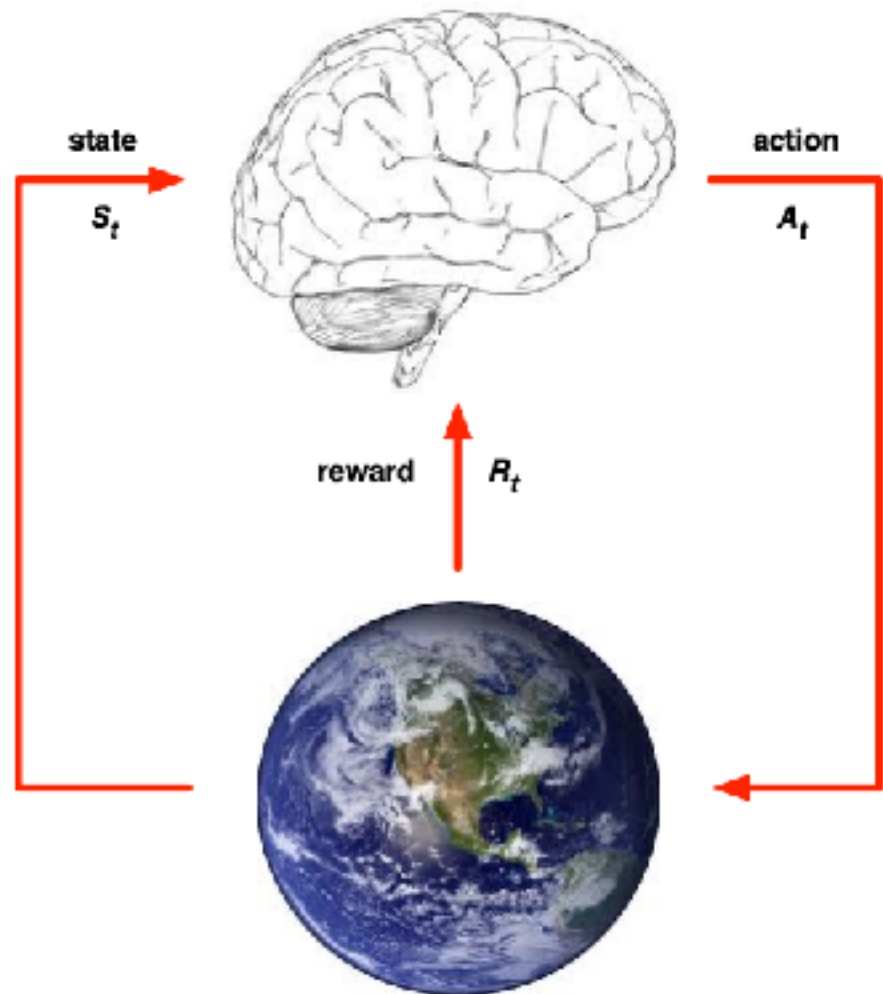


Model-Based OR Model Free?



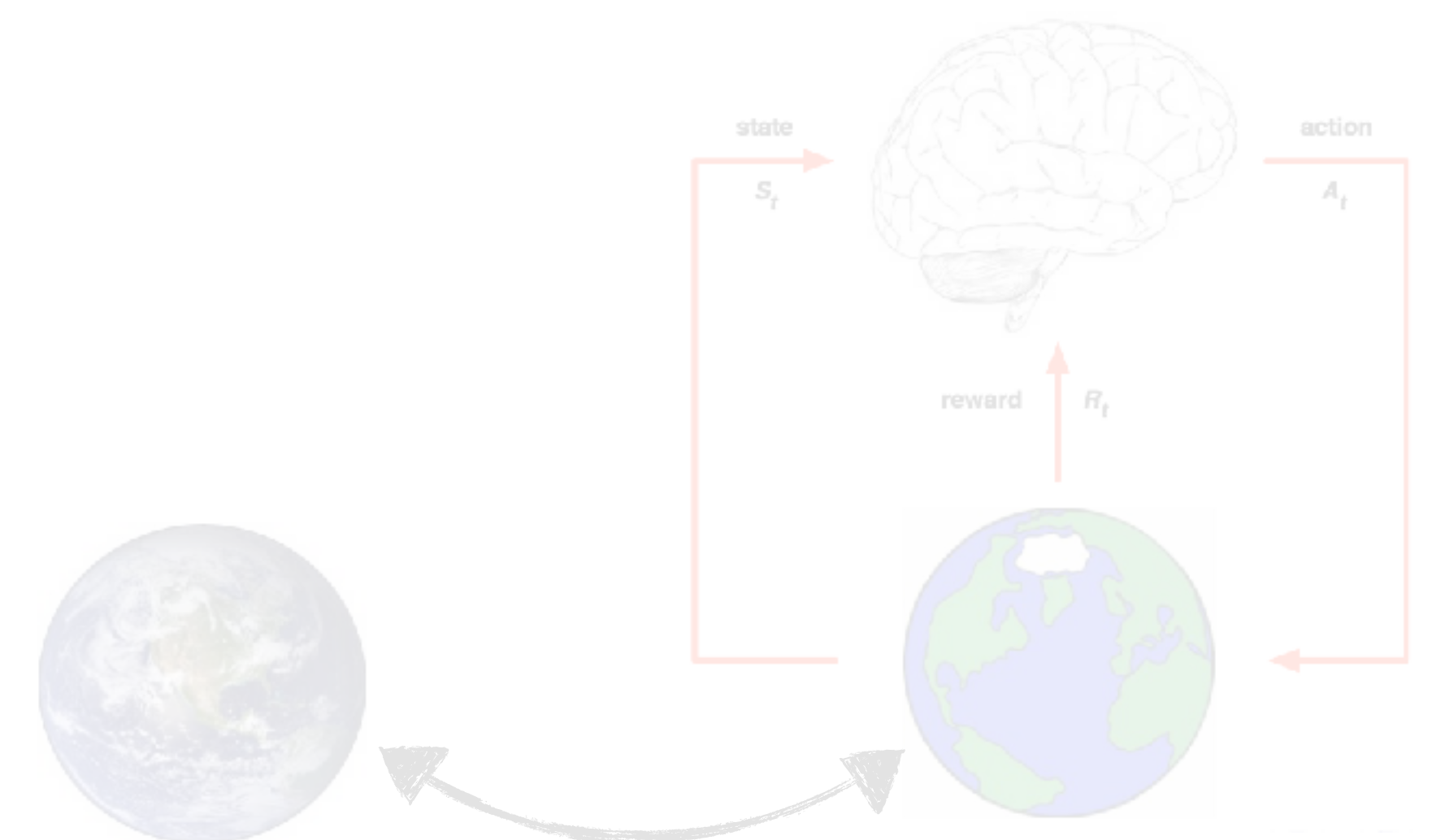
Model Free

Directly learn
 π or $Q(s, a)$

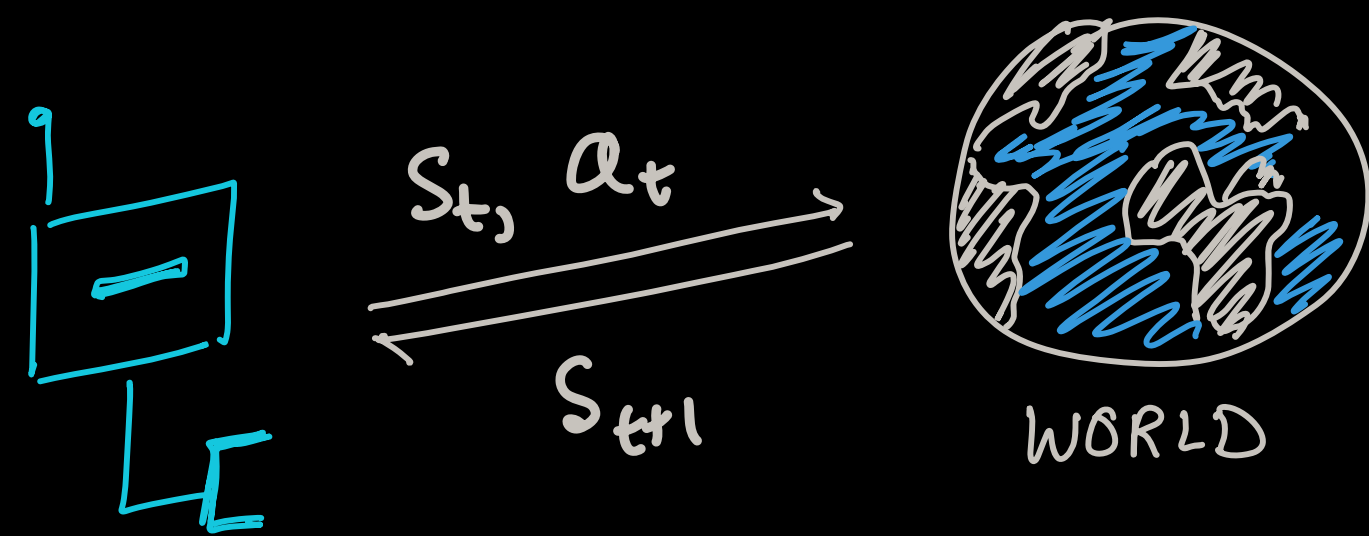


Model Based

Learn a model
 $P(s' | s, a)$, plan with
model to find π



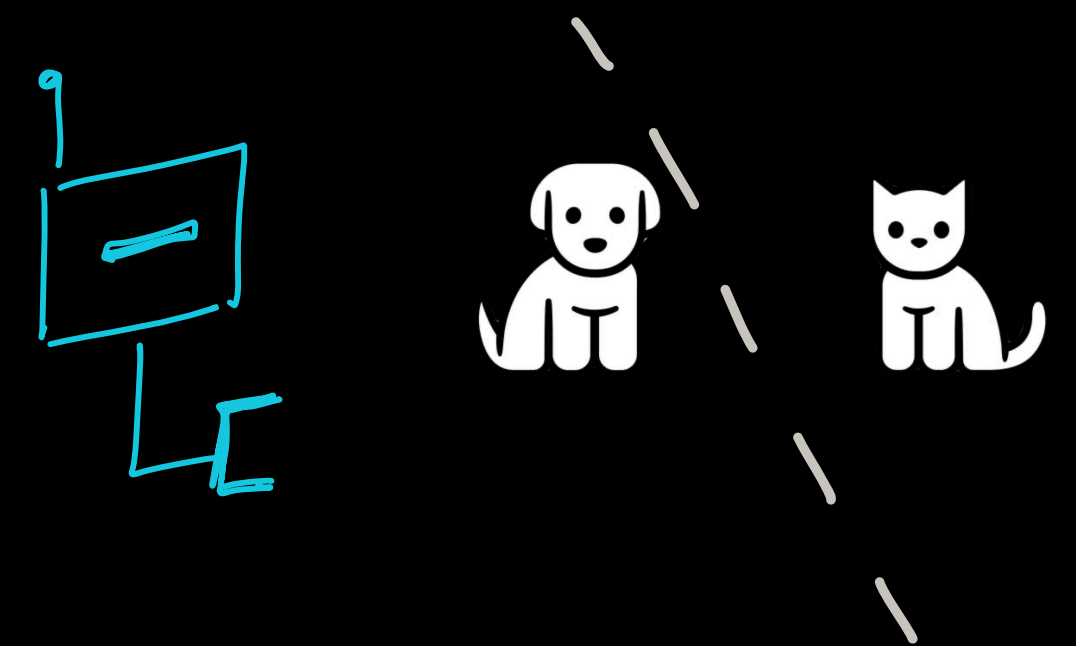
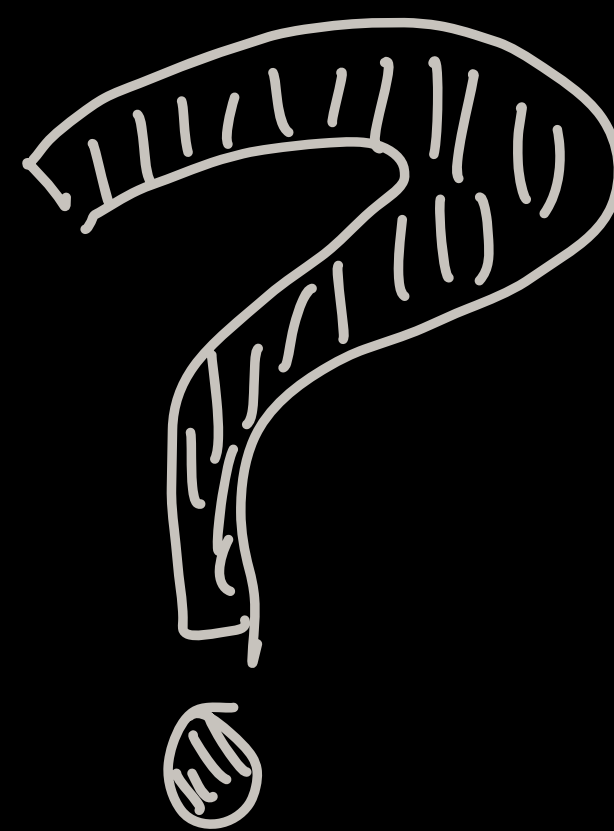
WHAT MAKES



REINFORCEMENT
LEARNING

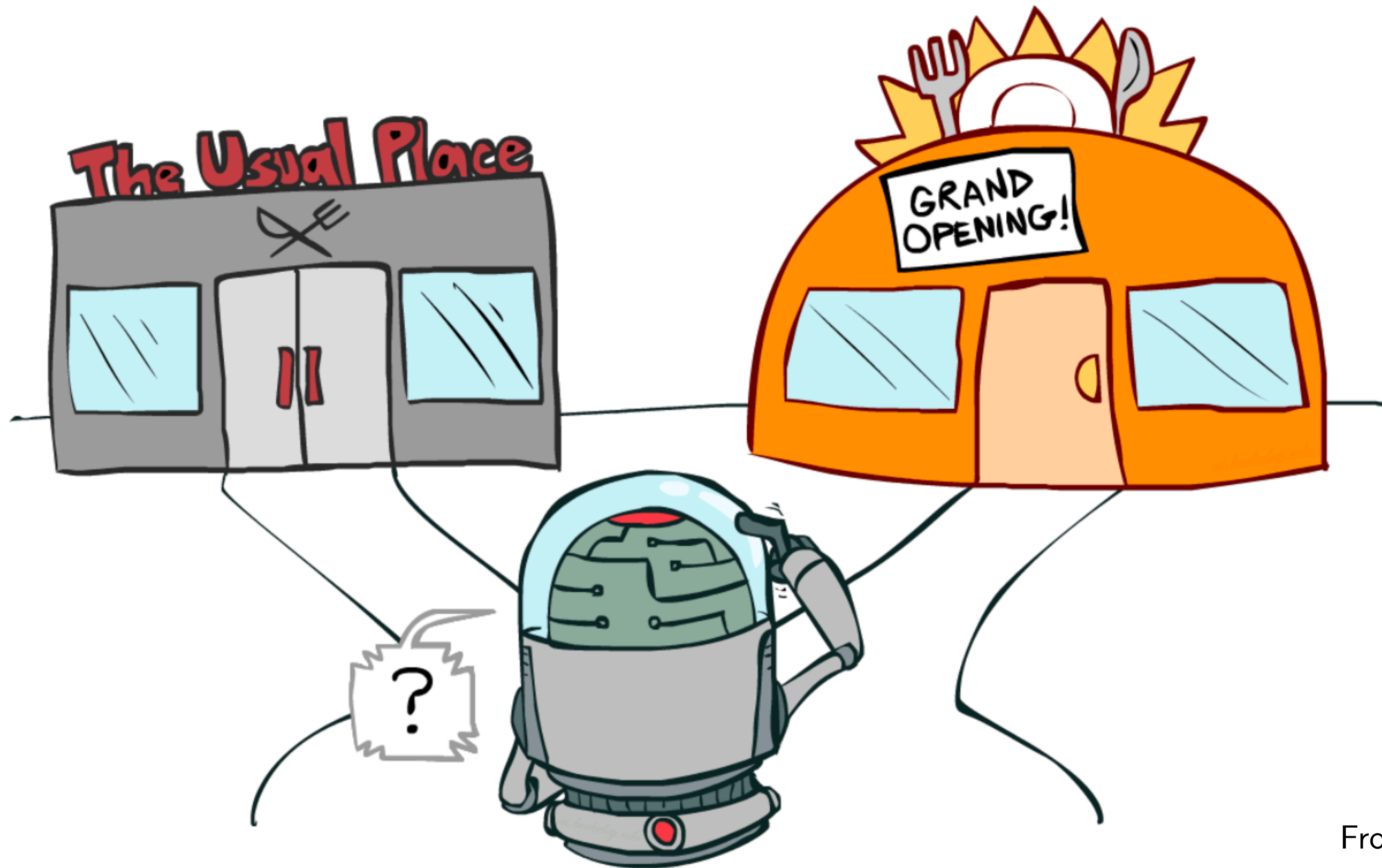
HARDER

THAN



SUPERVISED
LEARNING

Exploration vs Exploitation



From Dan Klein

Doors

a^1



?



