

# Welcome!

# The Plan for Today

1. Course goals & logistics
2. Why study deep reinforcement learning?
3. Intro to modeling behavior and reinforcement learning

Key learnings goals:

- what is deep reinforcement learning??
- how to represent behavior
- how to formulate a reinforcement learning problem

# What do we mean by deep reinforcement learning?

## Sequential decision-making problems

A system needs to make *multiple* decisions based on stream of information.

observe, take action, observe, take action, ...

## AND the solutions to such problems

- imitation learning
- model-free & model-based RL
- offline & online RL
- multi-task & meta RL
- RL for LLMs
- RL for robots

and more!

Emphasis on solutions that scale to deep neural networks

# How does deep RL differ from other ML topics?

## Supervised learning

Given labeled data:  $\{(x_i, y_i)\}$ , learn  $f(x) \approx y$

- directly told what to output
- inputs  $x$  are independently, identically distributed (i.i.d.)

## Reinforcement learning

Learn behavior  $\pi(a|s)$ .

- from experience, indirect feedback
- data **not** i.i.d.: actions affect the future observations.

Behavior can include:

motor control

chat bots

game playing

driving

web agents

# We can't cover everything in deep RL.

We'll focus on:

- core concepts behind deep RL methods
- implementation of algorithms
- examples in robotics, control, language models (but techniques generalize broadly)
- topics that we think are most useful & exciting!

For more theory & other applications, see CS234!

Core goal: Able to understand and implement existing and emerging methods.

# Why study deep reinforcement learning?

## 1. Going beyond supervised (x, y) examples

- AI model predictions have consequences! How can we take them into account?
- When direct supervision isn't available Learn from *any* objective.

## 2. Widely used and deployed for performant AI systems

## 3. Learning from experience seems fundamental to intelligence

- RL can **discover new solutions**

## 4. Plenty of exciting open research problems

# Why study deep reinforcement learning?

Beyond supervised learning from  $(x, y)$  examples

Decision-making problems are everywhere!

- a. Any sort of AI agent: robots, autonomous vehicles, web assistants
- b. What if you want your AI system to interact with people? chatbots, recommenders
- c. What if deploying your system affects future outcomes & observations?
- d. What if don't have labels or your objective isn't just accuracy? “feedback loops”  
(and isn't differentiable!)

# Why study deep reinforcement learning?

Widely used for performant AI systems

Learning complex physical tasks: legged robots

<https://www.youtube.com/watch?v=9G9E-TIKGM8>

<https://www.youtube.com/watch?v=xG7WkPU8tgs>

# Why study deep reinforcement learning?

Widely used for performant AI systems

Learning complex physical tasks: robot manipulation



<https://www.youtube.com/watch?v=gg9AYgtYoNs>

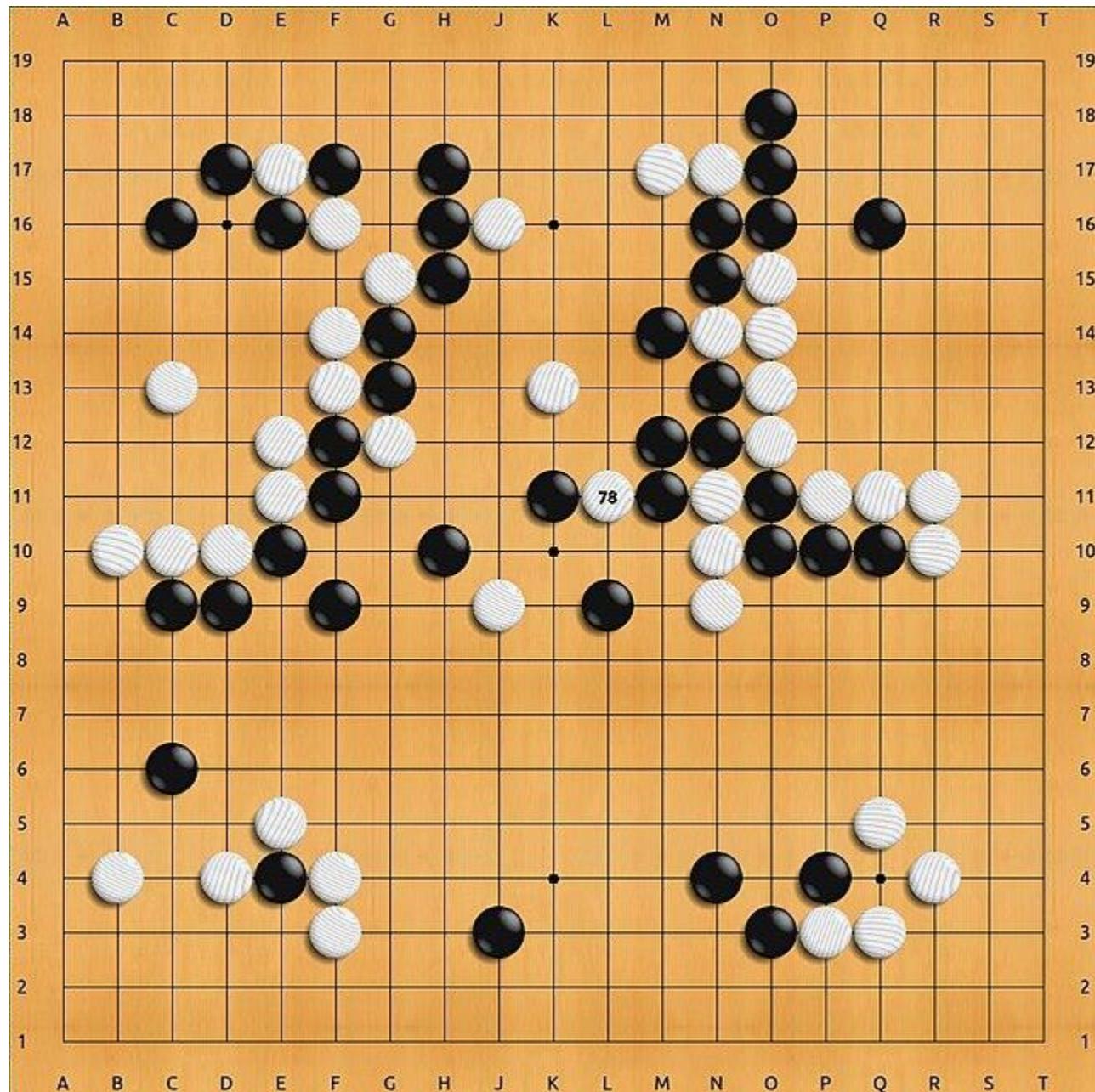
Google Gemini Robotics

Physical Intelligence  $\pi$

# Why study deep reinforcement learning?

Widely used for performant AI systems

Learning to play complex games



Ability to **discover** new solutions:  
“Move 37” in Lee Sedol AlphaGo  
match surprises everyone

# Why study deep reinforcement learning?

Widely used for performant AI systems

**Not just robots and games!**

Nearly all modern language models use some form of RL for post-training.



ChatGPT



deepseek



CURSOR

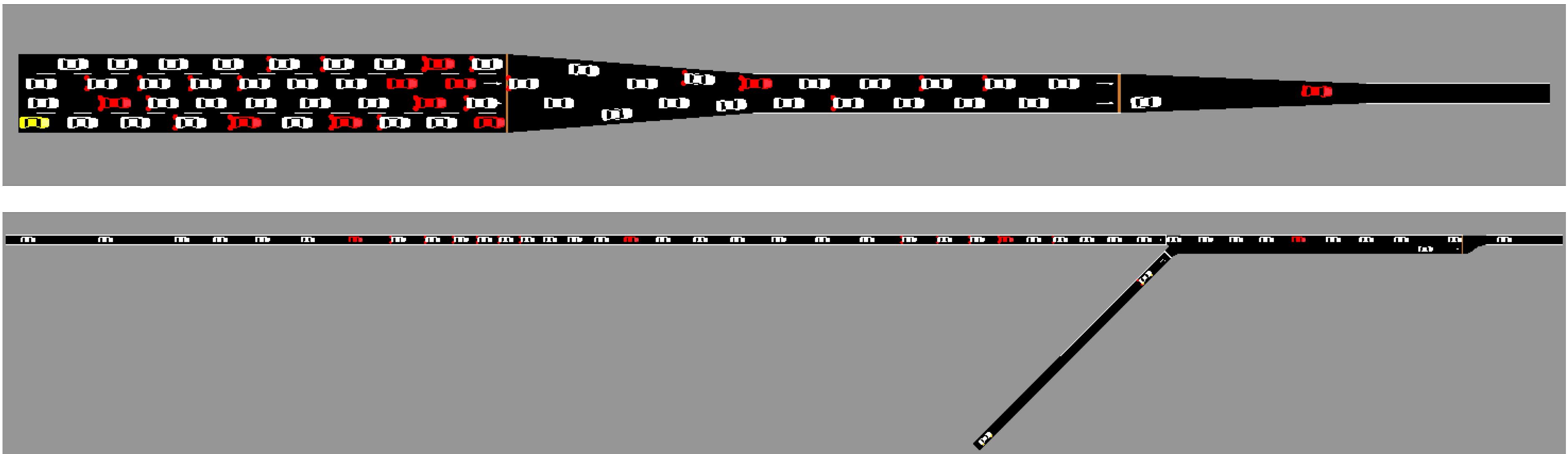
*especially* for more advanced reasoning.

# Why study deep reinforcement learning?

Widely used for performant AI systems

**Not just robots and games!**

Research on traffic control



# Why study deep reinforcement learning?

Widely used for performant AI systems

**Not just robots and games!**

Training generative image models to follow their prompt

— *a dolphin riding a bike* —→



— *an ant playing chess* —→

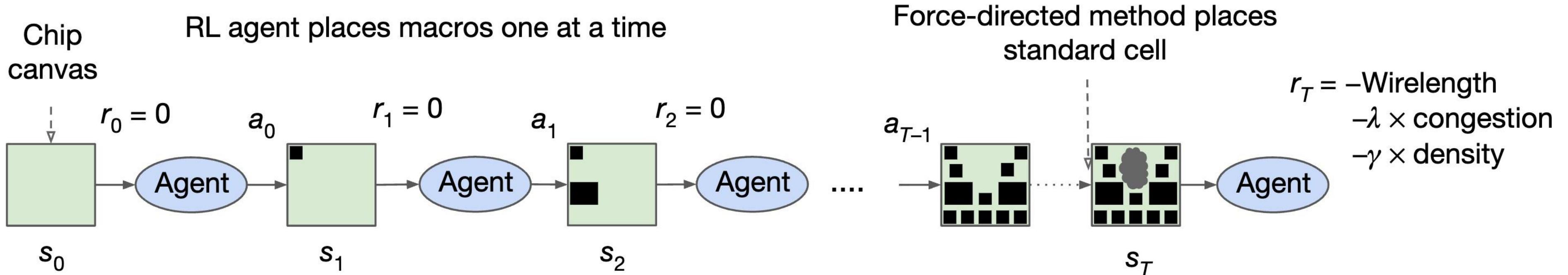


# Why study deep reinforcement learning?

Widely used for performant AI systems

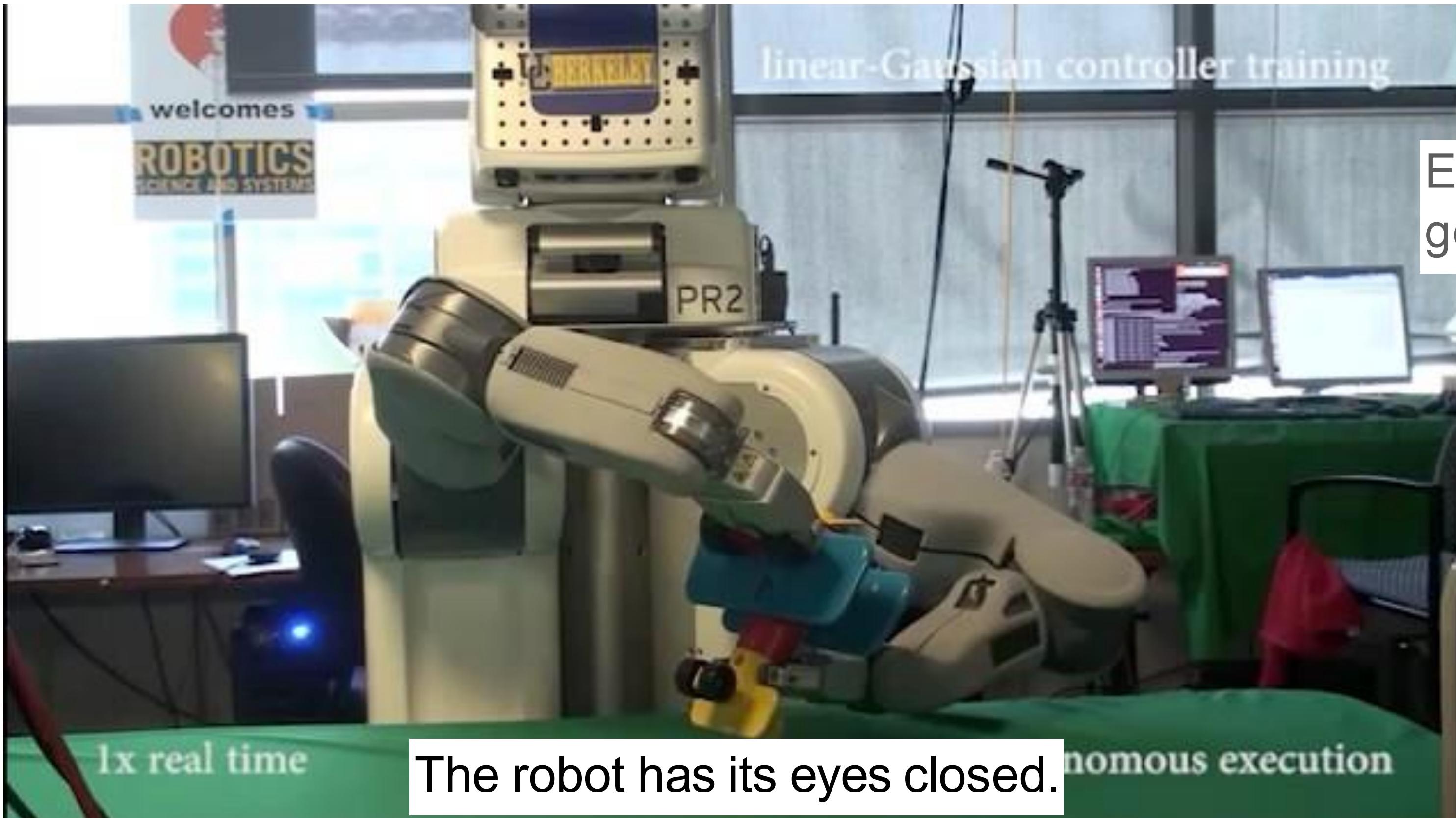
**Not just robots and games!**

Chip design, in Google's production TPU chips



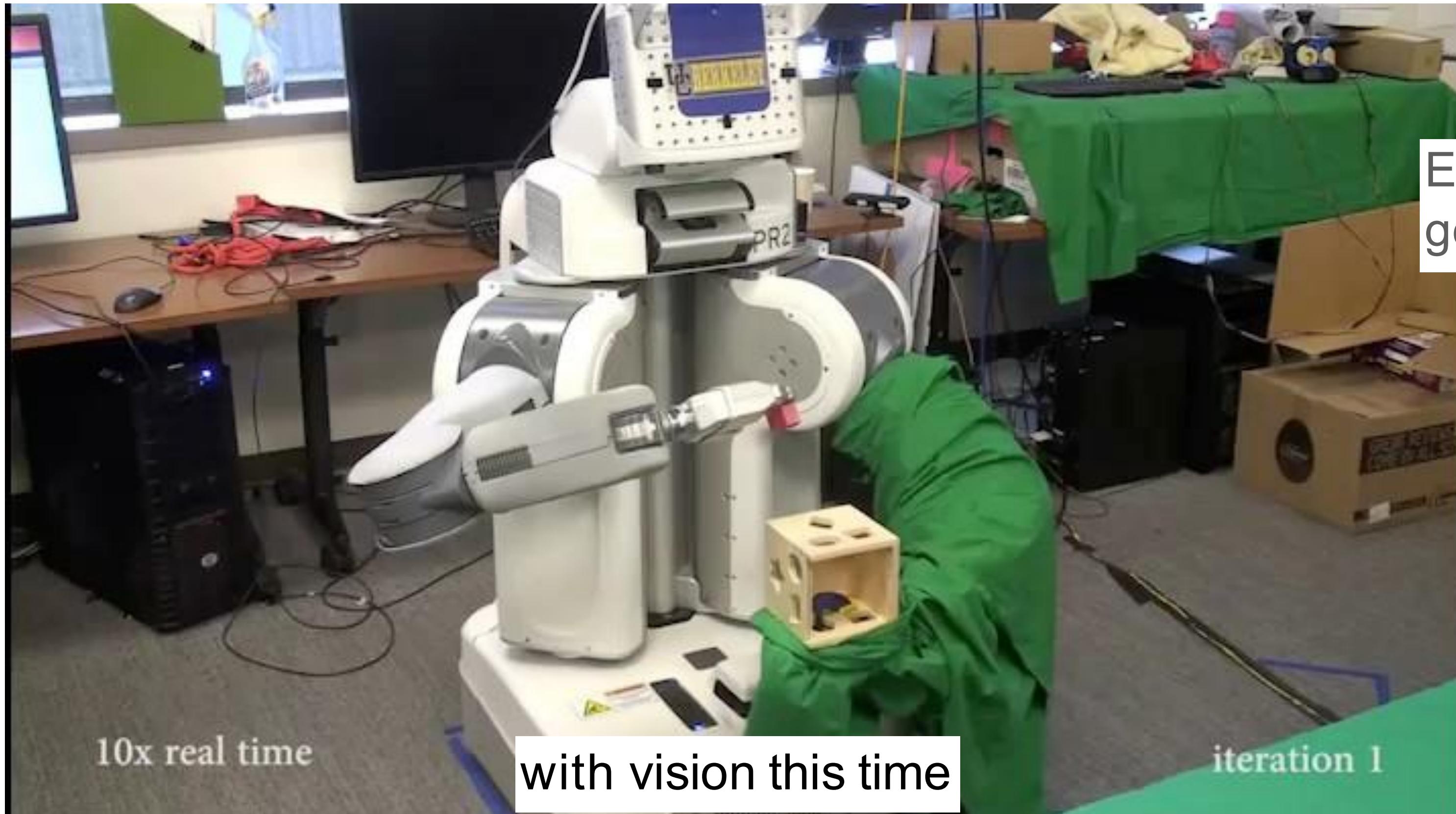
# Why study deep reinforcement learning?

Fundamental aspect of intelligence



# Why study deep reinforcement learning?

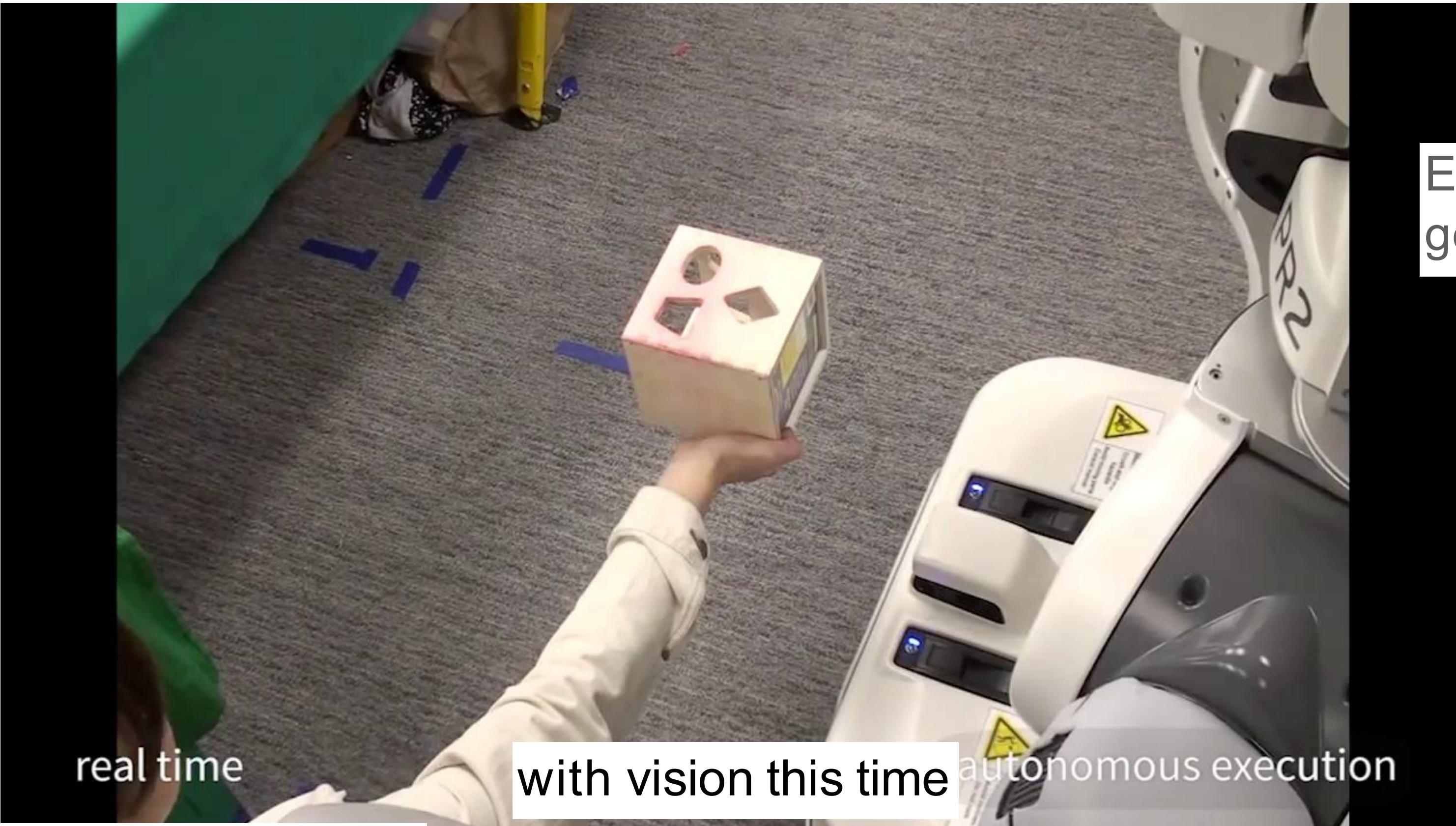
Fundamental aspect of intelligence



Enables the ability to  
get better with practice

# Why study deep reinforcement learning?

Fundamental aspect of intelligence



Enables the ability to  
get better with practice

# Why study deep reinforcement learning?

Still lots of exciting research problems!

How does robot learn to represent what is good or bad for the task? —> reward learning

How can an agent generalize its behavior to many different scenarios?

(Can we apply such a system at scale?)

Leverage large, diverse datasets —> offline RL

Transfer from other tasks, goals —> multitask RL, meta-RL

Can use RL to learn long-horizon tasks, like cooking a meal? —> hierarchy, reasoning

Can robots practice fully autonomously? —> reset-free RL

Behind the scenes of  
RL...



Yevge  
n

Yevgen is doing more work than the robot!  
It's not practical to collect a lot of data this way.

# The Plan for Today

1. Course goals & logistics
2. Why study deep reinforcement learning?
- 3. Intro to modeling behavior and reinforcement learning**

# How to represent experience as data?

*state*  $s_t$  - the state of the “world” at time  $t$

*observation*  $o_t$  - what the agent observes at time  $t$

(only used when missing information)

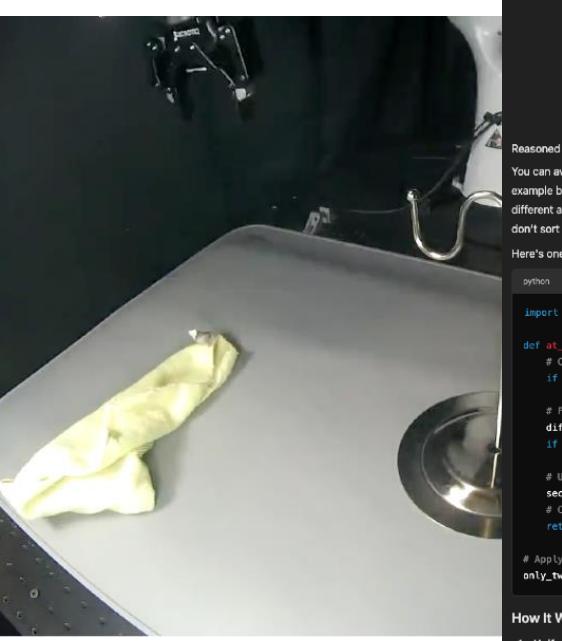
*action*  $a_t$  - the decision taken at time  $t$

*trajectory* - sequence of states/observations and actions

$(s_1, a_1, s_2, a_2, \dots, s_T, a_T)$

(could be length  $T=1!$ )

*reward function*  $r(s, a)$  - how good is  $s, a$ ?



I'm using this line of code to detect if there are at most 2 unique values in a dimension of an array, but the code is really slow (I think because unique() is doing a lot more work than needed). Can you write a version of the code that is faster?

```
only_two_unique_values = np.unique(x).size == 2  
len(np.unique(x)) >= 2
```

Reasoned about optimizing unique value check for 8 seconds

You can avoid the overhead of computing every unique element by “short-circuited” checking. In the example below, we first check if all values equal the first element; if not, we grab the first value that's different and then verify that every element is either the first value or that second one. This way we don't sort or scan the entire array for all uniques when more than two are present.

Here's one implementation:

```
python ⌂ Copy  
  
import numpy as np  
  
def fast_two_unique(values):  
    # values shape: (A, B, C)  
    # Get the first element from each slice.  
    first = values[..., 0] # shape: (A,)  
  
    # Create a mask that is True where elements differ from the first element.  
    mask = (values != first)[..., None]  
  
    # Find, for each slice, the index of the first differing element.  
    idx = np.argmax(mask, axis=0)  
  
    # Extract the candidate second unique value.  
    candidate2 = np.take_along_axis(values, idx[...], axis=0).squeeze(axis=0)  
  
    # Check that every element is either equal to 'first' or 'candidate2'.  
    second_val = candidate2[0]  
    return np.all((values == first) | (values == candidate2[...], None))  
  
# Example usage:  
# only_two_unique_values = np.apply_along_axis(fast_two_unique, arr=values)
```

How It Works

1. Uniform Check: If all values in the slice are identical, the function returns `True` immediately.
2. Candidate Selection: If not, it finds the first value that differs from the first element.
3. Verification: Finally, it verifies that every element in the slice is either equal to the first element or this second candidate.

Additional Notes

- This approach avoids the full overhead of `np.unique` and should be faster, especially when slices contain many elements or more than two unique values.
- For even greater speed (especially if your slices are large), you might consider using a just-in-time compiler (like Numba) to compile this function.

This version should offer a noticeable performance improvement compared to using `np.unique` on every slice.

Explanation

- Extracting the First Element:  
We set `first = values[..., 0]` so each  $(i, j)$  slice's first value becomes our reference.
- Finding a Candidate for the Second Unique Value:  
The mask `(values != first)[..., None]` identifies elements that differ from the first value. Using `np.argmax` finds the first index along axis 0 where the element is different. If all elements are the same, `np.argmax` returns 0 and `candidate2` will equal `first`.
- Verifying Uniqueness:  
Finally, we check that every element in each slice is either equal to `'first'` or `'candidate2'`. If a third unique value exists, the condition fails for that slice.

This vectorized solution should be considerably faster than using `apply_along_axis` with a Python loop.

That works great! Can you generalize the function to work with any axis, not just axis 0?

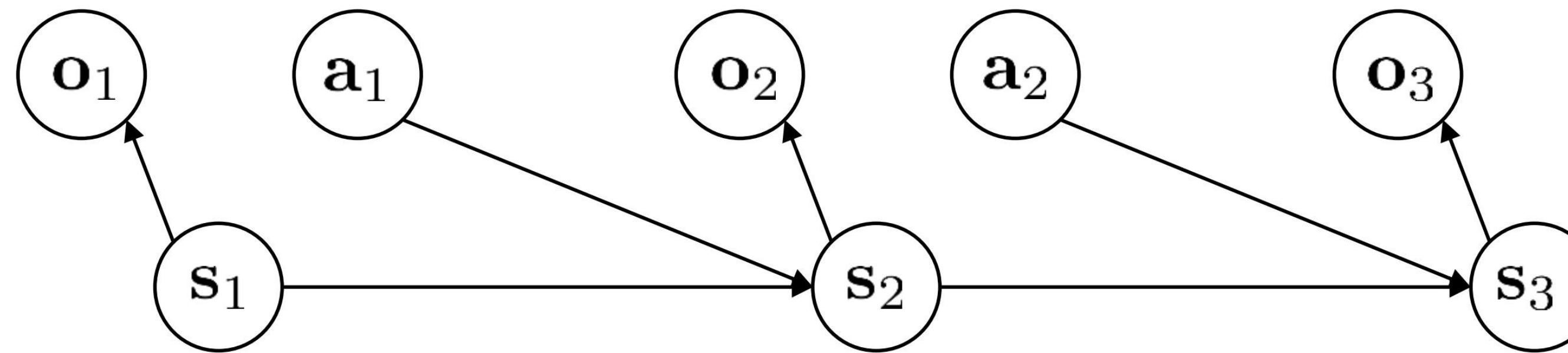
Reasoned about generalized vectorized solution for 5 seconds

Below is a generalized version that lets you choose the axis along which to check for at most two unique values. It works by moving the target axis to the end, applying the logic, and then reducing that axis:

```
python ⌂ Copy  
  
def fast_generalized_two_unique(values):  
    # values shape: (A, B, C)  
    # Get the first element from each slice.  
    first = values[..., 0] # shape: (A,)  
  
    # Create a mask that is True where elements differ from the first element.  
    mask = (values != first)[..., None]  
  
    # Find, for each slice, the index of the first differing element.  
    idx = np.argmax(mask, axis=-1)  
  
    # Extract the candidate second unique value.  
    candidate2 = np.take_along_axis(values, idx[...], axis=-1).squeeze(-1)  
  
    # Check that every element is either equal to 'first' or 'candidate2'.  
    second_val = candidate2[0]  
    return np.all((values == first) | (values == candidate2[...], None))
```

# States vs. observations

Next state is purely a function of the current state and action (and randomness)



unknown dynamics  $p(s_{t+1} | s_t, a_t)$

*independent of  $s_{t-1}$*

**“Markov property”**

# Examples



*state s* - RGB images, joint positions, joint velocities

*action a* - commanded next joint position

*trajectory*  $\tau$  - 10-sec sequence of camera, joint readings, controls at 20 Hz

$$(s_1, a_1, s_2, a_2, \dots, s_T, a_T), \quad = 200$$

*reward r(s, a)* = 1 if the towel is on the hook in state 0 otherwise

```

I'm using this line of code to detect if there are at most 2 unique values in a dimension of an array, but the code is really slow! What better implementation would you suggest than needed?
Can you write a version of the code that is faster?
only_two_unique_values = np.apply_along_axis(lambda x:
    only_x[~np.unique(x)], axis=0, arrvalues < 2)

Reasoned about optimizing unique value check for 8 seconds
You can avoid the overhead of computing every unique element by "short-circuited" checking. In the example below, we first check if all values equal the first element. If we grab the first value that's different and then verify that every element is either the first value or that second one. This way you don't sort or scan the entire array for all uniques when more than two are present.

Here's one implementation:
jupyter
import numpy as np

def at_most_two_unique(values):
    # Get the first element from each slice.
    first = values[:, 0] # shape (A, B)

    # Create a mask that is True where elements differ from the first element.
    mask = (values != first) ...

def at_most_two_unique(values):
    # Find, for each slice, the index of the first differing element.
    # np.argmax returns 0 if no True exists; that's fine because then candidate2 equals first.
    diff_idx = np.argmax(values, axis=1)

    # Extract the candidate second unique value.
    candidate2 = np.unique(values[:, diff_idx], axis=1).squeeze(axis=1)

    # Check that every element is either equal to 'first' or 'candidate'.
    return np.all((values == first) | (values == candidate2))

# Only two unique values is a boolean array indicating, for each (i, j) slice, if the
only_two_unique_values = fast_two_unique_values

# Apply along axis 2.
only_two_unique_values = np.apply_along_axis(at_most_two_unique, axis=0, arrvalues)

How it Works
1. Uniform Check: If all values in the slice are identical, the function returns 'first' immediately.
2. Candidate Selection: If not, it finds the first value that differs from the first element.
3. Verification: Finally, it verifies that every element in the slice is either equal to the first element or this second candidate.

Additional Notes
- This approach avoids the full overhead of np.unique and should be faster, especially when slices contain many elements or more than two unique values.
- For even greater speed (especially if your slices are large), you might consider using a just-in-time compiler (like NumPy) to compile this function.

This version should offer a noticeable performance improvement compared to using np.unique on every slice.

Hrm, that made the code slower. Can you make it faster instead of slower? (maybe using a different function than apply_along_axis)

Reasoned about optimizing unique value check for 18 seconds
jupyter
import numpy as np

def fast_two_unique(values):
    # We set 'first' = values[:, 0] so each (i, j) slice's first value becomes our reference.
    # The mask (values != first)[:, None] identifies elements that differ from the first value.
    # Using np.argmax finds the first index along axis 2 where the element is different. If all elements are the same, np.argmax returns 0 and candidate2 will equal 'first'.
    diff_idx = np.argmax(values, axis=1)

    # Check that every element in each slice is either equal to 'first' or 'candidate2'. If a third unique value exists, the condition fails for that slice.

Explanation
• Extracting the First Element:
  We set 'first' = values[:, 0] so each (i, j) slice's first value becomes our reference.
• Finding a Candidate for the Second Unique Value:
  The mask (values != first)[:, None] identifies elements that differ from the first value.
  Using np.argmax finds the first index along axis 2 where the element is different. If all elements are the same, np.argmax returns 0 and candidate2 will equal 'first'.
• Verification:
  Finally, we check that every element in each slice is either equal to 'first' or 'candidate2'. If a third unique value exists, the condition fails for that slice.

This vectorized solution should be considerably faster than using apply_along_axis with a Python loop.

That works great! Can you generalize the function to work with any axis, not just axis 2?

Reasoned about generalized vectorized solution for 5 seconds
Below is a generalized version that lets you choose the axis along which to check for at most two unique values. It works by moving the target axis to the end, applying the logic, and then reducing that axis.
jupyter

```

*observation o* - the user's most recent message

*action a* - chatbot's next message

*trajectory*  $\tau$  - variable length conversation trace

$$(\mathbf{o}_1, \mathbf{a}_1, \mathbf{o}_2, \mathbf{a}_2, \dots, \mathbf{o}_T, \mathbf{a}_T)$$

*reward r(s, a)* = 1 if the user gives upvote  
-10 if the user downvotes  
0 if no user feedback

# Think: how to represent another example?



autonomous driving



web agent



poker player



choose your own!

## Define

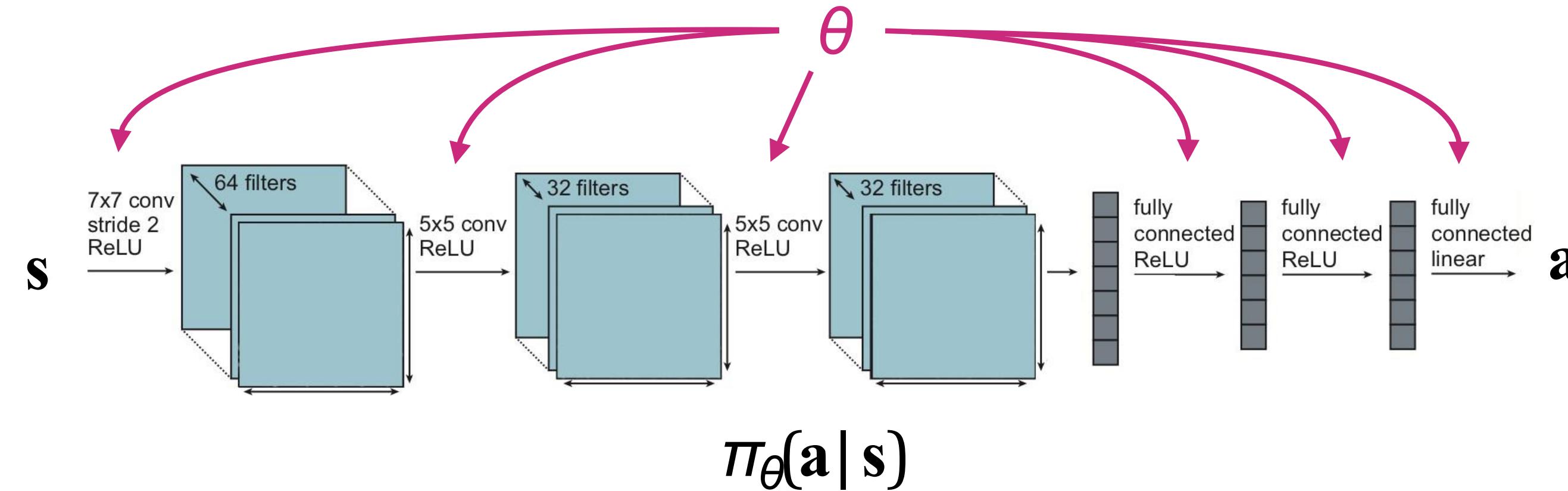
state **s** or observation **o**

action **a**

trajectory  $\tau$

reward  $r(s, a)$

# How to represent behavior with a neural network?



Observe state  $\mathbf{s}_t$

Take action  $\mathbf{a}_t$  (e.g. by sampling from policy  $\pi_\theta(\cdot | \mathbf{s}_t)$ )

Observe next state  $\mathbf{s}_{t+1}$  sampled from unknown world dynamics  $p(\cdot | \mathbf{s}_t, \mathbf{a}_t)$

Result: a *trajectory*  $\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T$ , also called a policy *roll-out* or an *episode*

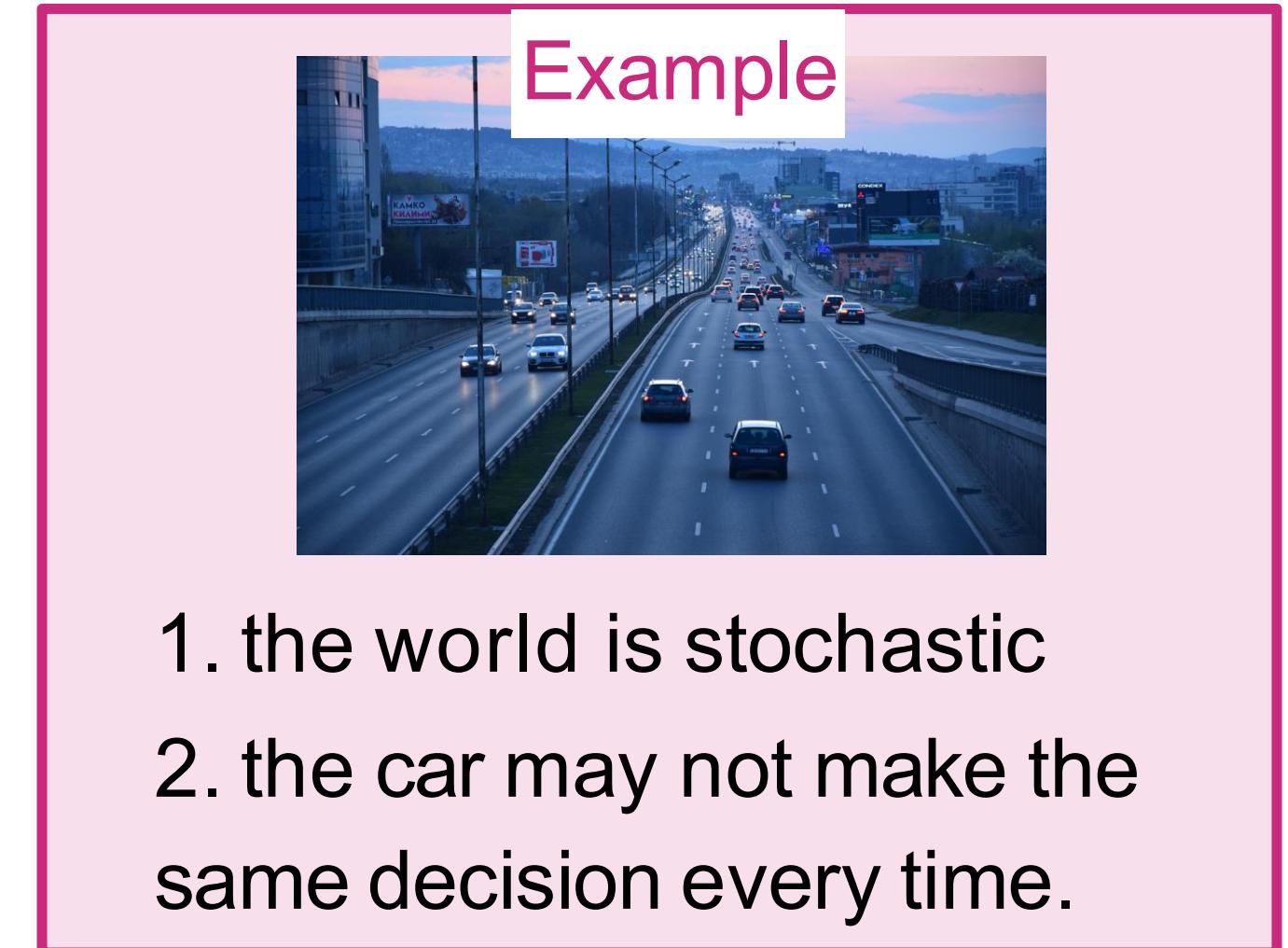
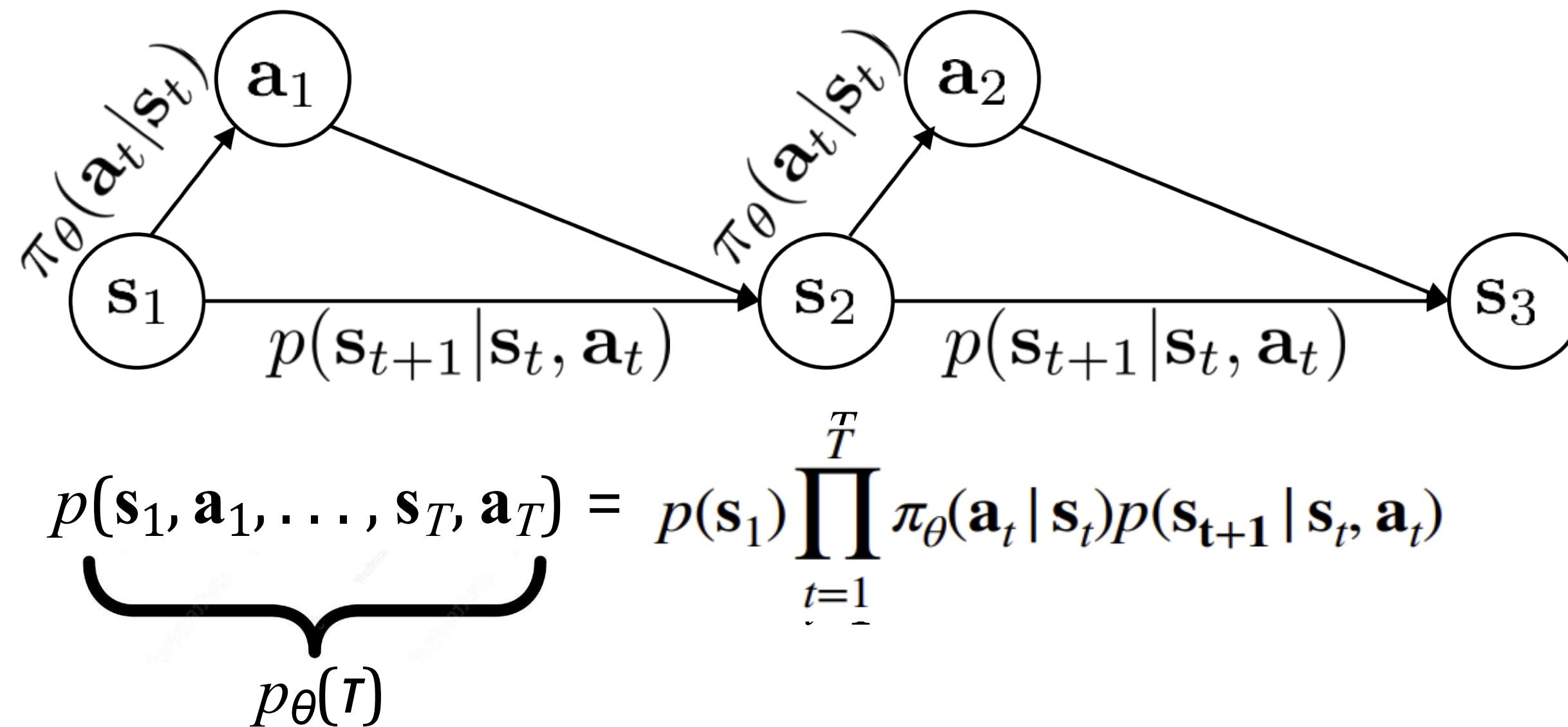
If you only have observations  $\mathbf{o}$ , give the policy memory:  $\pi_\theta(\mathbf{a}_t | \mathbf{o}_{t-m}, \dots, \mathbf{o}_t)$

# What is the goal of reinforcement learning?

maximize sum of rewards:  $\max \sum_t^T r(\mathbf{s}_t, \mathbf{a}_t)$

but this is not a deterministic quantity!

Question: what are the sources of variability?

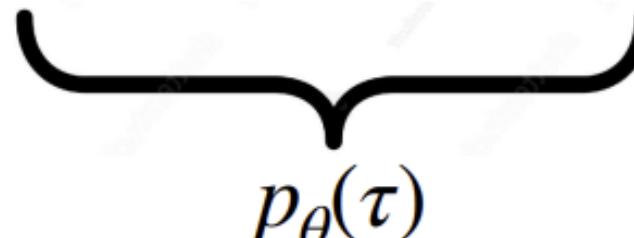


# What is the goal of reinforcement learning?

maximize sum of rewards:  $\max \sum_t r(\mathbf{s}_t, \mathbf{a}_t)$

maximize expected sum of rewards:  $\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$

$$p(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

  
 $p_{\theta}(\tau)$

## Aside: why stochastic policies?

1. Exploration: to learn from your own experience, must try different things.
2. Modeling stochastic behavior: existing data will exhibit varying behaviors

We can leverage tools from generative modeling!

—> generative model over *actions* given states/observations

# What is the goal of reinforcement learning?

maximize expected sum of rewards:  $\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t^T r(s_t, a_t) \right]$

$$p(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

  
 $p_{\theta}(\tau)$

**How good is a particular policy?**

*value function*  $V^{\pi}(s)$  - future expected reward starting at  $s$  and following

*Q-function*  $Q^{\pi}(s, a)$  - future expected reward starting at  $s$ , taking  $a$ , then following  $\pi$

# Types of algorithms

maximize expected sum of rewards:  $\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t^T r(s_t, a_t) \right]$

1. Imitation learning: mimic a policy that achieves high reward
2. Policy gradients: directly differentiate the above objective
3. Actor-critic: estimate value of the current policy and use it to make the policy better
4. Value-based: estimate value of the optimal policy
5. Model-based: learn to model the dynamics, and use it for planning or policy improvement

# Why so many algorithms?

Algorithms make different trade-offs, thrive under different assumptions.

- How easy / cheap is it to collect data with policy? (e.g. simulator vs. hand-written)
- How easy / cheap are different forms of supervision? (demos, detailed rewards)
- How important is stability and ease-of-use?
- Action space dimensionality, continuous vs. discrete
- Is it easy to learn the dynamics model?

# Recap of definitions

*state  $s_t$*  - the state of the “world” at time  $t$

(or *observation  $o_t$*  - what the agent observes at time  $t$ )

*action  $a_t$*  - the decision taken at time  $t$

*reward function  $r(s, a)$*  - how good is  $s, a$ ?

*initial state distr.  $p(s_1)$ , unknown dynamics  $p(s_{t+1} | s_t, a_t)$*

(partially-observed)  
Markov decision  
process

MDP, POMDP

# Recap of definitions

*trajectory*  $\tau$  - sequence of states/observations and actions  $(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T)$

*policy*  $\pi$  - represents behavior, selecting actions based on states or observations

**Goal:** learn policy  $\pi_\theta$  that maximizes expected sum of rewards:

$$\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

*value function*  $V^\pi(\mathbf{s})$  - future expected reward starting at  $\mathbf{s}$  and following  $\pi$

*Q-function*  $Q^\pi(\mathbf{s}, \mathbf{a})$  - future expected reward starting at  $\mathbf{s}$ , taking  $\mathbf{a}$ , then following  $\pi$

# Imitation Learning

CS 224R

# Partial Recap

*state  $s_t$*  - the state of the “world” at time

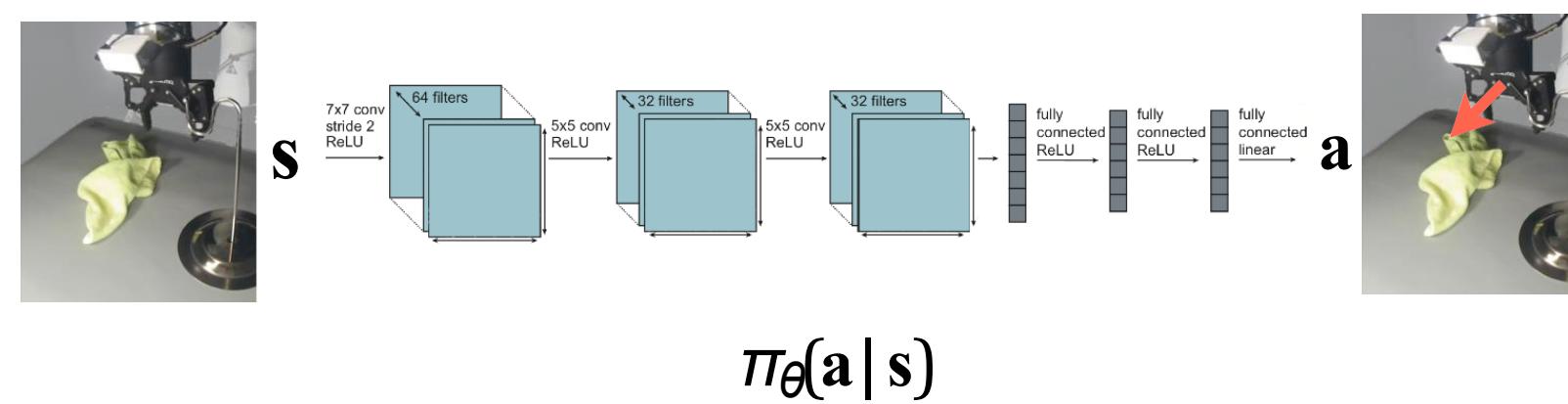
*observation  $\mathbf{o}_t$*  - what the agent observes at time

*action  $a_t$*  - the decision taken at time

*trajectory* - sequence of states/observations and actions

$$(s_1, a_1, s_2, a_2, \dots, s_T, a_T)$$

*reward function  $r(s, a)$*  - how good is  $s, a$ ?



*policy  $\pi(a_t|s_t)$  or  $\pi(a_t|\mathbf{o}_{t-m:t})$*  - behavior, usually what we are trying to learn

can be represented using a generative model

# The plan for today

## Imitation Learning

1. Imitation learning basics
2. Learning expressive policy distributions
3. Learning from online interventions
4. Time permitting: how to collect demonstrations

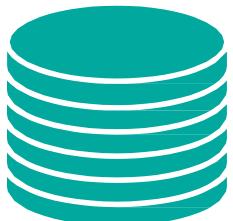
## Key learning goals:

- how to represent distributions with neural networks
- why expressive distributions matter for imitation learning
- what are compounding errors and how to address them

# What is the goal of imitation learning?

**Data:** Given trajectories collected by an expert

“demonstrations”  $\mathcal{D} := \{(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T)\}$



(sampled from some unknown policy  $\pi_{\text{expert}}$ )

**Goal:** Learn a policy  $\pi_\theta$  that performs at the level of the expert policy, by mimicking it.

Example



- Dataset from human drivers
- Sensor readings + steering commands

# Imitation learning - version 0

Deterministic policy

0. Given demonstrations collected by an expert  $\mathcal{D} := \{(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T)\}$

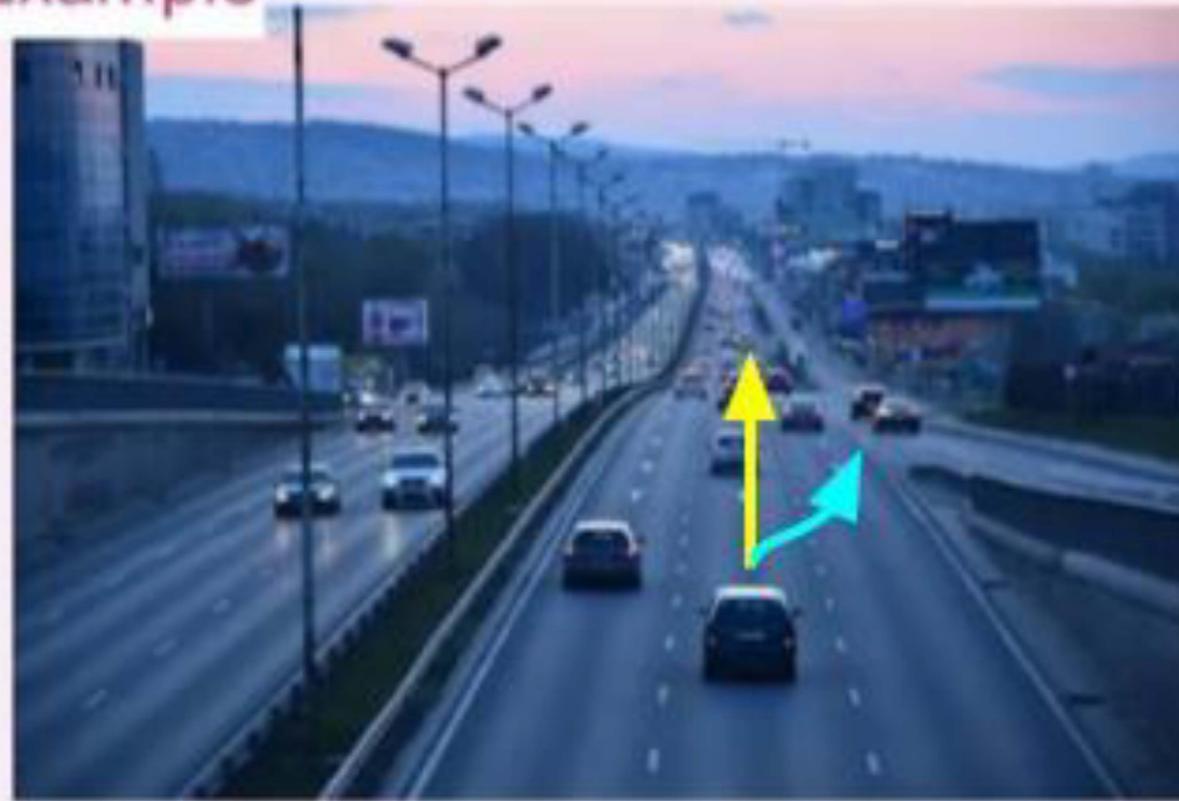
1. For deterministic policy, supervised regression to the expert's actions

$$\min_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{s}, \mathbf{a}) \in \mathcal{D}} \|\mathbf{a} - \hat{\mathbf{a}}\|^2 \quad \text{where } \hat{\mathbf{a}} = \pi_{\theta}(\mathbf{s})$$

2. Deploy learned policy  $\pi_{\theta}$

# What could go wrong?

## Example

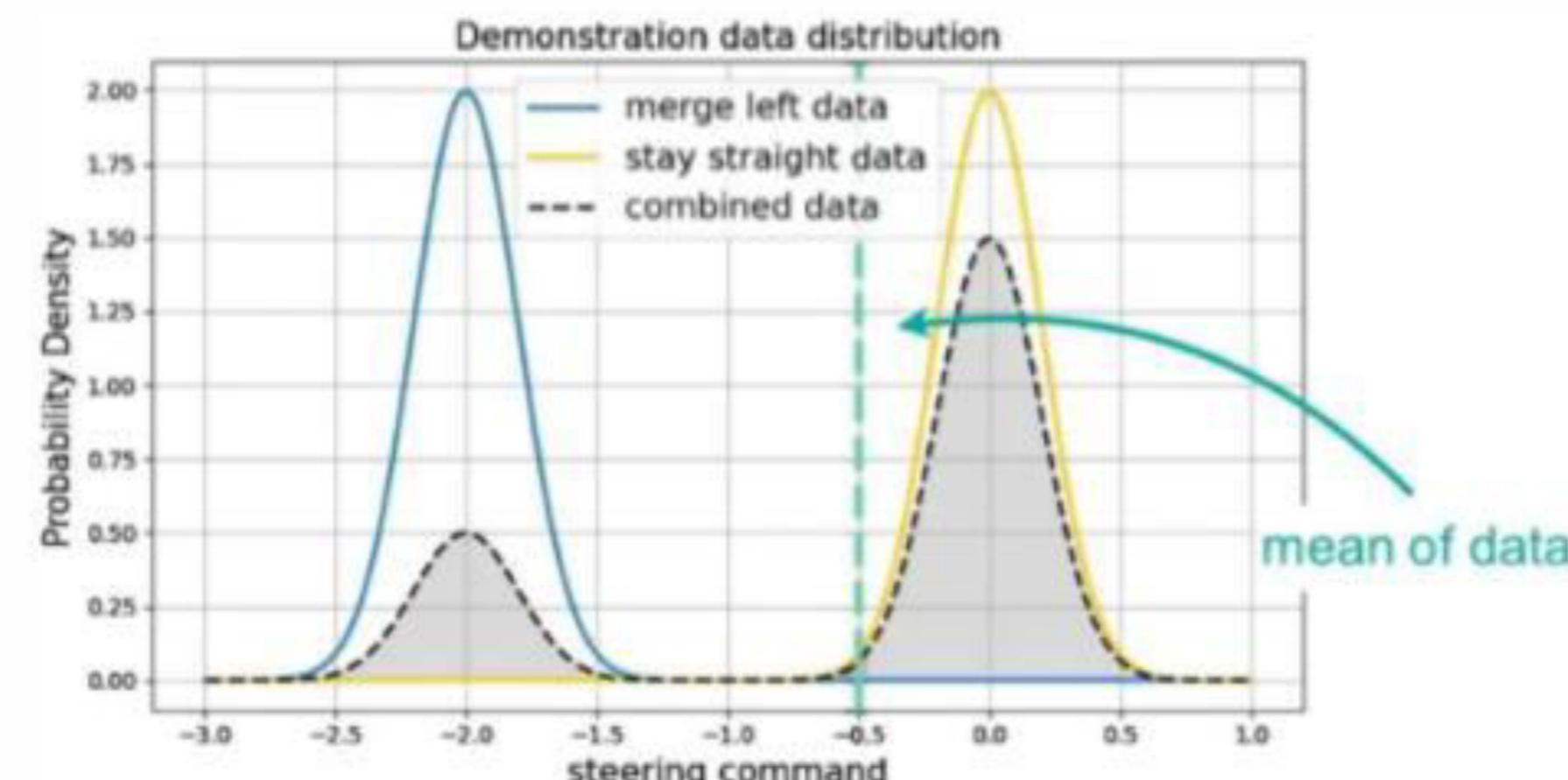


- Dataset from human drivers
- Sensor readings + steering commands

How often does this happen in practice? All the time!

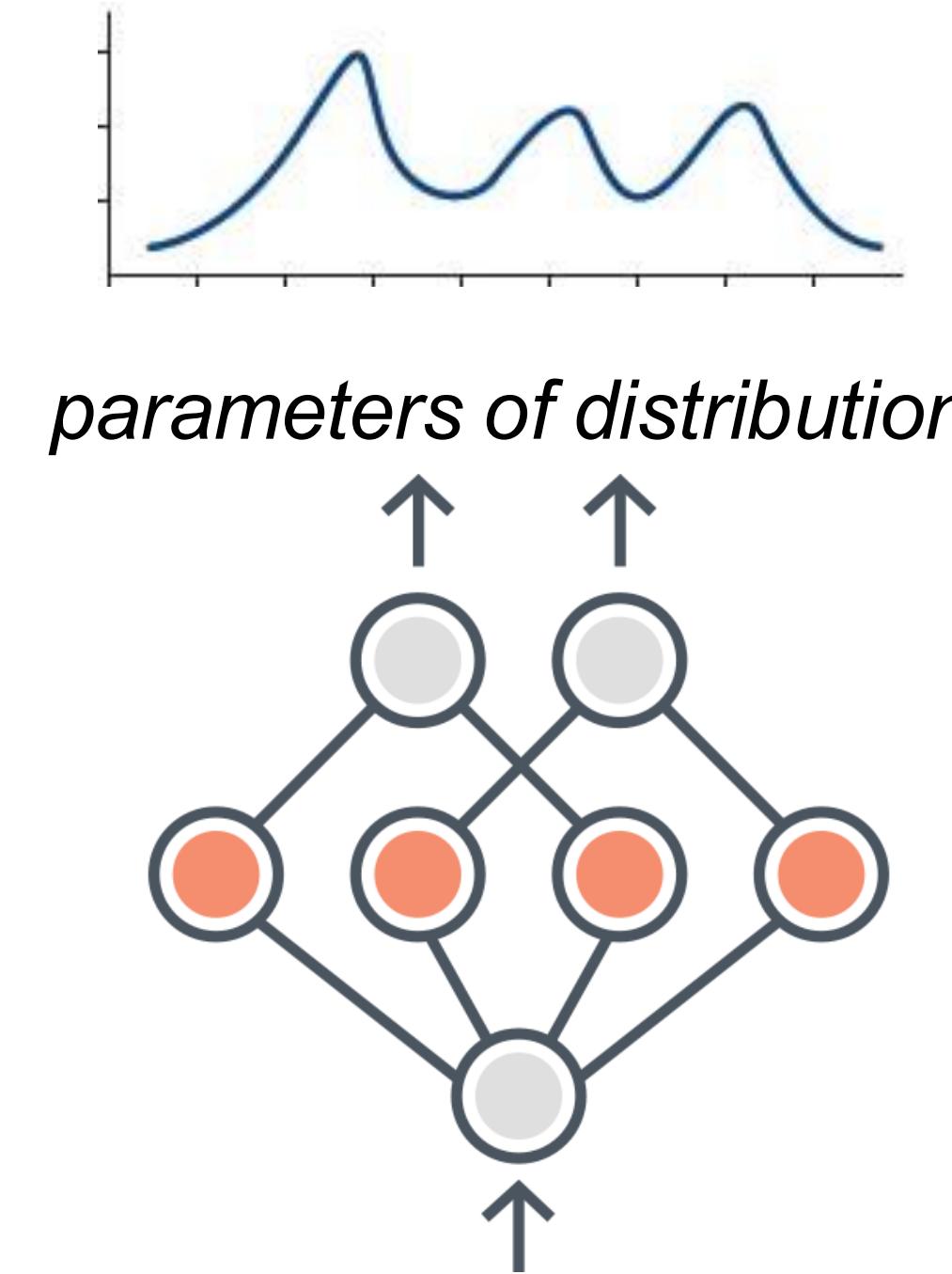
Esp. when data is collected by multiple people.

Question: what might policy trained with  $\ell_2$ -regression do?

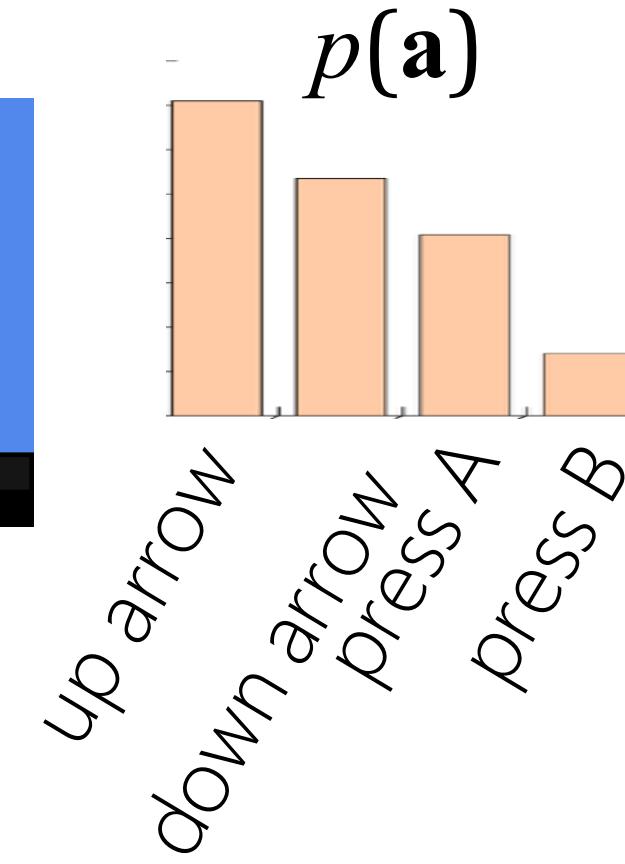
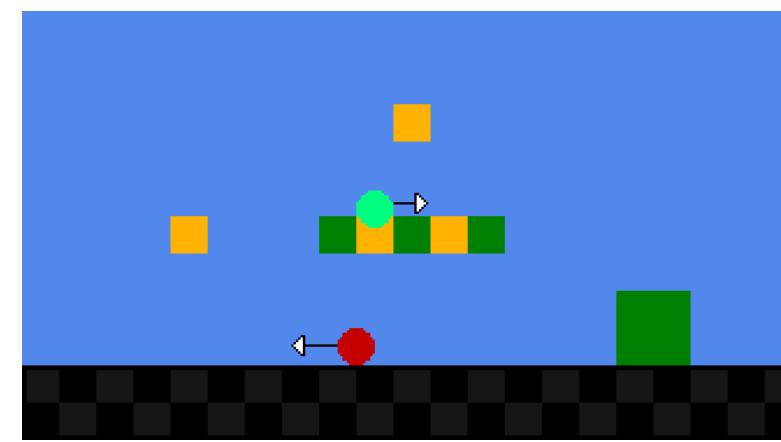


How can we  
represent more than  
the mean?

# Learning *distributions* with neural networks

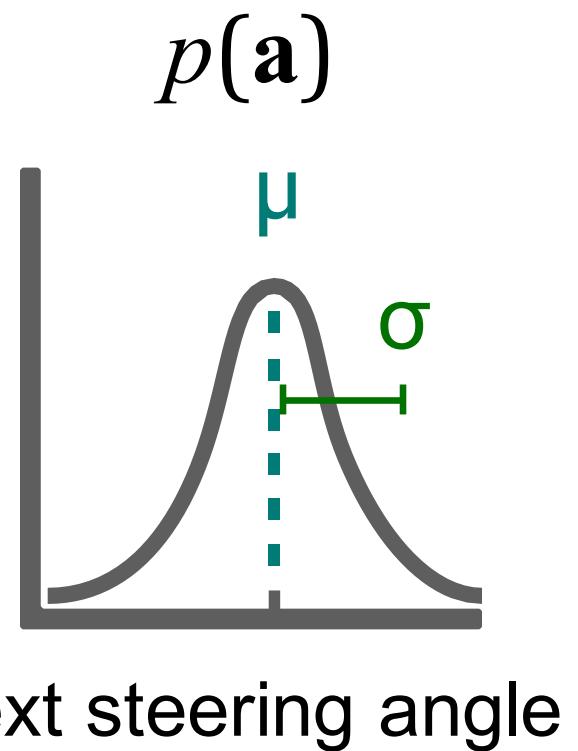


## 1D discrete actions



Neural net outputs  $p(\text{up})$ ,  $p(\text{down})$ , ...  
represent categorical distribution.  
**Maximally expressive**

## Continuous actions



Neural net outputs  $\mu$ ,  $\sigma$  to represent  
Gaussian distribution.  
**Not very expressive!**

# Learning *distributions* with neural networks

💡 Can we use generative modeling?

image diffusion models



autoregressive models

*Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat.*

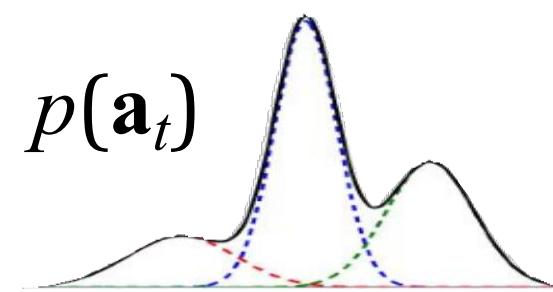
learning  $p(\text{image} \mid \text{text description})$

learning  $p(\text{next word} \mid \text{words so far})$

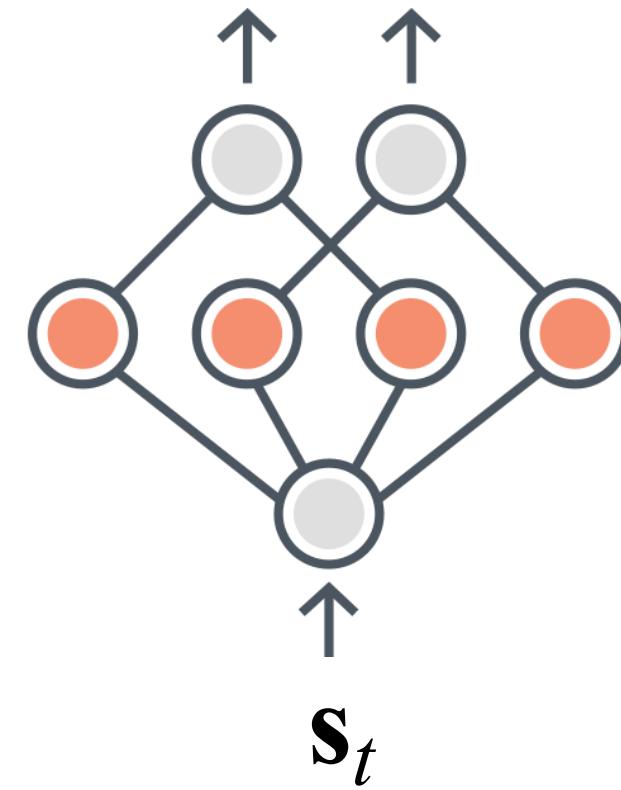
Our goal: learning  $p(\text{action} \mid \text{observations})$

# Generative models for policies (approximating $p(\mathbf{a} | \mathbf{s})$ )

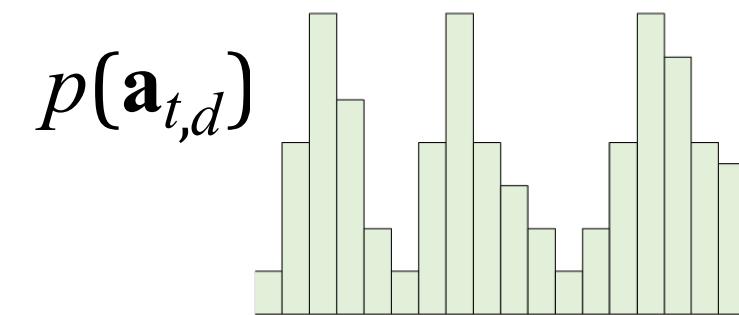
## Mixture of Gaussians



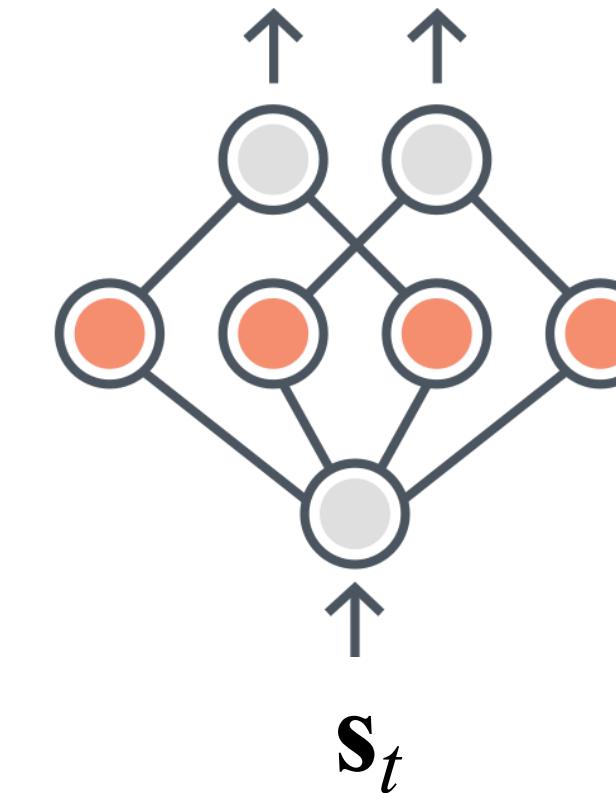
output  $\mu_1, \sigma_1, w_1, \mu_2, \sigma_2, w_2, \dots$



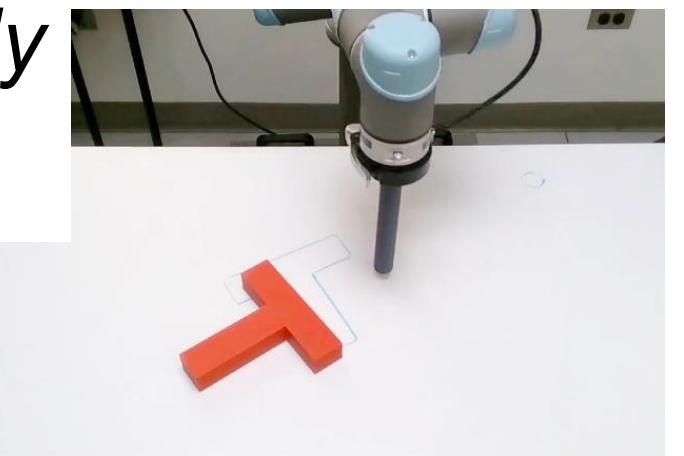
## Discretize + Autoregressive



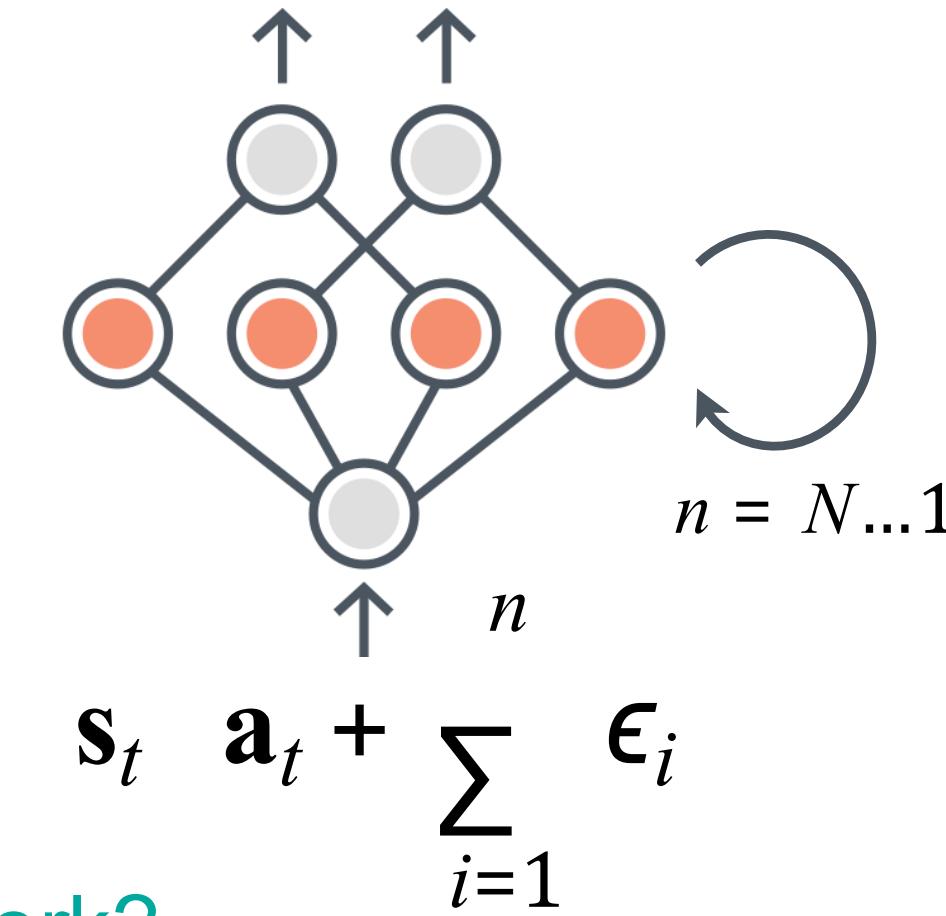
output  $p(\mathbf{a}_{t,1}), p(\mathbf{a}_{t,2} | \hat{\mathbf{a}}_{t,1}), p(\mathbf{a}_{t,3} | \hat{\mathbf{a}}_{t,1:2}), \dots$



*iteratively  
denoise*



output  $\mathbf{a}_t$



Question: how are these different from  $\ell_2$  loss with larger network?

**Important Note:** Neural network expressivity is often distinct from distribution expressivity!

# Imitation learning - version 1

Expressive policies

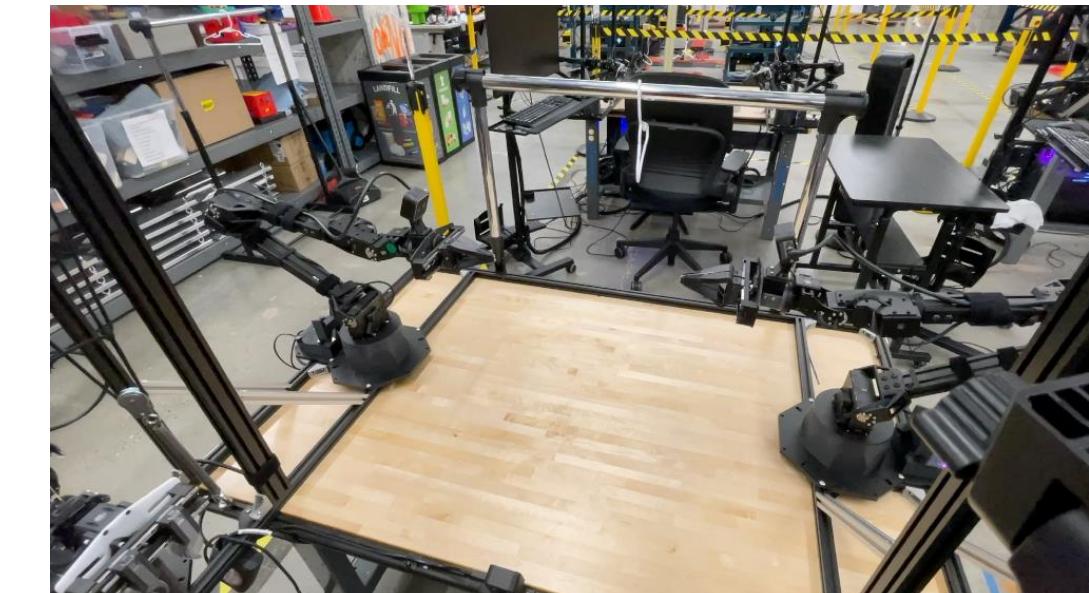
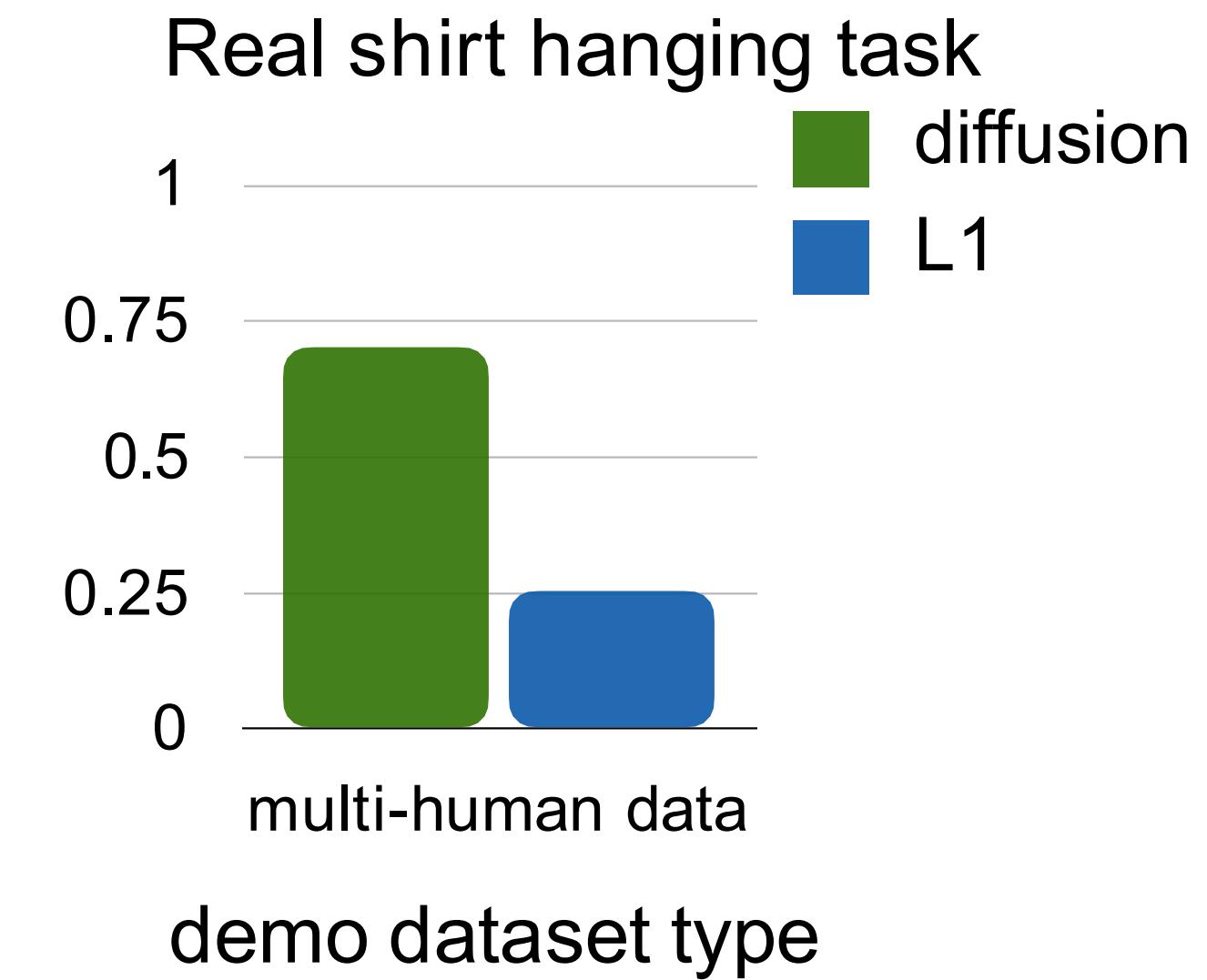
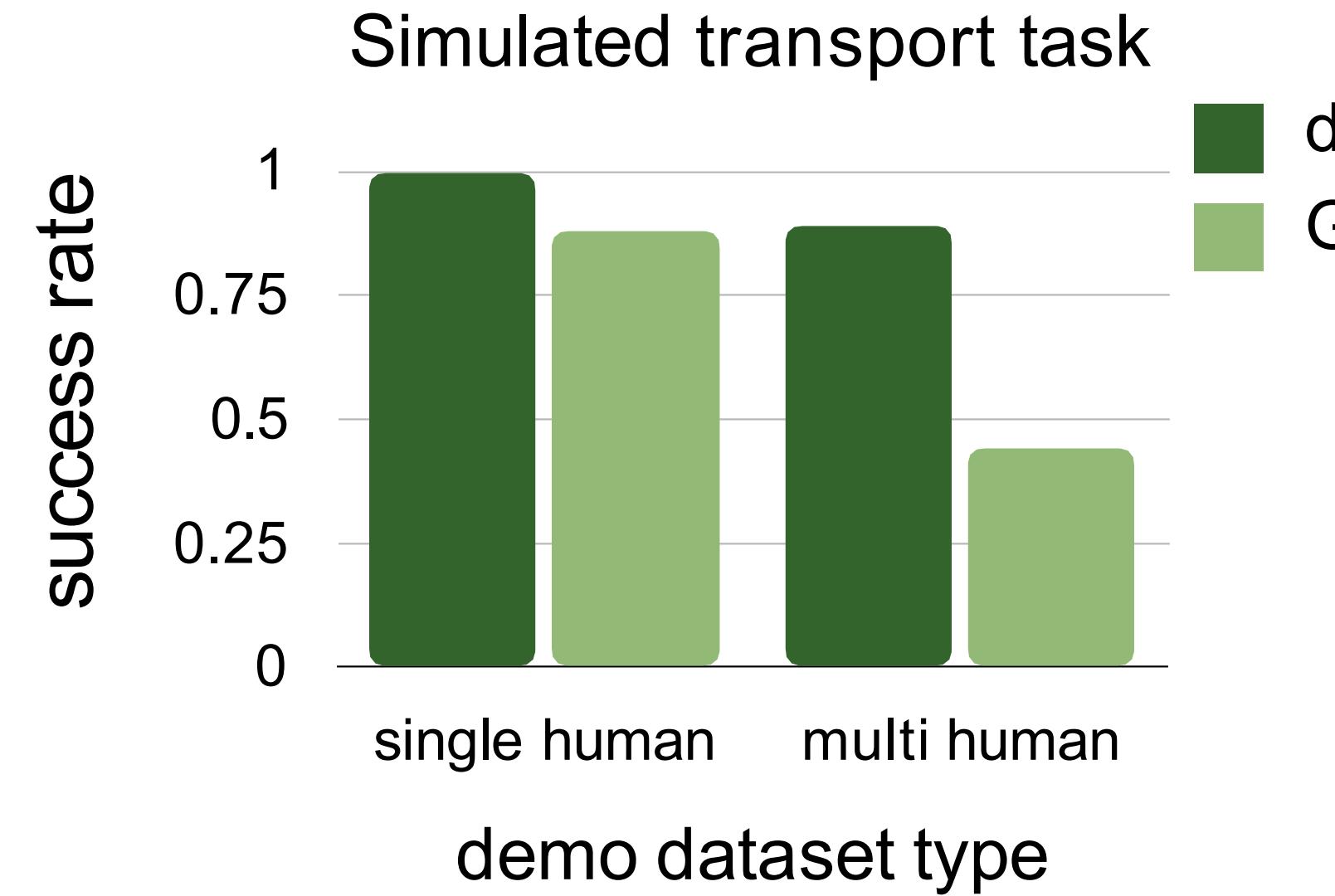
0. Given demonstrations collected by an expert  $\mathcal{D} := \{(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T)\}$

1. Train **generative model** of the expert's actions

$$\min_{\theta} - \mathbb{E}_{(s,a) \sim \mathcal{D}} [\log \pi_{\theta}(a | s)] \quad \text{with expressive distribution } \pi(\cdot | s)$$

2. Deploy learned policy  $\theta$       maximize the **log probability** of the  
   **demo actions** under the **policy**

# Imitation learning - version 1 vs version 0



# Robotics: Imitation learning + expressive policies

Physical Intelligence  $\pi_0$

NVIDIA Gr00t N1

[https://www.youtube.com/watch?v=YyXC\\_Mhnb\\_IU](https://www.youtube.com/watch?v=YyXC_Mhnb_IU)

<https://youtu.be/unnkmL6QJ78?feature=shared>

diffusion

diffusion

Figure Helix

OpenVLA

<https://youtu.be/Z3yQHYNXPws?feature=shared>

<https://openvla.github.io/>

diffusion

discretize + autoregressive prediction

# Autonomous driving: Imitation learning + expressive policies

Waymo EMMA

<https://waymo.com/blog/2024/10/introducing-emma/>

discretize + autoregressive prediction

Wayve LINGO-2

<https://wayve.ai/thinking/lingo-2-driving-with-language/>

discretize + autoregressive prediction

# Summary so far

Data from one consistent demonstrator

Unimodal policy distribution is enough

Multimodal data, e.g. from multiple demonstrators

Need expressive generative model for policy

# Summary so far

- Key idea: Train expressive policy class via generative modeling on dataset of demonstrations.