+100

a^2



?



+1

a^3

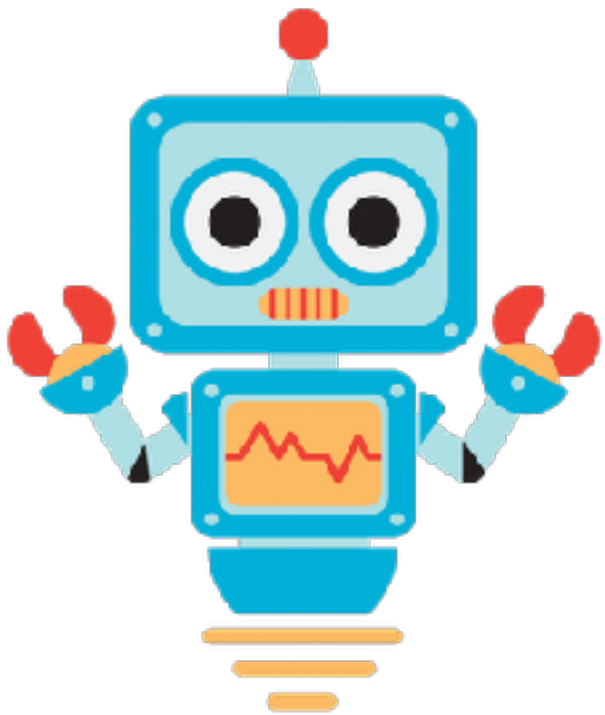


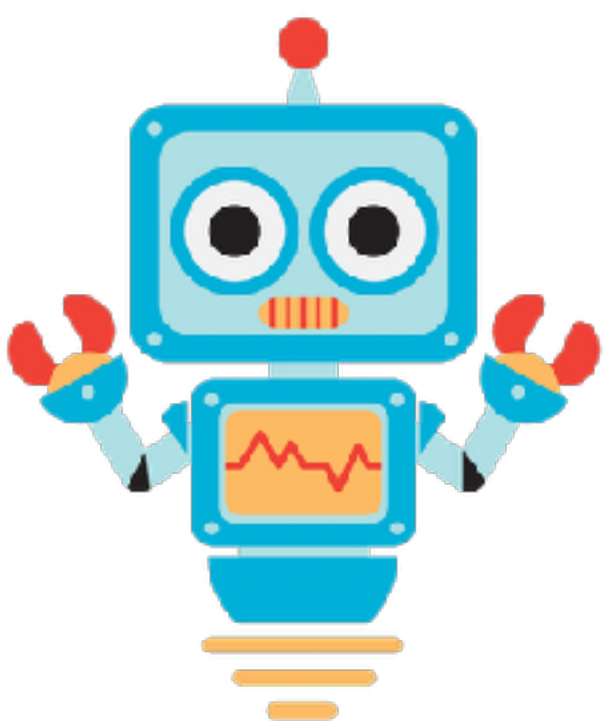
?



-1000

⋮





a^1



a^2



a^3



⋮

Doors

Round 1



Round 2



Round 3



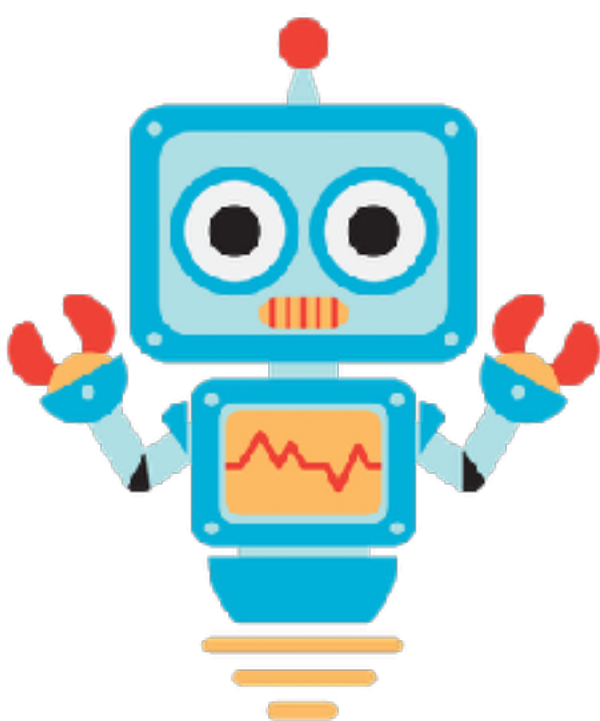
+100



+1



-1000



a^1



a^2



a^3



⋮

Doors

Round 1

Round 2

Round 3



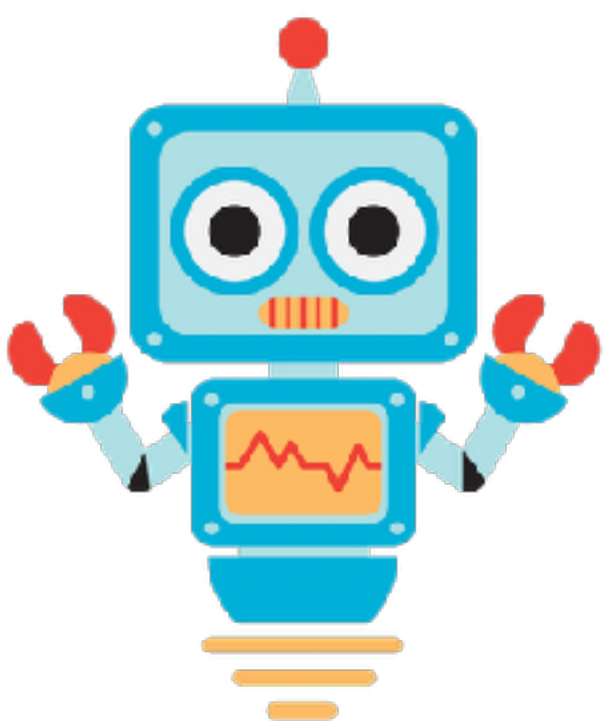
+100



+1



-1000



	Doors
a^1	
a^2	
a^3	
	⋮

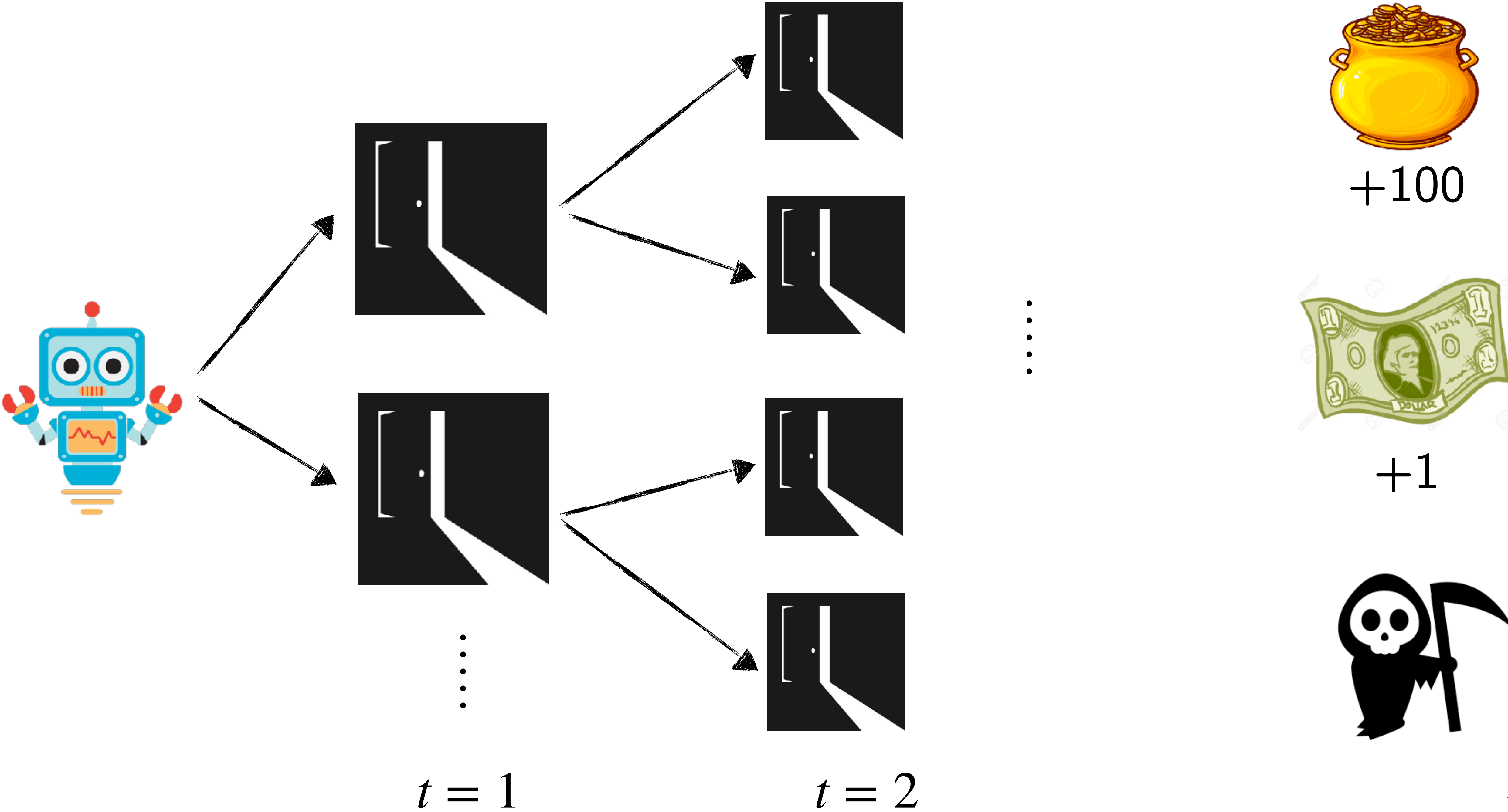
Round 1 Round 2 Round 3

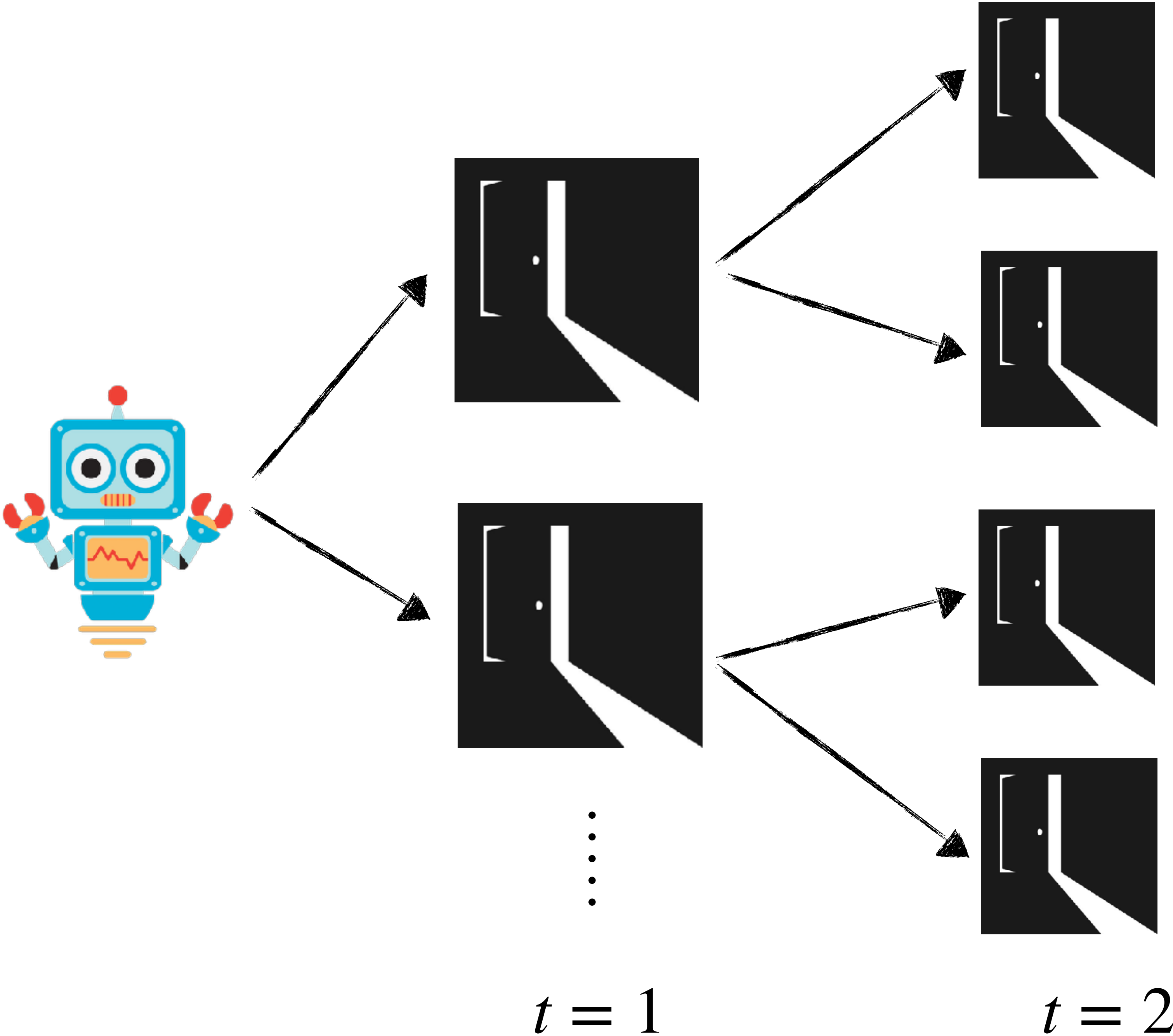


How do we explore/
exploit when picking
doors?

What if we played the
game over multiple time
steps?