- Algorithm is fully *offline*

Definitions.

*offline*: using only an existing dataset, no new data from learned policy

*online*: using new data from learned policy

- + no need for data from policy (online data can be unsafe, expensive to collect)
- + no need to define a reward function
- may need **a lot** of data for reliable performance

# The plan for today

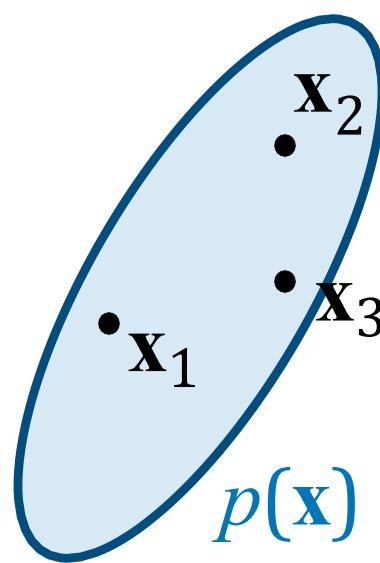
## Imitation Learning

1. Imitation learning basics
2. Learning expressive policy distributions
3. **Learning from online interventions**
4. Time permitting: how to collect demonstrations

# What can go wrong in imitation learning?

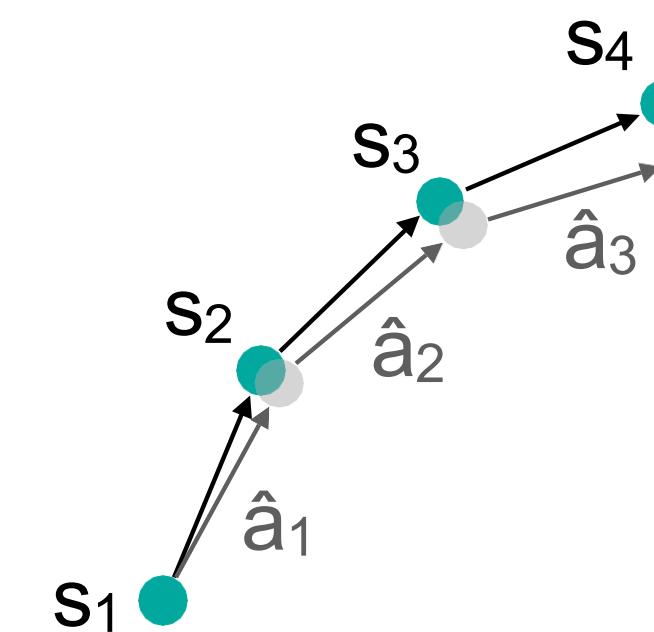
## Compounding errors

Supervised learning



Inputs independent  
of predicted labels  $\hat{y}$

Supervised learning of behavior



Predicted actions affect  
next state.

Errors can lead to drift  
away from the data  
distribution!

Errors can then compound!

$$\frac{p_{expert}(\mathbf{s})}{p_{\pi}(\mathbf{s})}$$

states visited  
by expert

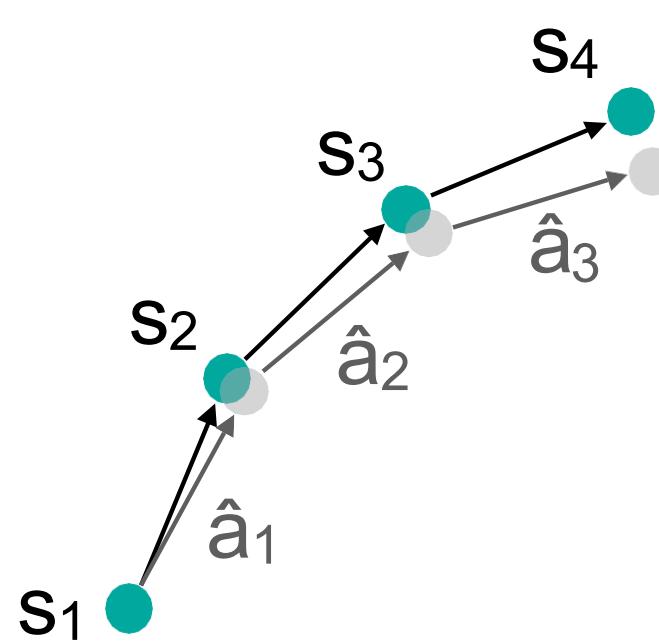
states visited by  
learned policy

“covariate shift”

# What can go wrong in imitation learning?

## Compounding errors

Supervised learning of behavior



Predicted actions affect next state.

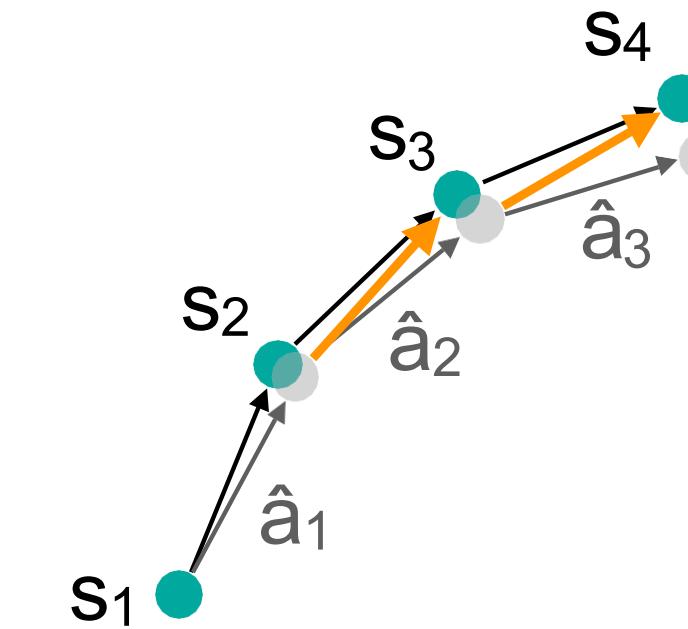
Errors can lead to drift away from the data distribution!

Errors can then compound!

$$p_{expert}(\mathbf{s}) \neq p_{\pi}(\mathbf{s})$$

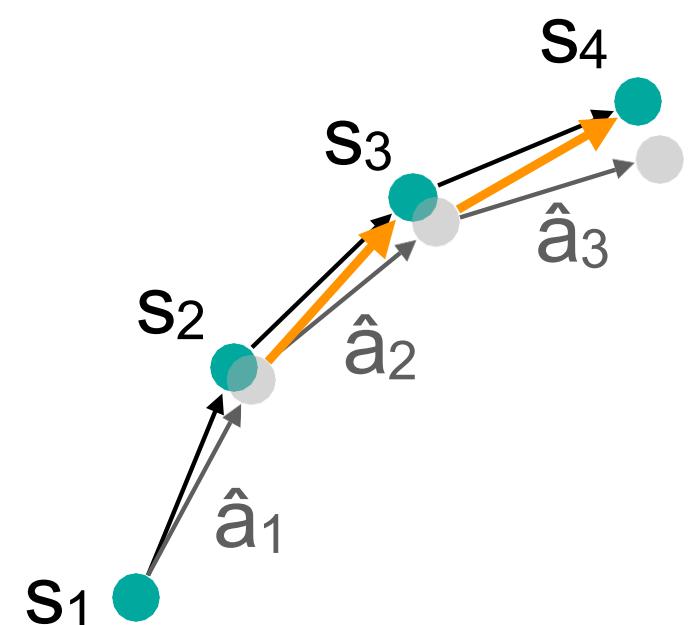
Solutions?

1. Collect A LOT of demo data & hope for the best.
2. Collect **corrective behavior data**



# Addressing Compounding Errors with Interventions

Collect corrective behavior data



1. Roll out learned policy  $\theta: \mathbf{s}'_1, \hat{\mathbf{a}}_1, \dots, \mathbf{s}'_T$
2. Query expert action at visited states  $\mathbf{a}^* \sim \pi_{expert}(\cdot | \mathbf{s}')$
3. Aggregate corrections with existing data  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}' \ \mathbf{a}^*)\}$
4. Update policy  $\min_{\theta} \mathcal{L}(\pi_{\theta}, \mathcal{D})$

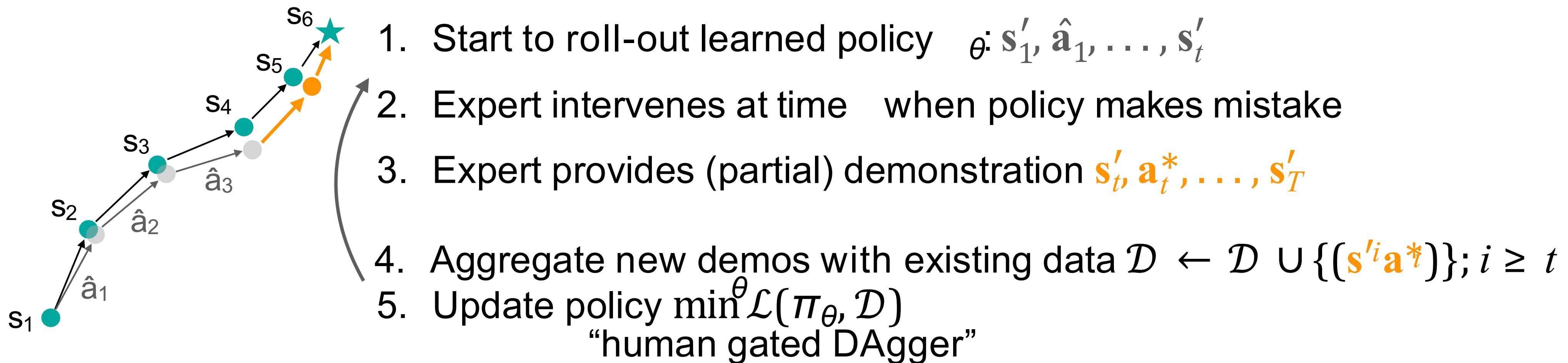
“dataset aggregation” (DAgger)

- + data-efficient way to learn from an expert
- can be challenging to query expert when agent has control

Is there another way to collect corrective data?

# Addressing Compounding Errors with Interventions

Collect corrective behavior data while *taking full control*



+ (much) more practical interface for providing corrections

- can be hard to catch mistakes quickly in some application domains

Question: could you automatically detect when intervention is needed?

# Policy Gradients

CS 224R

# Recap

*state*  $\mathbf{s}_t$  - the state of the “world” at time  $t$

*action*  $\mathbf{a}_t$  - the decision taken at time  $t$

*trajectory*  $\tau$  - sequence of states/observations and actions

$$(\mathbf{s}_1, \mathbf{a}_1, \mathbf{s}_2, \mathbf{a}_2, \dots, \mathbf{s}_T, \mathbf{a}_T)$$

*reward function*  $r(\mathbf{s}, \mathbf{a})$  - how good is  $\mathbf{s}, \mathbf{a}$ ?

*policy*  $\pi(\mathbf{a}|\mathbf{s})$  or  $\pi(\mathbf{a}|\mathbf{o})$  - behavior, usually what we are trying to learn

**Goal:** learn policy  $\pi_\theta$  that maximizes *expected sum of rewards*:

$$\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t^T r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

**Where:**  $p(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$

$p_{\theta}(\tau)$

Imitation learning from demonstrations  $\mathcal{D}$

$$\min_{\theta} - \mathbb{E}_{(s,a) \sim \mathcal{D}} [\log \pi_{\theta}(a|s)]$$

- + Simple, scalable approach for learning performant behavior
- Cannot outperform demonstrator, doesn't allow improvement from practice

**Definitions.**

***offline:* using only an existing dataset, no new data from learned policy**

***online:* using new data from learned policy<sup>3</sup>**

# The plan for today

the basis for:

Policy gradients: our first online RL algorithm

1. On-policy policy gradient
  - a. Derivation and intuition of policy gradients
  - b. Full algorithm
  - c. How to make it better - causality and baselines
2. Off-policy policy gradients
  - a. Importance sampling
  - b. KL constraints

Source: Unitree

<https://www.youtube.com/watch?v=xG7WkPU8tgs>

2. Creating a Boolean mask for elements that differ from that first element.

3. Using a vectorized argmax (via np.take\_along\_axis) to pick the first different value (which will equal the first element if none exist).

4. Verifying that every element in the slice is either equal to the first element or that candidate second unique value.

Try this code:

```
python ⌂ Copy
import numpy as np

def fast_two_unique(values):
    # values shape: (A, B, C)
    # Get the first element from each slice.
    first = values[..., 0] # shape: (A, B)

    # Create a mask that is True where elements differ from the first element.
    mask = (values != first)[..., None]

    # Find, for each slice, the index of the first differing element.
    # np.argmax returns 0 if no True exists; that's fine because then candidate2 equals None.
    idx = np.argmax(mask, axis=2)

    # Extract the candidate second unique value.
    candidate2 = np.take_along_axis(values, idx[...], None).squeeze(axis=2)

    # Check that every element is either equal to 'first' or 'candidate2'.
    return np.all((values == first)[..., None] | (values == candidate2)[..., None])

# Example usage:
# only_two_unique_values is a boolean array indicating, for each (i, j) slice, if the
only_two_unique_values = fast_two_unique(values)
```

Explanation

- Extracting the First Element:  
We set `first = values[..., 0]` so each (i, j) slice's first value becomes our reference.
- Finding a Candidate for the Second Unique Value:  
The mask `(values != first)[..., None]` identifies elements that differ from the first value. Using `np.argmax` finds the first index along axis 2 where the element is different. If all elements are the same, `np.argmax` returns 0 and `candidate2` will equal `first`.
- Verifying Uniqueness:  
Finally, we check that every element in each slice is either equal to `first` or `candidate2`. If a third unique value exists, the condition fails for that slice.

This vectorized solution should be considerably faster than using `apply_along_axis` with a Python loop.

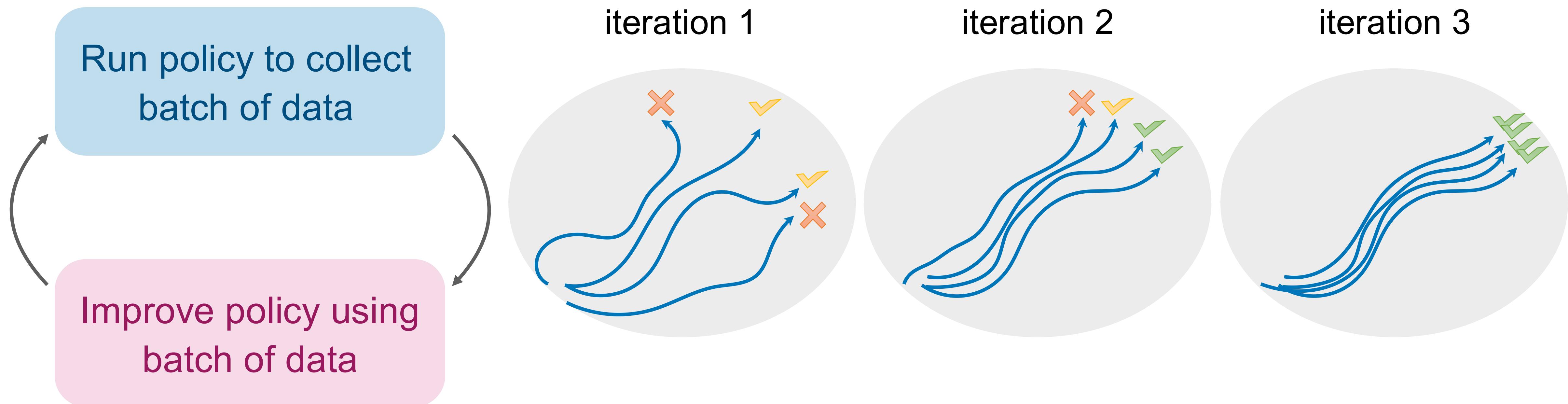
That works great! Can you generalize the function to work with any axis, not just axis 2?

Key learning goals:

- Key intuition behind policy gradients
- How to implement, when to use policy gradients

# Online RL Outline

First: Initialize the policy (randomly, with imitation learning, with heuristics)

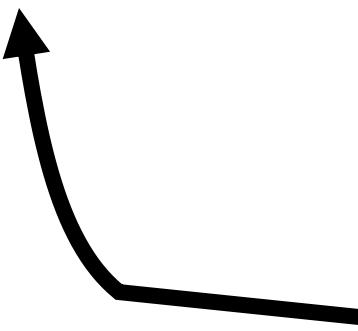


# Evaluating the RL objective

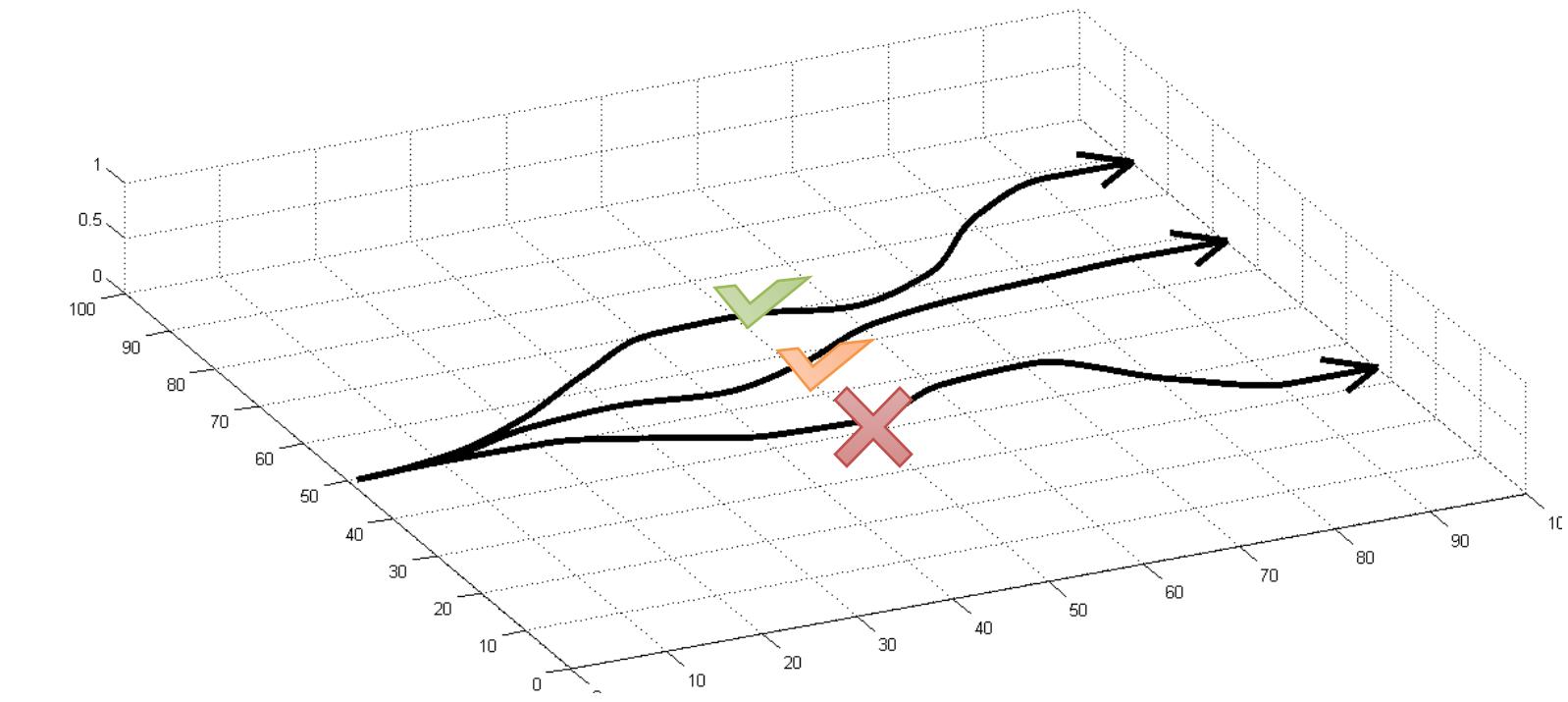
$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$J(\theta)$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$



sum over samples from  $\pi_{\theta}$



# Can we get the gradient of the RL objective?

Let's start in terms of trajectories.

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \underbrace{\sum_t r(\mathbf{s}_t, \mathbf{a}_t)}_{J(\theta)} \right]$$

a convenient identity

$$\underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} = p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} = \underline{\nabla_{\theta} p_{\theta}(\tau)}$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$
$$\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)$$

$$\nabla_{\theta} J(\theta) = \int \underline{\nabla_{\theta} p_{\theta}(\tau)} r(\tau) d\tau = \int \underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} r(\tau) d\tau = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

# Can we get the gradient of the RL objective?

From trajectories to final form.

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \underbrace{\log p_{\theta}(\tau)}_{\text{green}} r(\tau)]$$

$$\nabla_{\theta} \left[ \cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \cancel{\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \right]$$

$$p_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\log p_{\theta}(\tau) = \cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

# Estimating the gradient

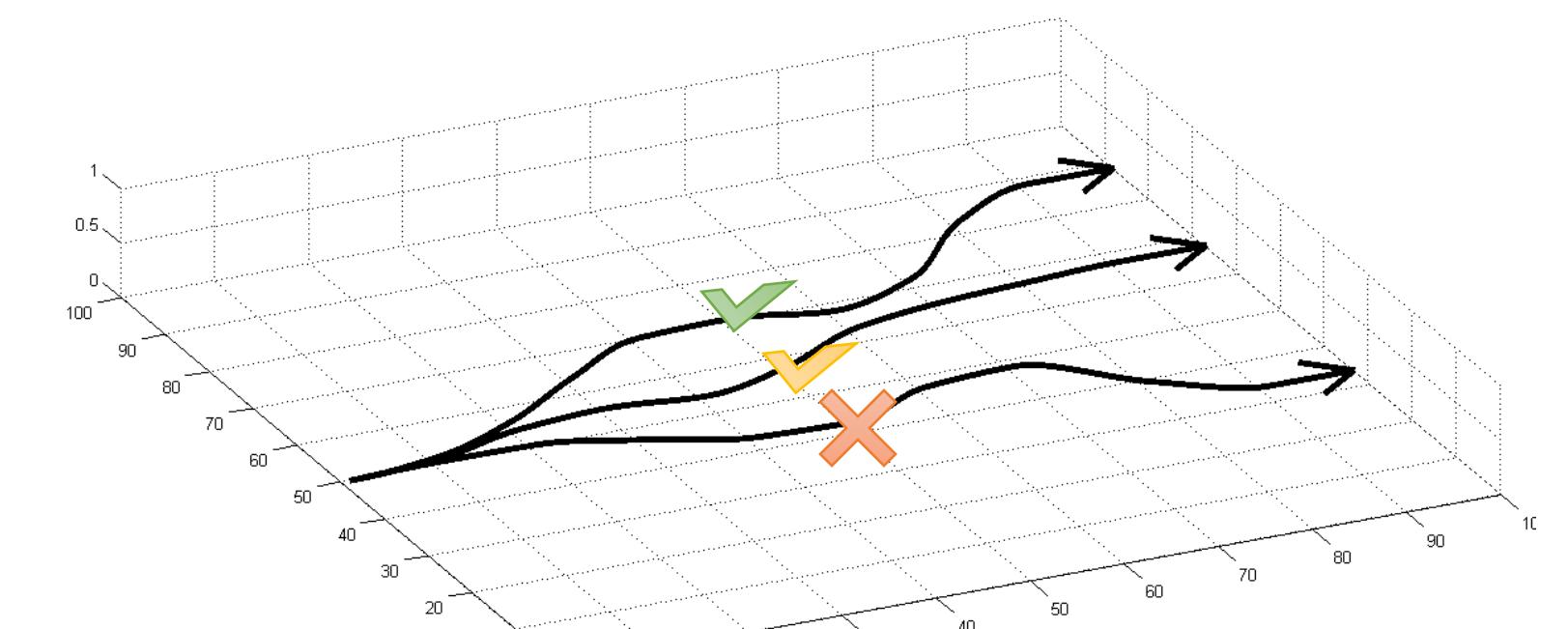
$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

Full algorithm:

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$
2.  $\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



Run policy to collect  
batch of data

Improve policy using  
batch of data

“REINFORCE algorithm”, vanilla policy gradient

# What does the gradient mean?

$$\nabla_{\theta} J(\theta) \approx \left( \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \left( \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right)$$

Imitation gradient, but weighted by reward

**Recall:** imitation learning

$$\min_{\theta} -\mathbb{E}_{(\mathbf{s}, \mathbf{a}) \sim \mathcal{D}} [\log \pi_{\theta}(\mathbf{a}|\mathbf{s})]$$

$$\nabla_{\theta} J_{\text{BC}}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$

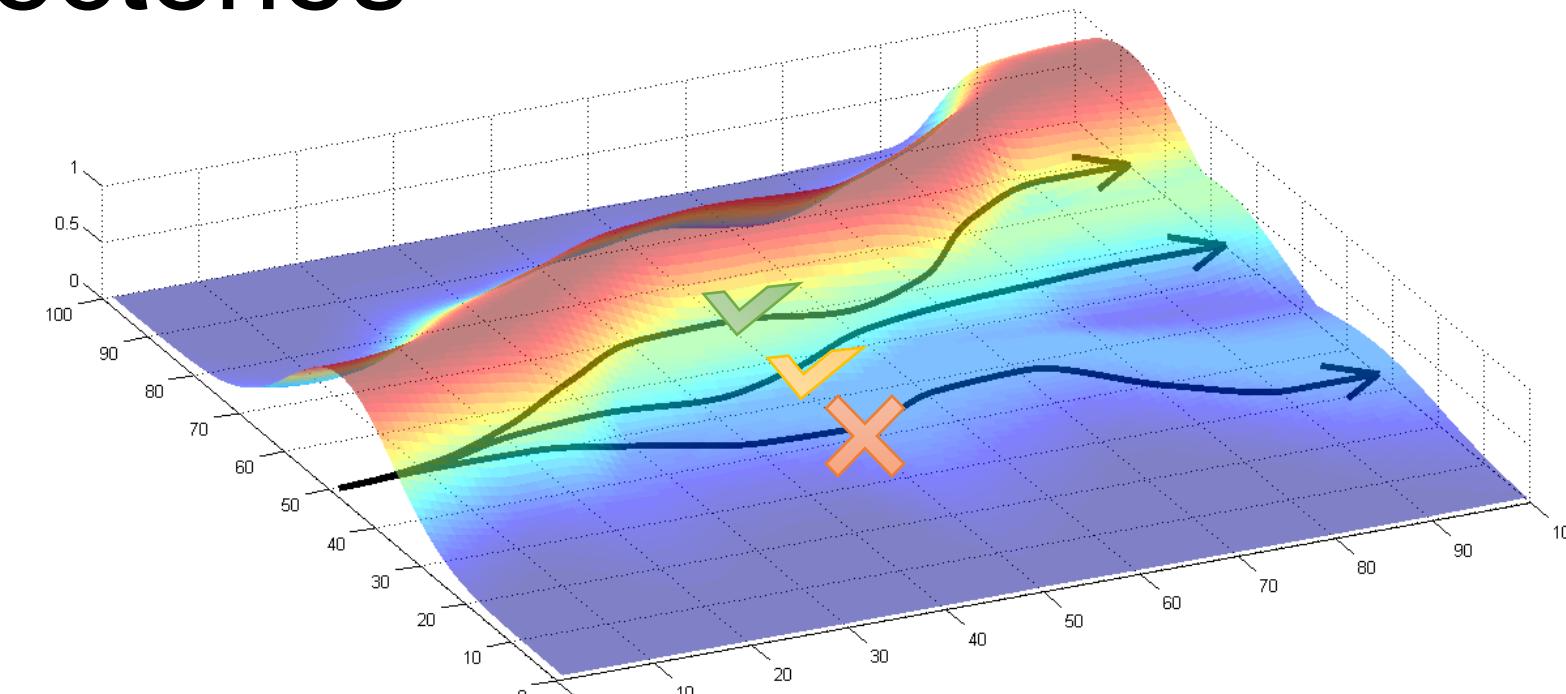
**Intuition:**

Increase likelihood of actions you took in high reward trajectories.

Decrease likelihood of actions you took in negative reward trajectories

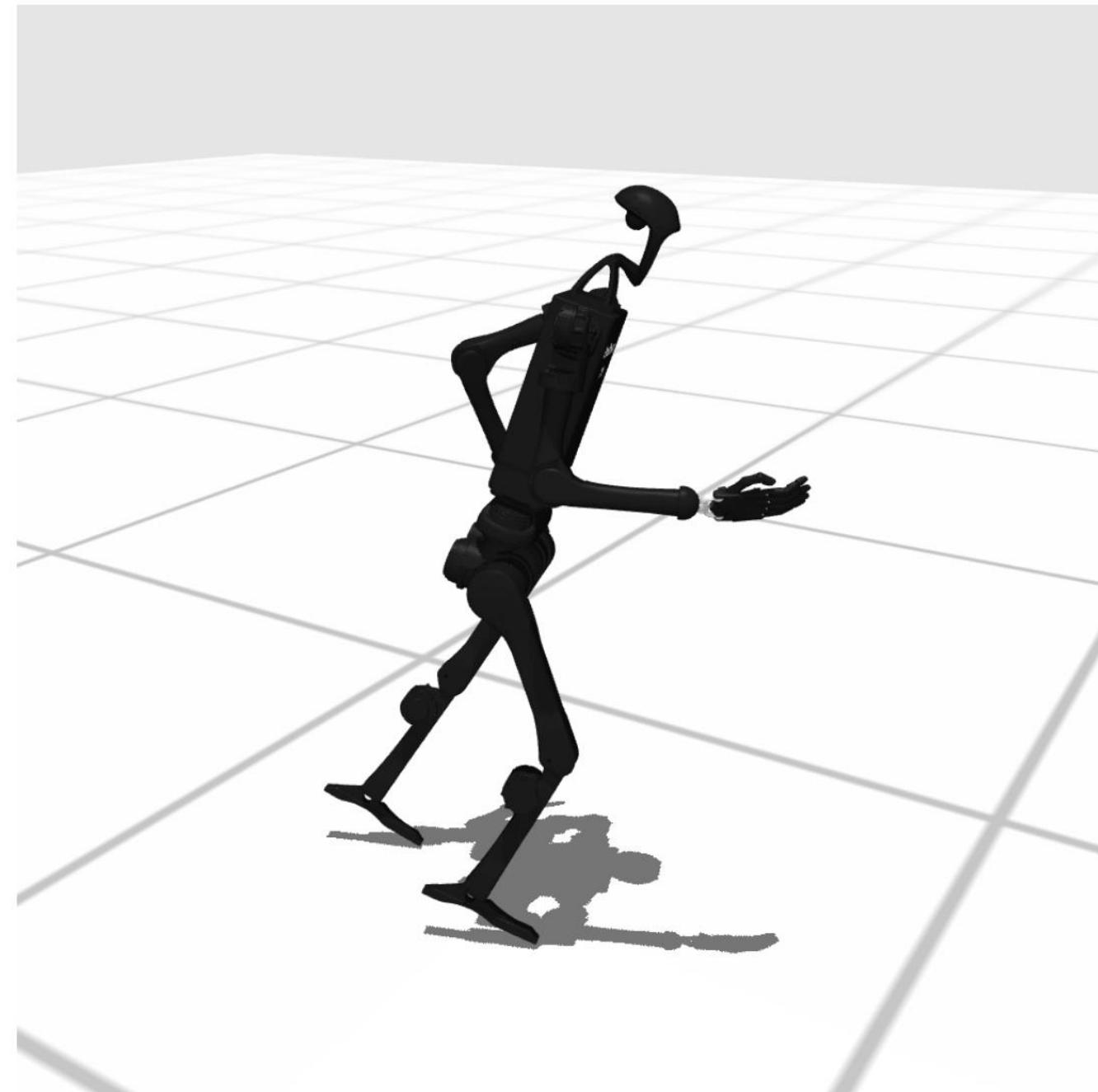
i.e. do more of the good stuff, less of the bad stuff.

formalization of “trial-and-error”



# What does the gradient do?

Example: learning humanoid walking in simulation



reward:  $r(\mathbf{s}, \mathbf{a}) = \text{forward velocity of robot}$   
*(can be negative if robot goes backwards)*

1. sample  $\{\tau^i\}$  from  $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$

$\tau^1$ : falls backwards ✗

$\tau^2$ : falls forwards ✓

$\tau^3$ : manages to stand still ✅

$\tau^4$ : one small step forward then falls backwards ✗

$\tau^5$ : one large step backwards then small step forwards ✗

2.  $\nabla_\theta J(\theta) \approx \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$

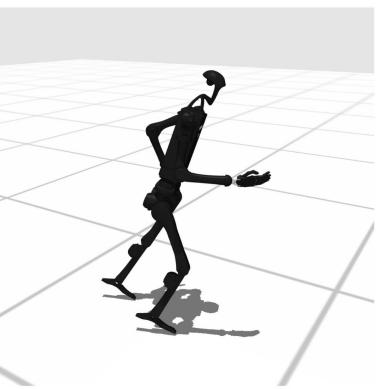
-> will encourage policy to fall forward, and not take step forward 😳

Policy gradient is **noisy / high-variance**

# Improving the gradient

Using causality

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$



$\tau^5$ : one large step backwards  
then small step forwards X

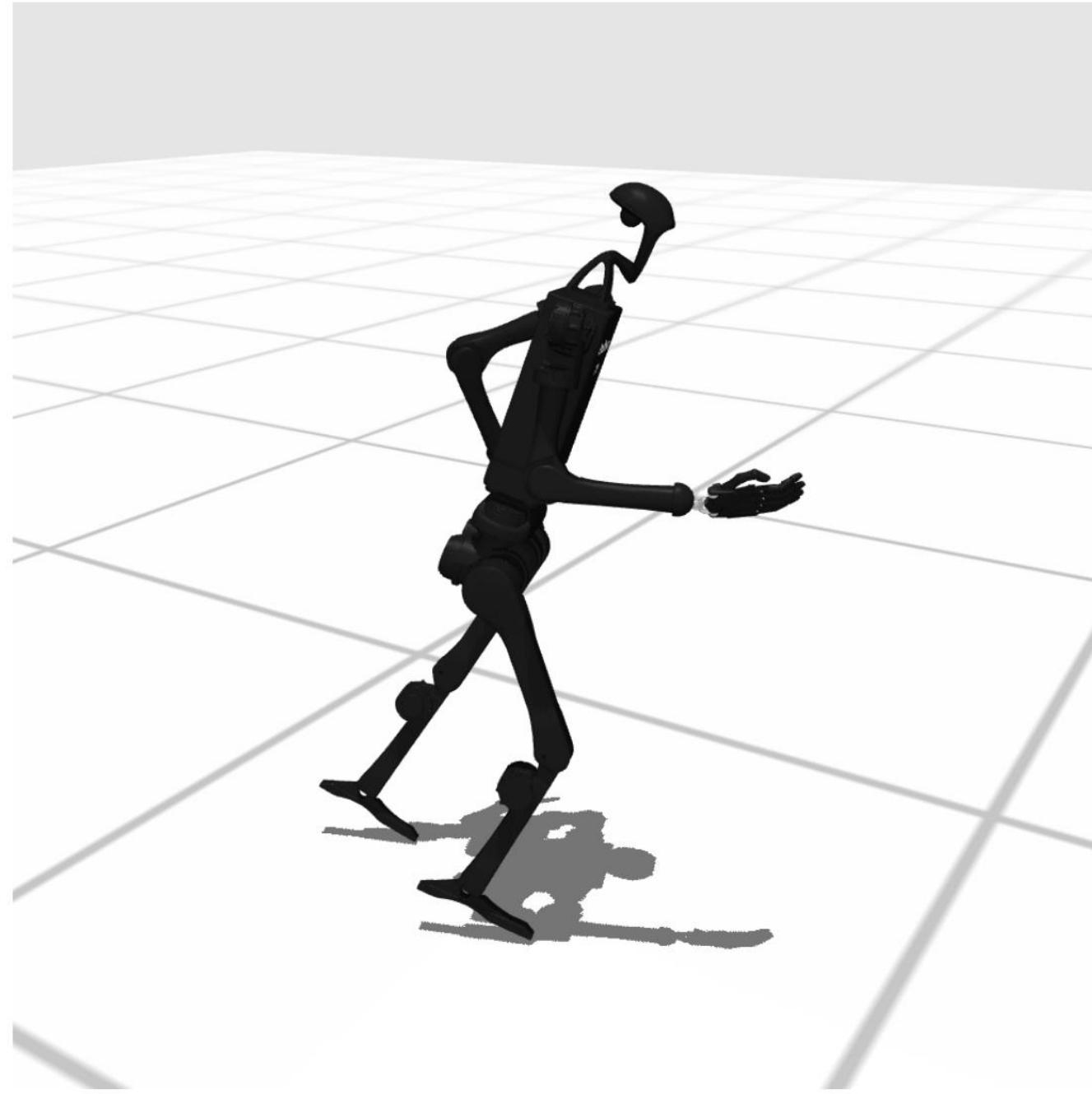
Policy behavior at time  $t$  does not affect rewards at time  $t' < t$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$

sum of future rewards

# What does the gradient do?

Example: learning humanoid walking in simulation



reward:  $r(\mathbf{s}, \mathbf{a}) = \text{forward velocity of robot}$   
*(can be negative if robot goes backwards)*

1. sample  $\{\tau^i\}$  from  $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$ 
  - $\tau^1$ : falls forwards ✓
  - $\tau^2$ : slowly stumbles forwards ✓
  - $\tau^3$ : steadily walks forwards ✓
  - $\tau^4$ : runs forwards ✓

$$2. \nabla_\theta J(\theta) \approx \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$$

-> encourages policy to fall, stumble forward some of the time 😱

Policy gradient is **noisy / high-variance**  
sensitive to reward scale

# Improving the gradient

Introducing baselines

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$



$\tau^1$ : falls forwards ✓

$\tau^4$ : runs forwards ✓

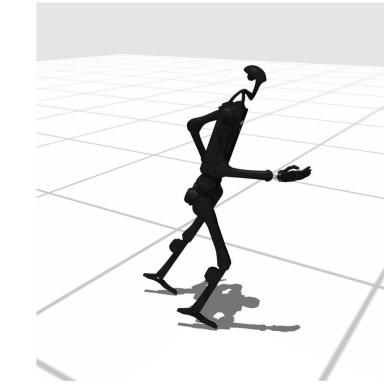
# Improving the gradient

Introducing baselines

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]$$

$$b = \frac{1}{N} \sum_{i=1}^N r(\tau)$$

subtracting a **constant**  
“baseline”



a convenient identity

$$p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} p_{\theta}(\tau)$$

$\tau^1$ : falls forwards ✓

$\tau^4$ : runs forwards ✓

If we subtract average reward, we get negative gradients for below-average behavior. 🎉

But, can we even do that? 🤔

$$E[\nabla_{\theta} \log p_{\theta}(\tau) b] = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) b d\tau = \int \nabla_{\theta} p_{\theta}(\tau) b d\tau = b \nabla_{\theta} \int p_{\theta}(\tau) d\tau = b \nabla_{\theta} 1 = 0$$

Subtracting a constant baseline is **unbiased** (change to the gradient is 0 in expectation)  
and can reduce **variance** of the gradient

Average reward is a pretty good baseline.

# What does the gradient do?

Example: learning to fold a jacket

<https://www.youtube.com/watch?v=e2aMSzMYUh0>

reward:  $r(\mathbf{s}, \mathbf{a}) = \begin{cases} 1 & \text{neatly folded} \\ 0.5 & \text{folded with some wrinkles} \\ 0 & \text{not folded} \end{cases}$

1. sample  $\{\tau^i\}$  from  $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$   
 $\tau^1$ : doesn't touch the jacket     $\tau^3$ : flattens the jacket but does not fold it  
 $\tau^2$ : folds only the sleeves     $\tau^4$ : folds the jacket
2.  $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) - b)$   
-> gradient is constant for all but one trajectory 

Policy gradient is still **noisy / high-variance**  
Best with dense rewards, large batches.

# How to implement policy gradients?

Our gradient

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left( \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) - b \right)$$



Computing these individually is inefficient.  
( $N*T$  backwards passes)

Can we use automatic differentiation on full objective?