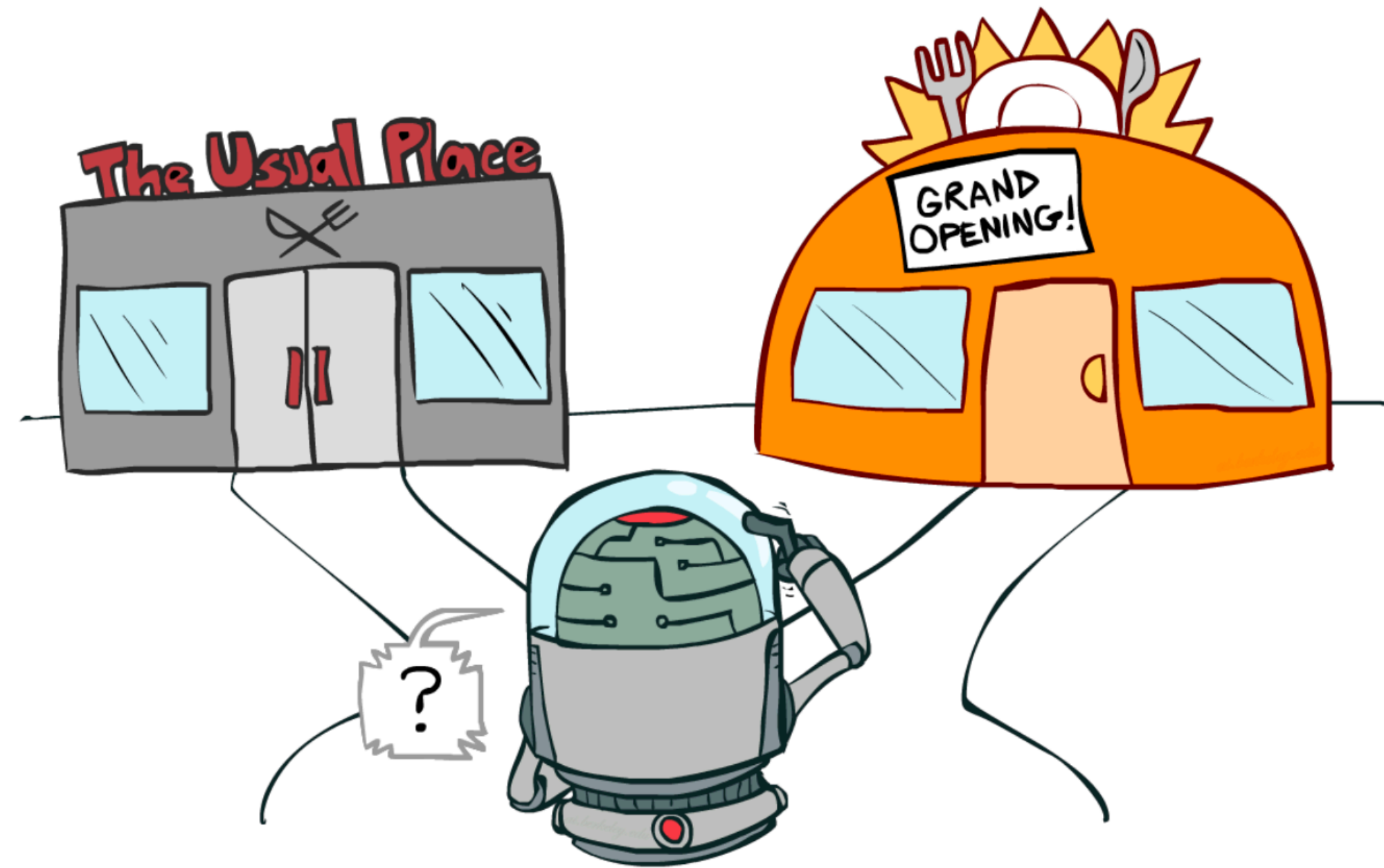




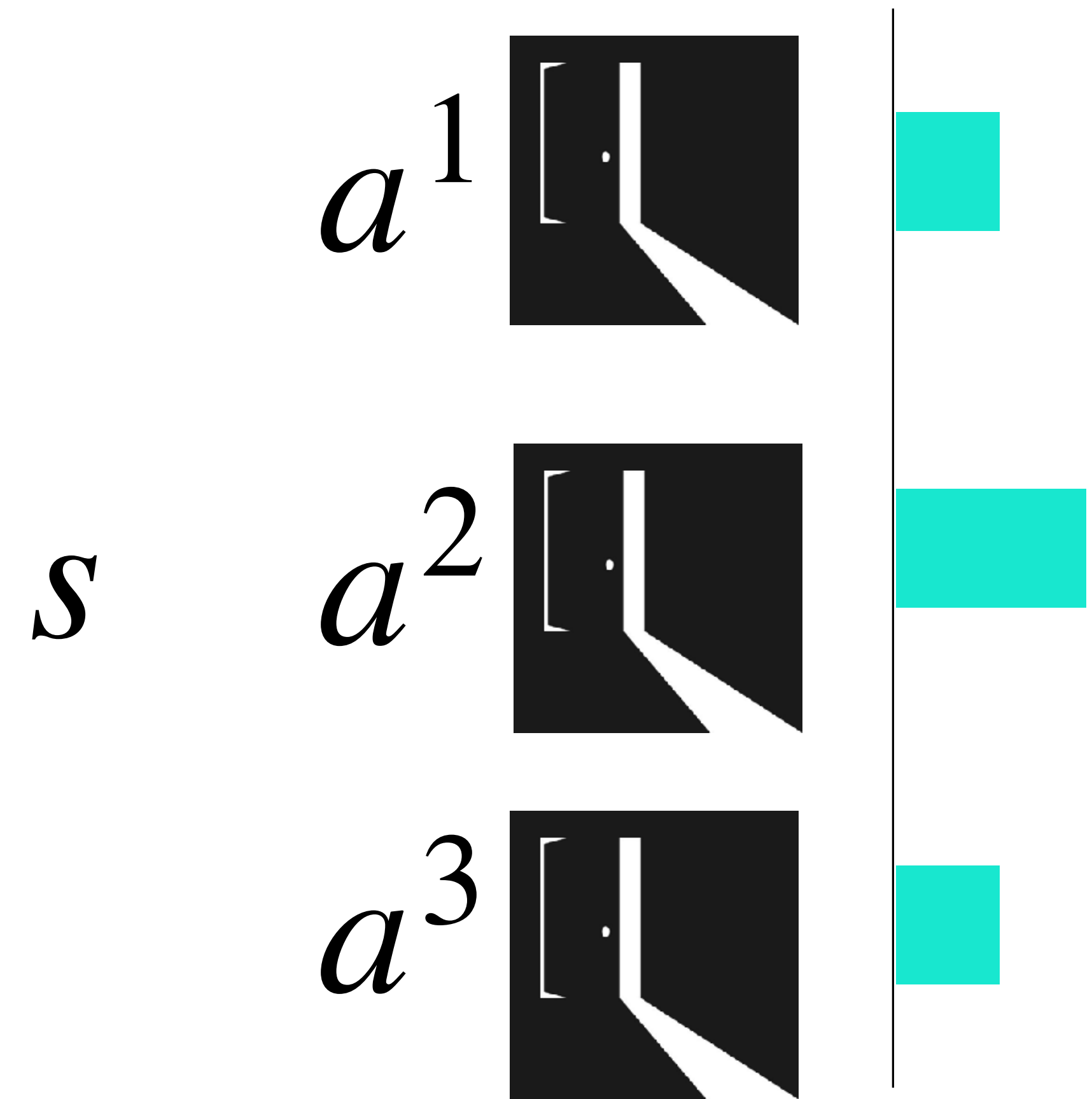


How do we
estimate
values of
each door?

Two Ingredients of RL

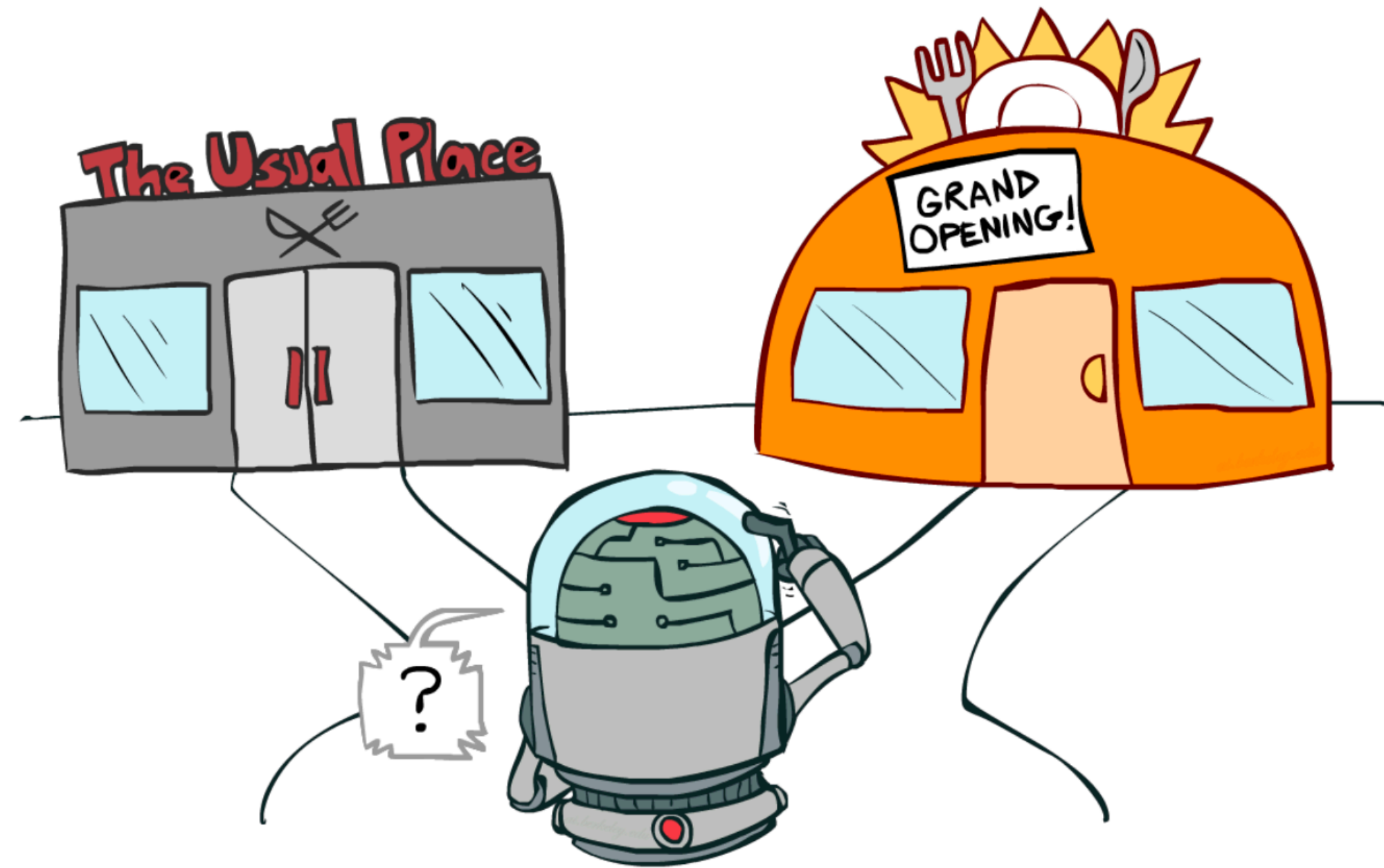


Exploration Exploitation

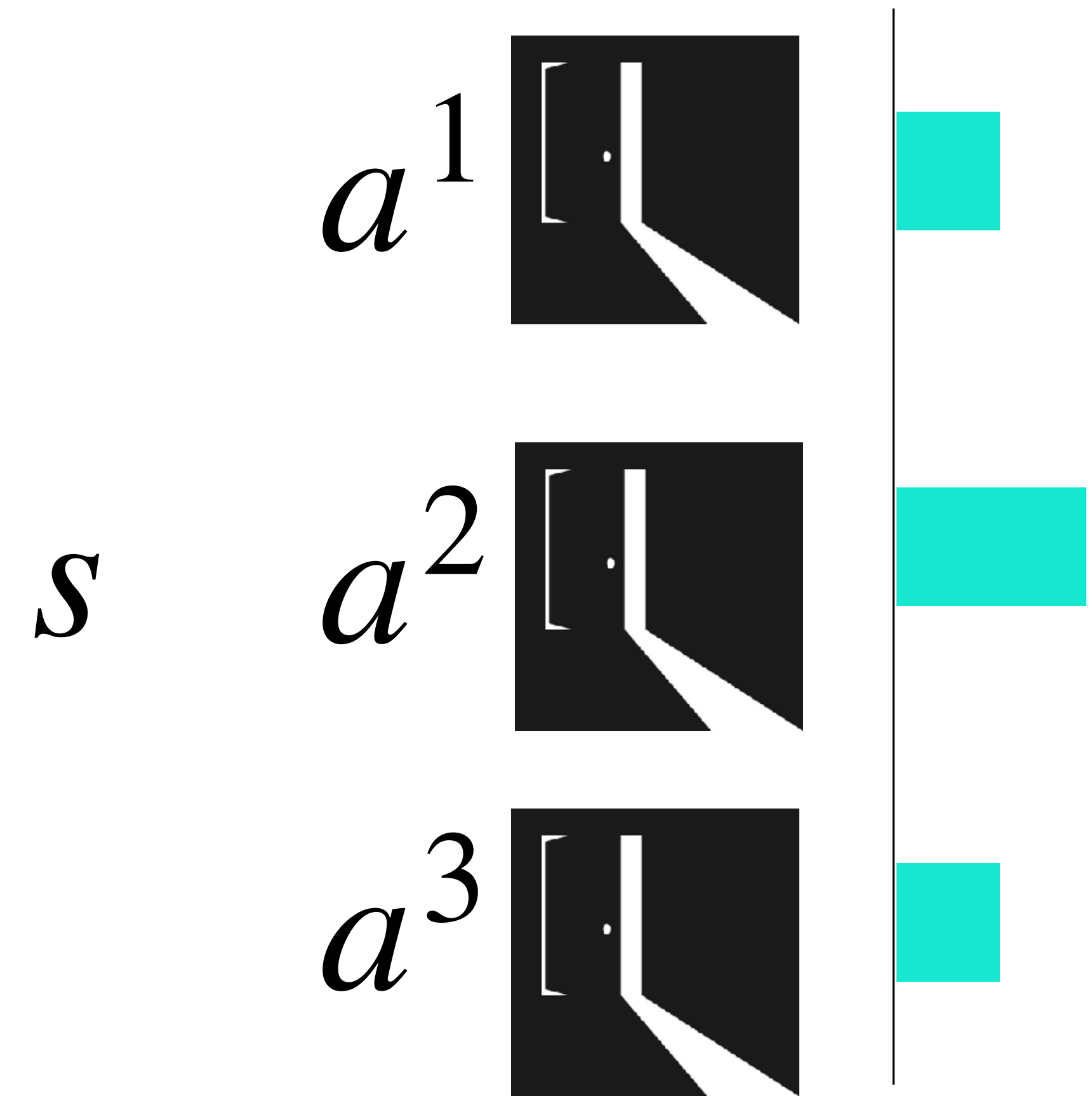


Estimate Values $Q(s, a)$

Two Ingredients of RL

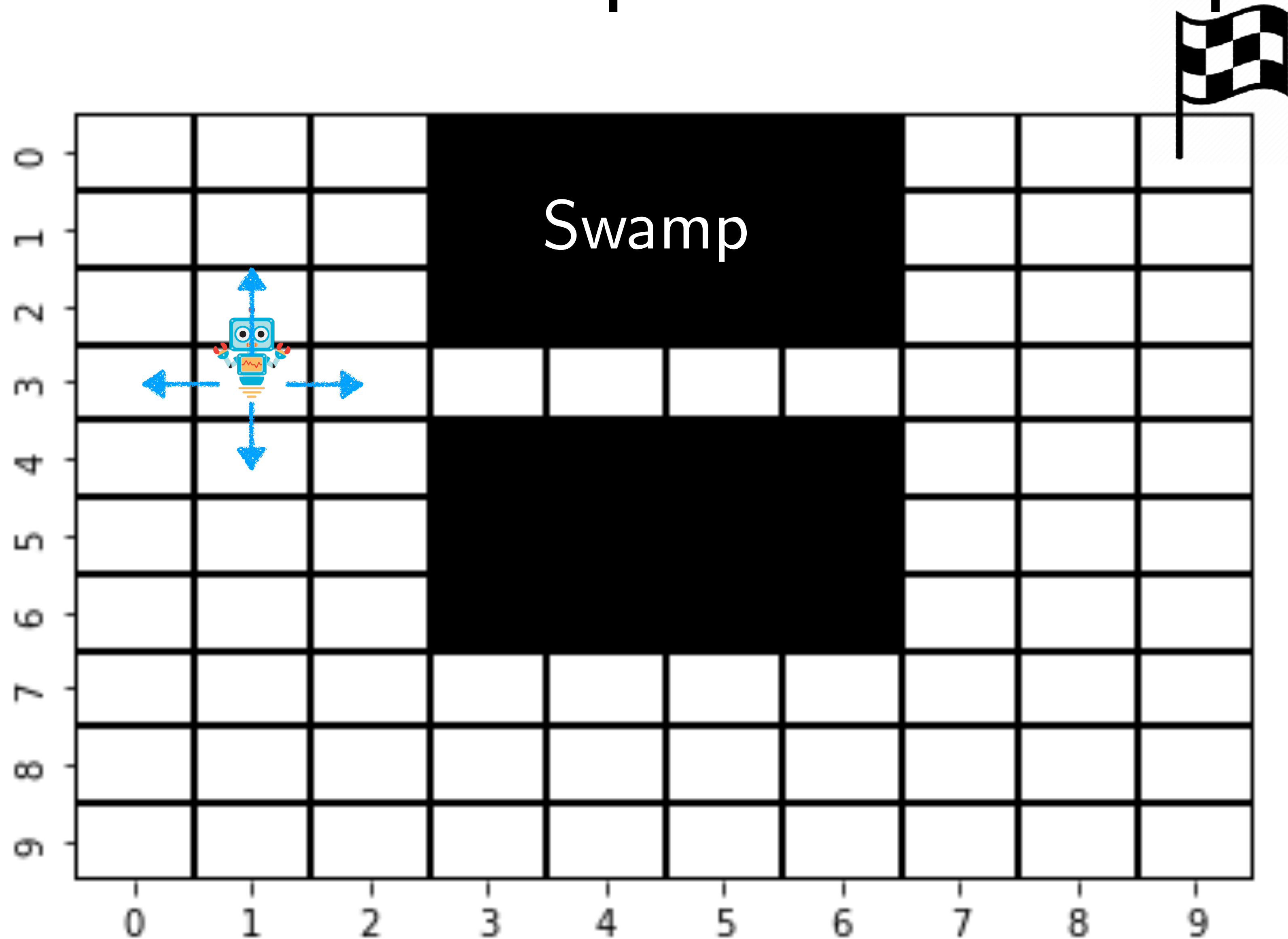


Exploration Exploitation



Estimate Values $Q(s, a)$

Recap: The Swamp MDP



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

- Two absorbing states: Goal and Swamp
- Cost of each state is 1 till you reach the goal
- Let's set $T = 30$

When the MDP is known!