Implement “surrogate objective” whose gradient is the same as  $\nabla J$

$$\tilde{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

*weighted maximum likelihood*

Cross-entropy for discrete action policy, squared error for Gaussian policy

# Summary so far

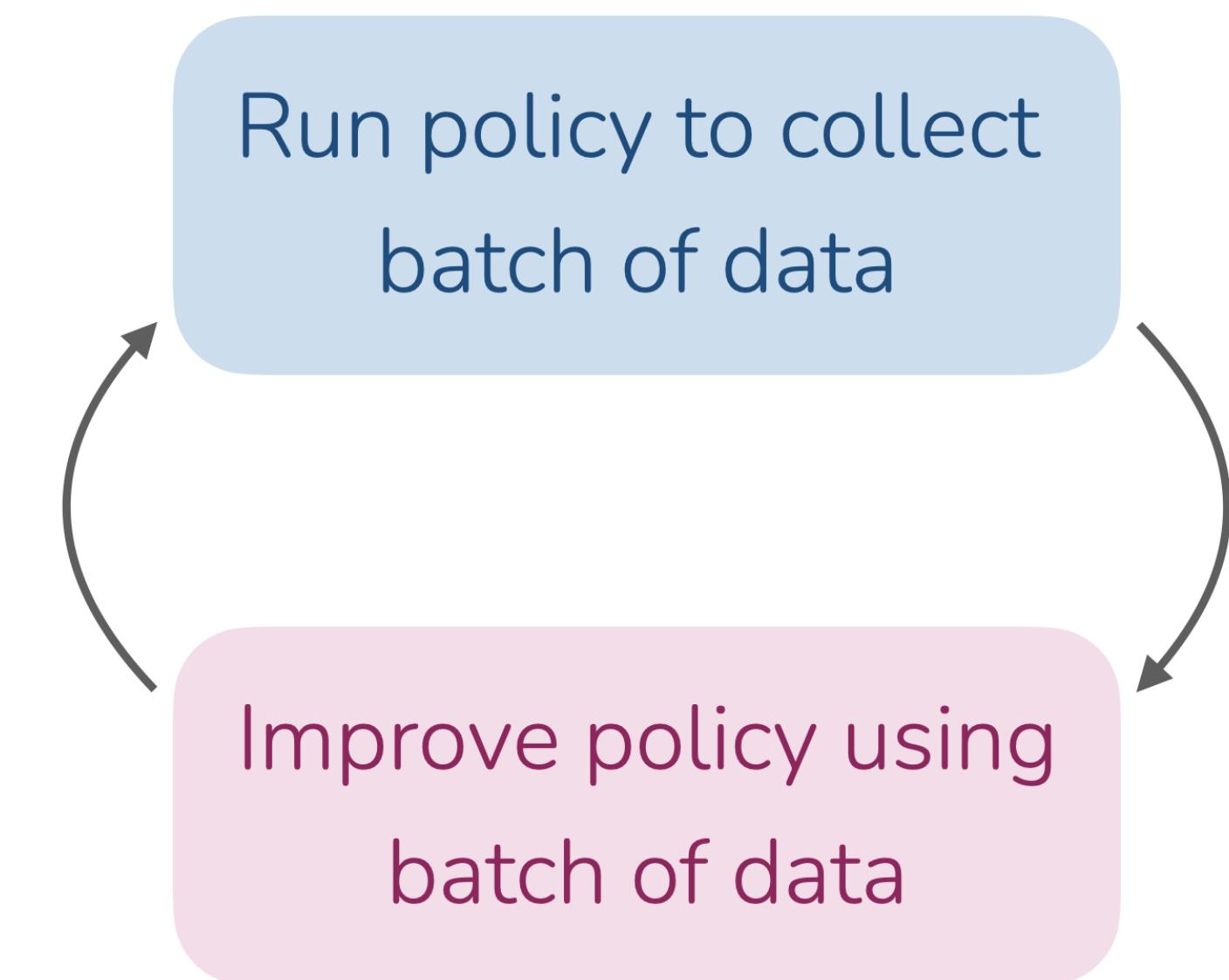
Estimating gradient of RL objective

- log gradient trick
- weigh policy likelihood by future rewards
- subtract baseline (e.g. average reward)
- even with tricks, gradient is **noisy**

First reinforcement learning algorithm

- collect batch of data, improve policy by applying gradient
- formalizes trial-and-error learning

Key intuition: do more high reward stuff, less low reward stuff



# What else is troublesome about policy gradients?

Latest version of our gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) - b \right) \right]$$

assumes samples from current policy  $\pi_{\theta}$

Full algorithm:

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$
2. compute  $\nabla_{\theta} J(\theta)$  using  $\{\tau^i\}$
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$       <- we change  $\theta$  right here

Need to recollect data every gradient step! 

Vanilla policy gradient is **on-policy**.

**attribute of online RL algorithms**

Definitions.

*on-policy*: update uses only data from current policy

*off-policy*: update can reuse data from other, past policies

# Off-policy version of policy gradient?

Importance sampling

$$J(\theta) = E_{\tau \sim p_\theta(\tau)}[r(\tau)]$$

What if we want to use samples from  $\bar{p}(\tau)$ ?  
(e.g. previous policy)

$$J(\theta) = E_{\tau \sim \bar{p}(\tau)} \left[ \frac{p_\theta(\tau)}{\bar{p}(\tau)} r(\tau) \right]$$

$$\frac{p_\theta(\tau)}{\bar{p}(\tau)} = \frac{\cancel{p(s_1)} \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \cancel{p(s_{t+1} | s_t, \mathbf{a}_t)}}{\cancel{p(s_1)} \prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t) \cancel{p(s_{t+1} | s_t, \mathbf{a}_t)}} = \frac{\prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t)}$$

## Importance sampling

Using proposal distribution  $q$

$$\begin{aligned} E_{x \sim p(x)}[f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{q(x)}{q(x)} p(x) f(x) dx \\ &= \int q(x) \frac{p(x)}{q(x)} f(x) dx \\ &= E_{x \sim q(x)} \left[ \frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

Note: Important for  $q$  to have non-zero support for high probability  $p(x)$

# Off-policy version of policy gradient?

Importance sampling

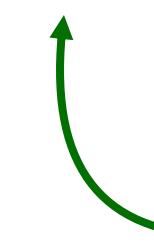
$$\frac{p_\theta(\tau)}{\bar{p}(\tau)} = \frac{\prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t)}$$

Say we want to update our latest policy  $\pi_{\theta'}$  but we want to use samples from  $\pi_\theta$

$$\nabla_{\theta'} J(\theta') = \mathbb{E}_{\tau \sim p_{\theta'}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) - b \right) \right]$$

$$\nabla_{\theta'} J(\theta') = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \frac{p_{\theta'}(\tau)}{p_\theta(\tau)} \left( \sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) - b \right) \right]$$

$$\nabla_{\theta'} J(\theta') = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \prod_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \left( \sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) - b \right) \right]$$



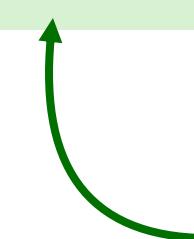
This can become very small or very large, for larger  $T$

# Off-policy version of policy gradient?

Importance sampling

Say we want to update our latest policy  $\pi_{\theta'}$  but we want to use samples from  $\pi_{\theta}$

$$\nabla_{\theta'} J(\theta') = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \prod_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)} \left( \sum_{t=1}^T \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) - b \right) \right]$$



This can become very small or very large, for larger  $T$

What if we consider the expectation over *timesteps* instead of trajectories?

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

Much less likely to explode/vanish ...but, hard to measure  $\frac{\pi_{\theta'}(\mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t})}$

Common final form

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

often approximated as 1

# Off-policy policy gradient

Common final form

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

Full algorithm:

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$
2. compute  $\nabla_{\theta} J(\theta)$  using  $\{\tau^i\}$
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Can take **multiple gradient steps** on the same batch

Run policy to collect  
batch of data

Improve policy using  
batch of data

What if our policy changes **a lot** before sampling new data?

Data no longer reflects states that policy will visit. Gradient estimate less accurate.

# Off-policy policy gradient

Common final form

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

What if our policy changes **a lot** before sampling new data?

Data no longer reflects states that policy will visit. Gradient estimate less accurate.

Can we constrain the policy to not stray too far during gradient updates?

Common choice:  $\mathbb{E}_{\mathbf{s} \sim \pi_{\theta}} [D_{KL}(\pi_{\theta}(\cdot | \mathbf{s}) \| \pi_{\theta'}(\cdot | \mathbf{s}))] \leq \delta$

# Review

## Online RL via policy gradients

- On-policy algorithm, differentiating the RL objective
- Baselines, causality for reducing gradient variance
- collect batch of data, improve policy by applying gradient

## Derived off-policy policy gradient

- importance sampling
- KL constraint on policy
- collect batch of data, apply multiple gradient updates

## Intuition

- Do more high reward stuff, less low reward stuff
- Gradient still very noisy, best with large batch sizes and dense rewards



# The plan for today

the basis for:

Policy gradients: our first online RL algorithm

<https://www.youtube.com/watch?v=xG7WkPU8tgs>

## 1. On-policy policy gradient

a. Derivation and intuition of policy gradients

b. Full algorithm

c. How to make it better - causality and baselines

Source: Unitree

## 2. Off-policy policy gradients

a. Importance sampling

b. KL constraints

Key learning goals:

- Key intuition behind policy gradients
- How to implement, when to use policy gradients

2. Creating a Boolean mask for elements that differ from that first element.  
3. Using a vectorized argmax (via np.take\_along\_axis) to pick the first different value (which will equal the first element if none exist).  
4. Verifying that every element in the slice is either equal to the first element or that candidate second unique value.

Try this code:

```
python ⌂ Copy
import numpy as np

def fast_two_unique(values):
    # values shape: (A, B, C)
    # Get the first element from each slice.
    first = values[..., 0] # shape: (A, B)

    # Create a mask that is True where elements differ from the first element.
    mask = (values != first)[..., None]

    # Find, for each slice, the index of the first differing element.
    # np.argmax returns 0 if no True exists; that's fine because then candidate2 equals first.
    idx = np.argmax(mask, axis=2)

    # Extract the candidate second unique value.
    candidate2 = np.take_along_axis(values, idx[...], None).squeeze(axis=2)

    # Check that every element is either equal to 'first' or 'candidate2'.
    return np.all((values == first)[..., None] | (values == candidate2)[..., None])

# Example usage:
# only_two_unique_values is a boolean array indicating, for each (i, j) slice, if the
only_two_unique_values = fast_two_unique(values)
```

Explanation

- Extracting the First Element:  
We set `first = values[..., 0]` so each (i, j) slice's first value becomes our reference.
- Finding a Candidate for the Second Unique Value:  
The mask `(values != first)[..., None]` identifies elements that differ from the first value. Using `np.argmax` finds the first index along axis 2 where the element is different. If all elements are the same, `np.argmax` returns 0 and `candidate2` will equal `first`.
- Verifying Uniqueness:  
Finally, we check that every element in each slice is either equal to `first` or `candidate2`. If a third unique value exists, the condition fails for that slice.

This vectorized solution should be considerably faster than using `apply_along_axis` with a Python loop.

That works great! Can you generalize the function to work with any axis, not just axis 2?

Reasoned about generalized vectorized solution for 5 seconds.  
Below is a generalized version that lets you choose the axis along which to check for at most two unique values. It works by moving the target axis to the end, applying the logic, and then reducing that axis.

python ⌂ Copy

# Next time

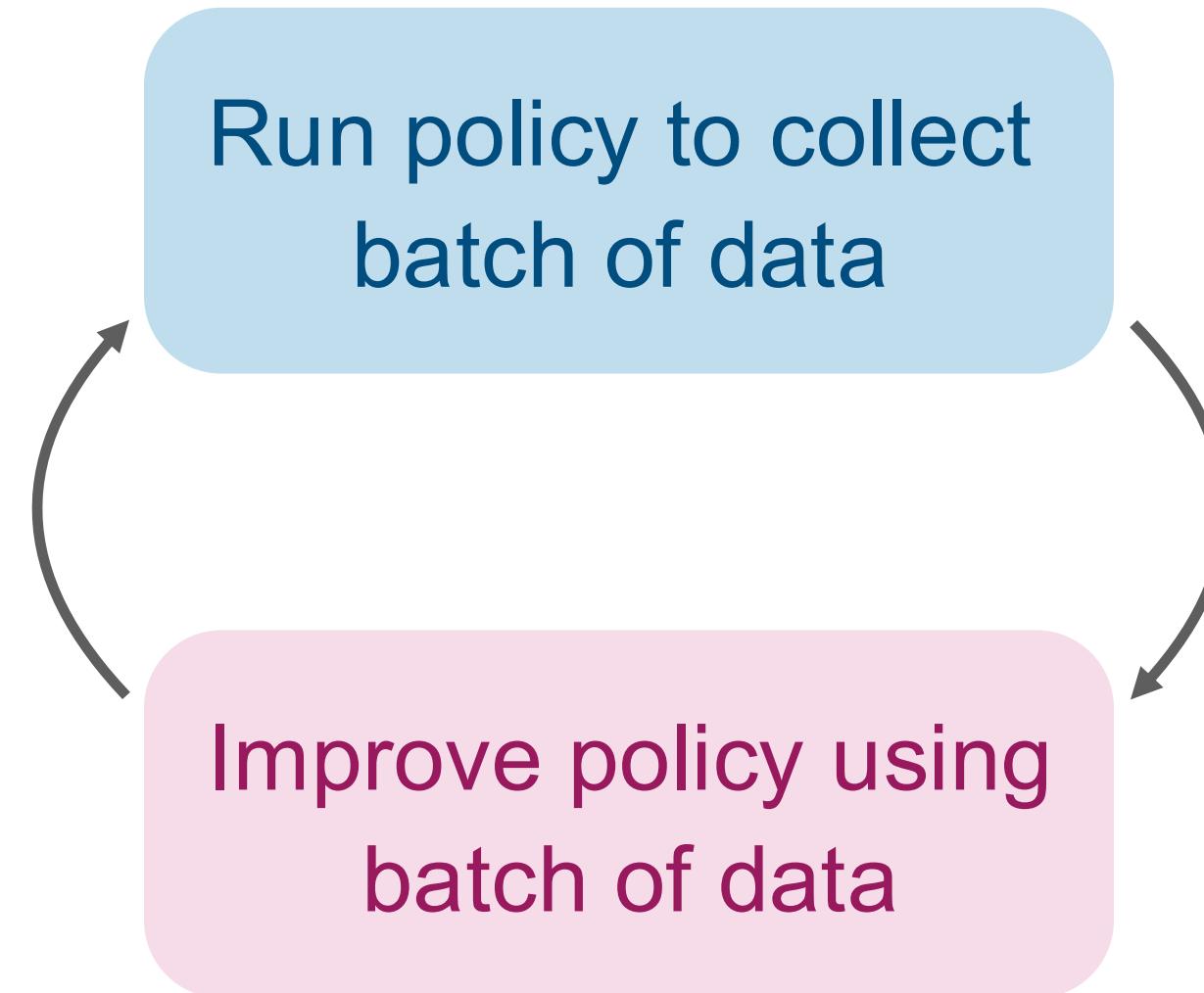
## *Actor critic methods*

- > build closely on policy gradients!
- > basis for popular algorithms like PPO

# Actor Critic Methods

CS 224R

# Recap of Last Time



Online reinforcement learning with policy gradients

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

samples from policy      policy log likelihood      reward to go      baseline

Do more of the above average stuff, less of the below average stuff.

# Recap of Last Time

Vanilla policy gradients is fully **on-policy**.

→ only one gradient step per batch of data

**attribute of online RL algorithms**

**Definitions.**

*on-policy*: update uses only data from current policy

*off-policy*: update can reuse data from other, past policies

**Importance weights** for *off-policy* policy gradient:

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \right) \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

Note: Not a silver bullet. Only suitable when policies are very similar.

# The plan for today

## Actor critic methods

1. Improving policy gradients
2. How to estimate the value of a policy
  - a. Sample & directly supervise
  - b. Use your own estimate
  - c. Something in-between?
3. Off-policy actor-critic
  - a. Importance weights & constraining step size
  - b. Full off-policy version with replay buffers

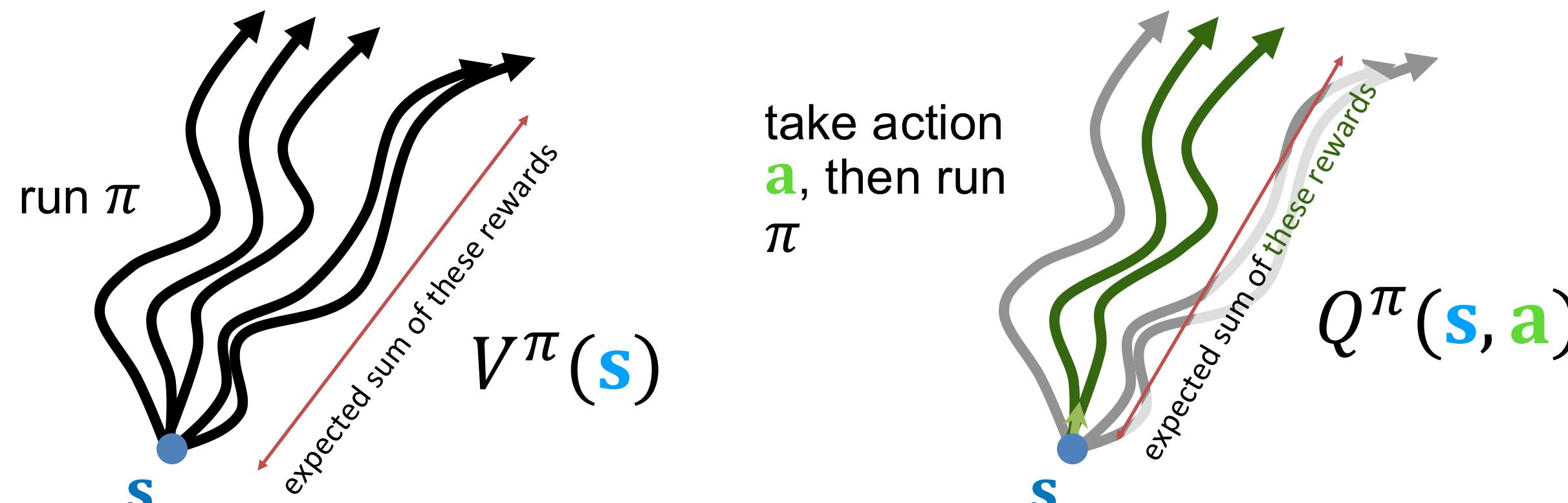
## Key learning goals:

- How to estimate how good a state and action is for a policy
- How to use those estimates to form a more efficient RL algorithm

# Revisiting Some Useful Objects

*value function*  $V^\pi(\mathbf{s})$  - future expected rewards starting at  $\mathbf{s}$  and following  $\pi$

*Q-function*  $Q^\pi(\mathbf{s}, \mathbf{a})$  - future expected rewards starting at  $\mathbf{s}$ , taking  $\mathbf{a}$ , then following  $\pi$

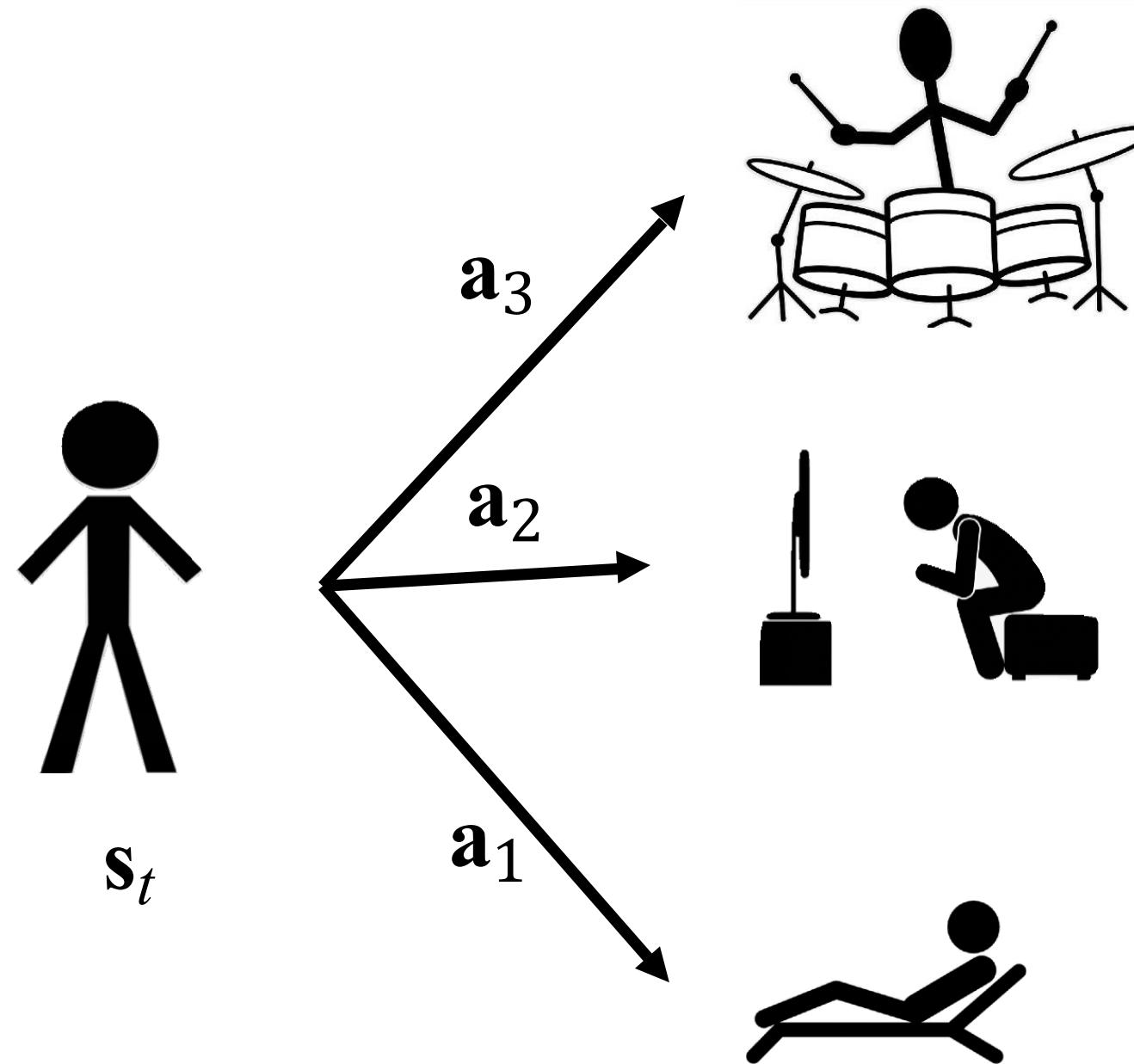


$$\text{Useful relation: } V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\cdot | \mathbf{s})}[Q^\pi(\mathbf{s}, \mathbf{a})]$$

*advantage*  $A^\pi(\mathbf{s}, \mathbf{a})$  - how much better it is to take  $\mathbf{a}$  than to follow policy  $\pi$  at state  $\mathbf{s}$

$$A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s})$$

# Let's look at an example



Reward = 1 if I can play it in  
a month, 0 otherwise



Current policy:  $\pi(a_1 | s) = 1 \forall s$

Question:

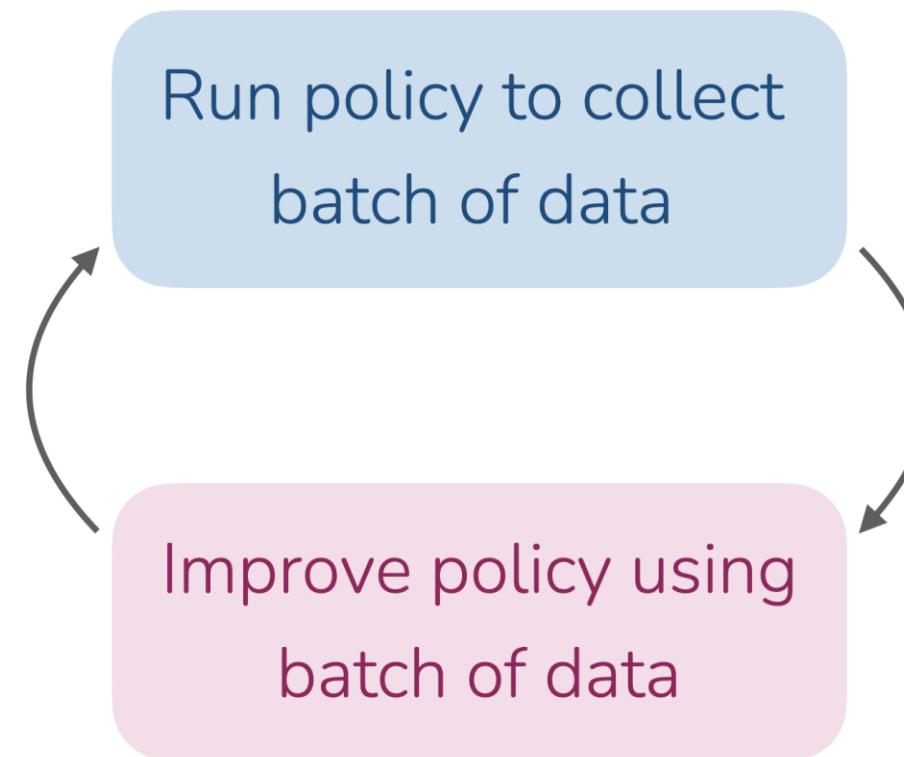
For each action, what are these?

Value function:  $V^\pi(s_t) = ?$

Q function:  $Q^\pi(s_t, a_t) = ?$

Advantage function:  $A^\pi(s_t, a_t) = ?$

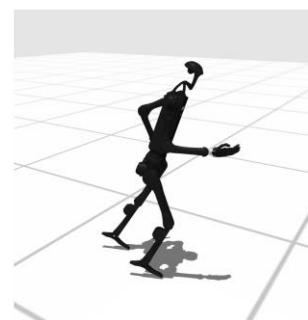
# What is dissatisfying about policy gradients?



$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

samples from policy      policy log likelihood      reward to go      baseline

<https://www.youtube.com/watch?v=e2aMSzMYUho>



$\tau^4$ : one small step forward then falls backwards

-> pushes down likelihood of step forward

$\tau^2$ : folds only the sleeves

$\tau^3$ : flattens the jacket but does not fold it

$\tau^4$ : folds the jacket

-> with sparse rewards, don't utilize  $\tau^2$  or  $\tau^3$

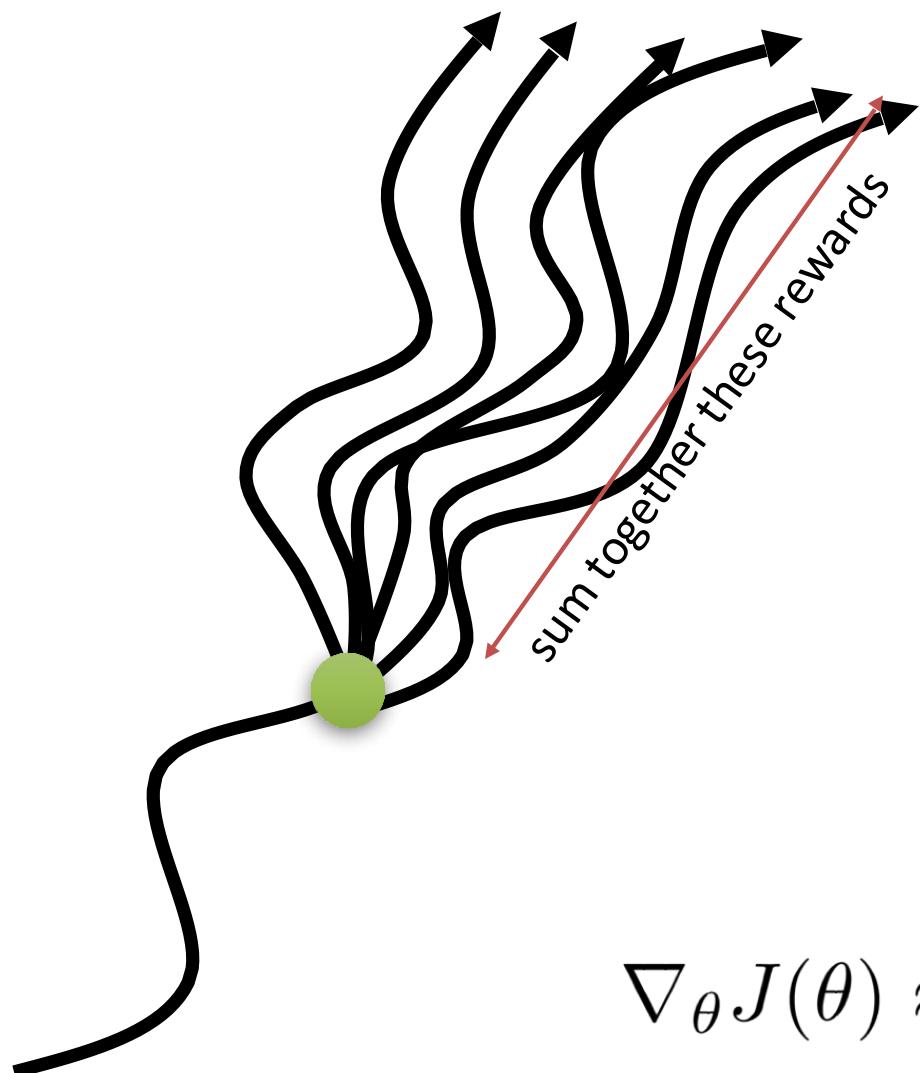
Policy gradients doesn't make efficient use of data! 💡 Can we learn what is good & bad?

# Improving policy gradients

Estimating reward to go

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$

“reward to go”



estimate of future rewards if we take  $\mathbf{a}_{i,t}$  in state  $\mathbf{s}_{i,t}$

$$\sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$$

Can we get a better estimate?

$$\sum_{t'=t}^T E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t] = Q(\mathbf{s}_t, \mathbf{a}_t)$$

true expected rewards to go

This would be way better!

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

# Improving policy gradients

What about baselines?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) (Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - \cancel{b}) \\ V^{\pi}(\mathbf{s}_t)$$

$$b = \text{average reward} = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$V(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_{\theta}(\cdot | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)]$$

Recall:  $A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi}(\mathbf{s}_t)$

With baseline:

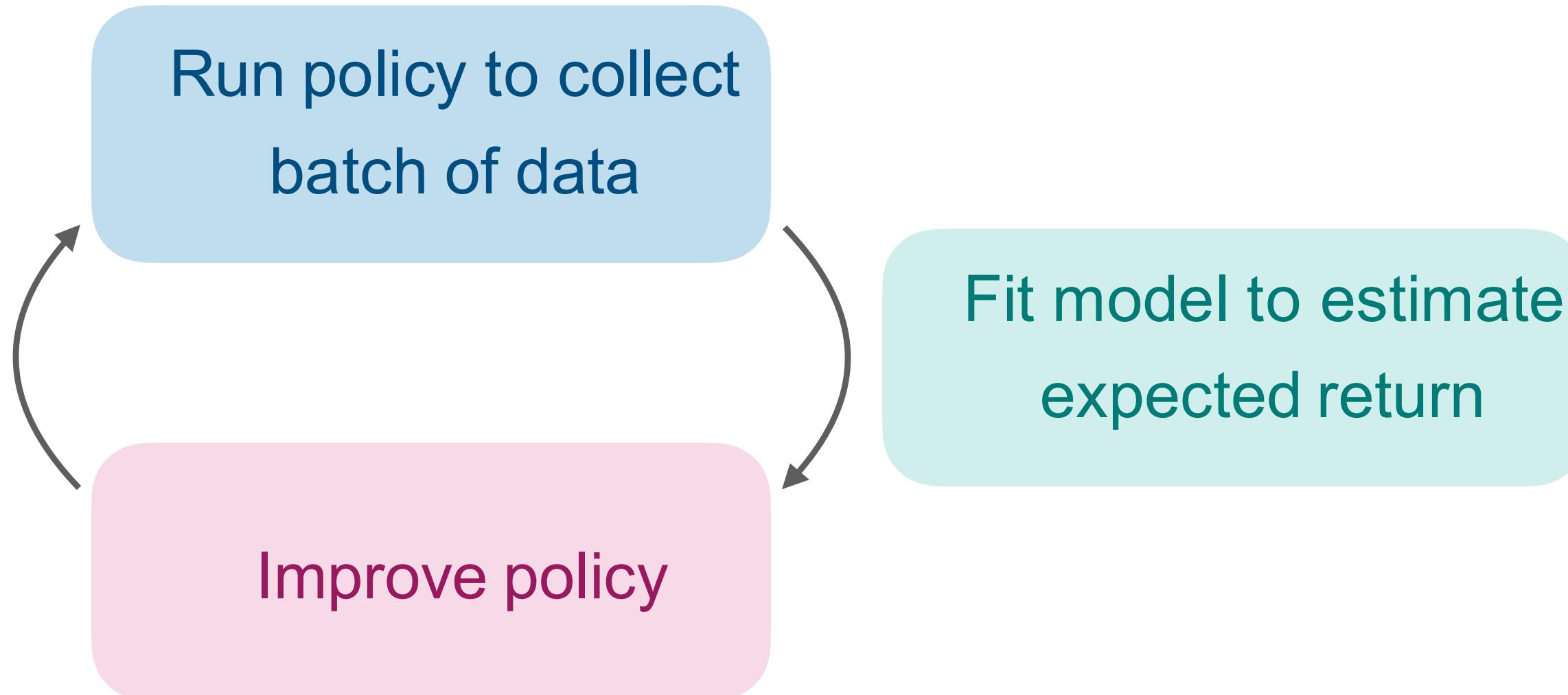
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

Better estimates of lead to less noisy gradients!

# Online RL Outline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

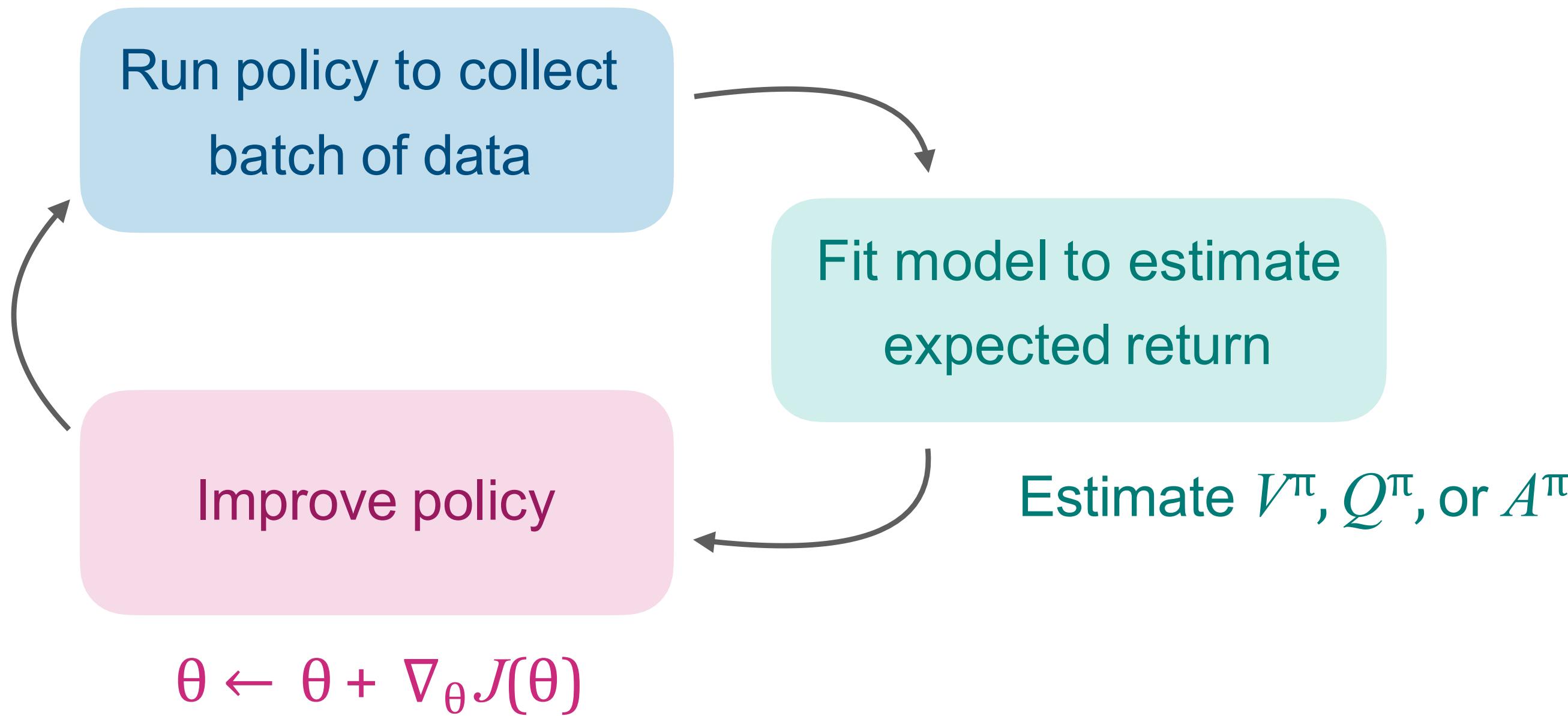
First: Initialize the policy (randomly, with imitation learning, with heuristics)



# Online RL Outline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

First: Initialize the policy (randomly, with imitation learning, with heuristics)



# The plan for today

## Actor critic methods

1. Improving policy gradients
2. **How to estimate the value of a policy**
  - a. Sample & directly supervise
  - b. Use your own estimate
  - c. Something in-between?
3. Off-policy actor-critic
  - a. Importance weights & constraining step size
  - b. Full off-policy version with replay buffers

## Key learning goals:

- How to estimate how good a state and action is for a policy
- How to use those estimates to form a better RL algorithm

# Estimating expected return

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

Should we fit  $V^\pi$ ,  $Q^\pi$ , or  $A^\pi$ ?

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t)$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

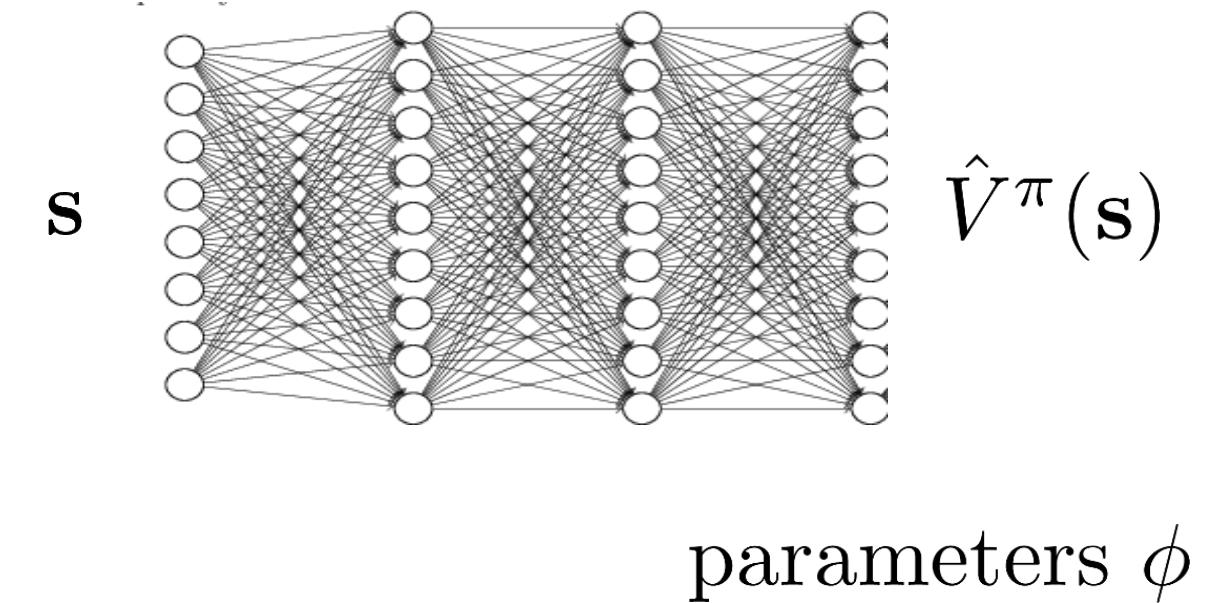
$$= r(\mathbf{s}_t, \mathbf{a}_t) + \sum_{t'=t+1}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$

$$= r(\mathbf{s}_t, \mathbf{a}_t) + E_{\mathbf{s}_{t+1} \sim p(\cdot | \mathbf{s}_t, \mathbf{a}_t)} [V^\pi(\mathbf{s}_{t+1})]$$

$$\approx r(\mathbf{s}_t, \mathbf{a}_t) + V^\pi(\mathbf{s}_{t+1}) \quad (\text{use the sampled } \mathbf{s}_{t+1})$$

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^\pi(\mathbf{s}_{t+1}) - V^\pi(\mathbf{s}_t)$$

Let's just fit  $V^\pi$ !

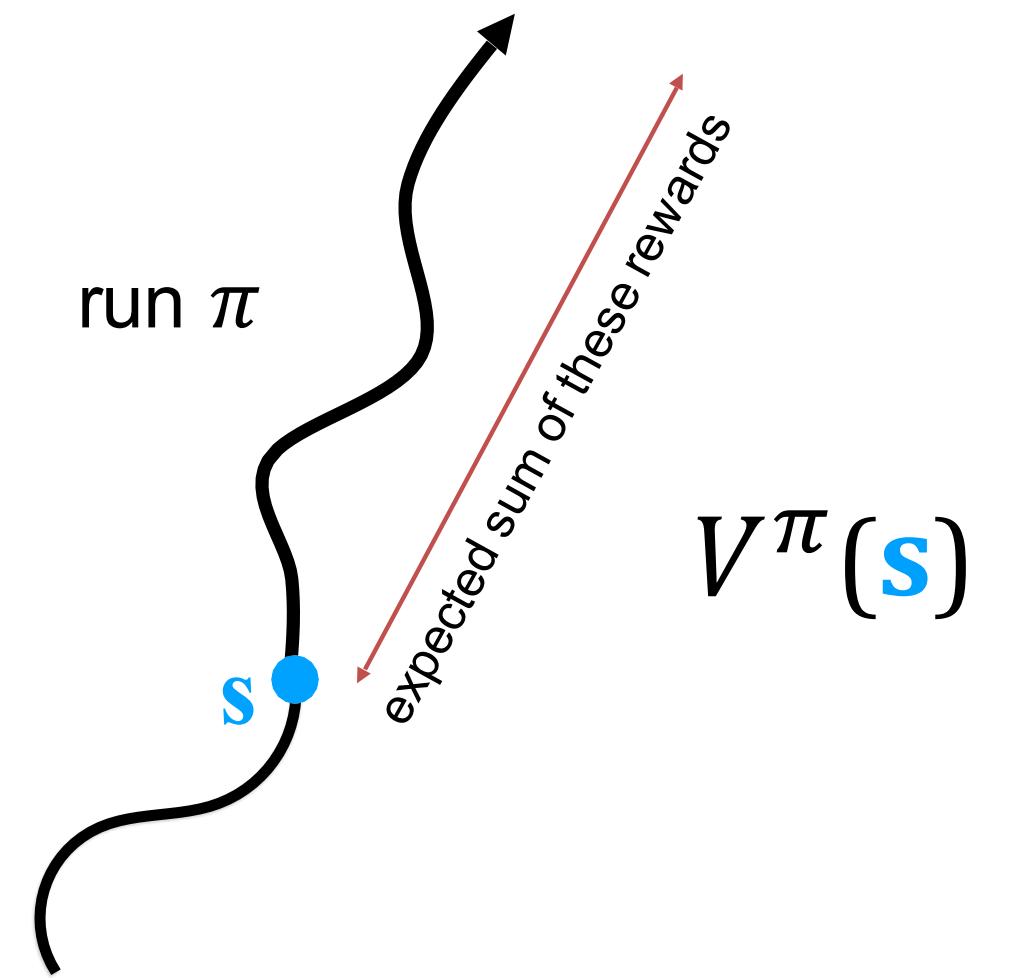


# Estimating $V^\pi$

Version 1: Monte Carlo estimation

$$V^\pi(s_t) \approx \sum_{t'=t}^T r(s_{t'}, a_{t'})$$

<- original single sample estimate

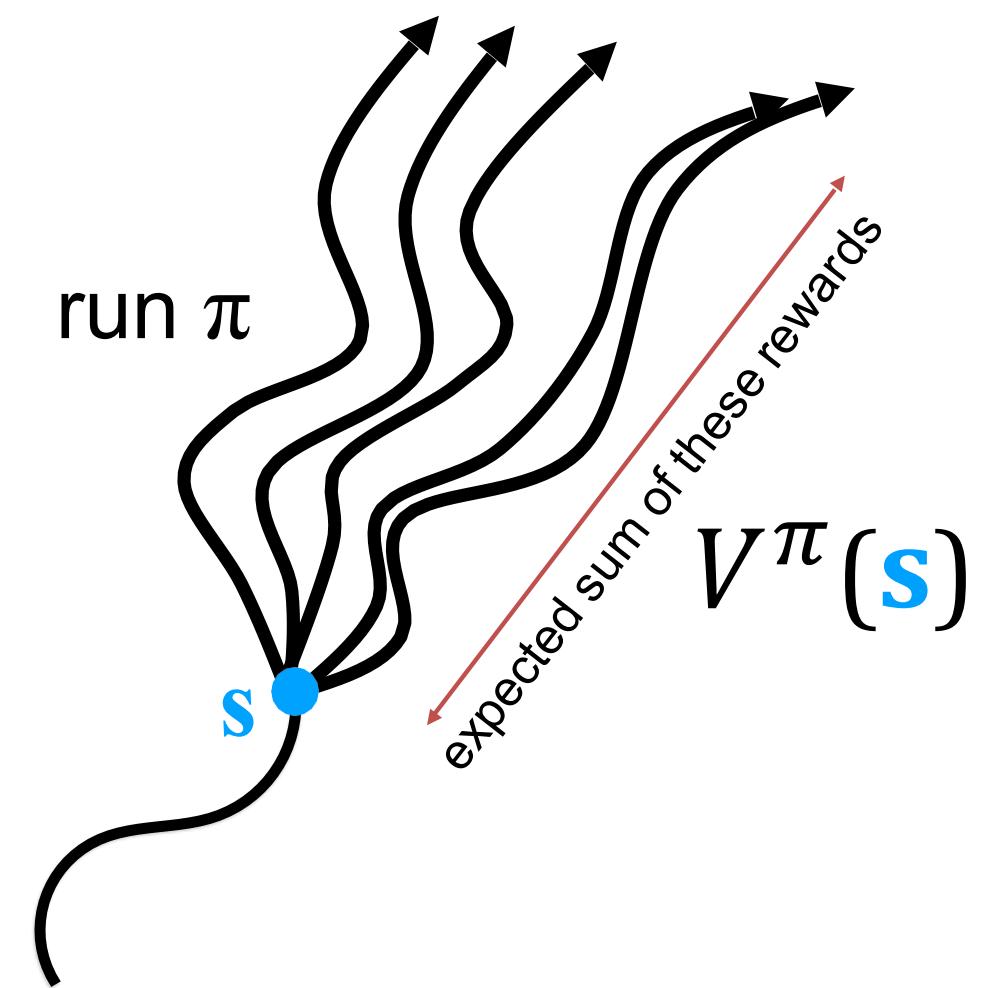


# Estimating $V^\pi$

Version 1: Monte Carlo estimation

$$V^\pi(s_t) \approx \sum_{t'=t}^T r(s_{t'}, a_{t'}) \quad <- \text{original single sample estimate}$$

$$V^\pi(s_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T r(s_{t'}, a_{t'}) \quad <- \text{multi-sample estimate} \\ (\text{but can't reset the world})$$



# Estimating $V^\pi$

Version 1: Monte Carlo estimation

$$V^\pi(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

<- original single sample estimate

$$V^\pi(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

<- multi-sample estimate

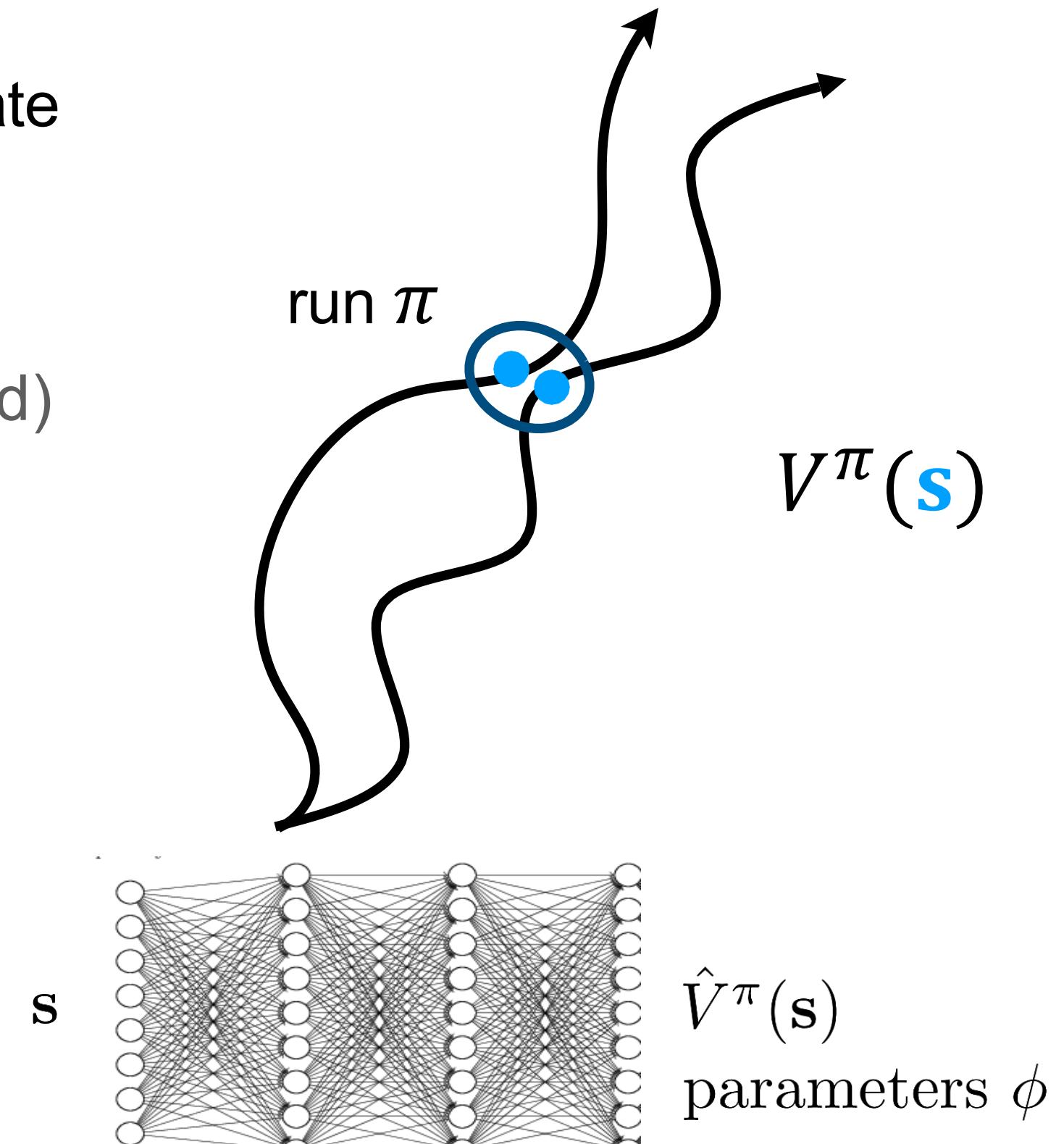
(but can't reset the world)

**Step 1:** Aggregate dataset of single sample estimates:

$$\left\{ \left( \mathbf{s}_{i,t}, \underbrace{\sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})}_{y_{i,t}} \right) \right\}$$

**Step 2:** Supervised learning to fit estimated value function

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$



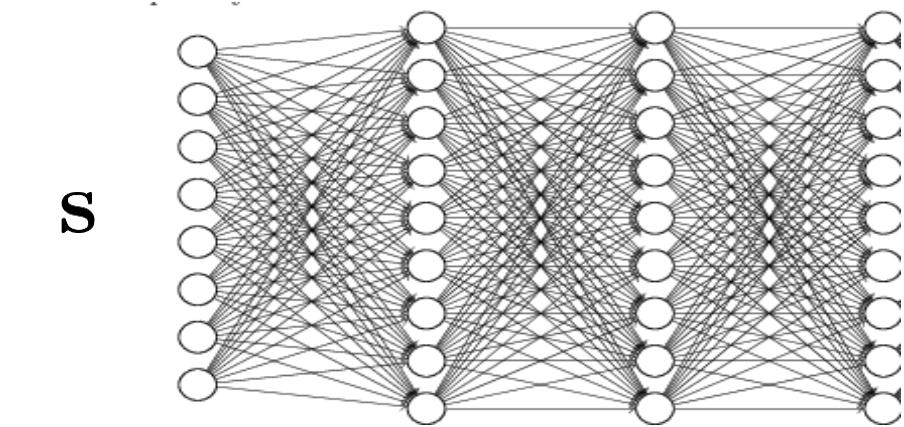
# Estimating $V^\pi$

Version 2: Bootstrapping

Monte Carlo target:  $y_{i,t} = \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$

ideal target:  $y_{i,t} = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t}] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \sum_{t'=t+1}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t+1}]$

$\approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + V^\pi(\mathbf{s}_{i,t+1}) \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$



directly use previous fitted value function!

**Step 1:** Aggregate dataset of “bootstrapped” estimates:

training data:  $\left\{ \left( \mathbf{s}_{i,t}, \underbrace{r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})}_{y_{i,t}} \right) \right\}$  <- update labels every gradient update!

**Step 2:** Supervised learning to fit estimated value function

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

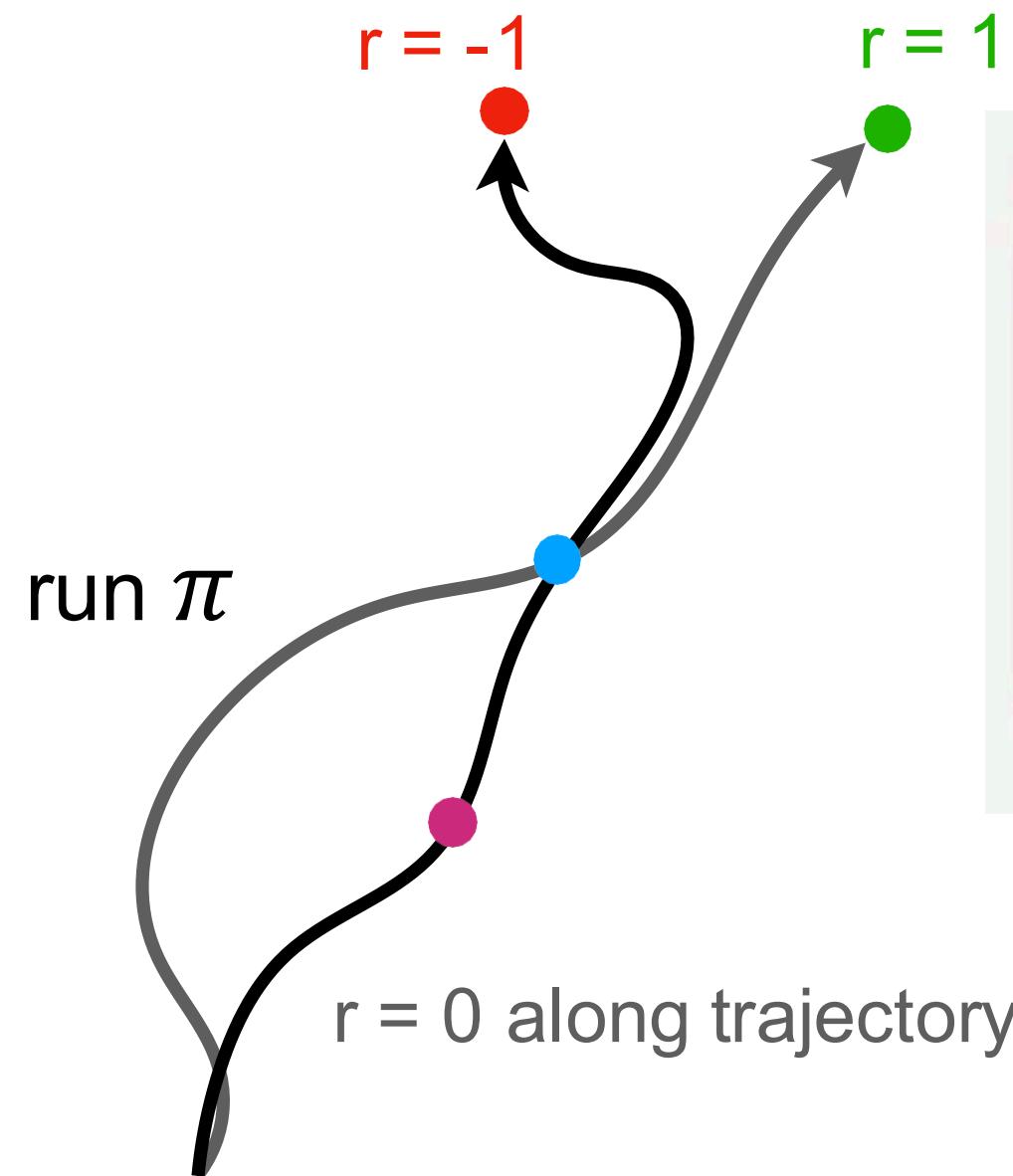
Also referred to as a form of  
**temporal difference learning**

# Estimating $V^\pi$ : Monte Carlo vs. Bootstrap

Monte Carlo       $y_{i,t} = \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$       supervise with roll-out's summed rewards

Bootstrapped       $y_{i,t} = r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$       supervise using reward and current value estimate

Let's look at an example:



Think-pair-share:

|   |   |
|---|---|
| $\hat{V}_{MC}^\pi(\textcolor{blue}{s}) \approx ?$ | $\hat{V}_{MC}^\pi(\textcolor{violet}{s}) \approx ?$ |
| $\hat{V}_{TD}^\pi(\textcolor{blue}{s}) \approx ?$ | $\hat{V}_{TD}^\pi(\textcolor{violet}{s}) \approx ?$ |

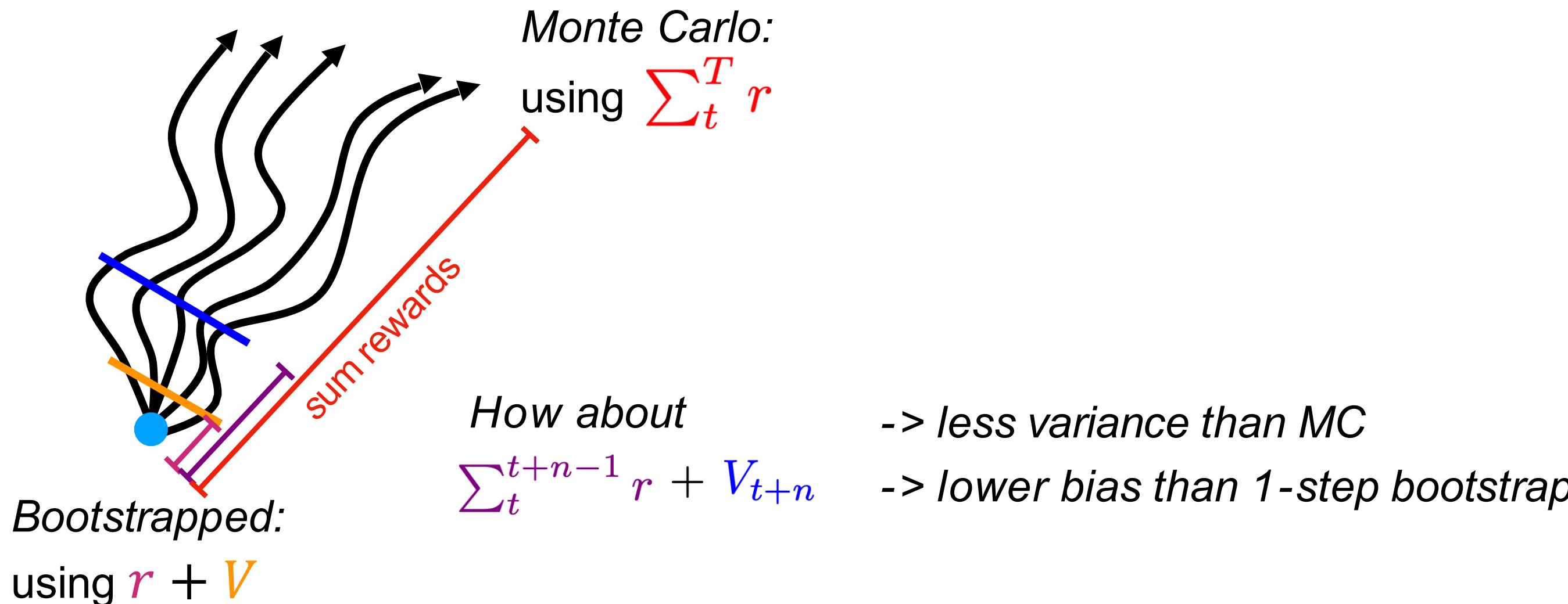
Is there middle ground?  
Can we balance bias and variance?

# Estimating $V^\pi$ : Monte Carlo vs. Bootstrap

Monte Carlo       $y_{i,t} = \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$       supervise with roll-out's summed rewards

Bootstrapped       $y_{i,t} = r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$       supervise using reward and current value estimate

N-step returns       $y_{i,t} = \sum_{t'=t}^{t+n-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+n})$



$n > 1, n < T$  often works the best in practice!

# Aside: discount factors

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$$

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

what if  $T$  (episode length) is  $\infty$ ?

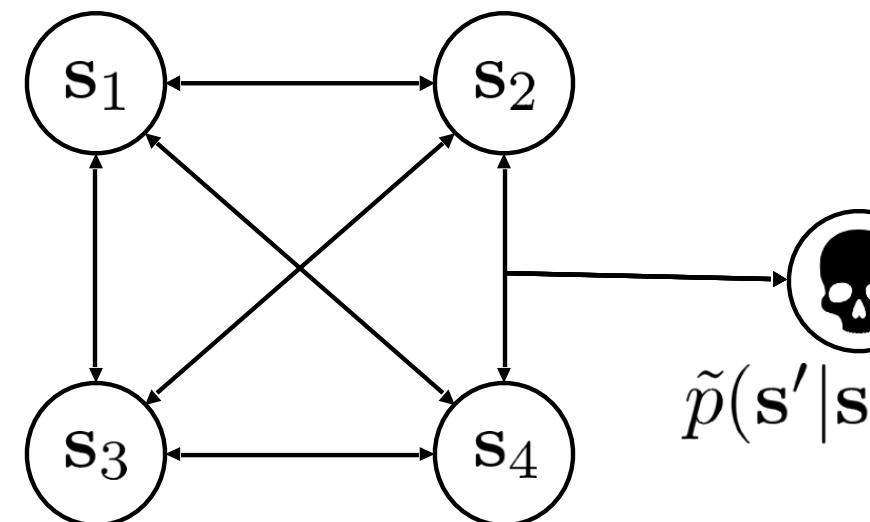
$\hat{V}_\phi^\pi$  can get infinitely large in many cases

simple trick: better to get rewards sooner than later

$\gamma$  changes the MDP:

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$$

↑  
discount factor  $\gamma \in [0, 1]$  (0.99 works well)



$$\tilde{p}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = (1 - \gamma)$$

$$\tilde{p}(\mathbf{s}'|\mathbf{s}, \mathbf{a}) = \gamma p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$$

# A Full Algorithm Walkthrough

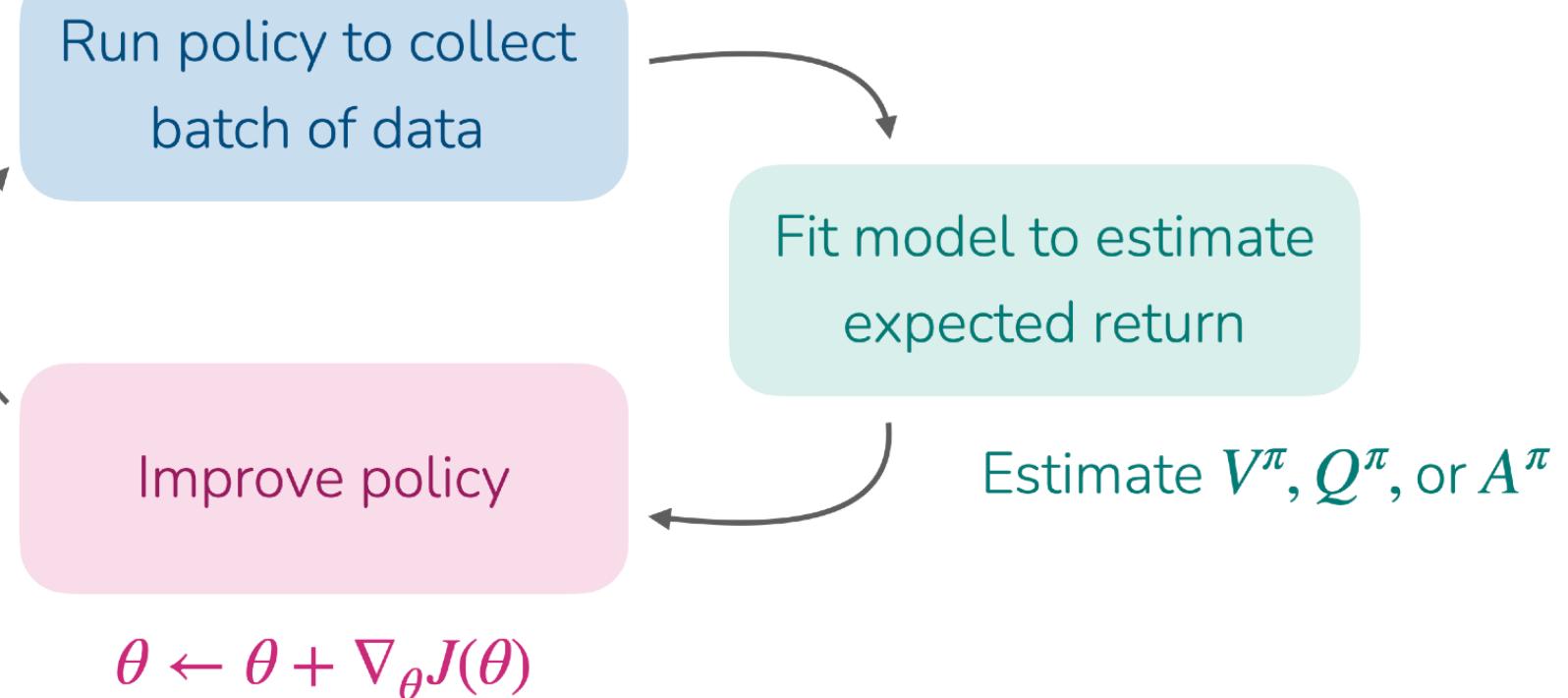
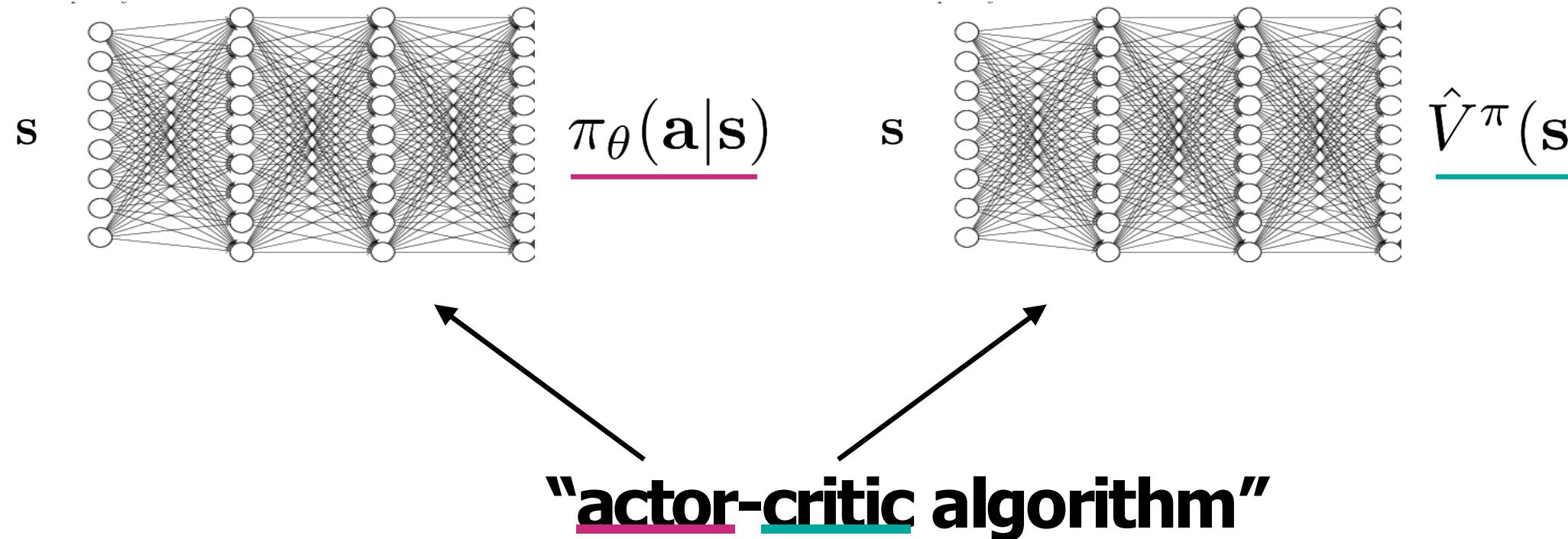
1. Sample batch of data  $\{(\mathbf{s}_{1,i}, \mathbf{a}_{1,i}, \dots, \mathbf{s}_{T,i}, \mathbf{a}_{T,i})\}$  from  $\pi_\theta$

2. Fit  $\hat{V}_\phi^{\pi_\theta}$  to summed rewards in data

3 Evaluate  $\hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) = r(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) + \gamma \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t+1,i}) - \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t,i}) \forall t, i$

4. Evaluate  $\nabla_\theta J(\theta) \approx \sum_{t,i} \nabla_\theta \log \pi_\theta(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$

5. Update  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$



$$y_{i,t} = \sum_{t'=t}^{t+n-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) + \gamma^n \hat{V}_\phi^\pi(\mathbf{s}_{i,t+n})$$

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

# Review so far

## Algorithms

*Policy gradients*: observe what is good vs. bad, then do more of good stuff

*Actor-critic*: learn to estimate what is good vs. bad, then do more of the good stuff

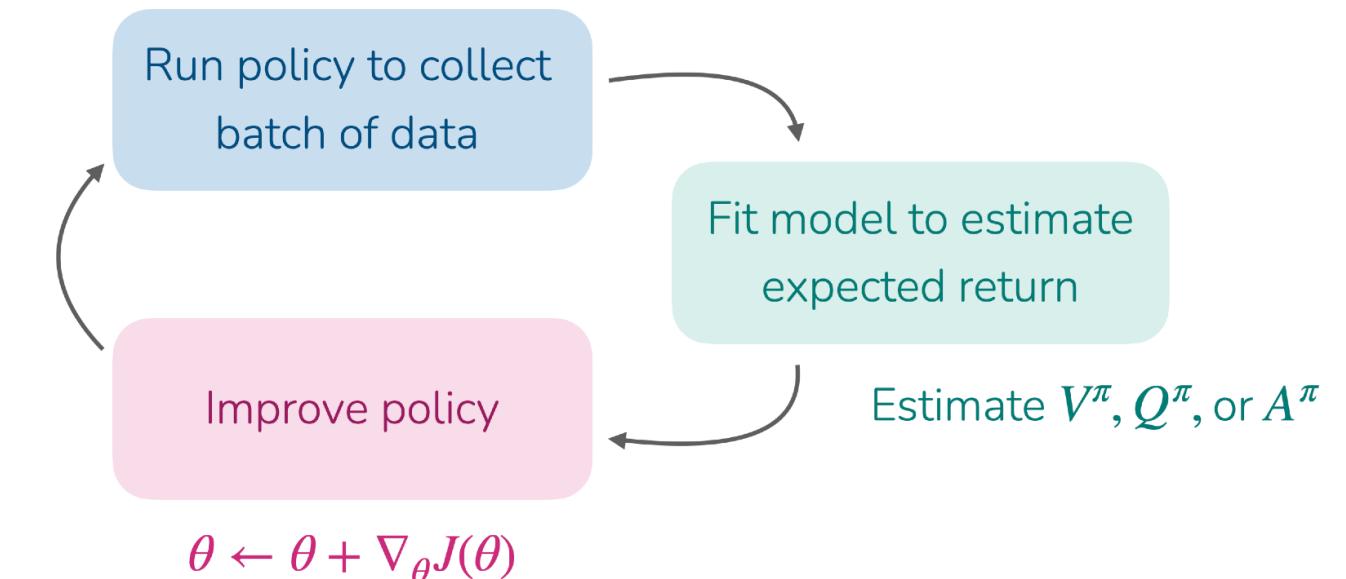
i.e. get a better policy gradient by **using neural network to estimate value function**

**How to estimate value:** “policy evaluation”

Supervised learning directly on observed sum of future rewards

Supervised learning on current reward + value estimate of next state

Hybrid: Supervised learning on sum of next rewards + value of state after that

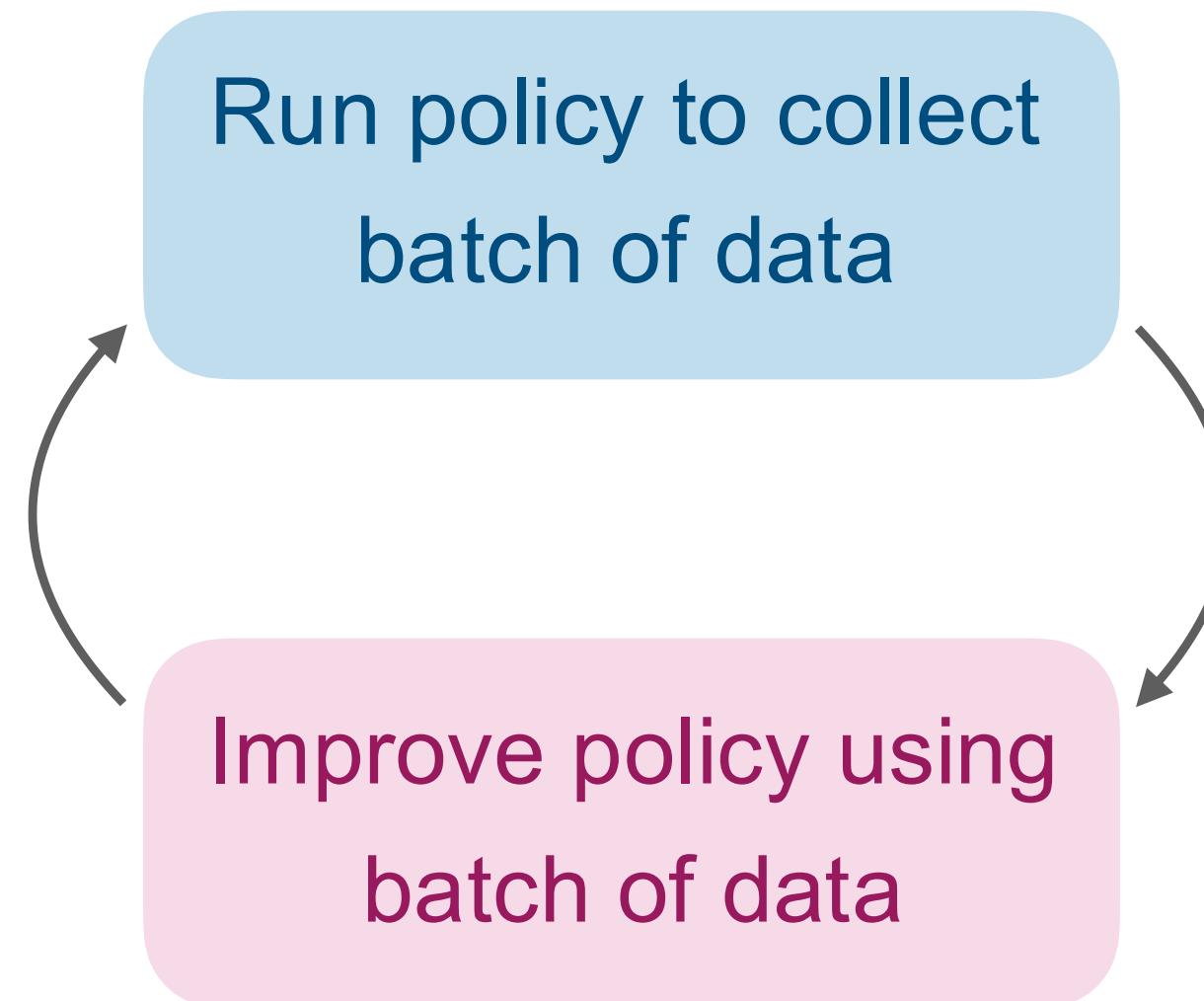


What about off-policy versions?

# Off-Policy Actor Critic Methods

CS 224R

# Recap: Policy Gradients



Online reinforcement learning with policy gradients

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

samples from policy      policy log likelihood      reward to go      baseline

Do more of the above average stuff, less of the below average stuff.

Importance weights for *off-policy* policy gradient:

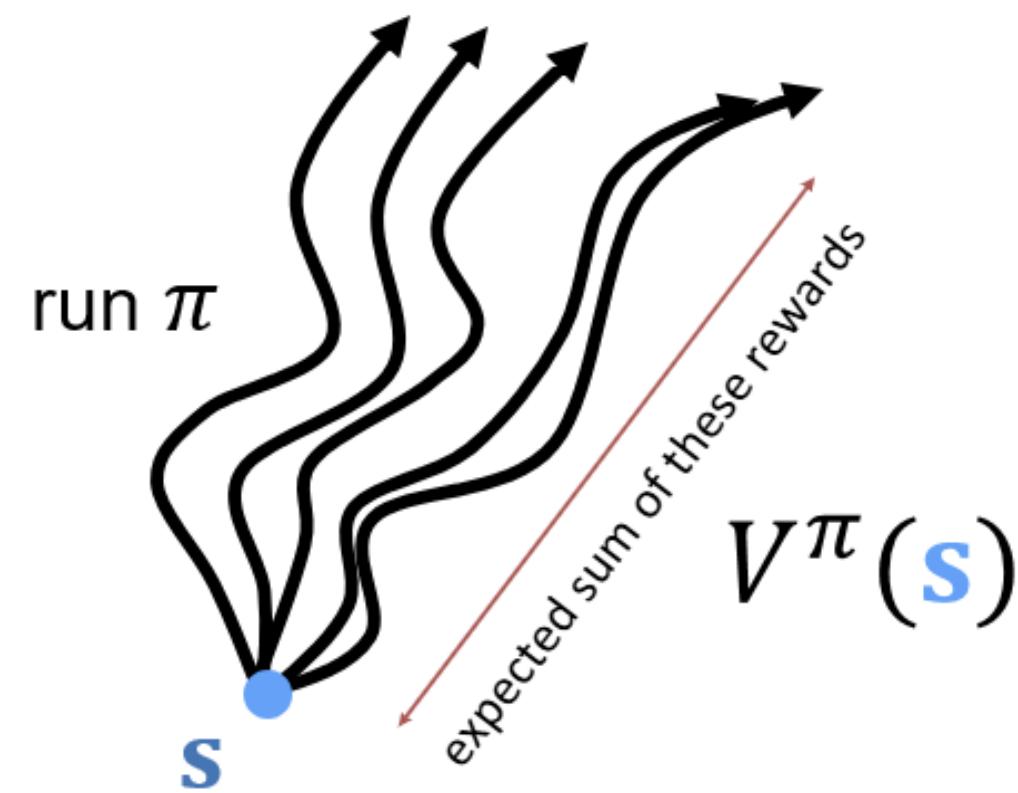
$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

# Recap: Some Useful Objects

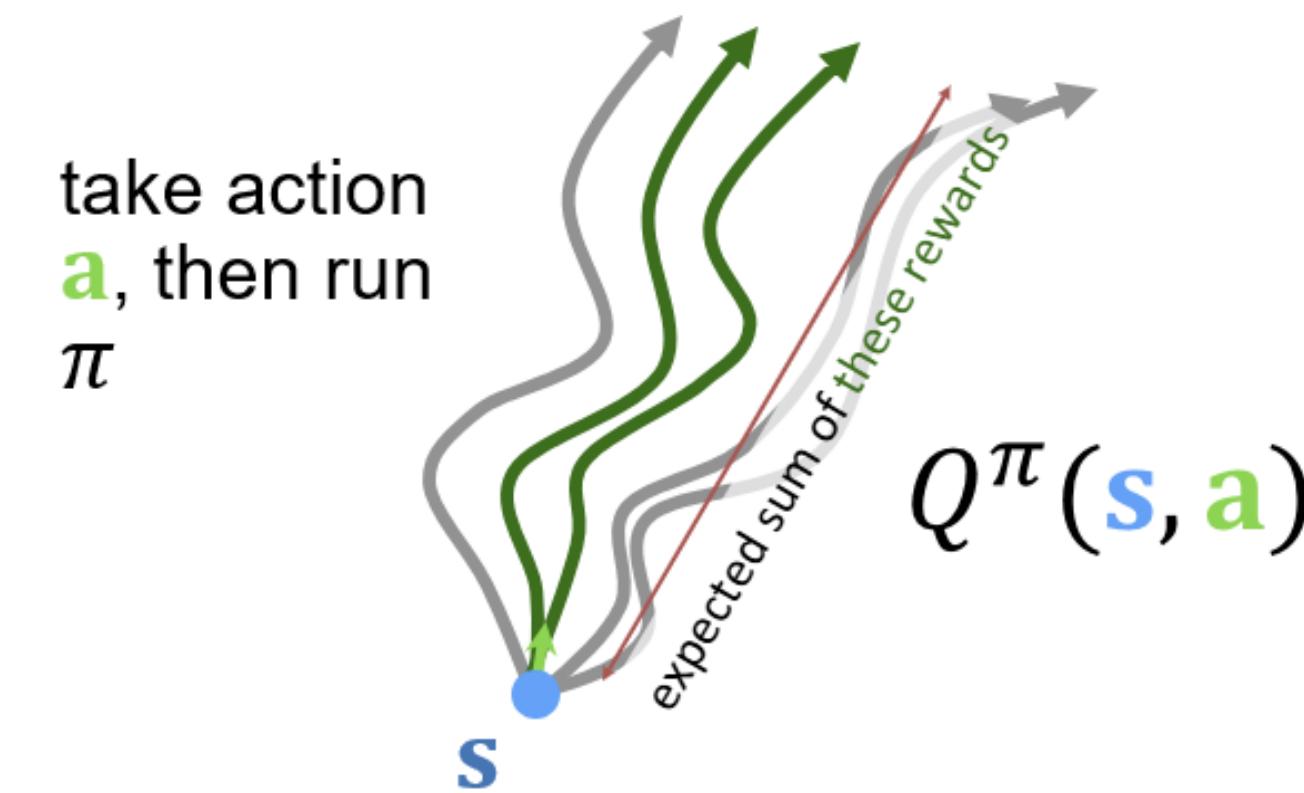
*value function*  $V^\pi(\mathbf{s})$  - future expected rewards starting at  $\mathbf{s}$  and following  $\pi$

*Q-function*  $Q^\pi(\mathbf{s}, \mathbf{a})$  - future expected rewards starting at  $\mathbf{s}$ , taking  $\mathbf{a}$ , then following  $\pi$

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$



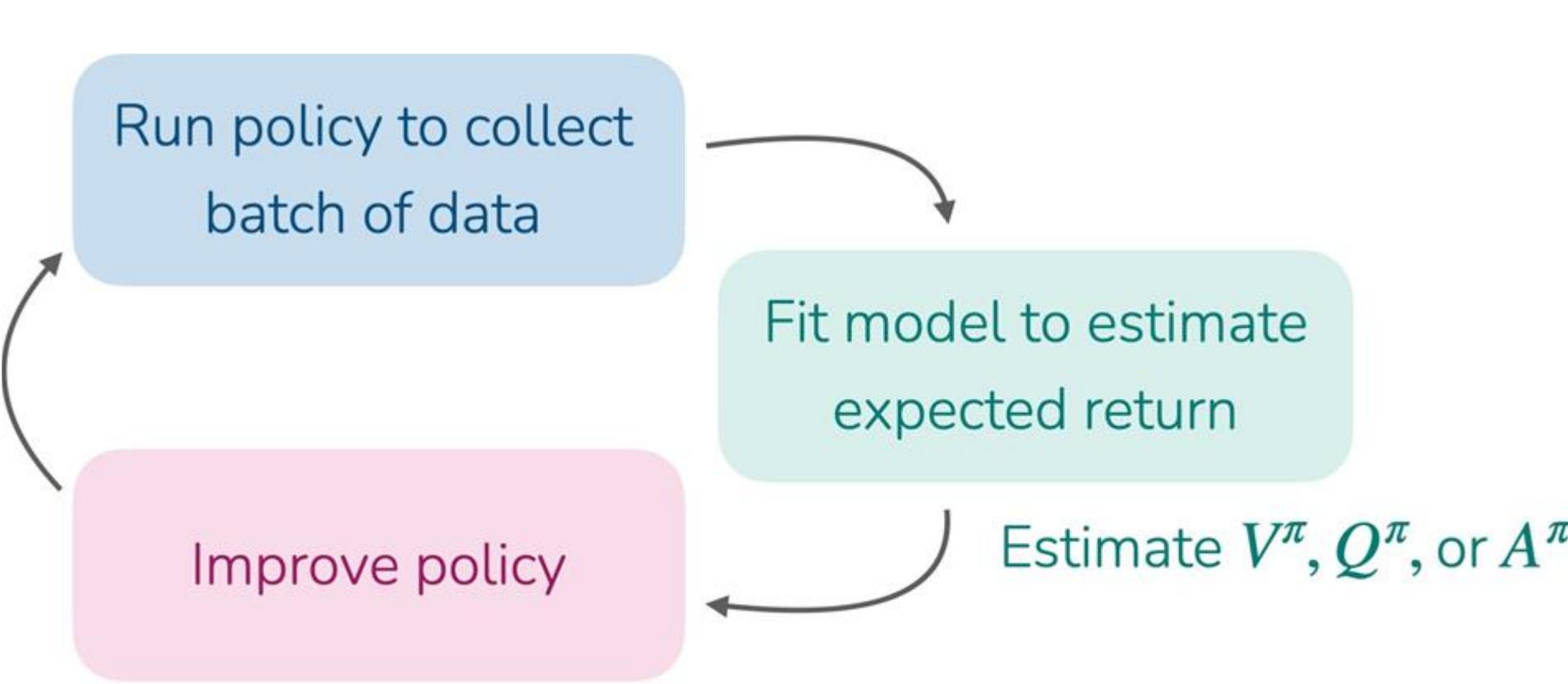
$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$



*advantage*  $A^\pi(\mathbf{s}, \mathbf{a})$  - how much better it is to take  $\mathbf{a}$  than to follow policy  $\pi$  at state  $\mathbf{s}$

$$A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s})$$

# Recap: Actor-Critic Methods



$$\theta \leftarrow \theta + \nabla_{\theta} J(\theta)$$

Online reinforcement learning with actor-critic

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

Estimate what is good and bad, then do more of the good stuff.