Run Value
/ Policy Iteration



When MDP is known: 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

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

Estimate value

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

Improve policy

What happens when the
MDP is *unknown*?



Need to *estimate the value* of policy



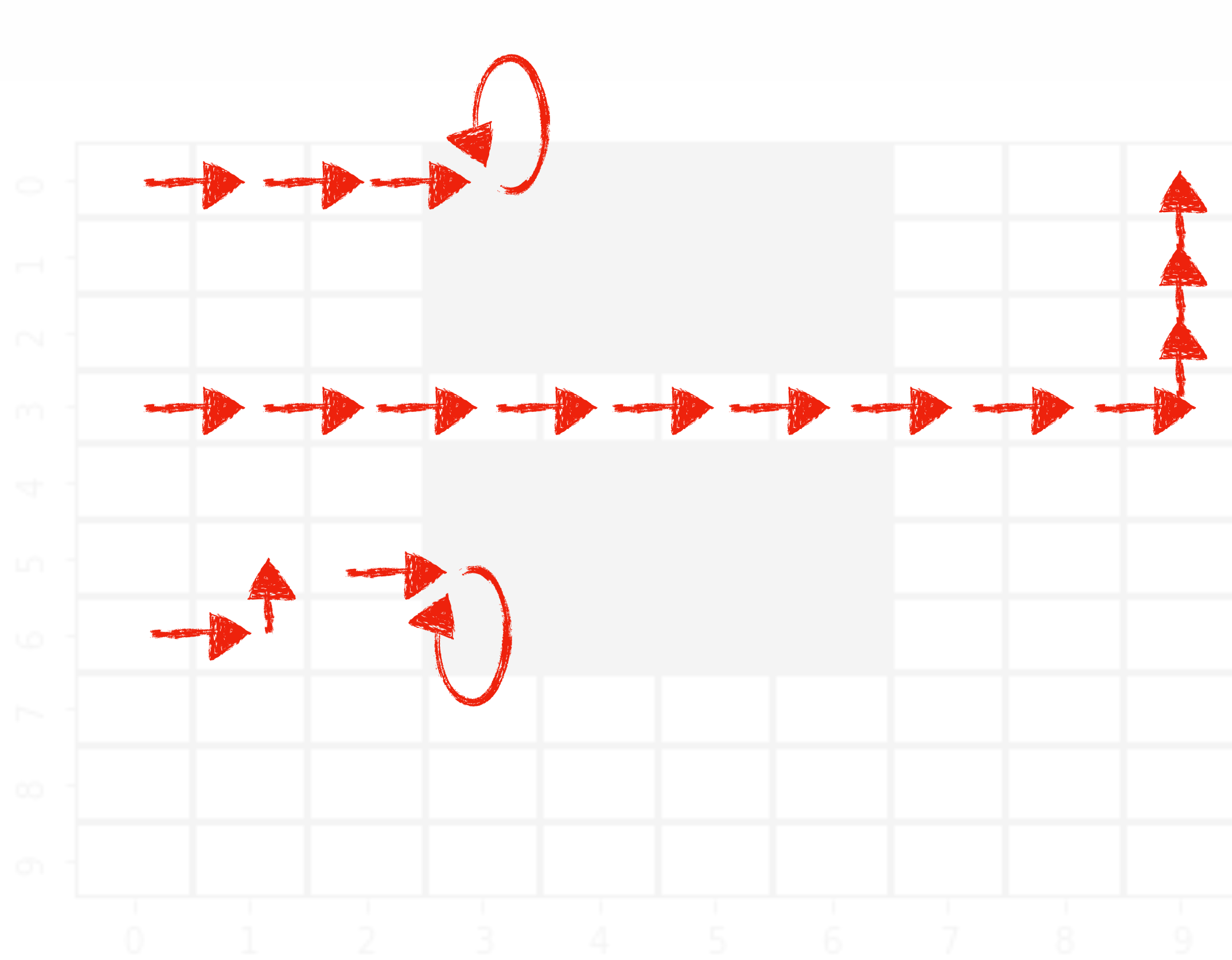
Iter: 0

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

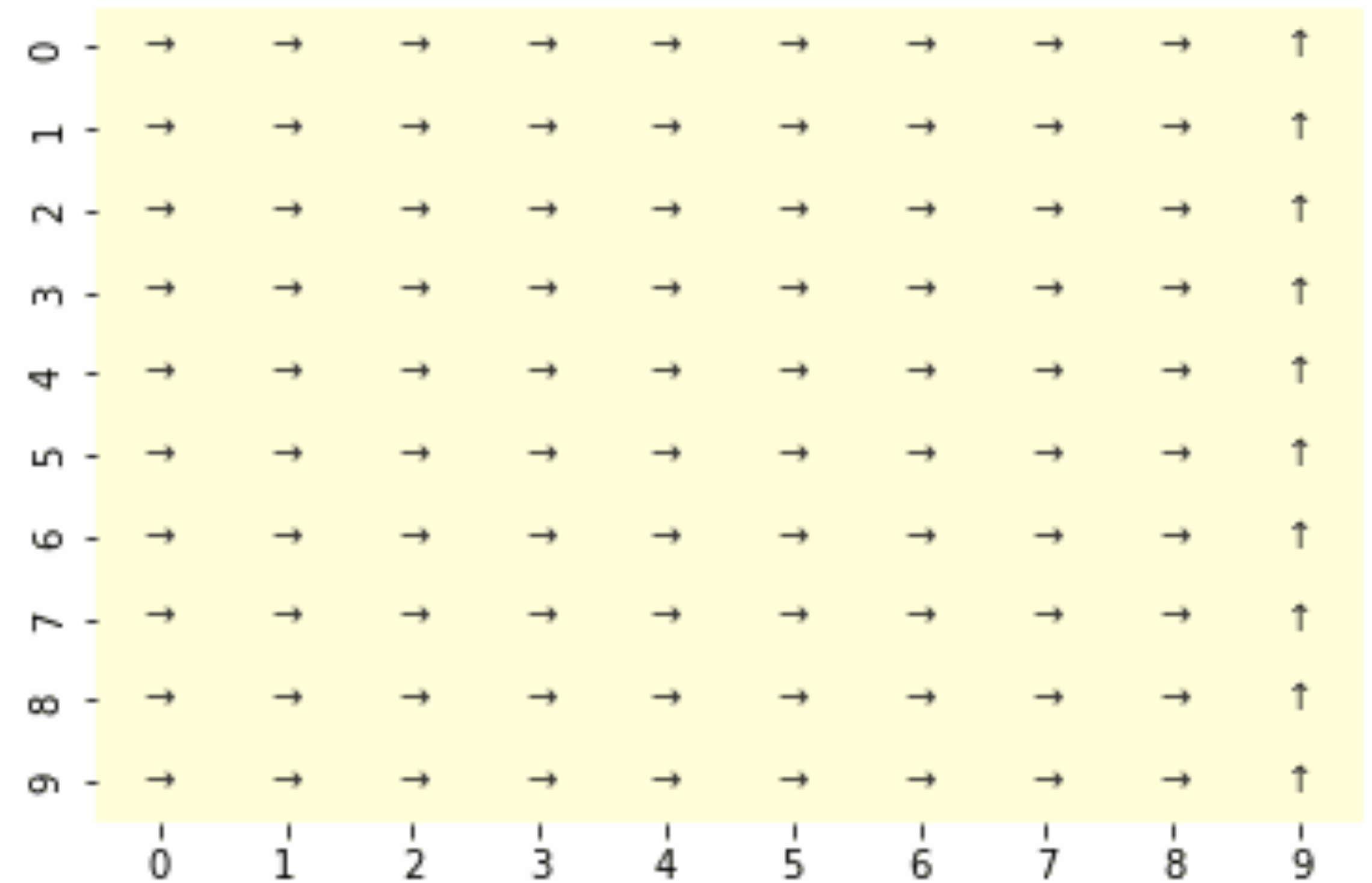
Value $V^\pi(s)$

Policy π

Estimate the value of policy from sample rollouts

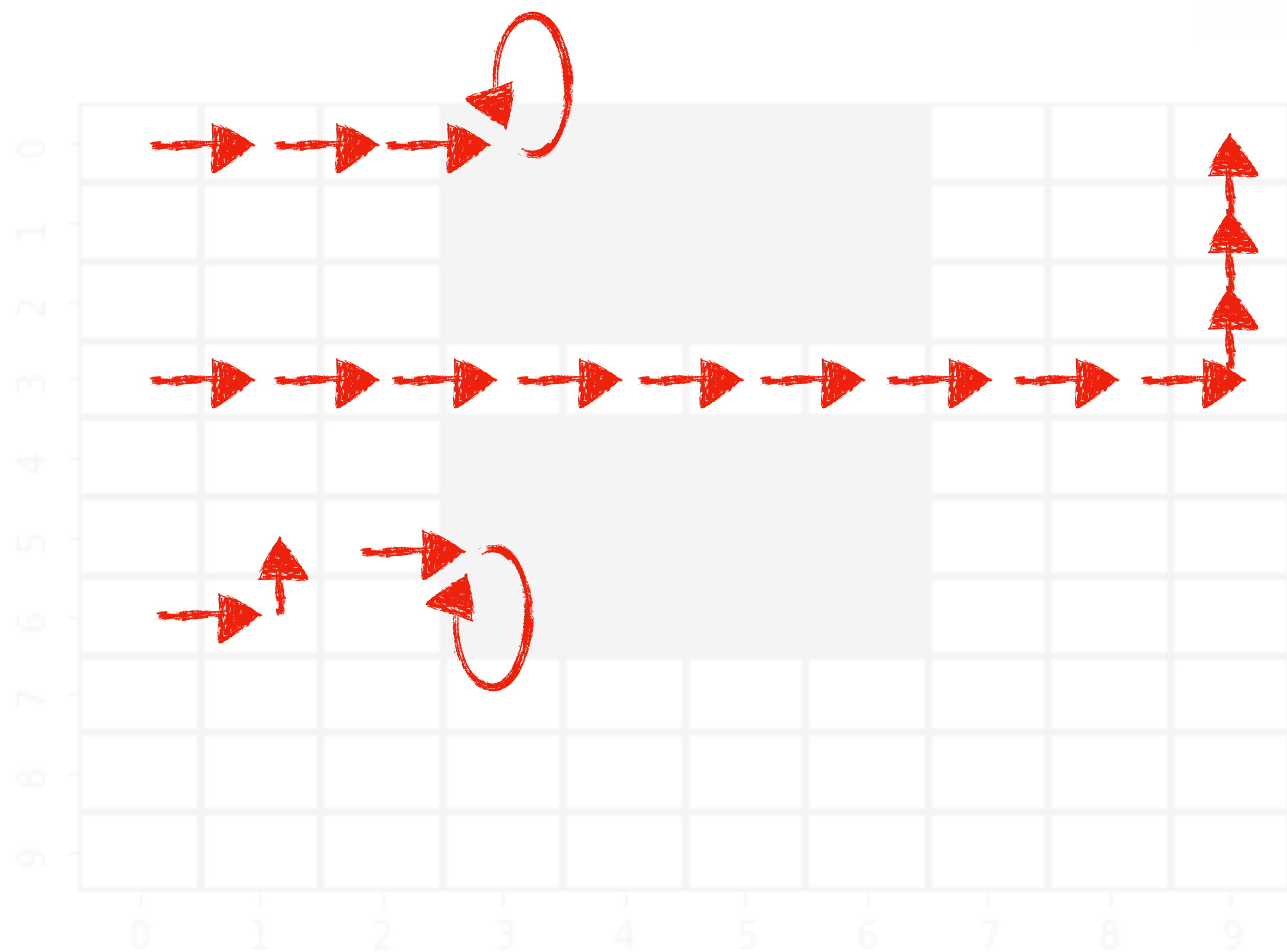


Roll outs

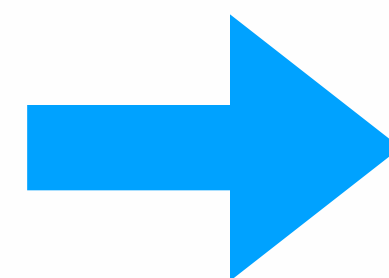


Policy π

Estimate the value of policy from sample rollouts



Roll outs



0	74	75	76	77	77	77	77	2	1	0
1	74	75	76	77	77	77	77	3	2	1
2	74	75	76	77	77	77	77	3.9	3	2
3	55	56	56	57	50	40	26	4.9	3.9	3
4	74	75	76	77	77	77	77	5.9	4.9	3.9
5	74	75	76	77	77	77	77	6.8	5.9	4.9
6	74	75	76	77	77	77	77	7.7	6.8	5.9
7	15	14	13	12	11	10	9.6	8.6	7.7	6.8
8	16	15	14	13	12	11	10	9.6	8.6	7.7
9	17	16	15	14	13	12	11	10	9.6	8.6
	0	1	2	3	4	5	6	7	8	9

Value $V^{\pi}(s)$

Activity!