To estimate value, fit  $V^{\pi}$ :

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_i) - y_i \right\|^2$$

where  $y_{i,t} = \sum_{t'=t}^{t+n-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) + \gamma^n \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+n})$

$$A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\pi}(\mathbf{s}_{t+1}) - V^{\pi}(\mathbf{s}_t)$$

# The plan for today

## Off-policy actor critic methods

1. Taking multiple gradient steps
  - a. Importance weights
  - b. Constraining step size with KL penalty or clipping
  - c. Practical PPO algorithm
- b. Even more off-policy algorithm
  - a. Fitting Q-value functions with data from other policies
  - b. Practical SAC algorithm

## Key learning goals:

- All of the key concepts for practical algorithms like PPO and SAC

# Off-Policy Actor-Critic Methods

## Version 1: Multiple Gradient Steps

1. Sample batch of data  $\{(\mathbf{s}_{1,i}, \mathbf{a}_{1,i}, \dots, \mathbf{s}_{T,i}, \mathbf{a}_{T,i})\}$  from  $\pi_\theta$
2. Fit  $\hat{V}_\phi^{\pi_\theta}$  to summed rewards in data
3. Evaluate  $A^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) = r(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) + \gamma \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t+1,i}) - \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t,i}) \forall t, i$
4. Evaluate  $\nabla_\theta J(\theta) \approx \sum_{t,i} \nabla_\theta \log \pi_\theta(\mathbf{a}_{t,i} | \mathbf{s}_{t,i}) A^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$  <- use importance weights here
5. Update  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

# Off-Policy Actor-Critic Methods

## Version 1: Multiple Gradient Steps

1. Sample batch of data  $\{(\mathbf{s}_{1,i}, \mathbf{a}_{1,i}, \dots, \mathbf{s}_{T,i}, \mathbf{a}_{T,i})\}$  from  $\pi_\theta$
2. Fit  $\hat{V}_\phi^{\pi_\theta}$  to summed rewards in data
3. Evaluate  $A^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) = r(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) + \gamma \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t+1,i}) - \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t,i}) \forall t, i$
4. Evaluate  $\nabla_{\theta'} J(\theta') \approx \sum_{t,i} \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{t,i}|\mathbf{s}_{t,i}) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$  <- use importance weights here
5. Update  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

What can go wrong?

# Off-Policy Actor-Critic Methods

## Version 1: Multiple Gradient

1. Sample batch of data  $\{(\mathbf{s}_{1,i}, \mathbf{a}_{1,i}, \dots, \mathbf{s}_{T,i}, \mathbf{a}_{T,i})\}$  from  $\pi_\theta$
2. Fit  $\hat{V}_\phi^{\pi_\theta}$  to summed rewards in data
3. Evaluate  $A^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) = r(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) + \gamma \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t+1,i}) - \hat{V}_\phi^{\pi_\theta}(\mathbf{s}_{t,i}) \forall t, i$
4. Evaluate  $\nabla_{\theta'} J(\theta') \approx \sum_{t,i} \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{t,i}|\mathbf{s}_{t,i}) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$  <- use importance weights here
5. Update  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

Let's look at surrogate objective:

$$\tilde{J}(\theta') \approx \sum_{t,i} \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

Advantages based  
on old policy.

Policy will increase probability on actions with high advantages.

a convenient identity

$$p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) = \nabla_\theta p_\theta(\tau)$$

**What can go wrong if you take a lot of gradient steps?**

Policy is incentivized to differ significantly from old policy, can easily overfit.

# Off-Policy Actor-Critic Methods

## Version 1: Multiple Gradient Steps

Let's look at surrogate objective:

$$\tilde{J}(\theta') \approx \sum_{t,i} \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

Advantages based  
on old policy.

Policy will increase probability on actions with high advantages.

a convenient identity

$$p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) = \nabla_\theta p_\theta(\tau)$$

### What can go wrong if you take a lot of gradient steps?

Policy is incentivized to differ significantly from old policy, can easily overfit.

💡 Idea 1: Use KL constraint on policy.

Very common. Will see in LLM preference optimization

$$\mathbb{E}_{\mathbf{s} \sim \pi_\theta} [D_{KL}(\pi_{\theta'}(\cdot | \mathbf{s}) \| \pi_\theta(\cdot | \mathbf{s}))] \leq \delta$$

💡 Idea 2: Can we bound the importance weights?

Doesn't directly constrain policy, but removes incentives

-> A key idea behind proximal policy optimization (PPO)

# Proximal Policy Optimization (PPO)

Off-policy actor-critic with a few tricks.

Surrogate objective:

$$\tilde{J}(\theta') \approx \sum_{t,i} \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

**Trick #1:** Clip the importance weights:

$$\tilde{J}(\theta') \approx \sum_{t,i} \text{clip}\left(\frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}, 1 - \epsilon, 1 + \epsilon\right) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i})$$

Policy no longer incentivized to deviate significantly

**Trick #2:** Take **minimum** w.r.t. original objective: (in rare event where clipping makes objective better)

$$\tilde{J}(\theta') \approx \sum_{t,i} \min \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}), \text{clip}\left(\frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}, 1 - \epsilon, 1 + \epsilon\right) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) \right)$$

This is the final PPO surrogate objective!

How do we estimate the advantage?

# Proximal Policy Optimization (PPO)

Off-policy actor-critic with a few tricks.

Surrogate objective:

$$\tilde{J}(\theta') \approx \sum_{t,i} \min \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})} \hat{A}^{\pi_{\theta}}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}), \text{clip} \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t}|\mathbf{s}_{i,t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}^{\pi_{\theta}}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) \right)$$

**Trick #3:** “Generalized advantage estimation”

(GAE) Fit  $\pi$  with Monte Carlo or bootstrapping.

Then, use varying horizon to estimate advantage:

$$\hat{A}_n^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{t+n} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t) + \gamma^n \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+n})$$

$$\hat{A}_{\text{GAE}}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{n=1}^{\infty} w_n \hat{A}_n^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$$

How to choose weights? cut earlier for less variance

$$w_n \propto \lambda^{n-1} \quad \text{for } n = 1, \dots, N, \text{ where } N \leq T$$

# Proximal Policy Optimization (PPO)

1. Sample batch of data  $\{(\mathbf{s}_{1,i}, \mathbf{a}_{1,i}, \dots, \mathbf{s}_{T,i}, \mathbf{a}_{T,i})\}$  from  $\pi_\theta$
2. Fit  $\hat{V}_\phi^{\pi_\theta}$  to summed rewards in data
3. Evaluate  $\hat{A}_{\text{GAE}}^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{n=1}^{\infty} w_n \left( \sum_{t'=t}^{t+n} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_\phi^\pi(\mathbf{s}_t) + \gamma^n \hat{V}_\phi^\pi(\mathbf{s}_{t+n}) \right)$
4. Update policy with M gradient steps on surrogate objective:

$$\tilde{J}(\theta') \approx \sum_{t,i} \min \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}), \text{clip} \left( \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}^{\pi_\theta}(\mathbf{s}_{t,i}, \mathbf{a}_{t,i}) \right)$$

Some example hyperparameters:

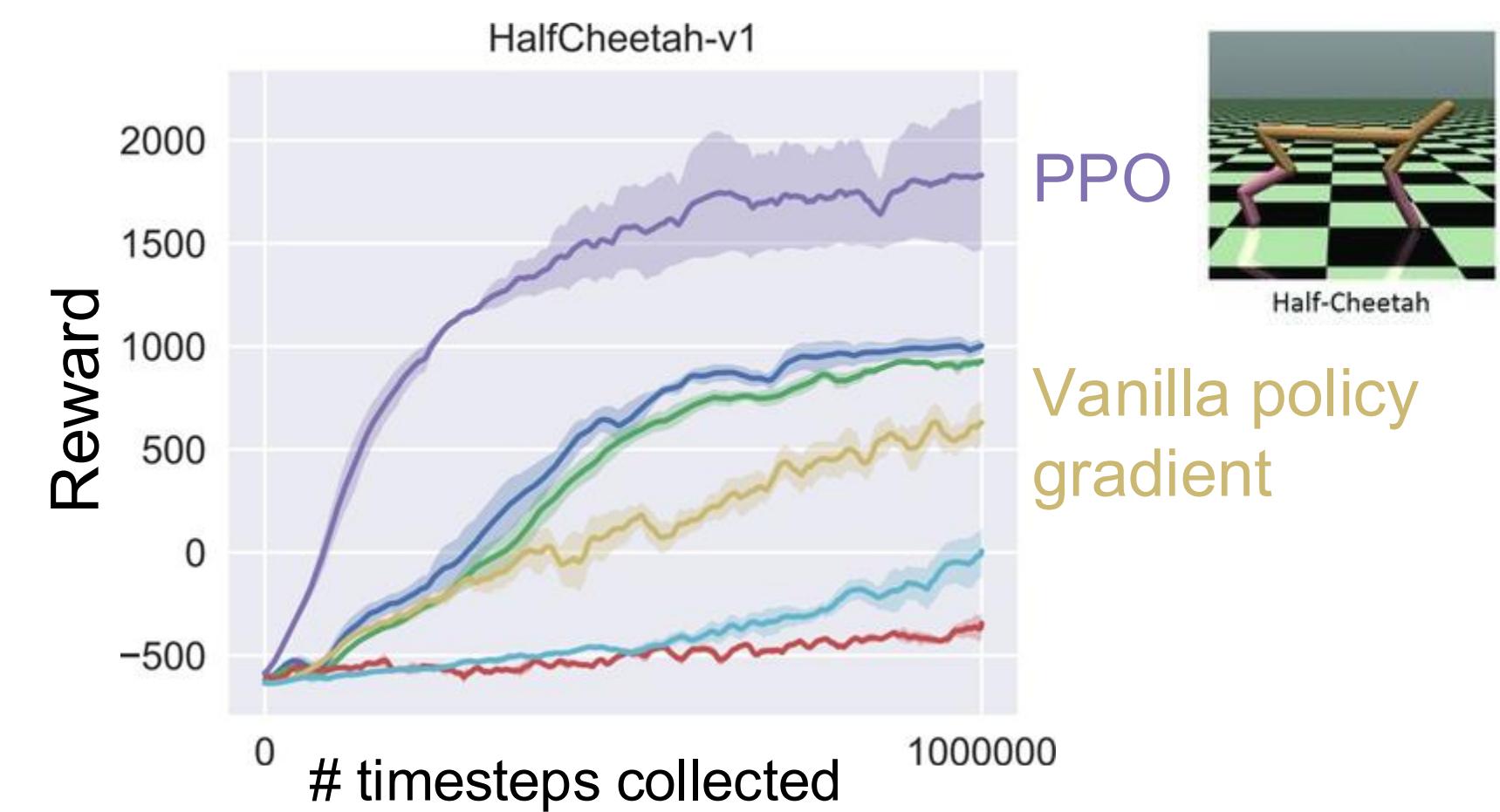
~2000 timesteps in batch of data

clipping range  
=0.2

~10 epochs when updating policy

( =~300 gradient steps with batch size 64)

~500 iterations —> 1M total timesteps of experience



# Off-Policy Actor-Critic Methods

**So far:**

- use one batch of policy data for one gradient step (fully on-policy)
- use one batch of policy data for multiple gradient steps (starting to be off-policy)

**Can we be even more off-policy?**

Can we reuse data from previous batches, i.e. all of the past trial-and-error data?

Two key ideas:

- maintain a buffer of all past data “*replay buffer*”
- adjust equations to remove on-policy assumptions

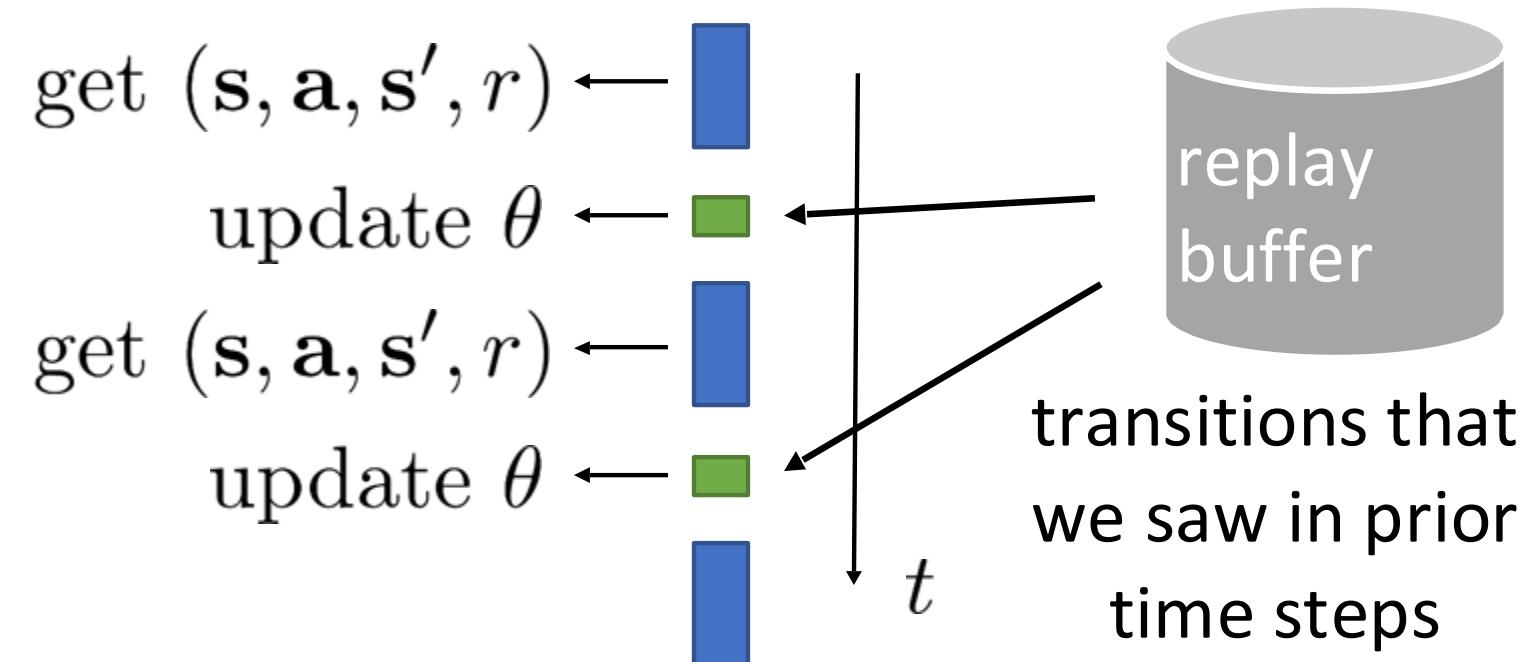
# Off-Policy Actor-Critic Methods

## Version 2: Replay buffers

actor-critic algorithm:

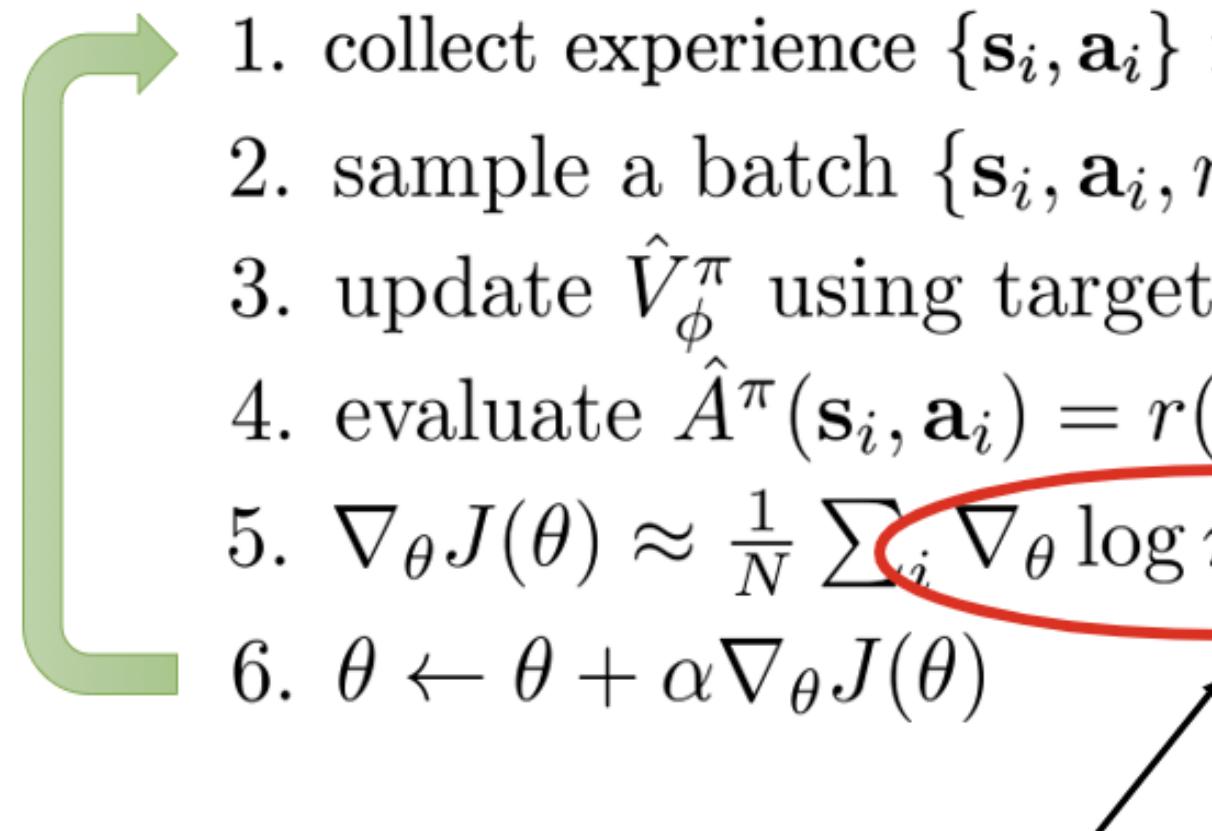
1. collect experience  $\{\mathbf{s}_i, \mathbf{a}_i\}$  from  $\pi_\theta(\mathbf{a}|\mathbf{s})$
  2. fit  $\hat{V}_\phi^\pi(\mathbf{s})$  to sampled reward sums
  3. evaluate  $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$
  4.  $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
  5.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
- Add this to replay buffer
- Do this on minibatch sampled from all previous data

off-policy actor-critic



# Off-Policy Actor-Critic Methods

## Version 2: Replay buffers

- 
1. collect experience  $\{\mathbf{s}_i, \mathbf{a}_i\}$  from  $\pi_\theta(\mathbf{a}|\mathbf{s})$  (& add to replay buffer)
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{V}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i)$  for each  $\mathbf{s}_i$
  4. evaluate  $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) \leftarrow \hat{V}_\phi^\pi(\mathbf{s}_i)$
  5.  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
  6.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

not the action  $\pi_\theta$  would have taken!

not the right target value



$$\mathcal{L}(\phi) = \frac{1}{N} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

↗  
mini batch size

**Question:** When we fit  $V^\pi$  on all data from replay buffer, what policy  $\pi$  is it learning a value function for?

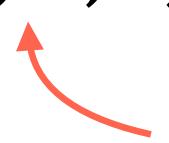
The algorithm is currently broken 😢

# Off-Policy Actor-Critic Methods

How do we fit a value function for  $\pi_\theta$  using replay buffer of data from past policies?

💡 What if we fit  $Q(s, a)$  instead of  $V(s)$ ?

The datapoints we have:  $(s, a, r, s')$  + future reward

 depends on past policy's actions  
 action from past policy

-> Don't want to use this.

-> If we fit  $Q(s, a)$ , we pass the action as input

(okay if it is a bit different from what policy would have done)

For any  $(s, a)$ :

$$Q^{\pi_\theta}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a), \bar{a}' \sim \pi_\theta(\cdot | s')} [Q^{\pi_\theta}(s', \bar{a}')]$$

Recall: Definition of Q-values

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$$

“total reward we get if we take  $a_t$  in  $s_t$ ...  
... and then follow the policy  $\pi$ ”

# Off-Policy Actor-Critic Methods

How to we fit a value function for  $\pi_\theta$  using replay buffer of data from past policies?

💡 What if we fit  $Q(s, a)$  instead of  $V(s)$ ?

The datapoints we have:  $(s, a, r, s')$  + future reward

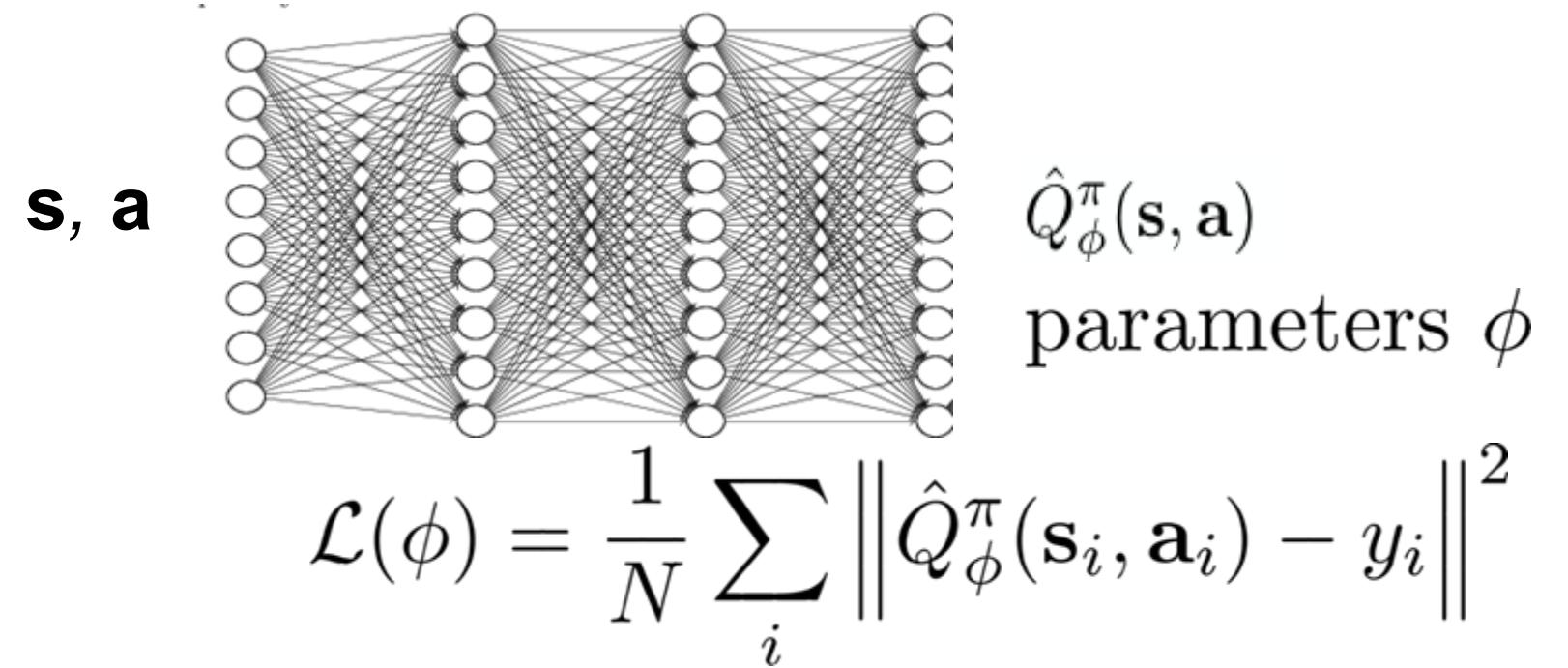
For any  $(s, a)$ :  $Q^{\pi_\theta}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a), \bar{a}' \sim \pi_\theta(\cdot | s')} [Q^{\pi_\theta}(s', \bar{a}')] \quad \text{Eqn 1}$

**Approach:**

1. Sample  $(s_i, a_i, s'_i)$  from the buffer
2. Sample  $\bar{a}'_i \sim \pi_\theta(\cdot | s'_i)$  from current policy

$$Q^{\pi_\theta}(s_i, a_i) \approx \underbrace{r(s_i, a_i) + \gamma Q^{\pi_\theta}(s'_i, \bar{a}'_i)}_{y_i} \quad \text{Eqn 2}$$

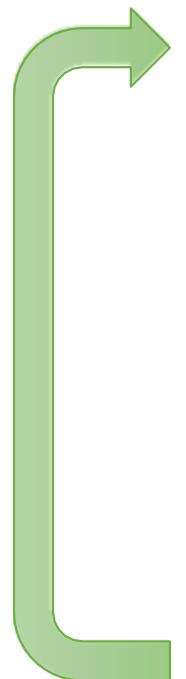
**Note:** To be accurate, targets  $y_i$  need sufficient action coverage



# Off-Policy Actor-Critic Methods

## Version 2: Fixing the value function

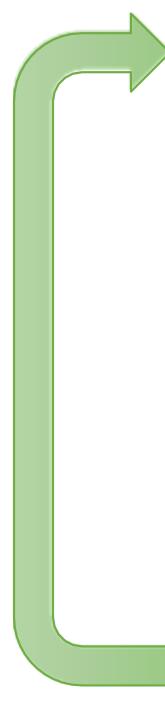
online actor-critic algorithm:

- 
1. take action  $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{V}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i)$  for each  $\mathbf{s}_i$
  4. evaluate  $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) \leftarrow \hat{V}_\phi^\pi(\mathbf{s}_i)$
  5.  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
  6.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
- not the right target value
- 
3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$
- ↑  
**not** from replay buffer  $\mathcal{R}!$
- $$\mathbf{a}'_i \sim \pi_\theta(\mathbf{a}'_i|\mathbf{s}'_i)$$

# Off-Policy Actor-Critic Methods

## Version 2: Fixing the value function

online actor-critic algorithm:

- 
1. take action  $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$
  4. evaluate  $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$
  5.  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
  6.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

**First:** convenient to use  $Q^\pi$  instead of  $A^\pi$ . (i.e. no average reward baseline)

higher variance, but okay b/c we're now using a lot more data (all data in buffer)

**Second:** current policy's actions likely better than past policy's actions

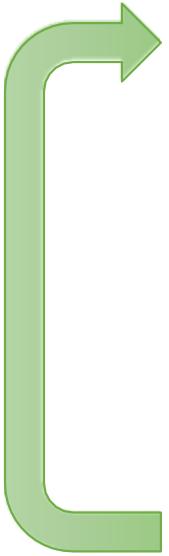
$$\text{sample } \mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a}|\mathbf{s}_i) \quad \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi|\mathbf{s}_i) \cdot \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$$

↑                      ↑  
**not** from replay buffer  $\mathcal{R}!$

# Off-Policy Actor-Critic Methods

Version 2: Anything else?

online actor-critic algorithm:

- 
1. take action  $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$
  4.  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$  where  $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a}|\mathbf{s}_i)$
  5.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

**Any remaining problems?**

$\mathbf{s}_i$  didn't come from  $p_\theta(\mathbf{s})$

nothing we can do here, just accept it

**intuition:** we want optimal policy on  $p_\theta(\mathbf{s})$   
but we get optimal policy on a *broader* distribution

# Off-Policy Actor-Critic Methods

Version 2: Some implementation details

online actor-critic algorithm:

1. take action  $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$  for each  $\mathbf{s}_i, \mathbf{a}_i$
4.  $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$  where  $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a}|\mathbf{s}_i)$
5.  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

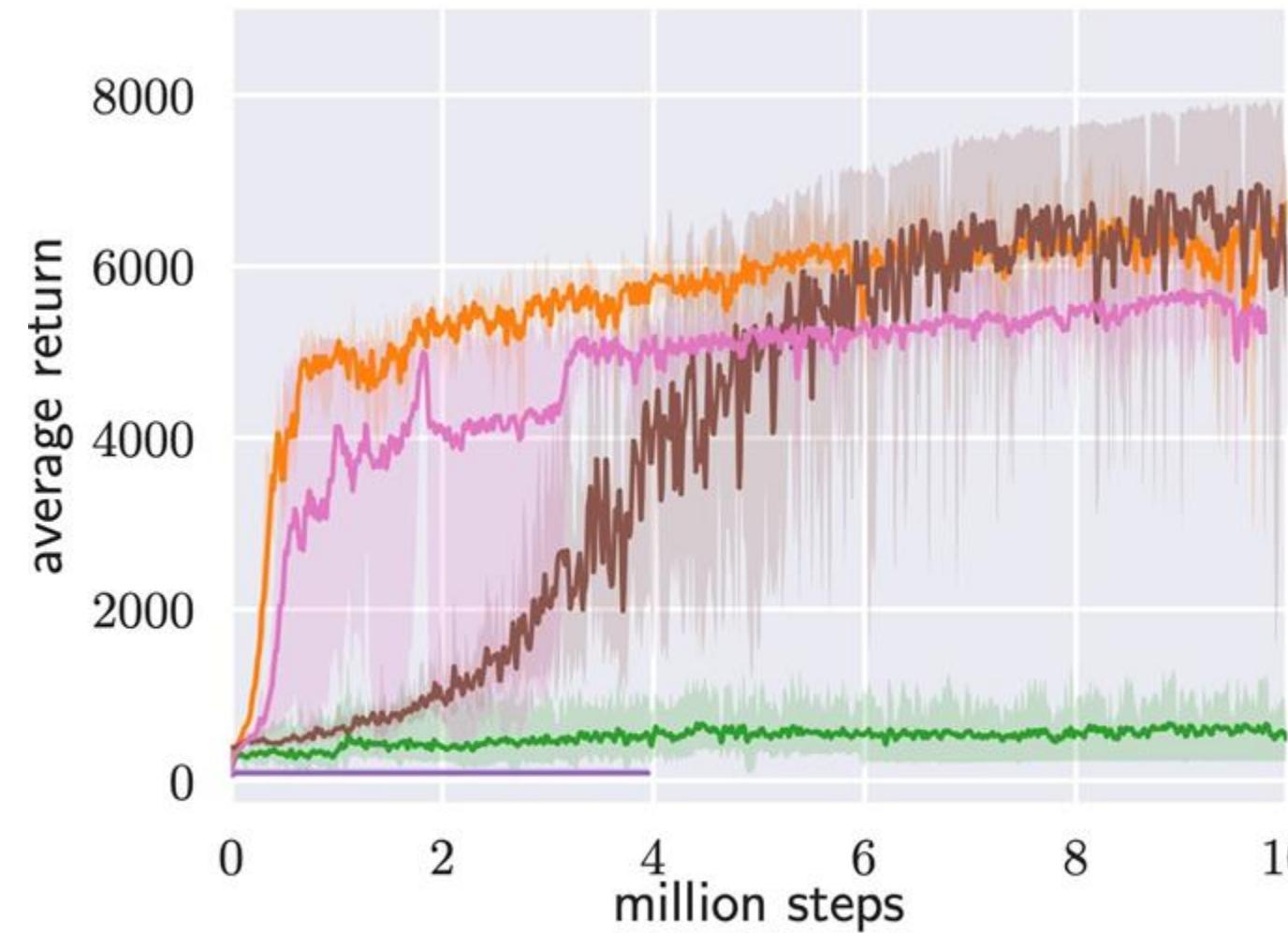
also fancier ways to fit Q-functions  
(more on this in next two lectures)

can also use **reparameterization trick** to  
better estimate the gradient (for Gaussian  
policy)

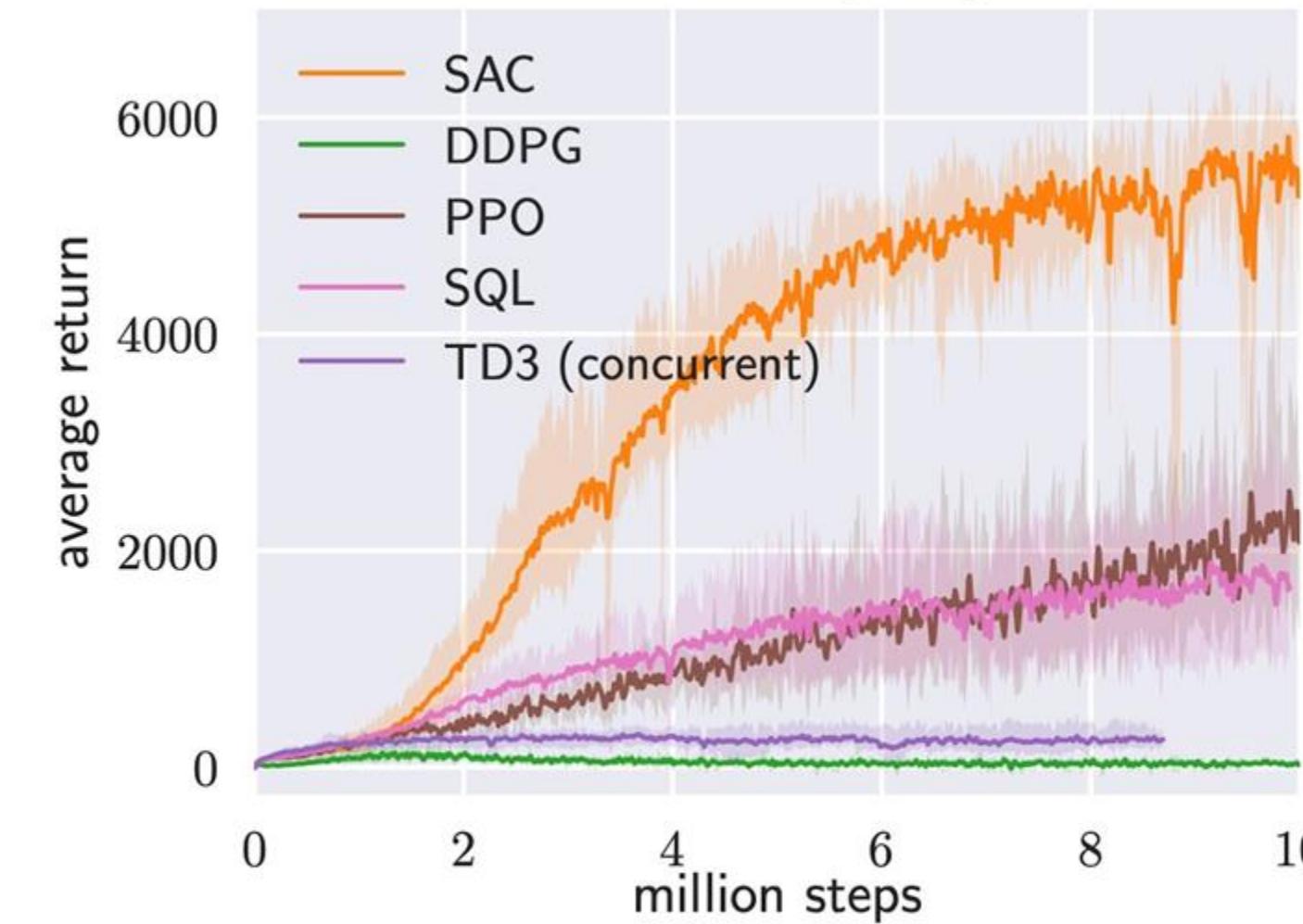
Example practical algorithm:

Haarnoja, Zhou, Abbeel, Levine. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. 2018.

# More off-policy vs. less off-policy actor critic?



(e) Humanoid-v1



(f) Humanoid (rllab)

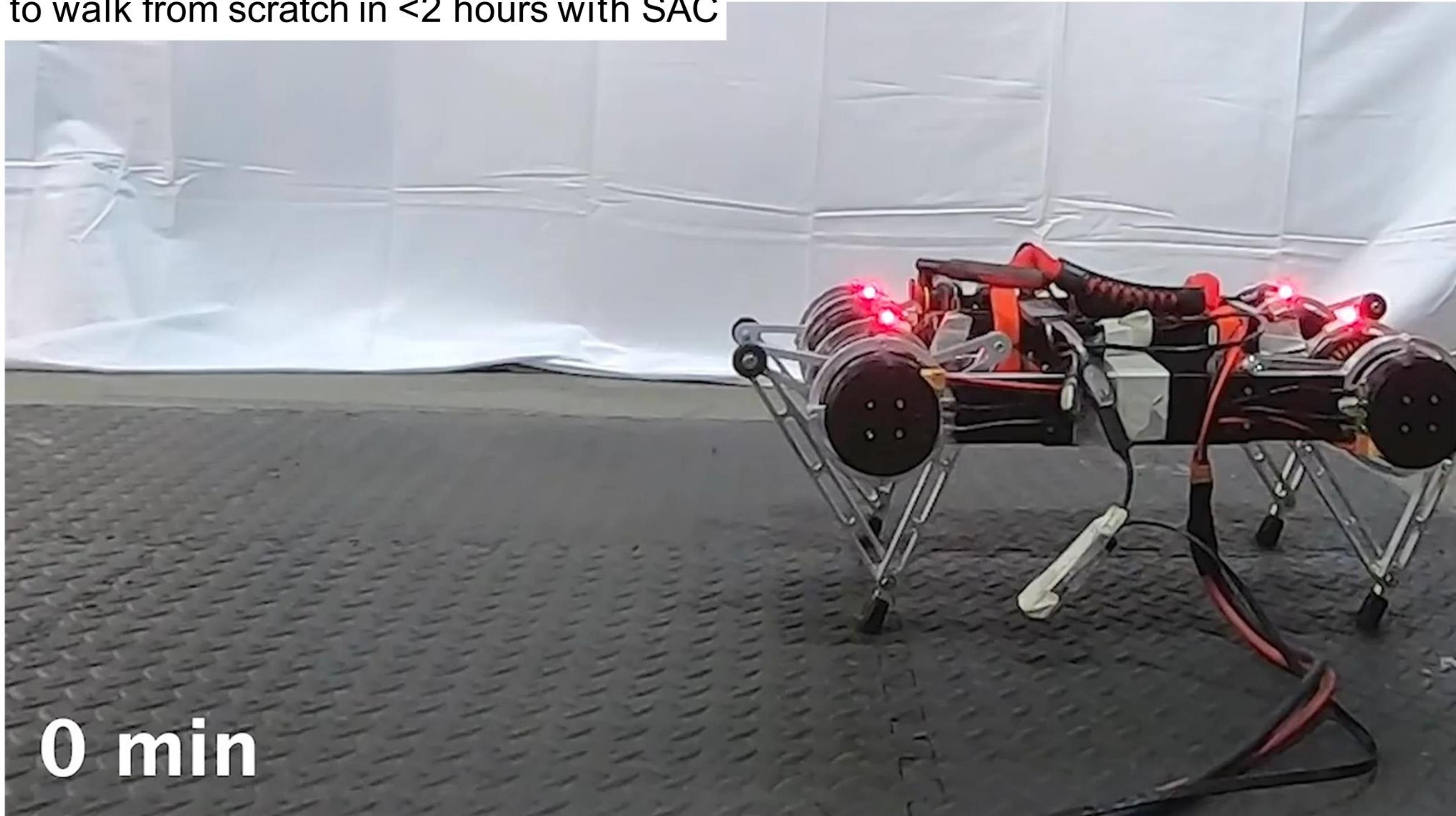
<- more off-policy  
(i.e. with replay buffer)

<- less off-policy  
(i.e. no replay buffer)

- + Off-policy with replay buffer (e.g. soft actor-critic) can be far more data efficient
- They can also generally be a lot harder to tune hyperparameters, less stable (than e.g. PPO)

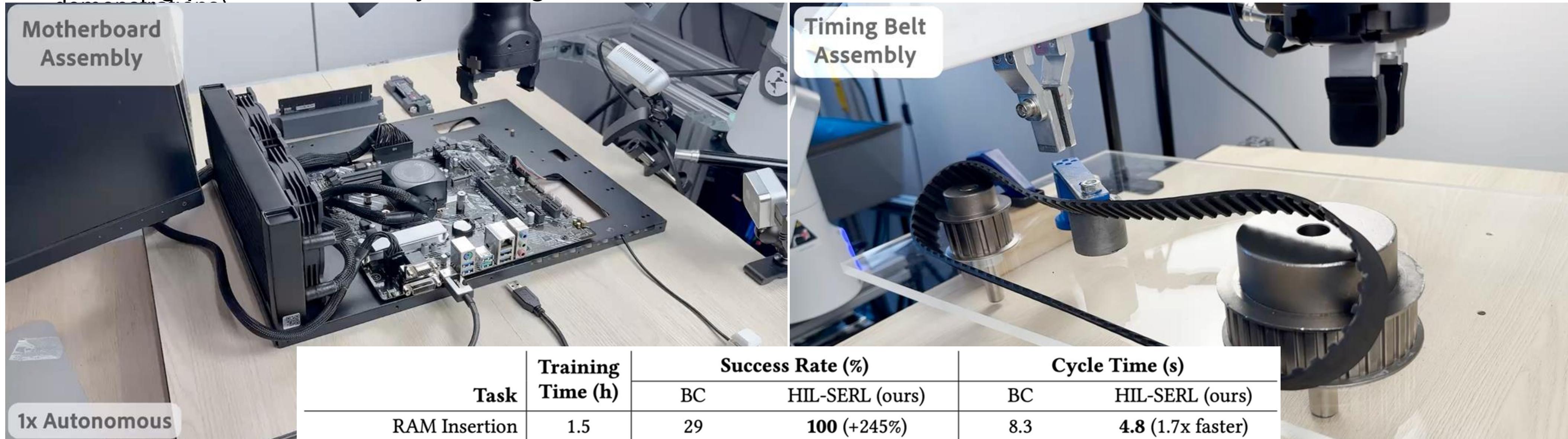
# Fully Off-Policy = Can be efficient enough for real-world RL

Learning to walk from scratch in <2 hours with SAC



# Fully Off-Policy = Can be efficient enough for real-world RL

Learning precise assembly with algorithm based on SAC (seeded with



| Task                   | Training Time (h) | Success Rate (%) |                      | Cycle Time (s) |                          |
|------------------------|-------------------|------------------|----------------------|----------------|--------------------------|
|                        |                   | BC               | HIL-SERL (ours)      | BC             | HIL-SERL (ours)          |
| RAM Insertion          | 1.5               | 29               | <b>100 (+245%)</b>   | 8.3            | <b>4.8 (1.7x faster)</b> |
| SSD Assembly           | 1                 | 79               | <b>100 (+27%)</b>    | 6.7            | <b>3.3 (2x faster)</b>   |
| IKEA - Side Panel 1    | 2                 | 77               | <b>100 (+30%)</b>    | 6.5            | <b>2.7 (2.4x faster)</b> |
| IKEA - Side Panel 2    | 1.75              | 79               | <b>100 (+27%)</b>    | 5.0            | <b>2.4 (2.1x faster)</b> |
| IKEA - Top Panel       | 1                 | 35               | <b>100 (+186%)</b>   | 8.9            | <b>2.4 (3.7x faster)</b> |
| IKEA - Whole Assembly  | -                 | 1/10             | <b>10/10 (+900%)</b> | -              | -                        |
| Car Dashboard Assembly | 2                 | 41               | <b>100 (+144%)</b>   | 20.3           | <b>8.8 (2.3x faster)</b> |
| Timing Belt Assembly   | 6                 | 2                | <b>100 (+4900%)</b>  | 9.1            | <b>7.2 (1.3x faster)</b> |

<- better, faster than imitation learning

# PPO = Common algorithm of choice for stable, less efficient learning

Learning to walk in simulation with PPO -> transfer to real world (with careful simulation choices)

[https://www.youtube.com/watch?v=xf\\_ UXK0OTIk](https://www.youtube.com/watch?v=xf_ UXK0OTIk)

# PPO = Common algorithm of choice for stable, less efficient learning

Learning to solve Rubik's cube in simulation -> transfer to real hand

(though, transfer from simulation generally much harder, less successful for robot manipulation)

<https://www.youtube.com/watch?v=jm-ihc7CASY>

# What about RL for language models?

PPO is a common choice!

Will start to discuss them when talking about reward modeling next week.