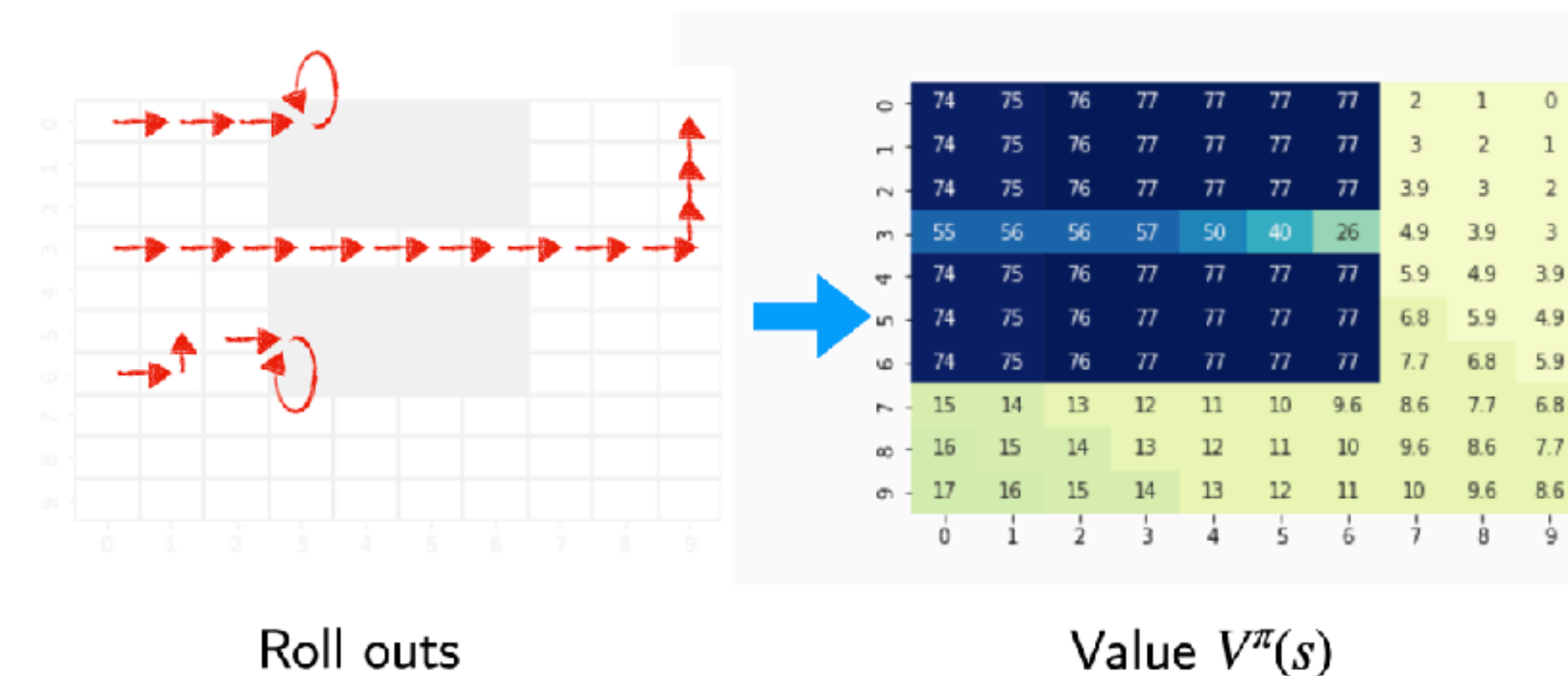


Think-Pair-Share

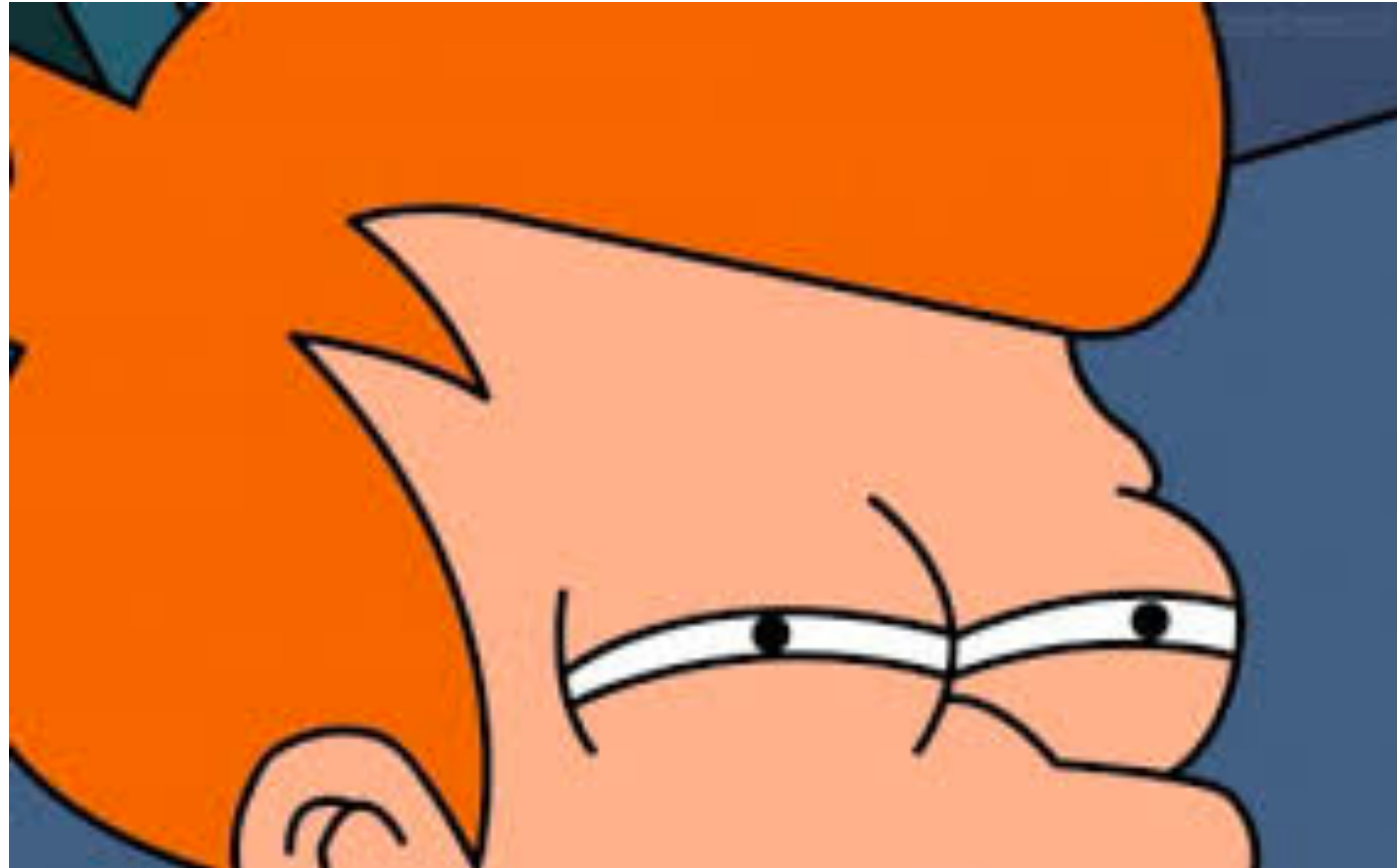
Think (30 sec): Given a bunch of roll-outs, how can you estimate value of a state? (Hint: More than one way!)

Pair: Find a partner

Share (45 sec): Partners exchange ideas



Option 1: Trust only the actual returns!



Monte Carlo Evaluation

Goal: Learn $V^\pi(s)$ from complete rollout $s_1, a_1, c_1, s_2, a_2, c_2, \dots \sim \pi$

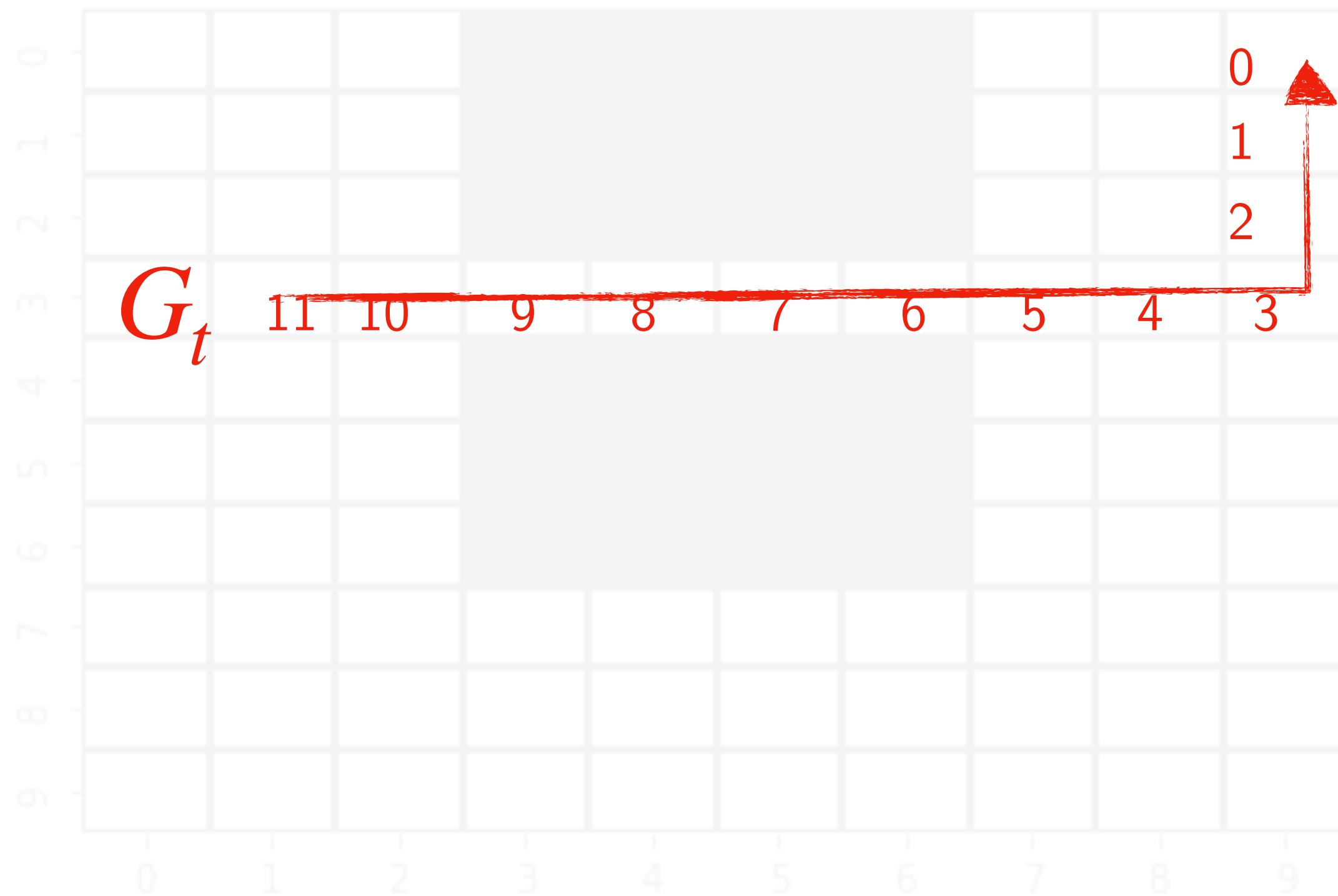
Define: *Return* is the total discounted cost

$$G_t = c_{t+1} + \gamma c_{t+2} + \gamma^2 c_{t+3} + \dots$$

Value function is the expected return

$$V^\pi(s) = \mathbb{E}_\pi[G_t \mid s_t = s]$$

First Visit Monte Carlo



For episode in rollouts:

If state s is visited for *first* time t

Update $V(s) \leftarrow V(s) + \alpha(G_t - V(s))$

Law of large numbers: $V(s) \rightarrow V^\pi(s)$

Can we do better than Monte Carlo?

What if we want quick updates?
(No patience to wait till end)

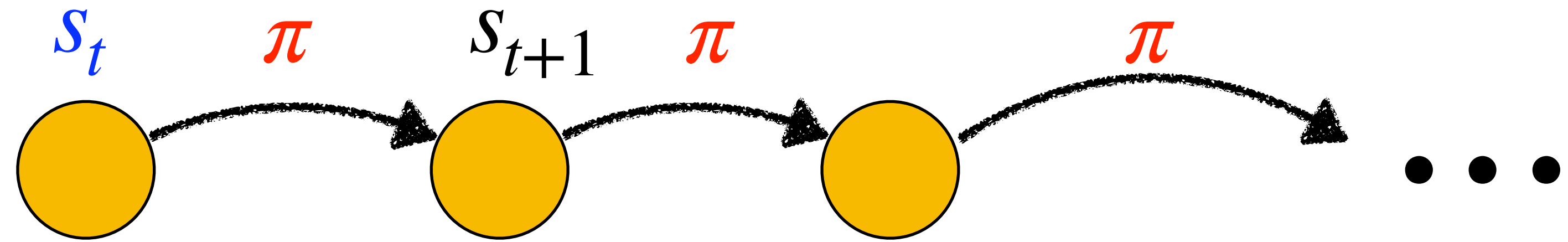
What if we don't have complete episodes?



Option 2: Believe in Bellman



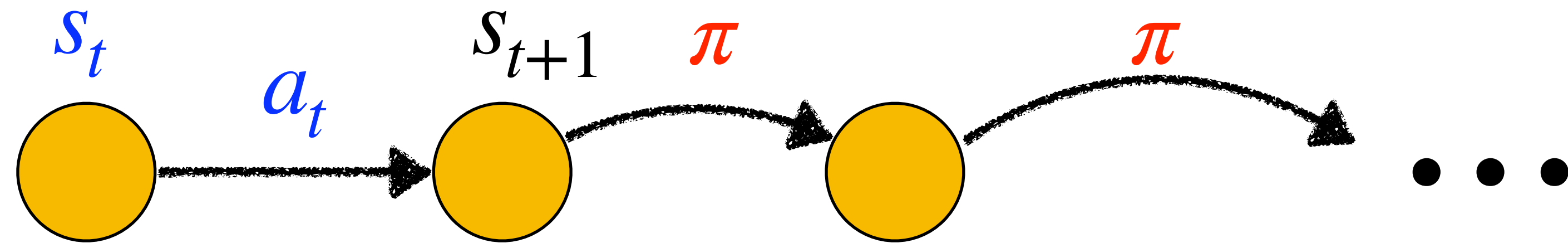
Recap: Value of a state



$$V^{\pi}(s_t) = c_t + \gamma c_{t+1} + \gamma^2 c_{t+2} +$$

*Expected discounted sum of cost from
starting at a state
and following a policy from then on*

Recap: 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

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

Temporal Difference (TD) learning

Goal: Learn $V^\pi(s)$ from traces

$$(s_t, a_t, c_t, s_{t+1}) \quad (s_t, a_t, c_t, s_{t+1}) \quad (s_t, a_t, c_t, s_{t+1}) \quad (s_t, a_t, c_t, s_{t+1})$$

Recall value function $V^\pi(s)$ satisfies

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s'} V^\pi(s')$$

TD Idea: Update value using estimate of next state value

$$V(s_t) \leftarrow V(s_t) + \alpha \left(c_t + \gamma V(s_{t+1}) - V(s_t) \right)$$

Temporal Difference Error

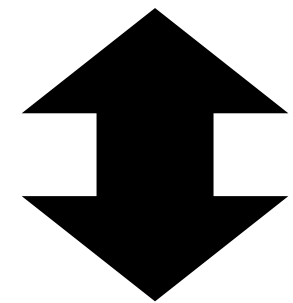
TD Learning

For every (s_t, a_t, c_t, s_{t+1})

$$V(s_t) \leftarrow V(s_t) + \alpha(c_t + \gamma V(s_{t+1}) - V(s_t))$$

Did you spot the trick?

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s'} V^\pi(s')$$



$$V(s_t) \leftarrow V(s_t) + \alpha(c_t + \gamma V(s_{t+1}) - V(s_t))$$





Monte-Carlo

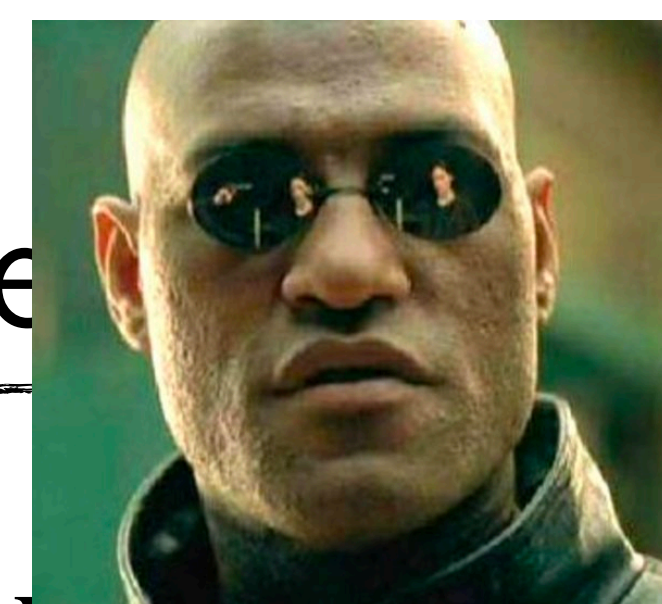
$$V(s) \leftarrow V(s) + \alpha(G_t - V(s))$$

Zero Bias

High Variance

Always convergence

(Just have to wait till heat death of the universe)



Temporal Difference

$$V(s) \leftarrow V(s) + \alpha(c + \gamma V(s') - V(s))$$

Can have bias

Low Variance

May *not* converge if
using function approximation

We have been talking about
trying to learn the value of a

$$\text{given policy } \pi \\ V^\pi(s) / Q^\pi(s, a)$$

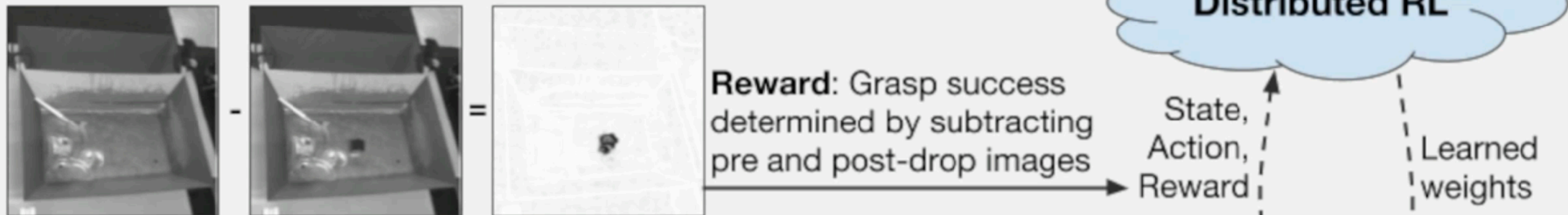
What if we wanted to learn
the optimal value function

$$V^*(s) / Q^*(s, a)$$

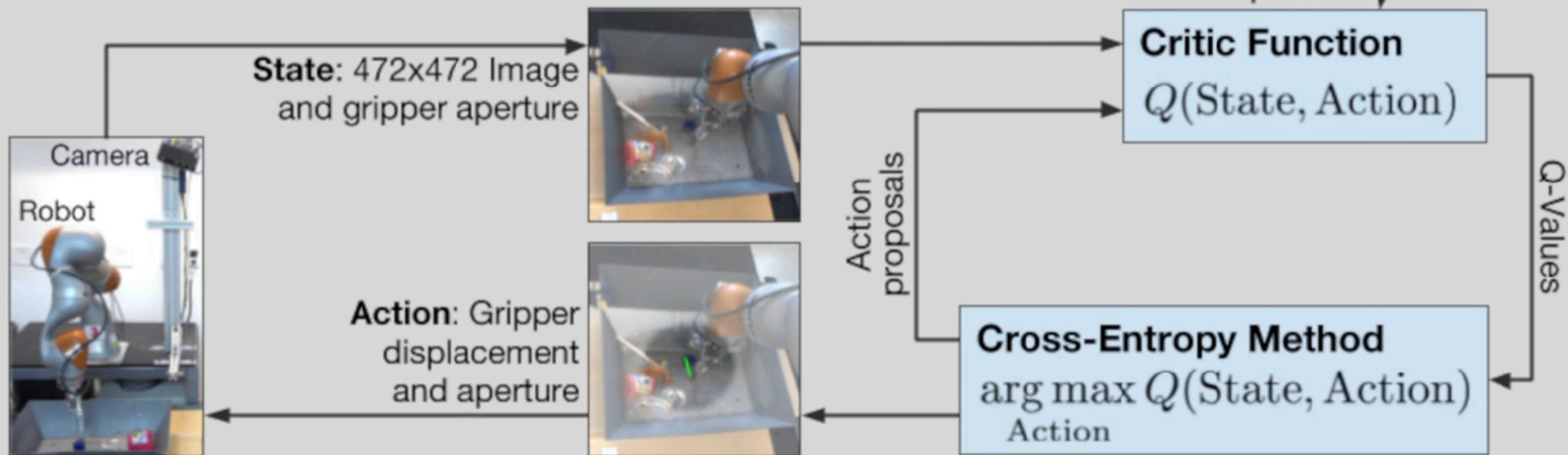


QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation

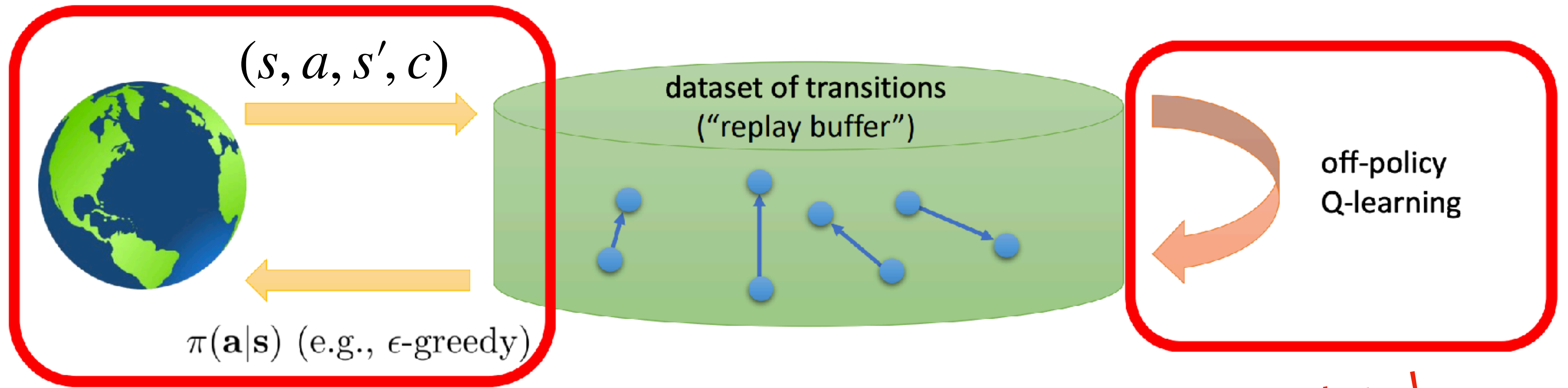
Training time



Inference time



Q-learning: Learning off-policy

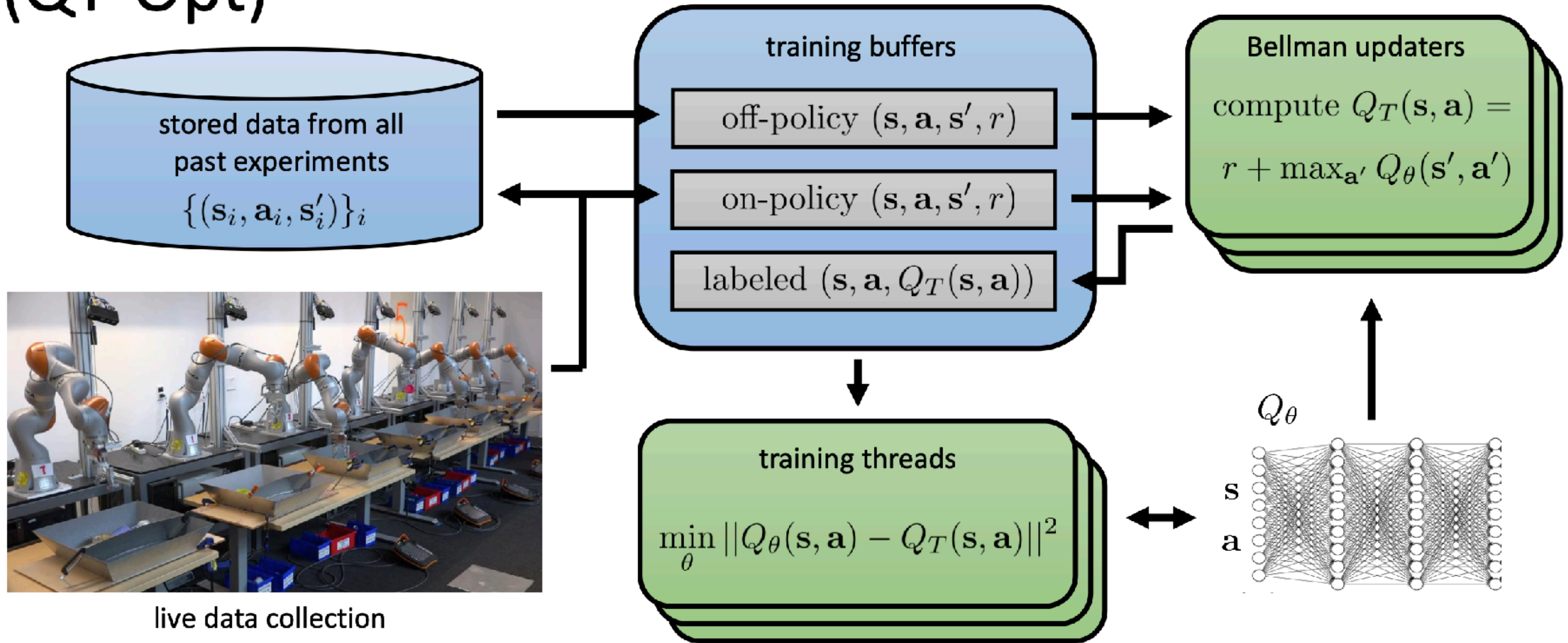


For every (s_t, a_t, c_t, s_{t+1})

Can learn from any data!

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

Large-scale Q-learning with continuous actions (QT-Opt)



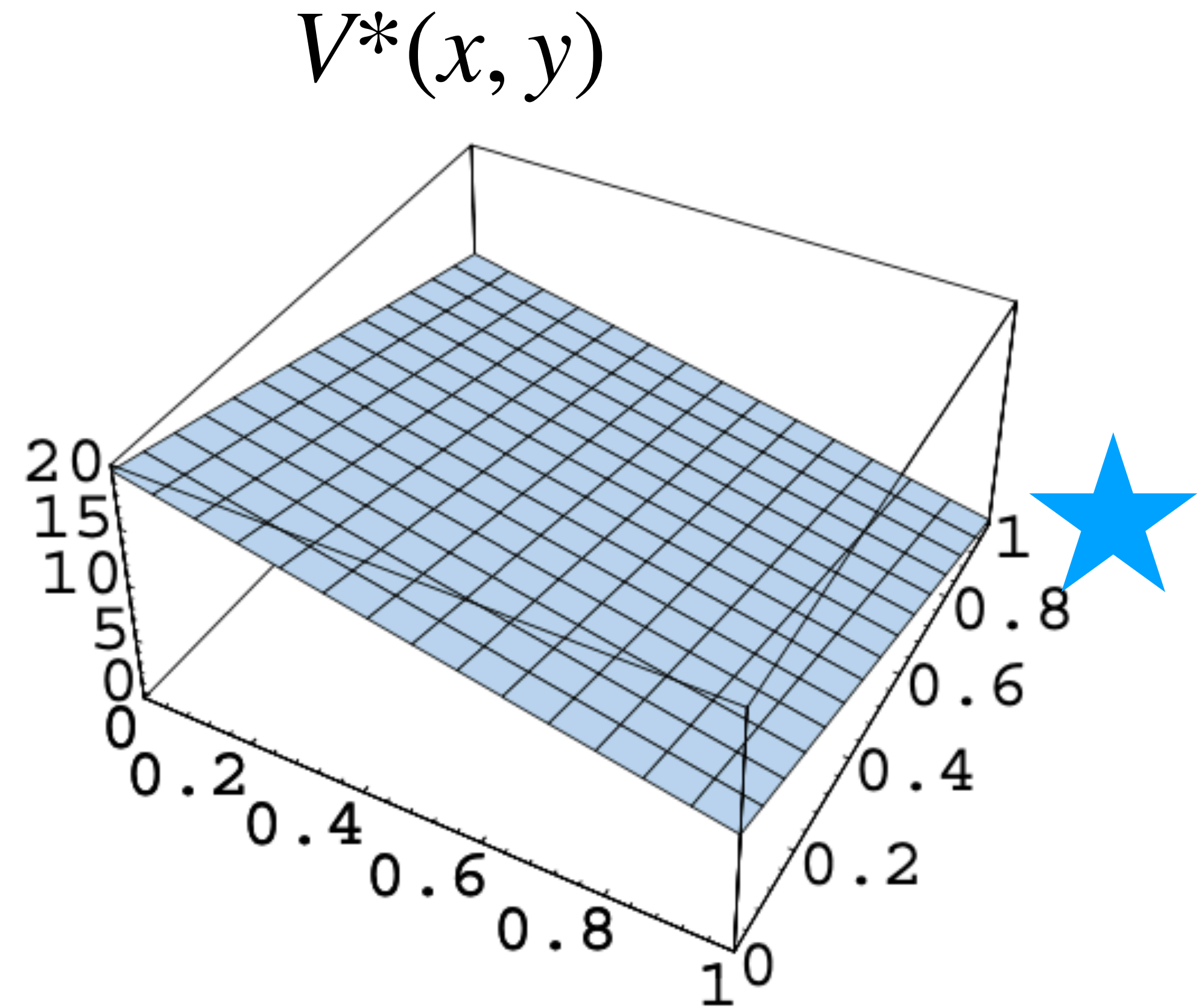
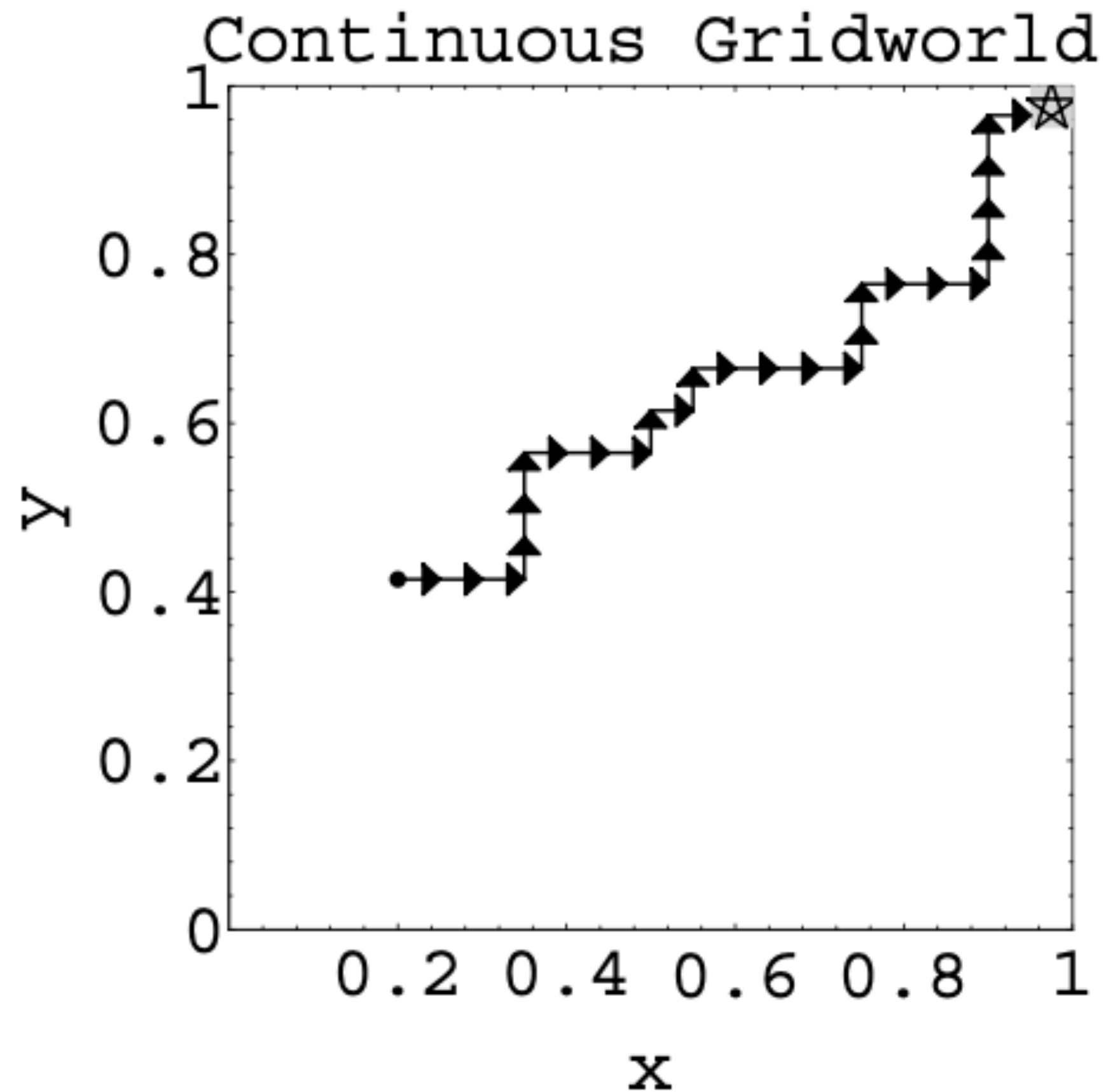
Is this ... magic?