# Q-Learning

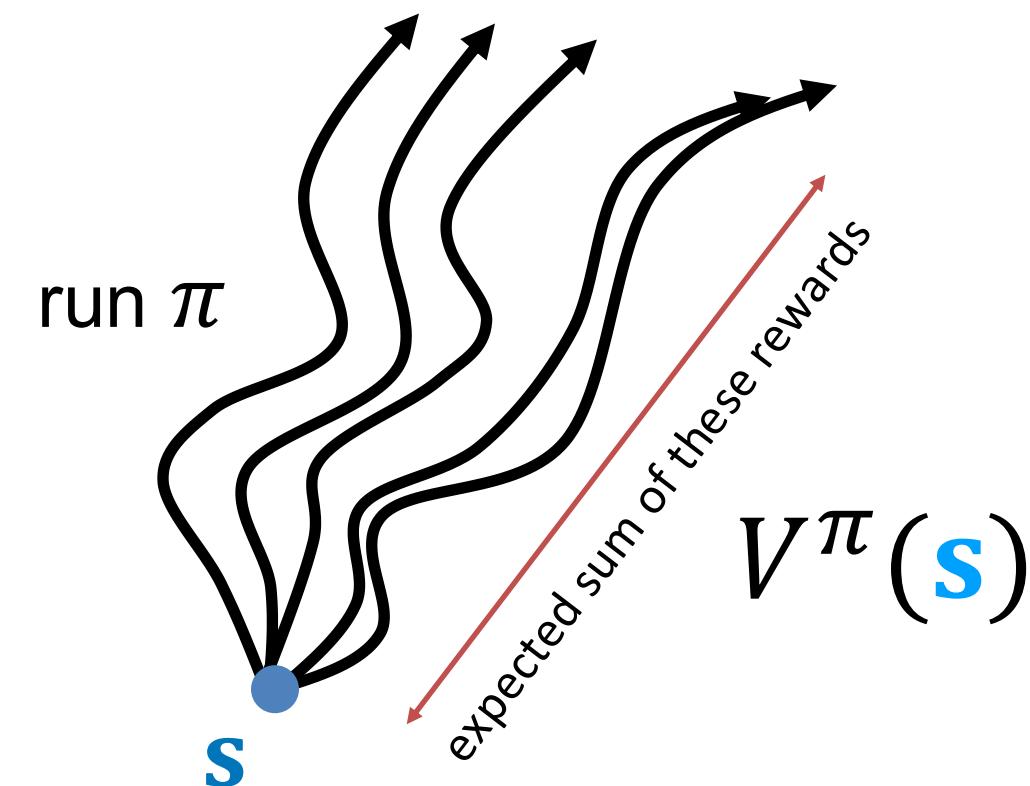
CS 224R

# Recap: Some Useful Objects

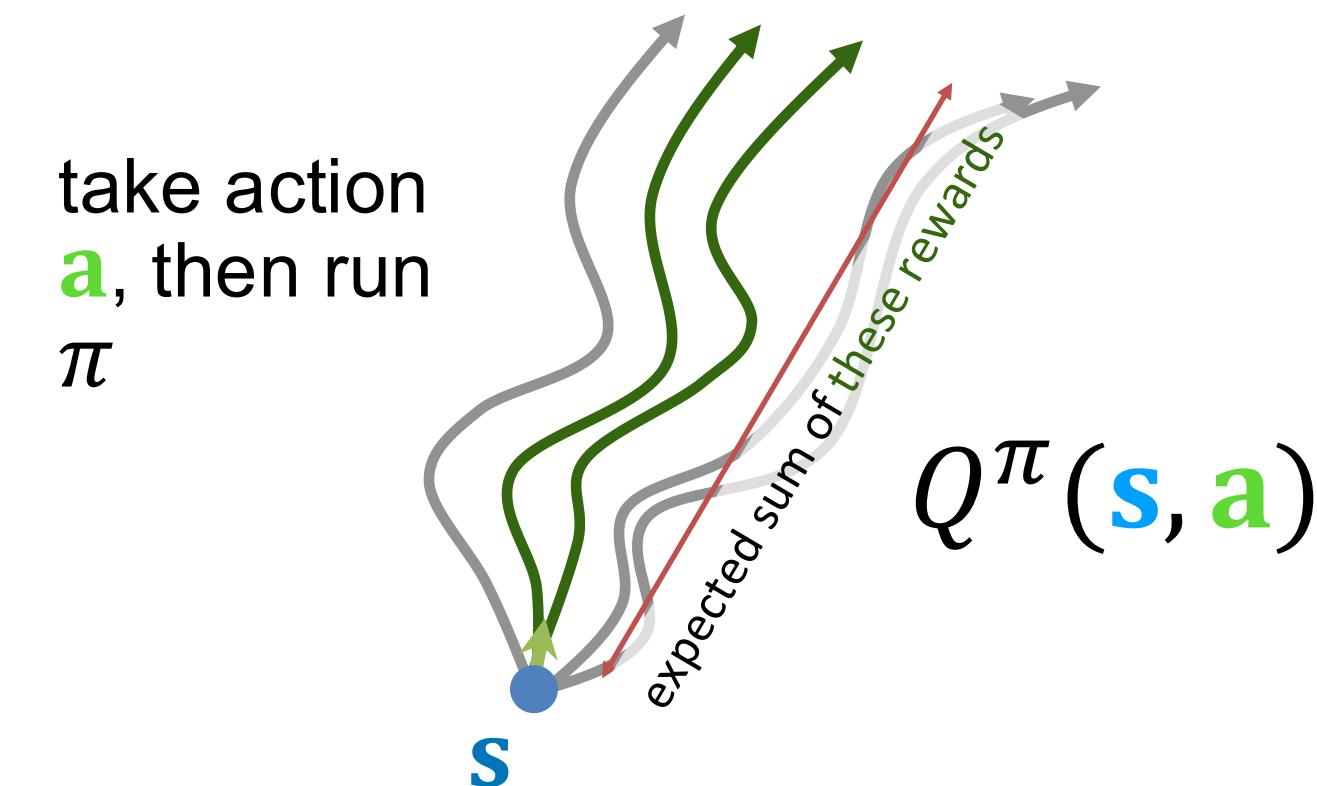
*value function*  $V^\pi(\mathbf{s})$  - future expected rewards starting at  $\mathbf{s}$  and following  $\pi$

*Q-function*  $Q^\pi(\mathbf{s}, \mathbf{a})$  - future expected rewards starting at  $\mathbf{s}$ , taking  $\mathbf{a}$ , then following  $\pi$

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

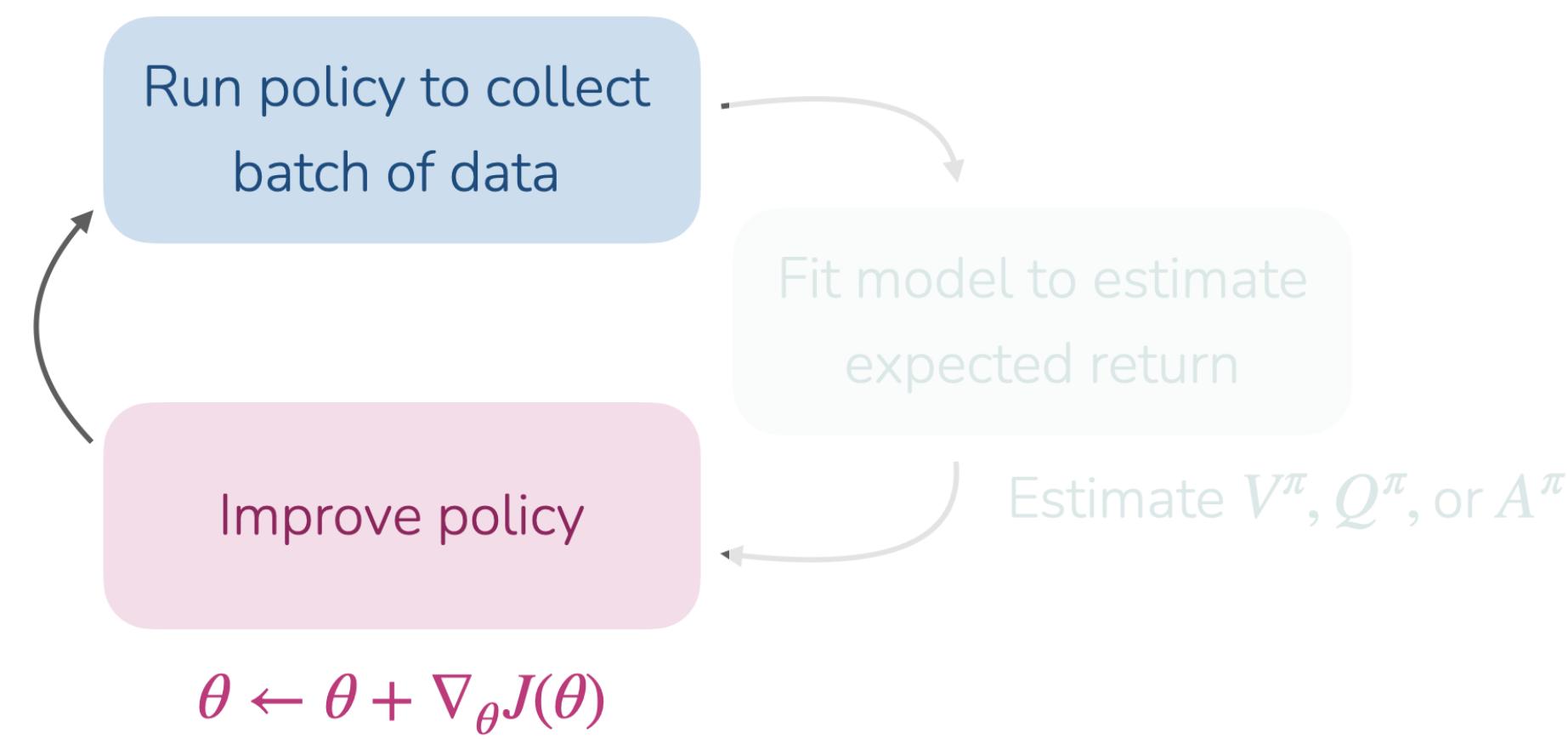


$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$$



# Recap: Methods

## Online RL with policy gradients



$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \left( \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

samples from policy      policy log likelihood      reward to go baseline

Do more of the **above average** stuff,  
less of the **below average** stuff.

## Online RL with actor-critic

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

Estimate what is good and bad, then  
do more of the good stuff.

# Recap: Off-Policy Policy Evaluation

## Online RL with actor-critic

Run policy to collect batch of data

Fit model to estimate expected return

Improve policy

$$\theta \leftarrow \theta + \nabla_{\theta} J(\theta)$$

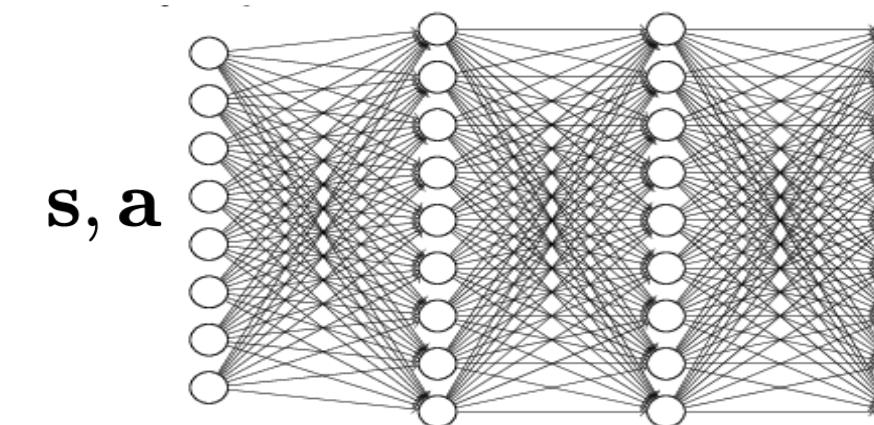
Estimate what is good and bad, then do more of the good stuff.

1. Sample  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i)$  from the buffer
2. Sample  $\bar{\mathbf{a}}'_i \sim \pi_{\theta}(\cdot | \mathbf{s}'_i)$  from *current* policy

$$Q^{\pi_{\theta}}(\mathbf{s}_i, \mathbf{a}_i) \approx r(\mathbf{s}_i, \mathbf{a}_i) + \gamma Q^{\pi_{\theta}}(\mathbf{s}'_i, \bar{\mathbf{a}}'_i)$$

$$y_i$$

**For off-policy version:** can we estimate **Q-function** using past policy data?



For any  $(\mathbf{s}, \mathbf{a})$ :

$$Q^{\pi_{\theta}}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a}), \bar{\mathbf{a}}' \sim \pi_{\theta}(\cdot | \mathbf{s}')} [Q^{\pi_{\theta}}(\mathbf{s}', \bar{\mathbf{a}}')]$$

$$\mathcal{L}(\phi) = \frac{1}{N} \sum_i \left\| \hat{Q}_{\phi}^{\pi}(\mathbf{s}_i, \mathbf{a}_i) - y_i \right\|^2$$

**Note:** To be accurate, targets  $y_i$  need sufficient action coverage

# The plan for today

## Value-based RL methods

1. Q-learning RL method
  - a. Policy iteration
  - b. Bellman optimality equation
  - c. How to collect data for Q-learning methods
2. Q-learning in practice
  - d. Target networks
  - e. Double DQN
  - f. N-step returns

## Key learning goals:

- How Q-functions relate to policies
- How to do RL without learning an explicit policy
- How to stabilize Q-learning in practice

# A thought exercise

Recall: Definition of Q-values

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_\pi [r(s_{t'}, a_{t'}) | s_t, a_t]$$

“total reward we get if we take  $a_t$  in  $s_t$ ...  
... and then follow the policy  $\pi$ ”

For some policy  $\pi$ , say you have an accurate estimate  $\hat{Q}^\pi(s, a)$  for all  $s, a$

Define a new policy  $\pi_{\text{new}}(a_t | s_t) = \begin{cases} 1 & \text{if } a_t = \arg \max_a \hat{Q}_\phi^\pi(s_t, a) \\ 0 & \text{otherwise} \end{cases}$

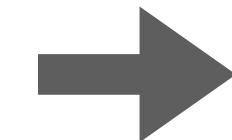
2D navigation example

Initial state distribution  
 $p(s_1)$



Reward is 1 at  $\star$ , 0 elsewhere

Current policy  $\pi$  goes right in all states



takes 1 timestep, 2 timesteps to travel these distances

**Questions:** Is the new policy better, worse, or the same as the original policy?  
Is it the optimal policy? Why or why not?

# Can we omit policy gradient completely?

$Q^\pi(s_t, a_t)$ : expected reward from taking  $a_t$  and subsequently following  $\pi$

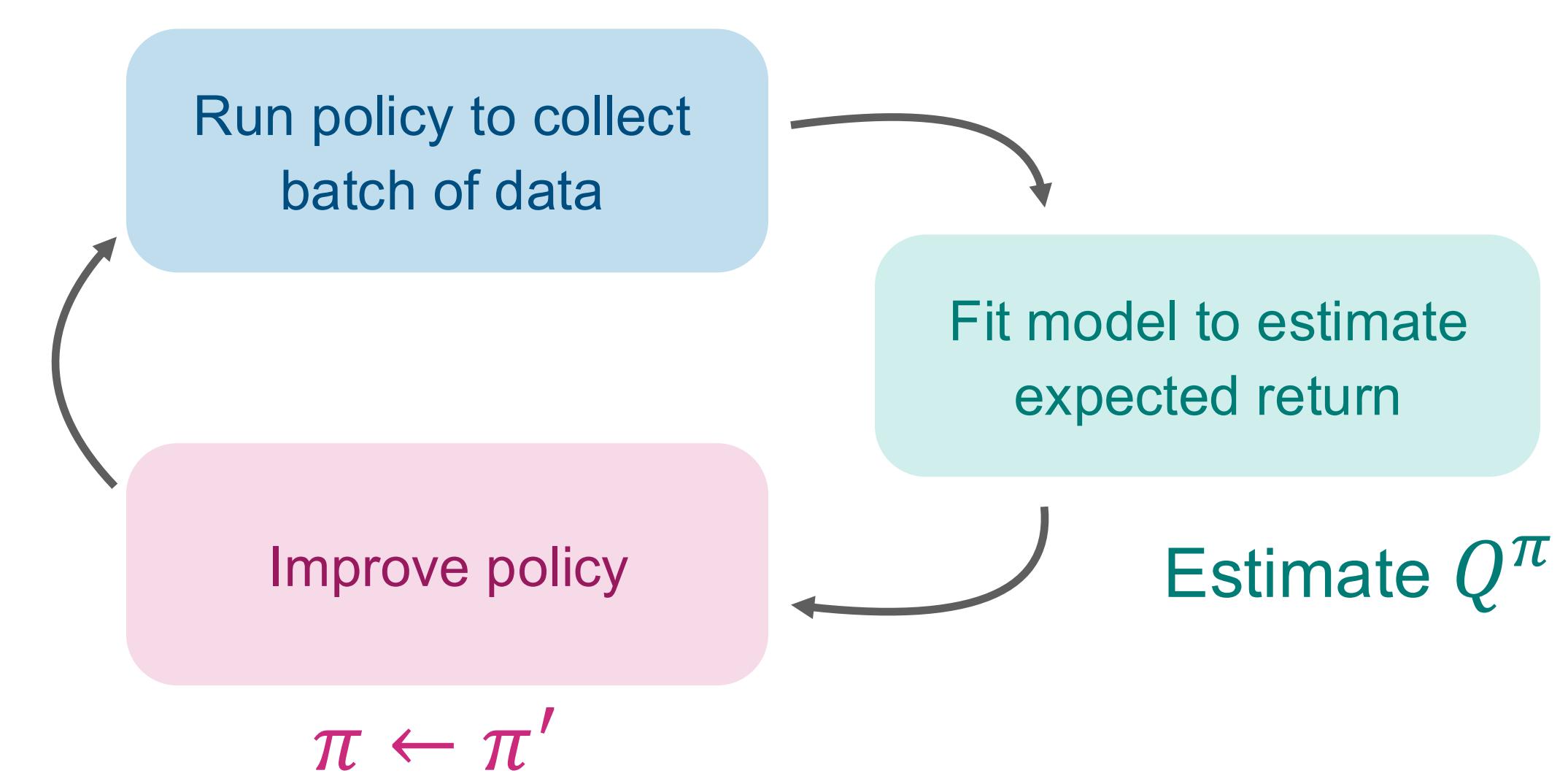
$\arg \max_{a_t} Q^\pi(s_t, a_t)$ : best action from  $s_t$ , if we then follow  $\pi$  afterwards

at *least* as good as any  $a_t \sim \pi(a_t|s_t)$   
regardless of what  $\pi(a_t|s_t)$  is!

forget policies, let's just do this!

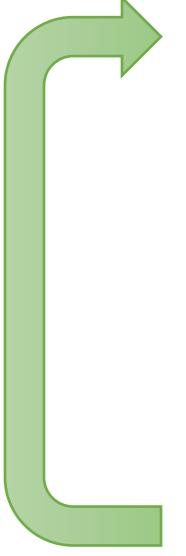
$$\pi'(a_t|s_t) = \begin{cases} 1 & \text{if } a_t = \arg \max_a Q^\pi(s_t, a) \\ 0 & \text{otherwise} \end{cases}$$

as good as  $\pi$   
(probably better)



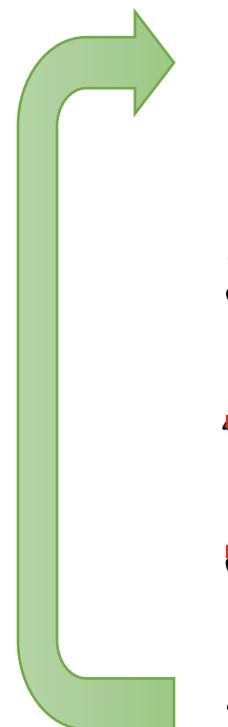
# From actor-critic to critic only

~~Off policy actor critic with replay buffer~~ Q-learning

- 
1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$  where  $\mathbf{a}'_i \sim \pi(\cdot|\mathbf{s}'_i)$
  4.  ~~$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$  where  $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a} | \mathbf{s}_i)$~~
  5.  ~~$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$~~
4. define new policy  $\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$

# From actor-critic to critic only

~~Off policy actor critic with replay buffer~~



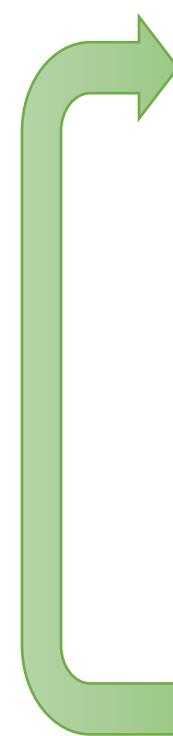
1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
3. update  $\hat{Q}_\phi^\pi$  using targets  $y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$  where  $\mathbf{a}'_i \sim \pi(\cdot|\mathbf{s}'_i)$       <- policy evaluation
4.  ~~$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$  where  $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a}|\mathbf{s}_i)$~~       (**Note:** can do multiple gradient steps here)
5.  ~~$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$~~
4. define new policy  $\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$       <- policy improvement

“Policy iteration”

Can we improve the policy in the Q-function update?

# From actor-critic to critic only

~~Off policy actor critic with replay buffer~~ Q-learning

- 
1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi^\pi$  using targets  ~~$y_i = r_i + \gamma \hat{Q}_\phi^\pi(\mathbf{s}'_i, \mathbf{a}'_i)$  where  $\mathbf{a}'_i \sim \pi(\cdot|\mathbf{s}'_i)$~~   $y_i = r_i + \gamma \max_{\mathbf{a}'} \hat{Q}_\phi(\mathbf{s}'_i, \mathbf{a}')$
  4.  ~~$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i^\pi | \mathbf{s}_i) \hat{Q}^\pi(\mathbf{s}_i, \mathbf{a}_i^\pi)$  where  $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a} | \mathbf{s}_i)$~~
  5.  ~~$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$~~
  4. define new policy  $\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$

Q-values for new policy!

# From actor-critic to critic only

## Q-learning

1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
3. update  $\hat{Q}_\phi$  using targets  $y_i = r_i + \gamma \max_{\mathbf{a}'} \hat{Q}_\phi(\mathbf{s}'_i, \mathbf{a}')$
4. define new policy  $\pi(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$

**Why does make sense?**

Recall:  $\underline{Q^\pi(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\cdot|\mathbf{s}, \mathbf{a}), \bar{\mathbf{a}}' \sim \pi(\cdot|\mathbf{s}')} [Q^\pi(\mathbf{s}', \bar{\mathbf{a}}')]} \text{ for all } (\mathbf{s}, \mathbf{a})$

-> This holds for any policy  $\pi$  (including the optimal policy  $\pi^*$ )

If  $\pi^*$  is the optimal policy, we also get:

$$\underline{Q^{\pi^*}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\cdot|\mathbf{s}, \mathbf{a})} \left[ \max_{\bar{\mathbf{a}}'} Q^{\pi^*}(\mathbf{s}', \bar{\mathbf{a}}') \right]} \text{ for all } (\mathbf{s}, \mathbf{a})$$

the action  $\pi^*$  would take in state  $\mathbf{s}'$

When we optimize step 3, we are trying to make this equation hold!

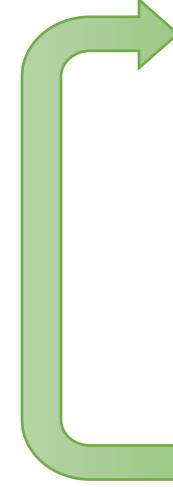
**Terminology**

**Bellman equation**

**Bellman optimality equation**

# From actor-critic to critic only

## Q-learning

- 
1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi$  using targets  $y_i = r_i + \gamma \max_{\mathbf{a}'} \hat{Q}_\phi(\mathbf{s}'_i, \mathbf{a}')$
  4. define new policy  $\pi(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$

**Note:** Q-learning is **off-policy**

## Why does make sense?

If  $\pi^*$  is the optimal policy, we also get:

$$Q^{\pi^*}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{s}' \sim p(\cdot | \mathbf{s}, \mathbf{a})} \left[ \max_{\bar{\mathbf{a}}'} Q^{\pi^*}(\mathbf{s}', \bar{\mathbf{a}}') \right] \text{ for all } (\mathbf{s}, \mathbf{a})$$

When we optimize step 3, we are trying to make this equation hold!

## Will this algorithm converge to optimal $Q^{\pi^*}$ ?

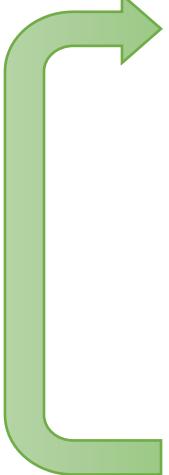
Yes, if you maintain a table of Q-values for every state and action. More generally, no. 😢

Can construct scenarios where it diverges, even with linear Q. *But*, it can be made to work well

# What kind of data should we collect?

Q-learning

from some exploration policy

- 
1. take action  $\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$ , store in  $\mathcal{R}$
  2. sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i\}$  from buffer  $\mathcal{R}$
  3. update  $\hat{Q}_\phi$  using targets  $y_i = r_i + \gamma \max_{\mathbf{a}'} \hat{Q}_\phi(\mathbf{s}'_i, \mathbf{a}')$
  4. define new policy  $\pi(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} \hat{Q}_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$

Note: Q-learning is off-policy

want coverage for many actions  $\mathbf{a}$

## What are good choices for data collection?

1. With probability  $\epsilon$ , take uniformly random action

$$\pi(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 - \epsilon & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ \epsilon / (|\mathcal{A}| - 1) & \text{otherwise} \end{cases} \quad \text{"epsilon-greedy"}$$

Often start with larger  $\epsilon$ , decrease over training.

2. Take actions with probability proportional to their Q-value.

$$\pi(\mathbf{a}_t|\mathbf{s}_t) \propto \exp(Q_\phi(\mathbf{s}_t, \mathbf{a}_t))$$

*"Boltzmann exploration"*

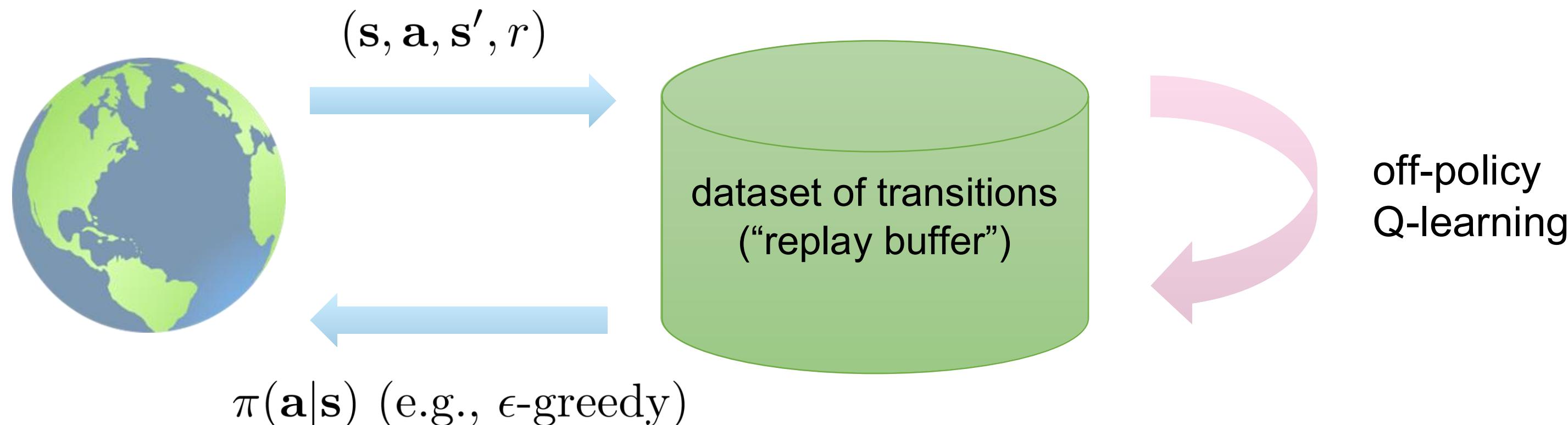
# Putting it together

full Q-learning with replay buffer:

1. collect data  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy, add it to  $\mathcal{R}$
2. sample a batch  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$  from  $\mathcal{R}$
3.  $\phi \leftarrow \phi - \alpha \sum_i \frac{dQ_\phi}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_\phi(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_\phi(\mathbf{s}'_i, \mathbf{a}')])$

$K = 1$  is common, though  
larger  $K$  more efficient

result:  $Q_\phi$       final policy  $\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \arg \max_{\mathbf{a}} Q_\phi(\mathbf{s}_t, \mathbf{a}) \\ 0 & \text{otherwise} \end{cases}$



# The plan for today

## Value-based RL methods

1. Q-learning RL method
  - a. Policy iteration
  - b. Bellman optimality equation
  - c. How to collect data for Q-learning methods
2. **Q-learning in practice**
  - d. Target networks
  - e. Double DQN
  - f. N-step returns

# How to make Q-learning stable?

full Q-learning with replay buffer:

- 
1. collect data  $\{(s_i, a_i, s'_i, r_i)\}$  using some policy, add it to  $\mathcal{R}$
  2. sample a batch  $(s_i, a_i, s'_i, r_i)$  from  $\mathcal{R}$
  3.  $\phi \leftarrow \phi - \alpha \sum_i \frac{dQ_\phi}{d\phi}(s_i, a_i)(Q_\phi(s_i, a_i) - [r(s_i, a_i) + \gamma \max_{a'} Q_\phi(s'_i, a')])$

this is a moving target

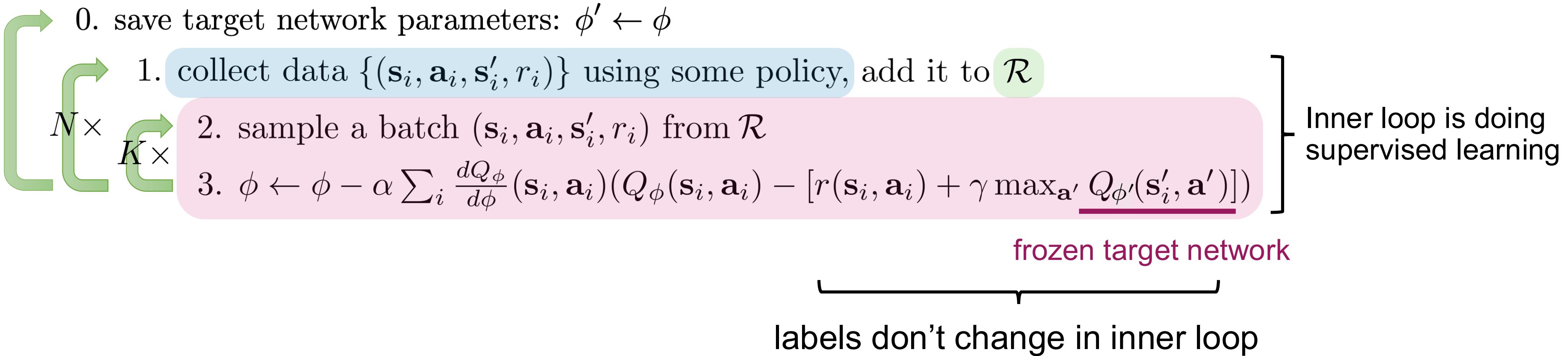
*can lead to unstable optimization*

Can we change the target Q-values more slowly?

💡 Simple idea: freeze parameters used for the target Q-values, update periodically

# How to make Q-learning stable?

Q-learning with replay buffer and target network:

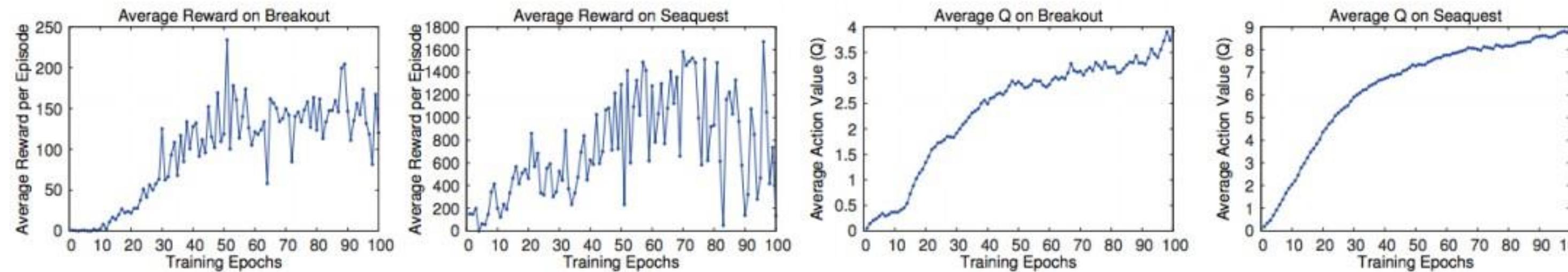


Corresponds to DQN algorithm (“deep Q network”)

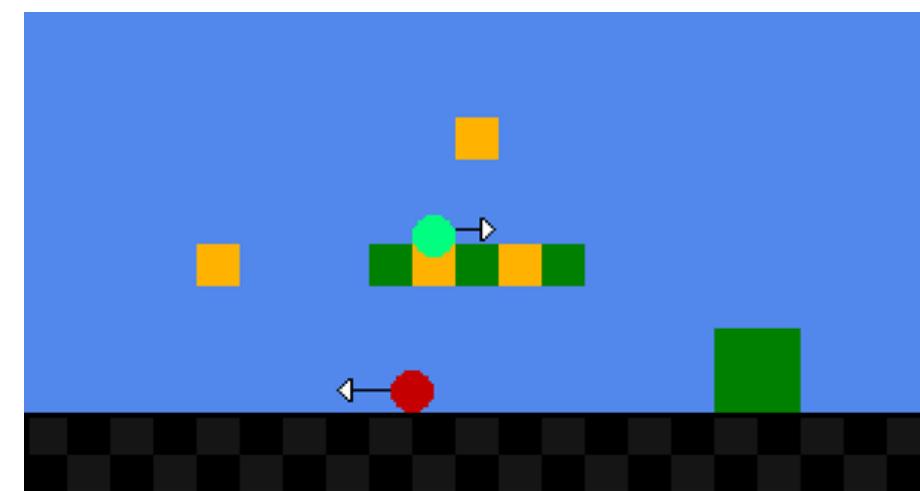
*Mnih et al. Playing Atari with Deep Reinforcement Learning. 2013.*

Along with an actor-critic method,  
you'll implement this algorithm in HW 2!

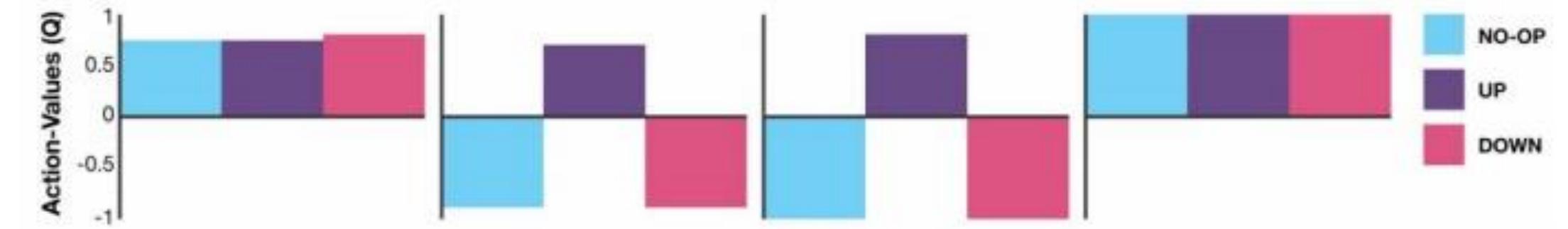
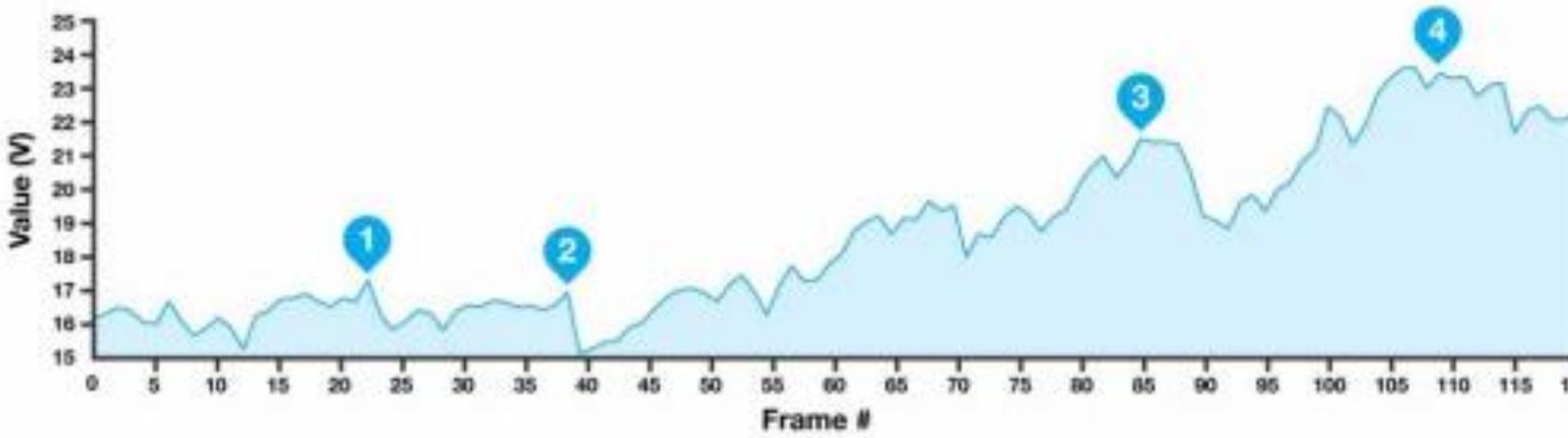
# Are the Q-values accurate?



As predicted Q increases,  
so does the return

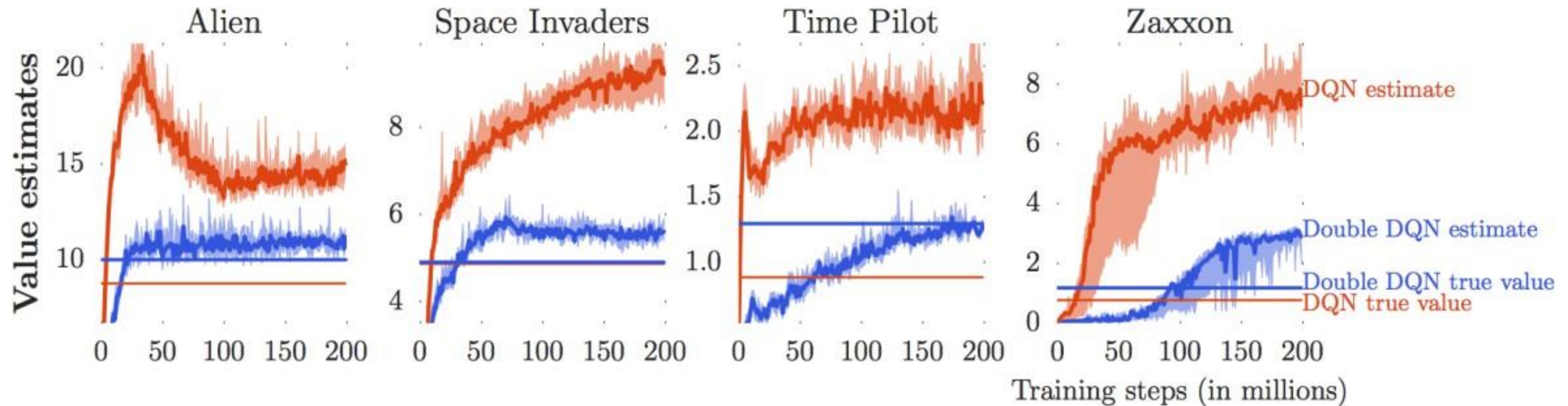


Game of Pong  
<https://www.youtube.com/watch?v=jjFCgAcM2tU>



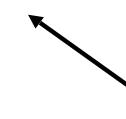
Atari Logo Image: Atari, Inc., Public domain, via  
Wikimedia Commons

# Are the Q-values accurate?



# Overestimation in Q-learning

target value  $y_j = r_j + \gamma \max_{\mathbf{a}'_j} Q_{\phi'}(\mathbf{s}'_j, \mathbf{a}'_j)$



this last term is the problem

$Q_{\phi'}(\mathbf{s}', \mathbf{a}')$  is not perfect – it looks “noisy”

hence  $\max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}')$  overestimates the next value!

note that  $\max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}') = \underline{Q_{\phi'}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}'))}$

value also comes from  $Q_{\phi'}$  action selected according to  $Q_{\phi'}$

# Double Q-learning

note that  $\max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}') = \underline{Q_{\phi'}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}'))}$

value *also* comes from  $Q_{\phi'}$  action selected according to  $Q_{\phi'}$

  
if the noise in these is decorrelated, the problem goes away!

idea: don't use the same network to choose the action and evaluate value!

“double” Q-learning: use two networks:

$$Q_{\phi_A}(\mathbf{s}, \mathbf{a}) \leftarrow r + \gamma Q_{\phi_B}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi_A}(\mathbf{s}', \mathbf{a}'))$$

$$Q_{\phi_B}(\mathbf{s}, \mathbf{a}) \leftarrow r + \gamma Q_{\phi_A}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi_B}(\mathbf{s}', \mathbf{a}'))$$

  
if the two Q's are noisy in *different* ways, there is no problem

# Double Q-learning in practice

where to get two Q-functions?

just use the current and target networks!

standard Q-learning:  $y = r + \gamma Q_{\phi'}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}', \mathbf{a}'))$

double Q-learning:  $y = r + \gamma Q_{\phi'}(\mathbf{s}', \arg \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}', \mathbf{a}'))$

just use current network (not target network) to evaluate action

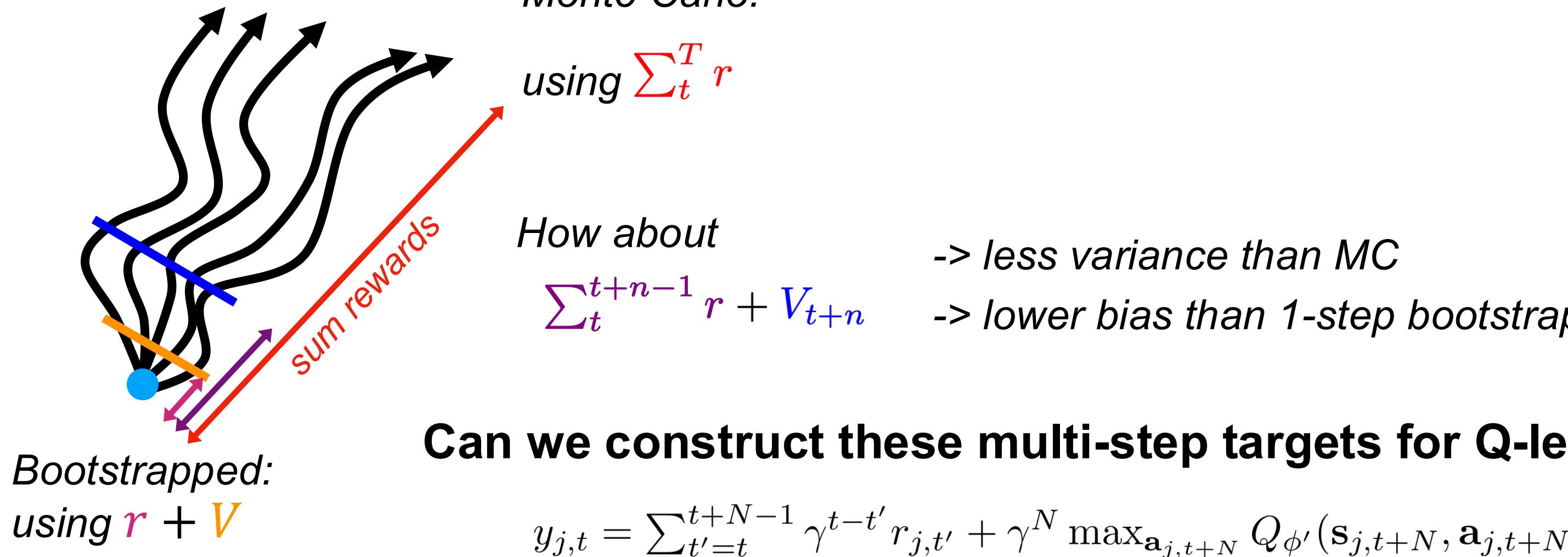
still use target network to evaluate value!