We just learned in IL how distribution shift is a big deal ...

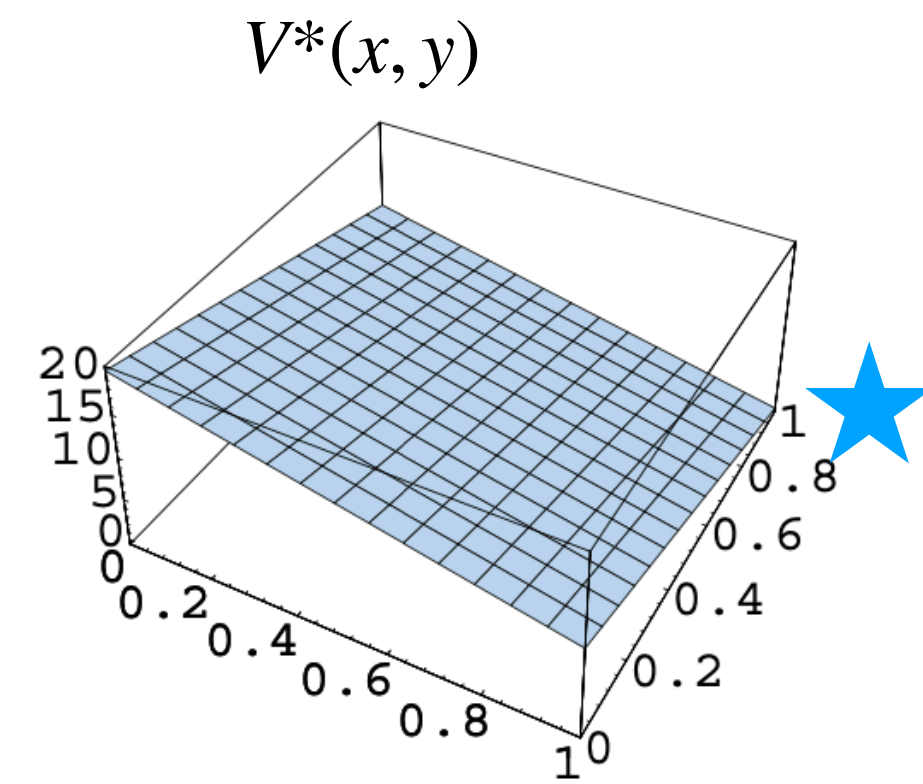
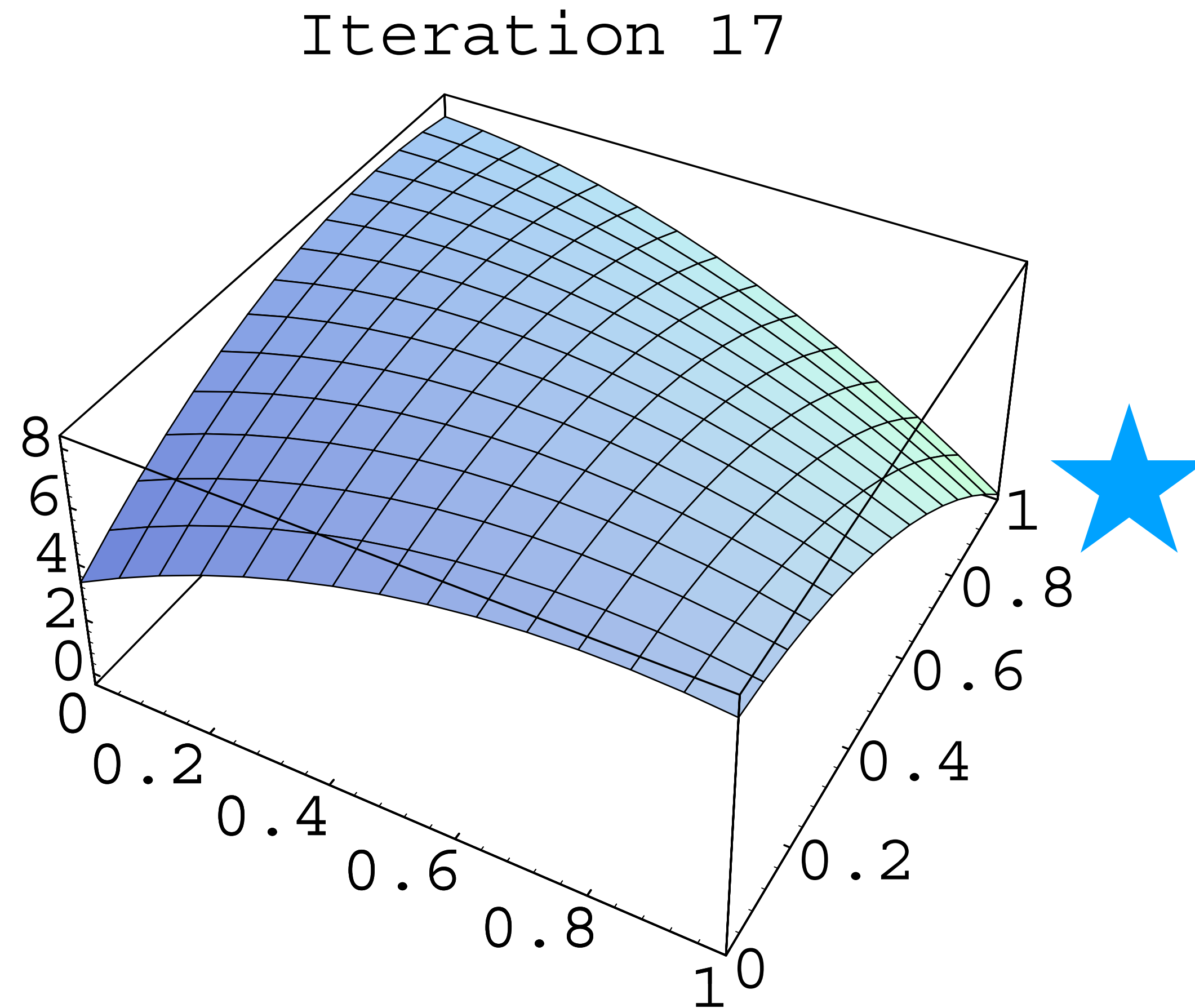
It's not magic. Q-learning relies on a set of assumptions:

1. Each state-action is visited *infinite* times
2. Learning rate α must be annealed over time

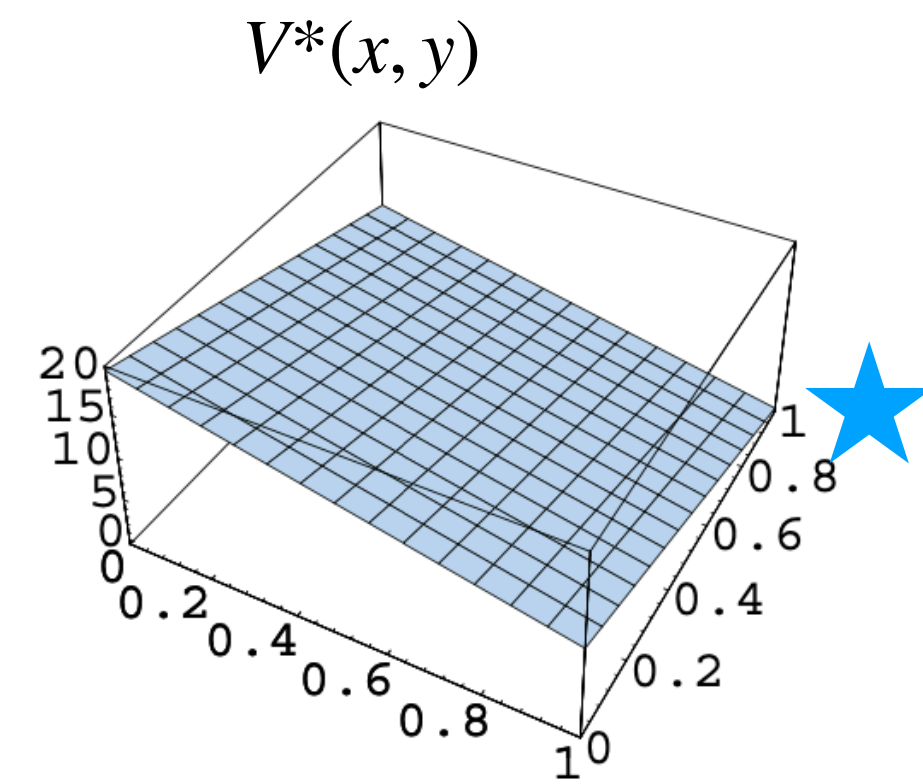
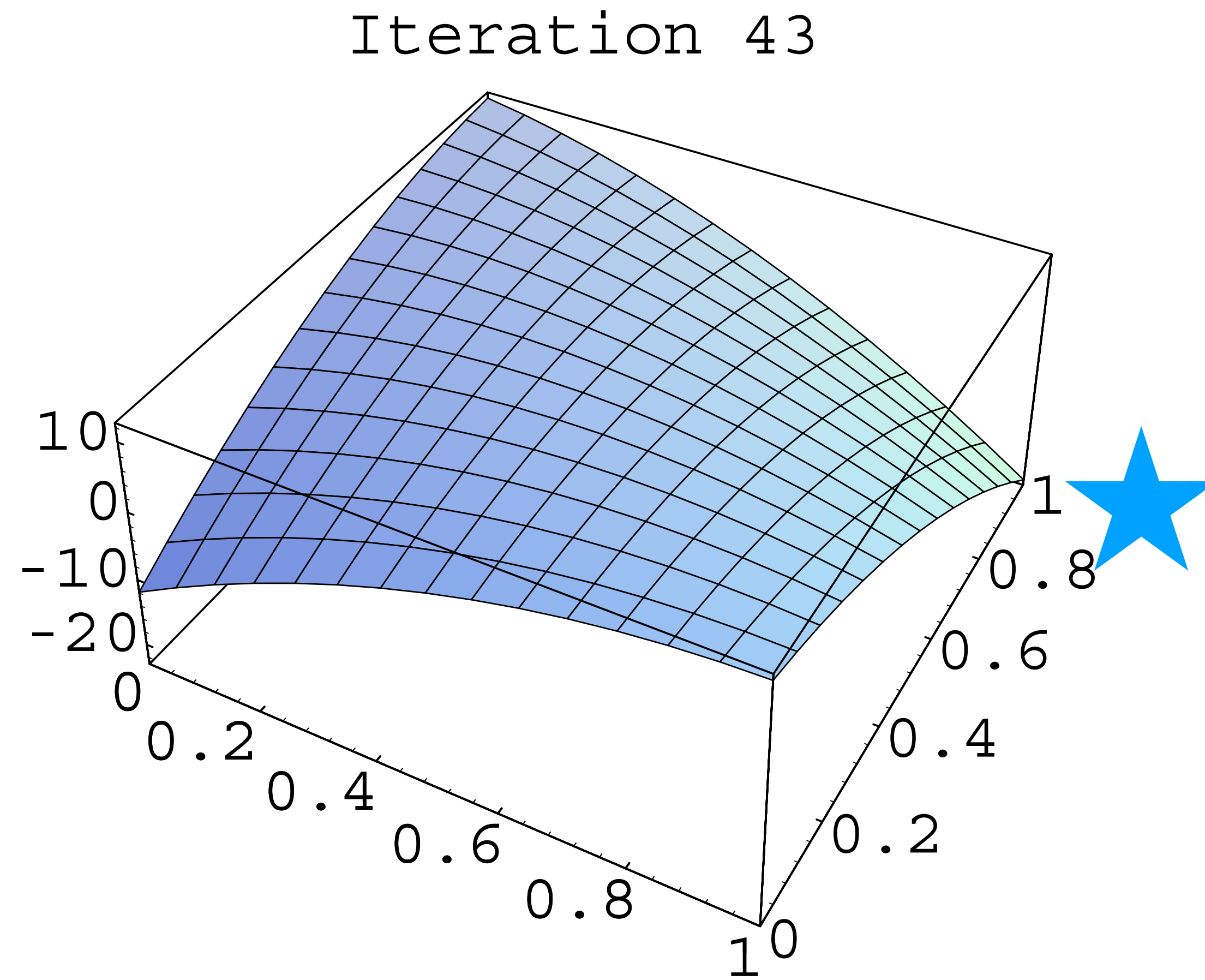
When things fail!



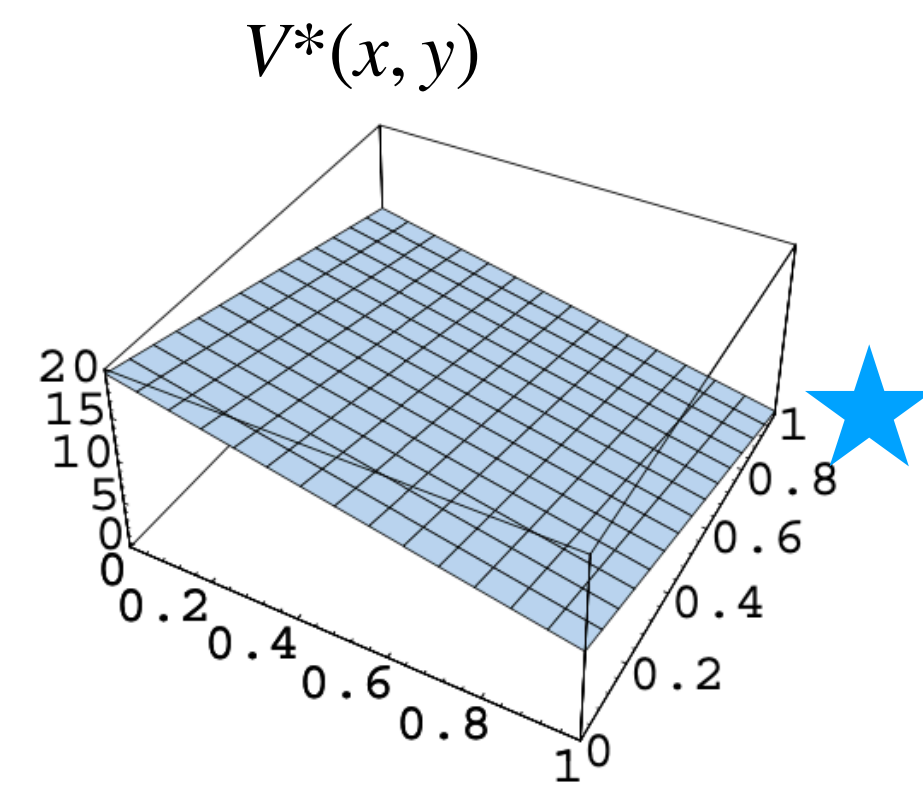
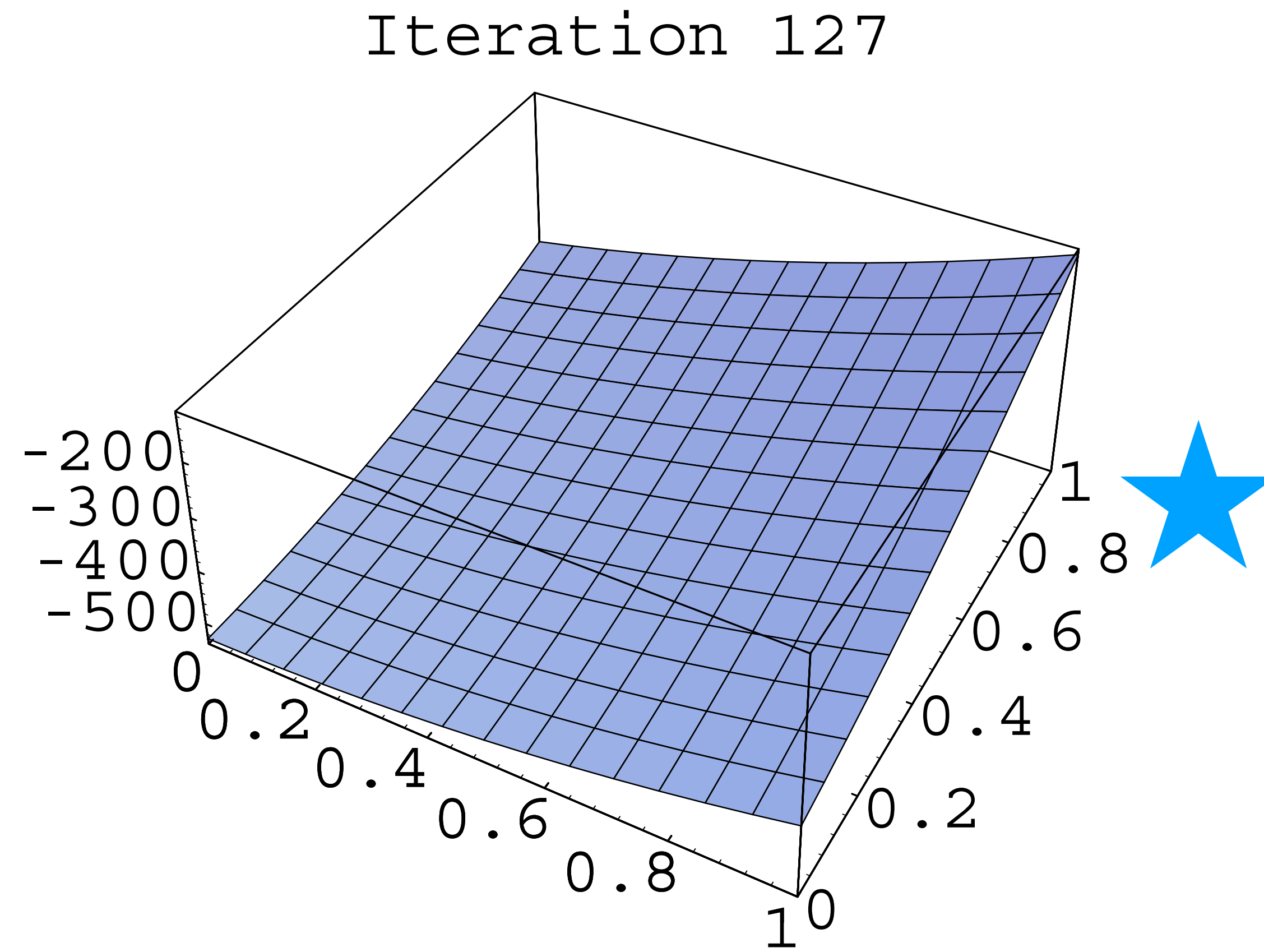
What happens when we run value iteration with a
quadratic?



What happens when we run value iteration with a
quadratic?



What happens when we run value iteration with a
quadratic?

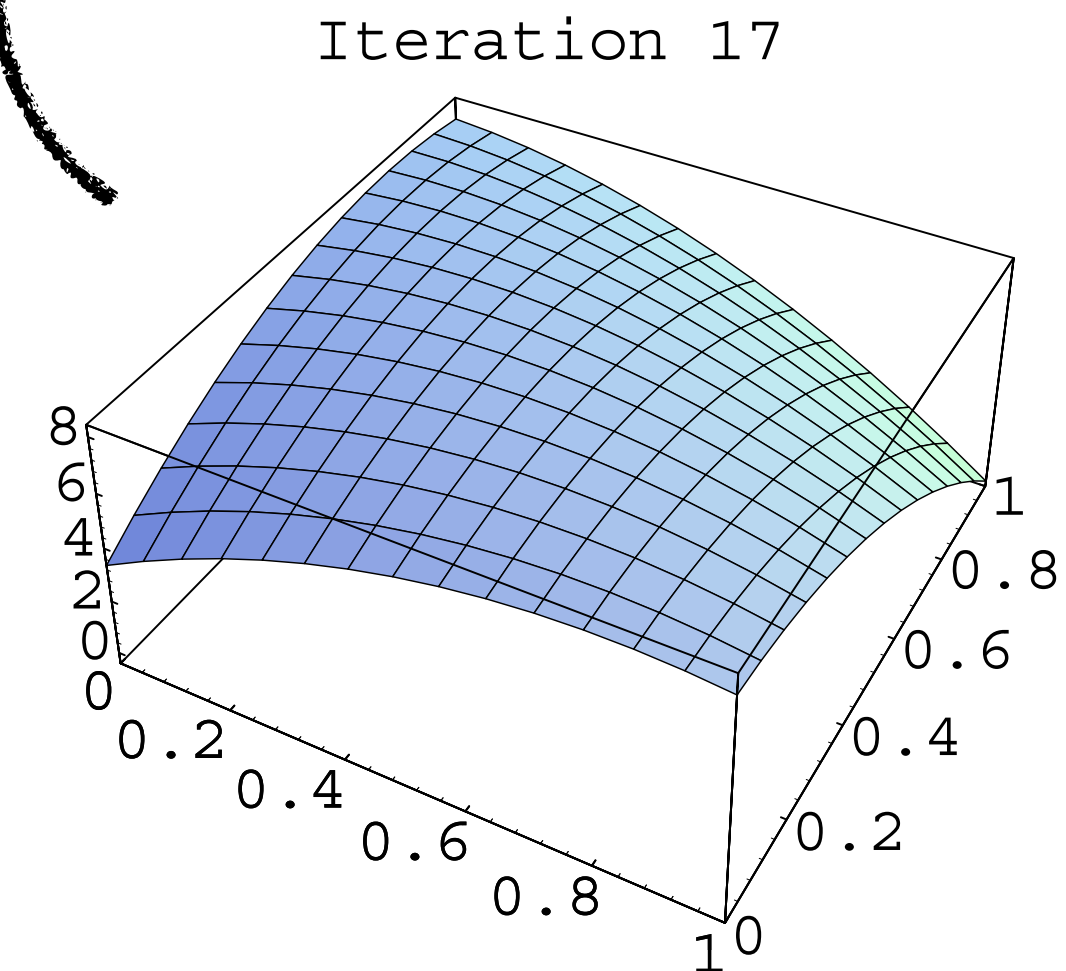
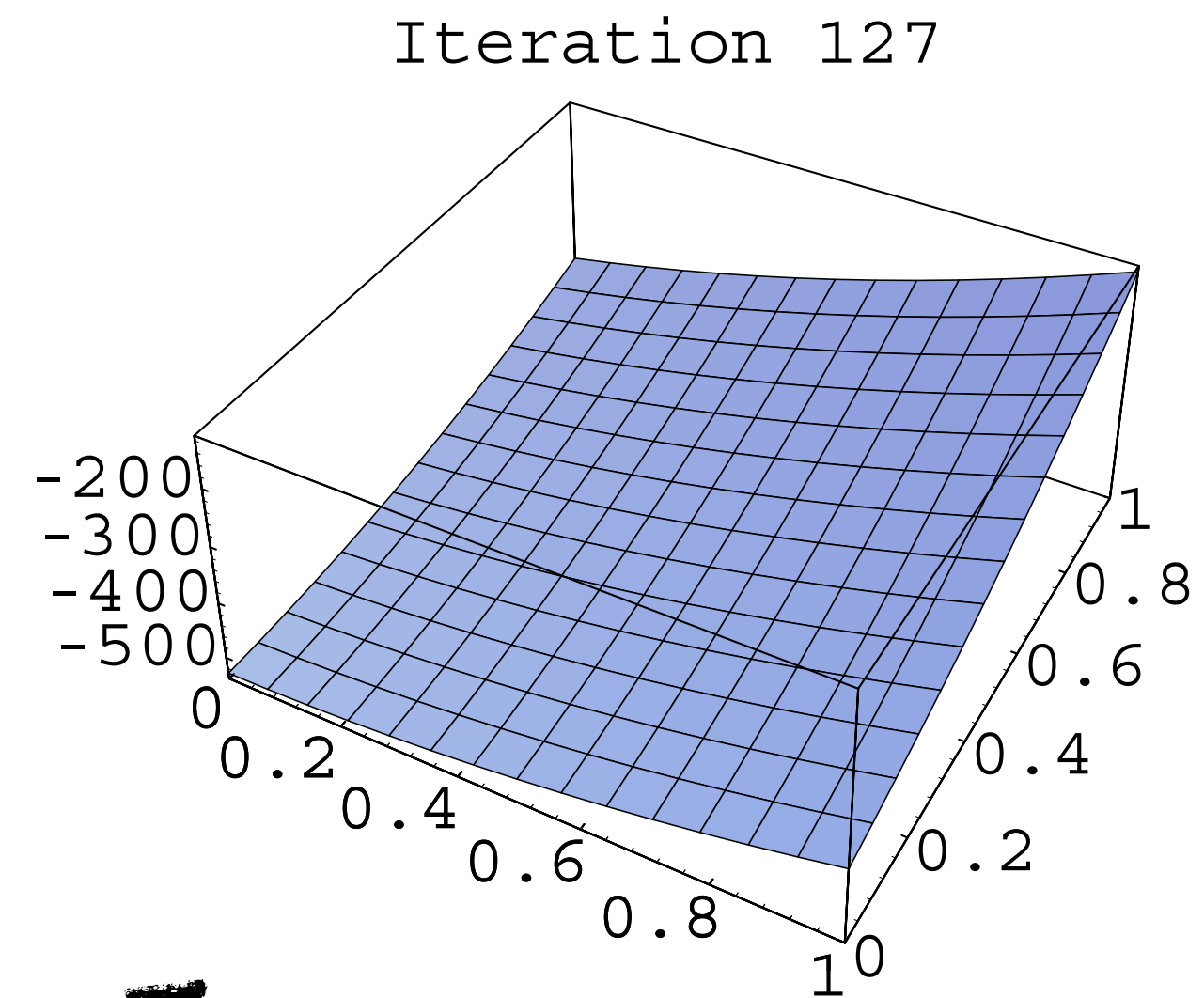


The problem of Bootstrapping!

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$



min()

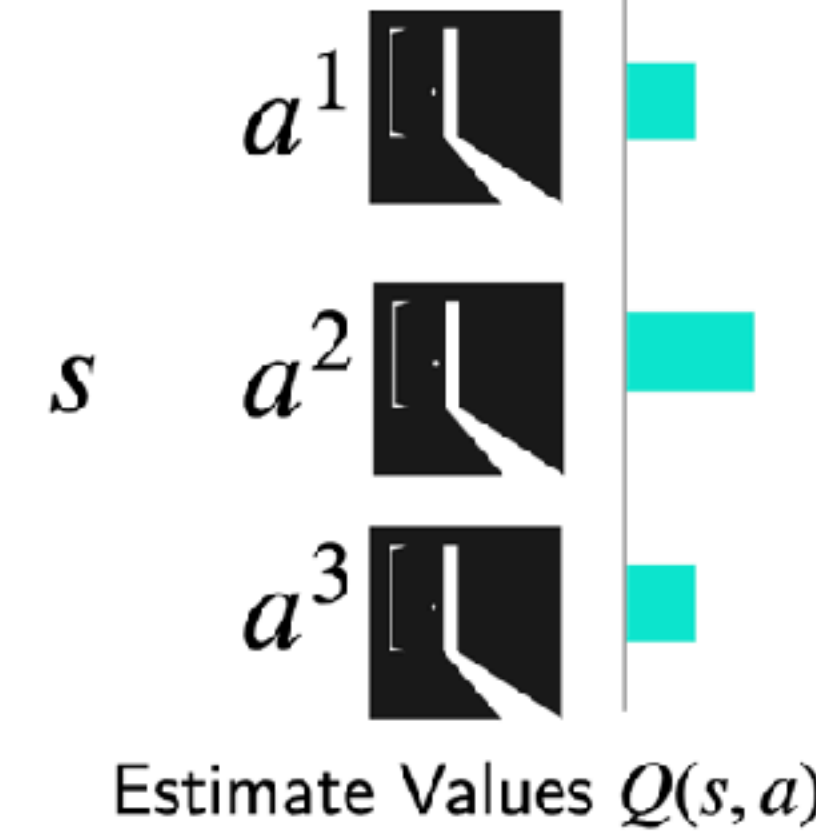


tl;dr

Two Ingredients of RL



Exploration Exploitation



Monte-Carlo

$$V(s) \leftarrow V(s) + \alpha(G_T - V(s))$$

Zero Bias

High Variance

Always convergence

(Just have to wait till heat death of the universe)

Temporal Difference



$$V(s) \leftarrow V(s) + \alpha(c + \gamma V(s') - V(s))$$

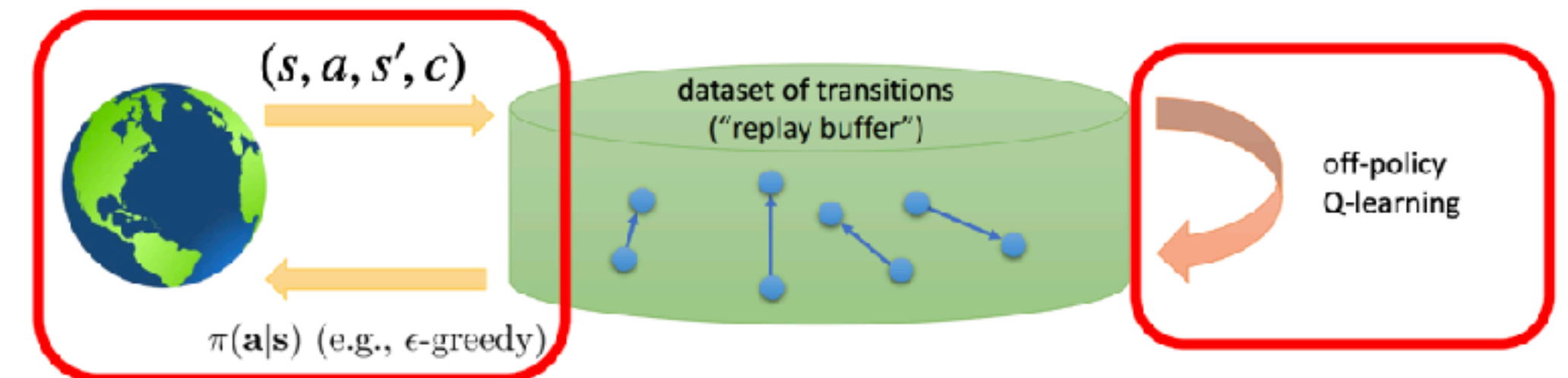
Can have bias

Low Variance

May *not* converge if
using function approximation

38

Q-learning: Learning off-policy



For every (s_t, a_t, c_t, s_{t+1})

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(c(s_t, a_t) + \gamma \min_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

41