# N-step returns?

Q-learning target:  $y_{j,t} = r_{j,t} + \gamma \max_{\mathbf{a}_{j,t+1}} Q_{\phi'}(\mathbf{s}_{j,t+1}, \mathbf{a}_{j,t+1})$

these are the only values that matter if  $Q_{\phi'}$  is bad!      these values are important if  $Q_{\phi'}$  is good

**Recall:** fitting value functions with rewards + bootstrapped  $V$



# N-step returns

Q-learning target:  $y_{j,t} = r_{j,t} + \gamma \max_{\mathbf{a}_{j,t+1}} Q_{\phi'}(\mathbf{s}_{j,t+1}, \mathbf{a}_{j,t+1})$

these are the only values that matter if  $Q_{\phi'}$  is bad!      these values are important if  $Q_{\phi'}$  is good

**Can we construct these multi-step targets for Q-learning?**

N-step target:  $y_{j,t} = \sum_{t'=t}^{t+N-1} \gamma^{t-t'} r_{j,t'} + \gamma^N \max_{\mathbf{a}_{j,t+N}} Q_{\phi'}(\mathbf{s}_{j,t+N}, \mathbf{a}_{j,t+N})$

these are the rewards for  
policy that collected the data

- + (far) less biased target values when Q-values are inaccurate
- + typically faster learning, especially early on
- only actually correct when learning on-policy  
(not an issue when N=1!)

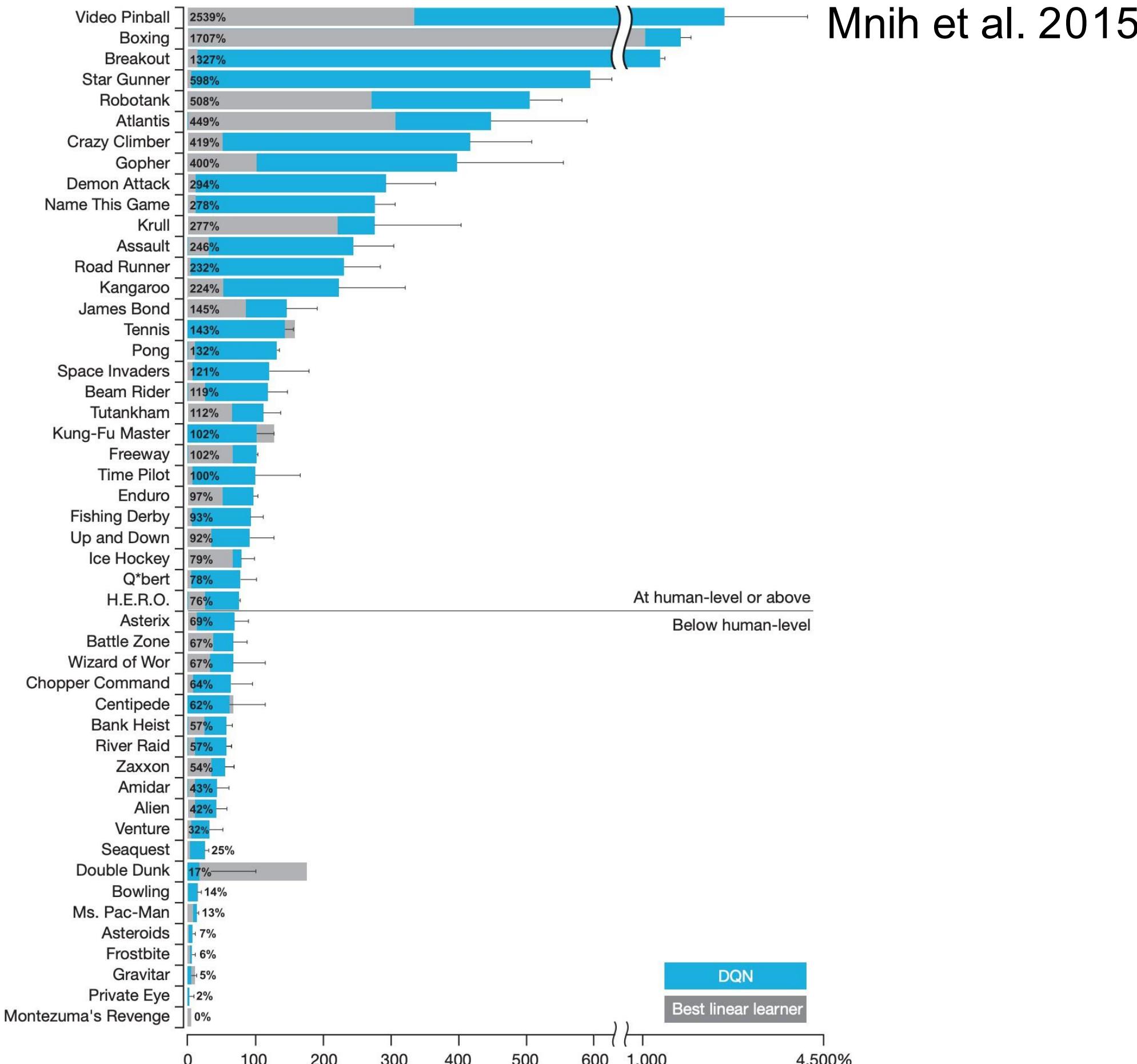
**Ways to fix?** Most commonly: ignore the problem & still use  $N > 1$

Can also:

- dynamically choose N to only use data that follows current policy  
(if data mostly on-policy, action space is small)
- use importance sampling

# Q-learning in practice?

Q-learning learns to outperform people on most Atari games



Mnih et al. 2015

Q-learning trains robot grasping system

Kalashnikov et al. 2018



| Method             | Dataset                             | Test       |
|--------------------|-------------------------------------|------------|
| QT-Opt (ours)      | 580k off-policy + 28k on-policy     | <b>96%</b> |
| Levine et al. [27] | 900k grasps from Levine et al. [27] | 78%        |
| QT-Opt (ours)      | 580k off-policy grasps only         | 87%        |
| Levine et al. [27] | 400k grasps from our dataset        | 67%        |

# When to use one online RL algorithm vs. another?

Chelsea's advice

**PPO & variants** When you care about stability, ease-of-use  
When you don't care about data efficiency

**DQN & variants** When you have discrete actions or low-dimensional continuous actions

**SAC & variants** When you care most about data efficiency  
When you are okay with tuning hyperparameters, less stability

# The plan for today

## Value-based RL methods

1. Q-learning RL method
  - a. Policy iteration
  - b. Bellman optimality equation
  - c. How to collect data for Q-learning methods
2. Q-learning in practice
  - d. Target networks
  - e. Double DQN
  - f. N-step returns

### Key learning goals:

- How Q-functions relate to policies
- How to do RL without learning an explicit policy
- How to stabilize Q-learning in practice

# Next time

Done with online RL methods!

# **CS 224R Tutorial**

## **Review of Q-Learning**

**Anikait Singh**

# Outline of Tutorial

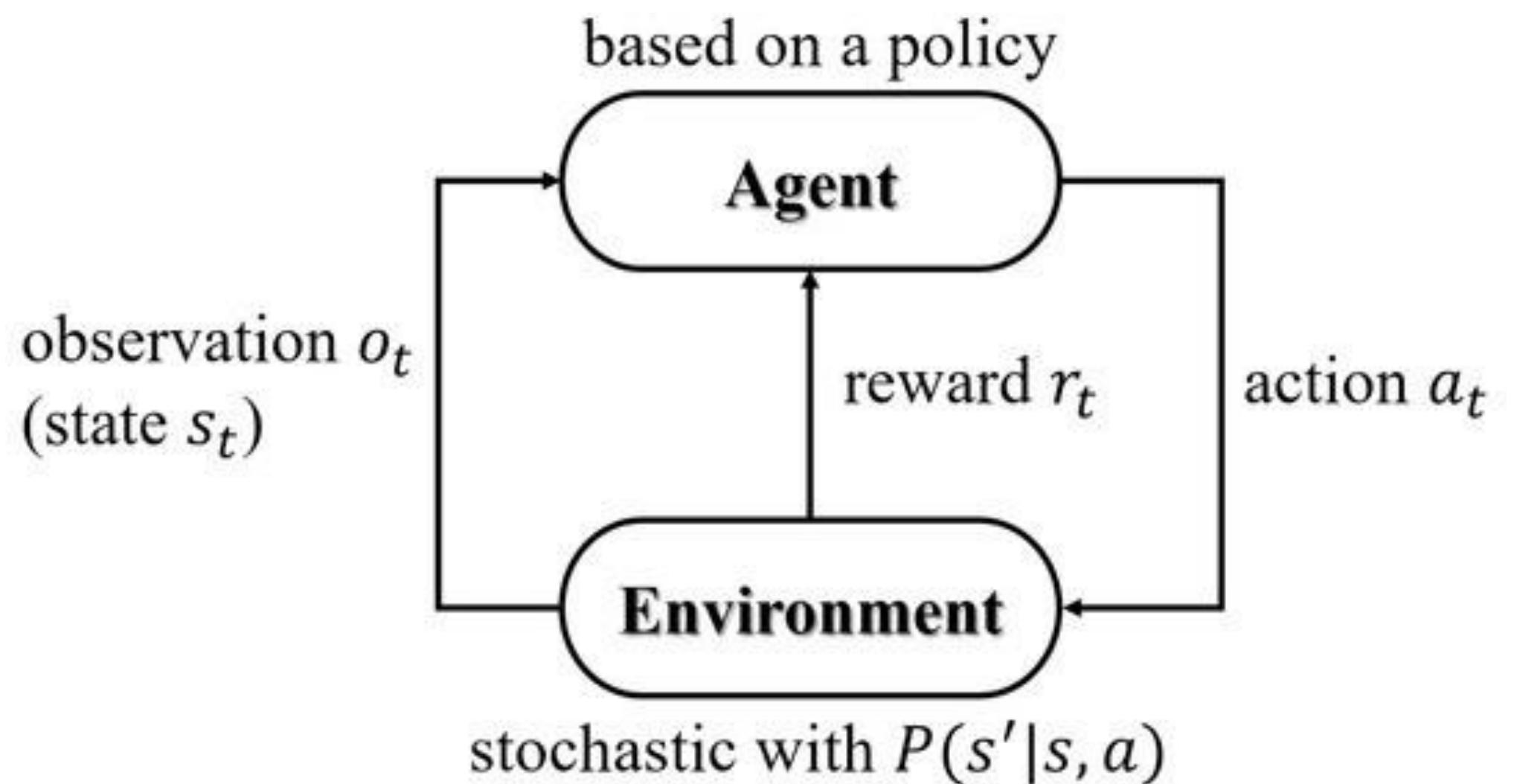
- Review of Markov Decision Processes (MDP)
- Exact MDPs
  - Fitted Q Iteration
- Parametric Q-Learning
  - Bias Variance Tradeoff (TD vs MC)
- Practical Details
  - Replay Buffer, Overestimation, TD Gradients
- Q-Learning/Actor Critic Algorithm Walkthrough

Many thanks to Pieter Abbeel, David Silver, Aviral Kumar, and Justin Fu!

# MDP Formulation

An MDP is defined by:

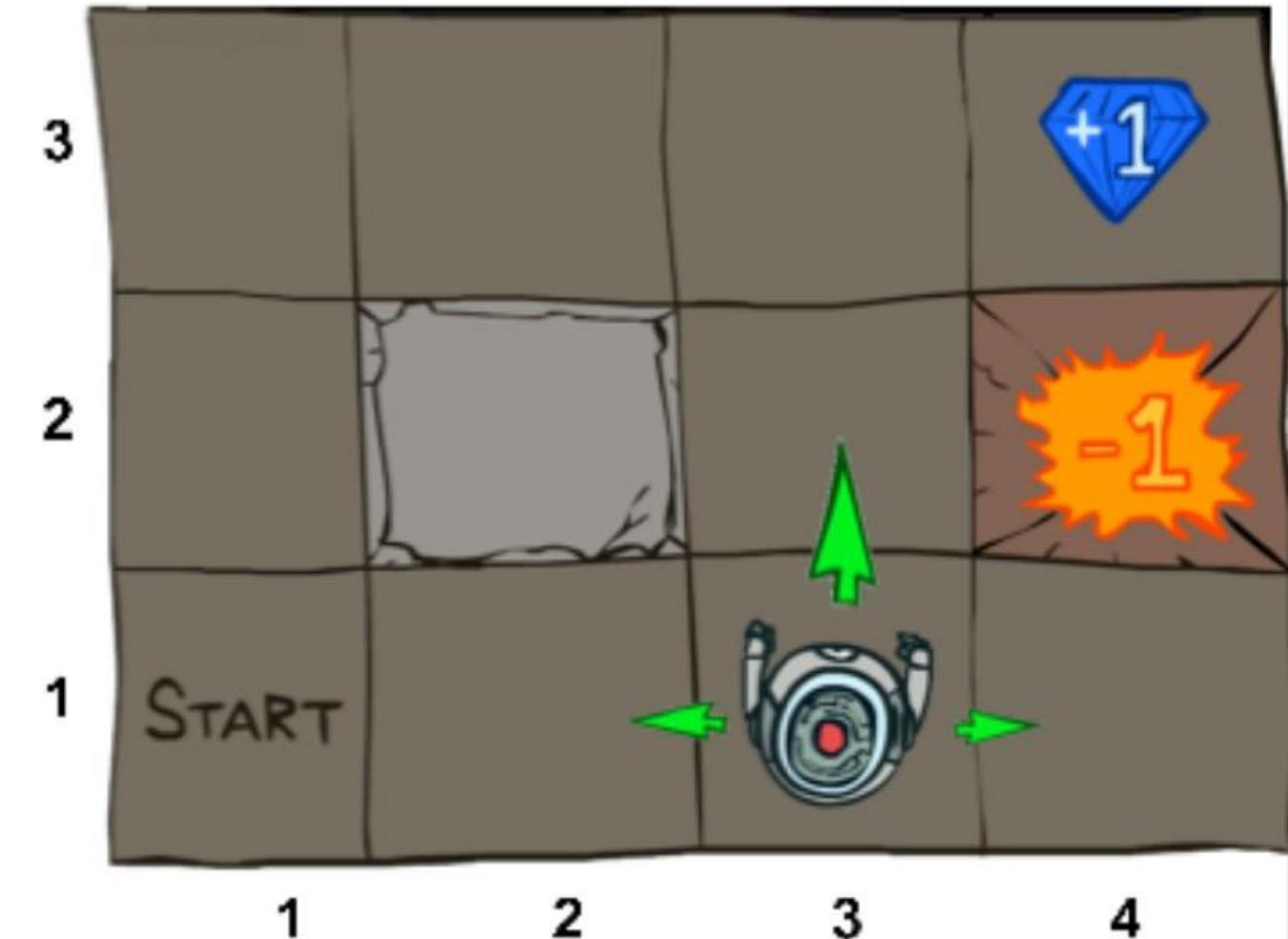
- Set of states  $s_t$
- Set of actions  $a_t$
- Transition function  $p(s_{t+1} | s_t, a_t)$
- Reward Function  $r(s_t, a_t)$
- Start state  $s_1$
- Discount Factor
- Horizon



# MDP Formulation

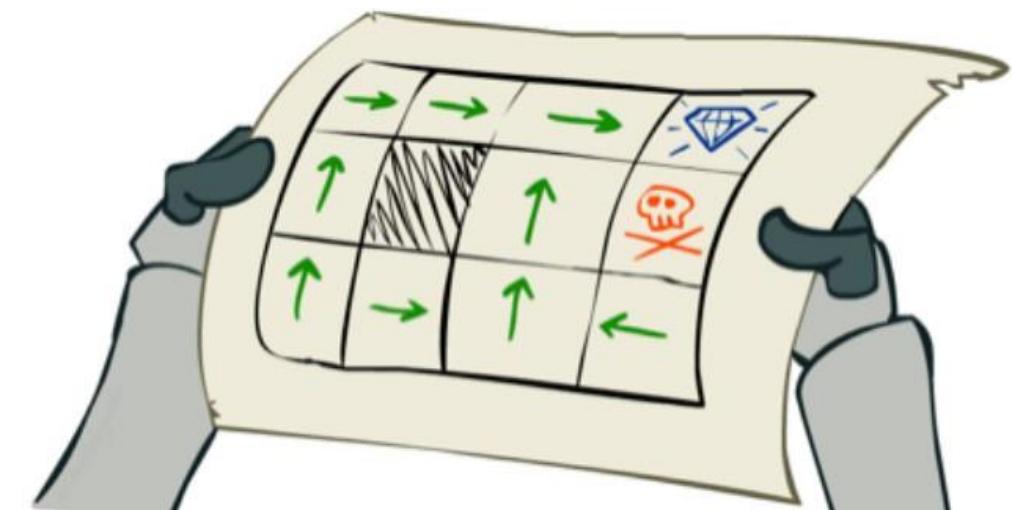
An MDP is defined by:

- Set of states  $s_t$
- Set of actions  $a_t$
- Transition function  $p(s_{t+1} | s_t, a_t)$
- Reward Function  $r(s_t, a_t)$
- Start state  $s_1$
- Discount Factor
- Horizon



Goal:  $\max_{\pi} \mathbb{E} \left[ \sum_{t=1}^T \gamma^t r(s_t, a_t) \mid \pi \right]$

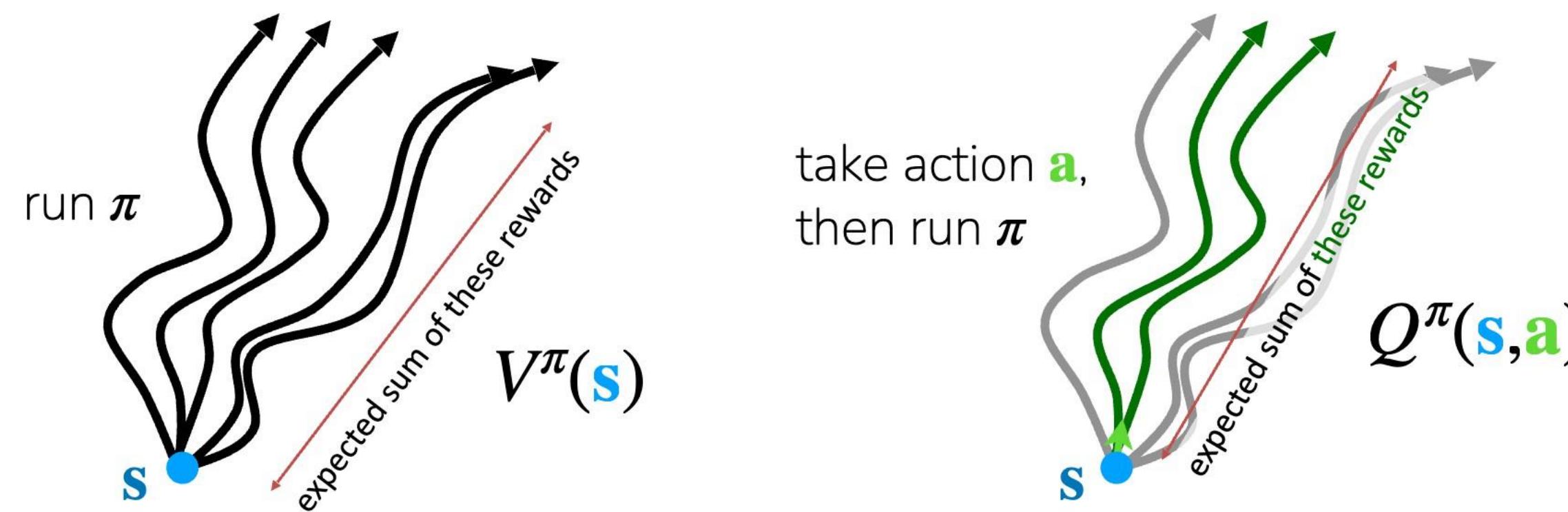
$\pi$ :



# Revisiting Some Useful Objects

value function  $V^\pi(\mathbf{s})$  - future expected rewards starting at  $\mathbf{s}$  and following  $\pi$

Q-function  $Q^\pi(\mathbf{s}, \mathbf{a})$  - future expected rewards starting at  $\mathbf{s}$ , taking  $\mathbf{a}$ , then following  $\pi$



$$\text{Useful relation: } V^\pi(\mathbf{s}) = \mathbb{E}_{\mathbf{a} \sim \pi(\cdot | \mathbf{s})} [Q^\pi(\mathbf{s}, \mathbf{a})]$$

advantage  $A^\pi(\mathbf{s}, \mathbf{a})$  - how much better it is to take  $\mathbf{a}$  than to follow policy  $\pi$  at state  $\mathbf{s}$

$$A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s})$$

# Tabular Learning - Fitted Q Iteration

Q-Function  $Q^\pi(s, a)$  - future expected rewards starting at state  $s$ , taking action  $a$ , then following  $\pi$

$$Q^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q^*(s_{t+1}, a_{t+1}) \right)$$

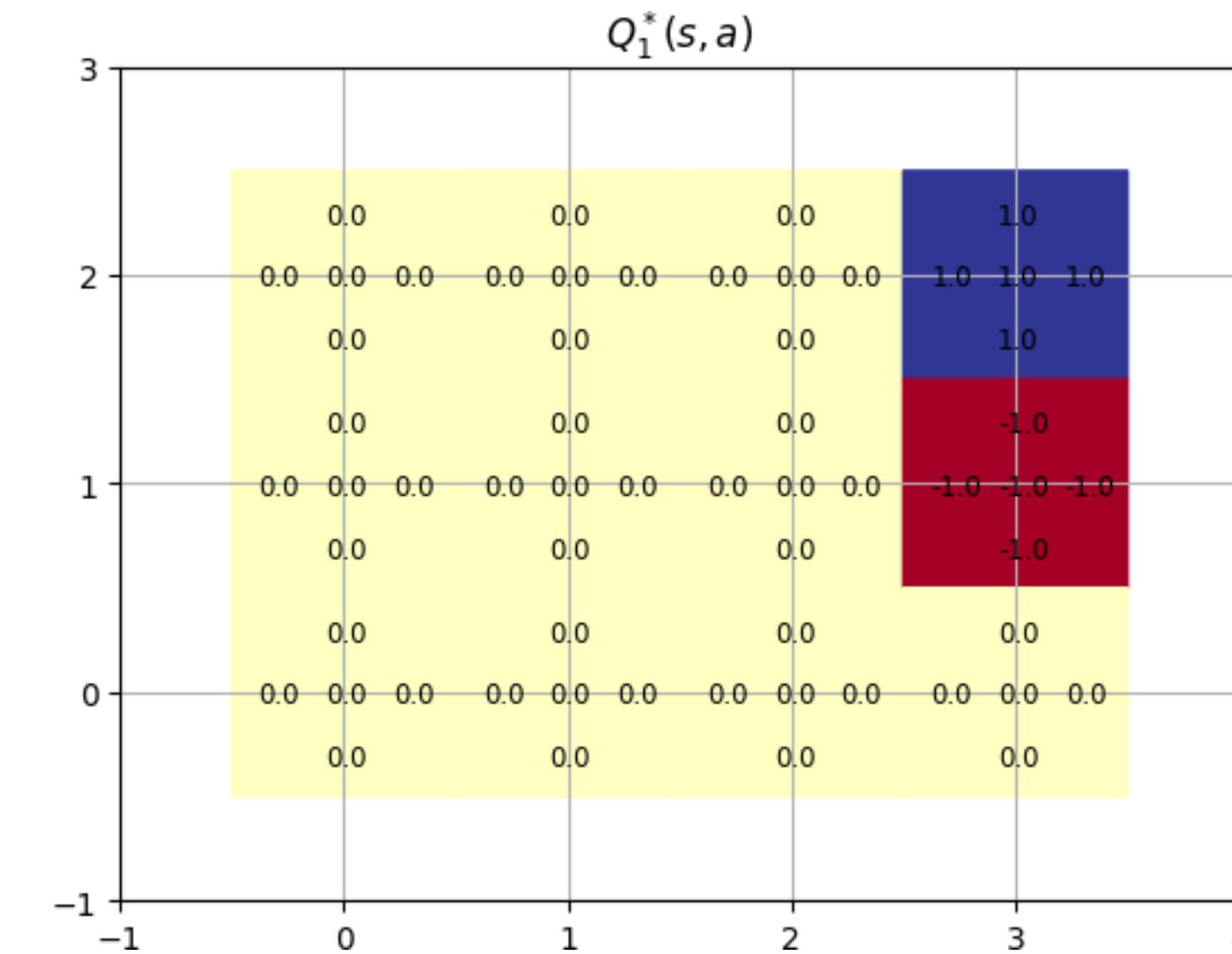
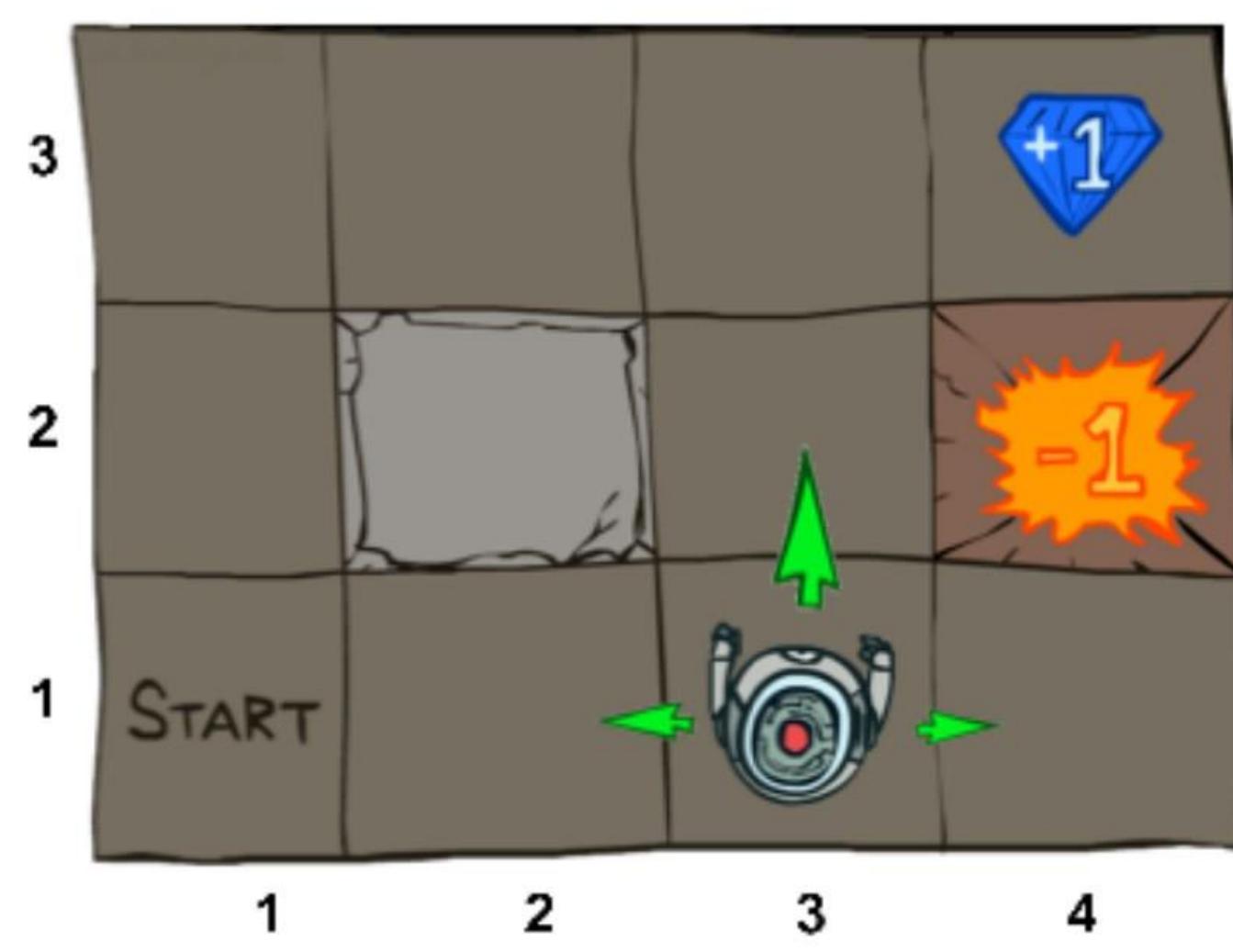
Q-Value Iteration

$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$

# Fitted Q Iteration with the MDP

Discount = 0.9

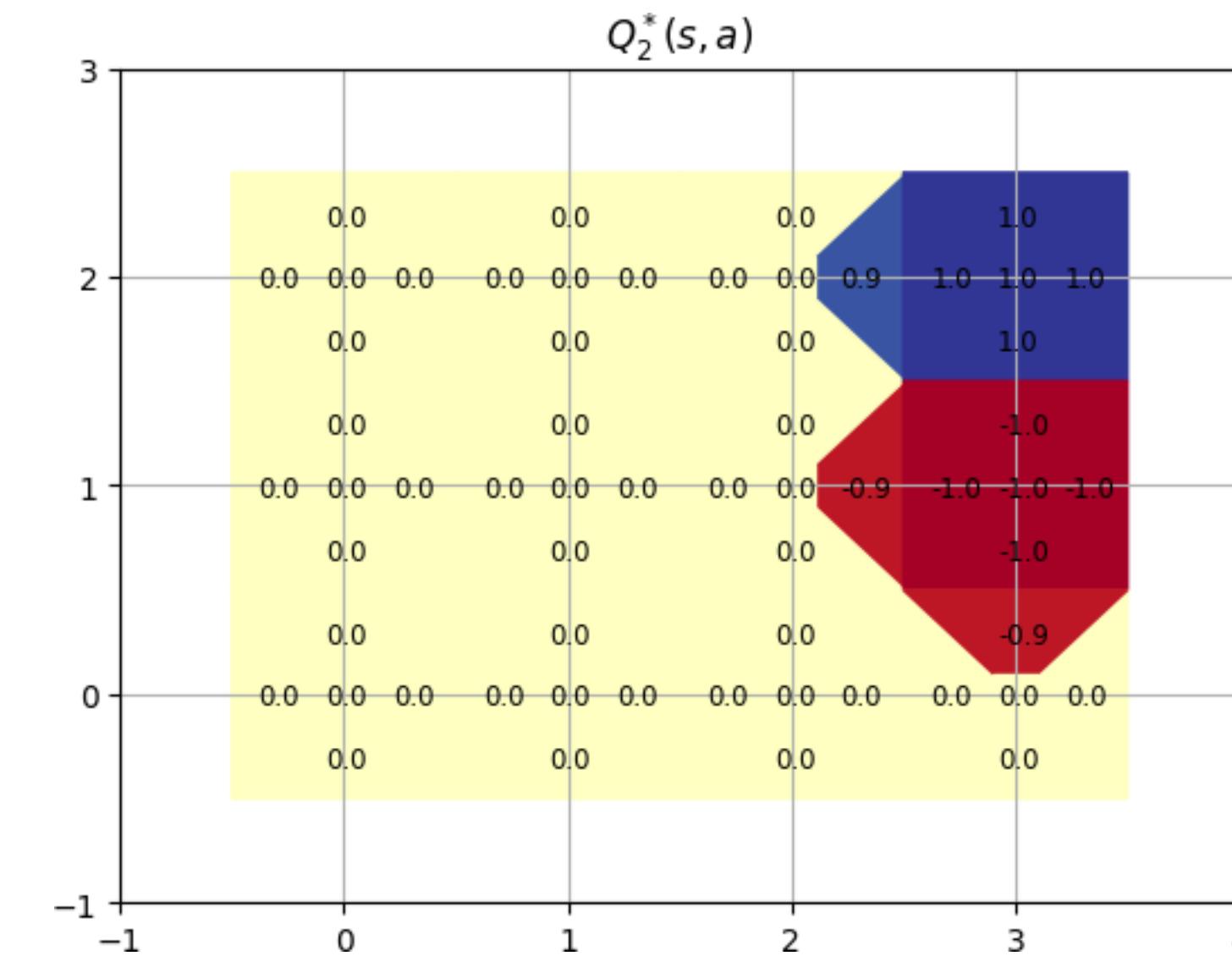
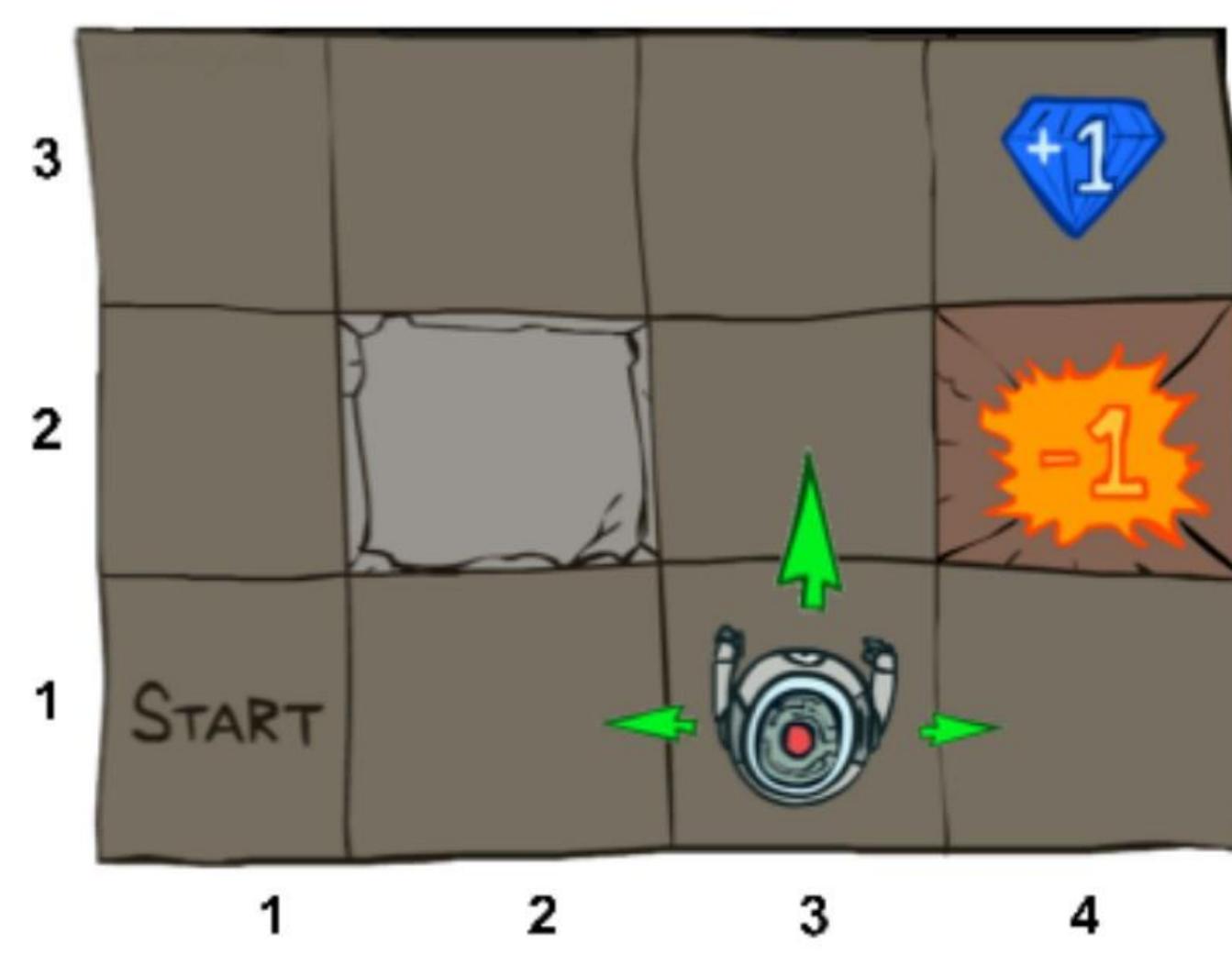
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



# Fitted Q Iteration with the MDP

# Discount = 0.9

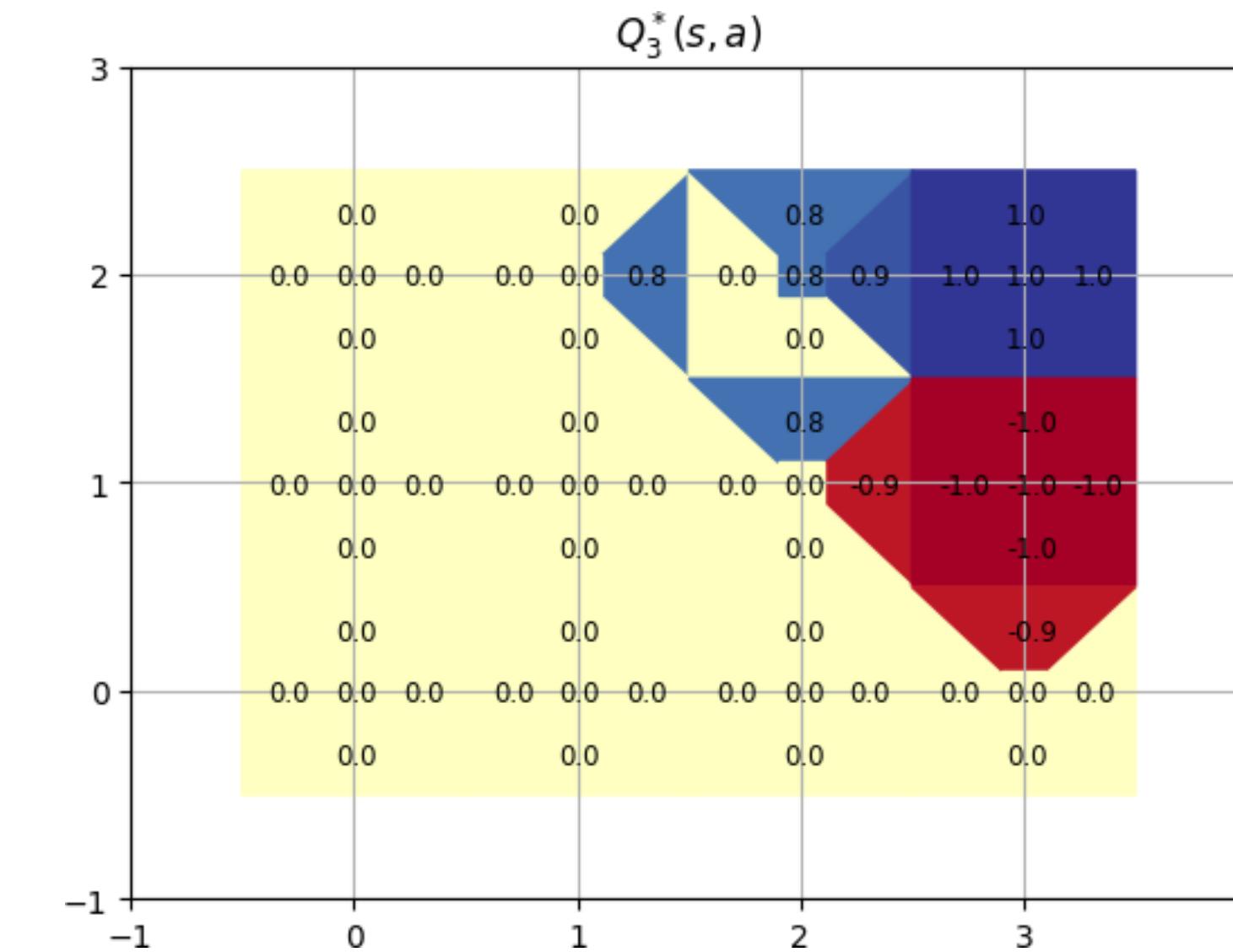
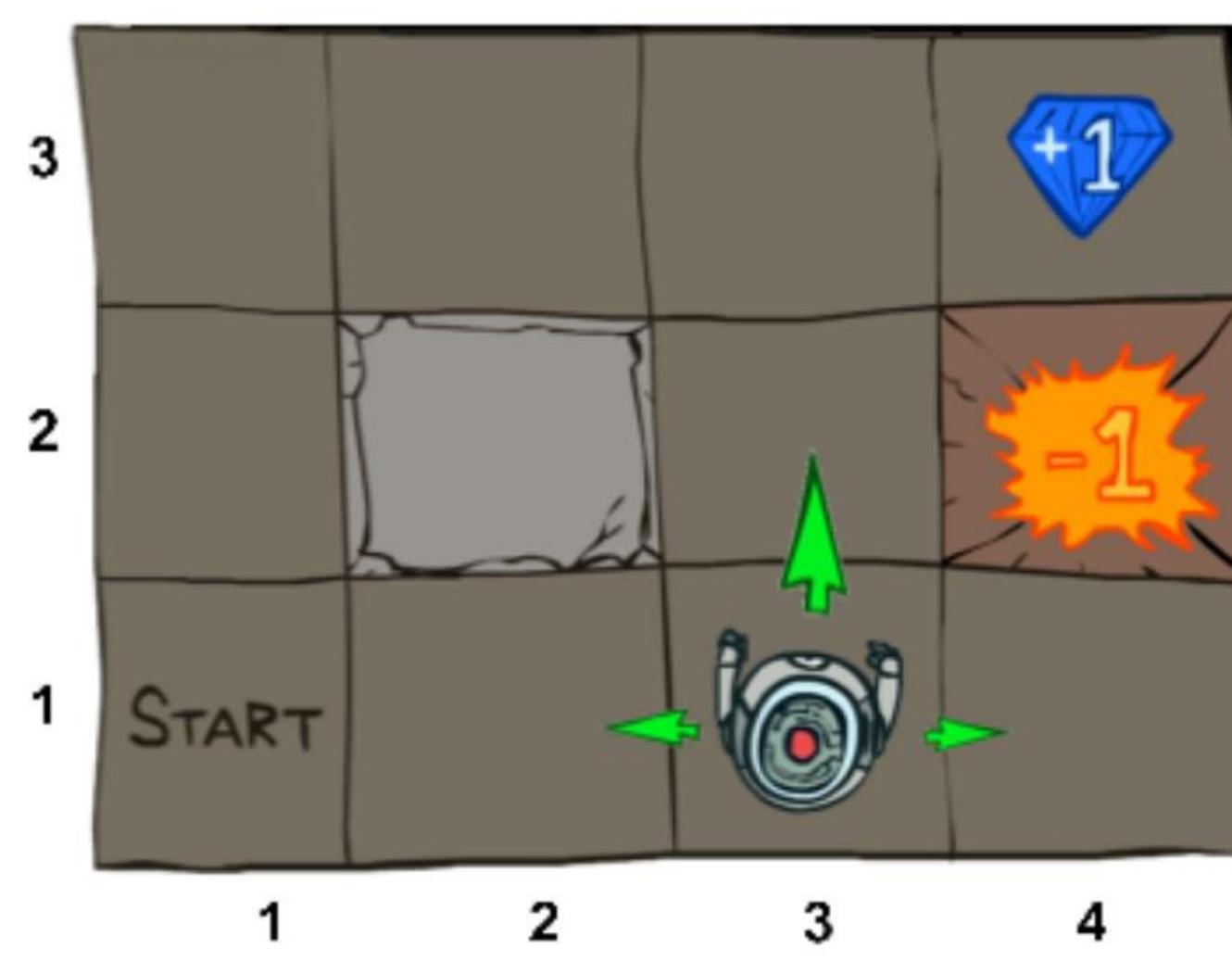
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



# Fitted Q Iteration with the MDP

Discount = 0.9

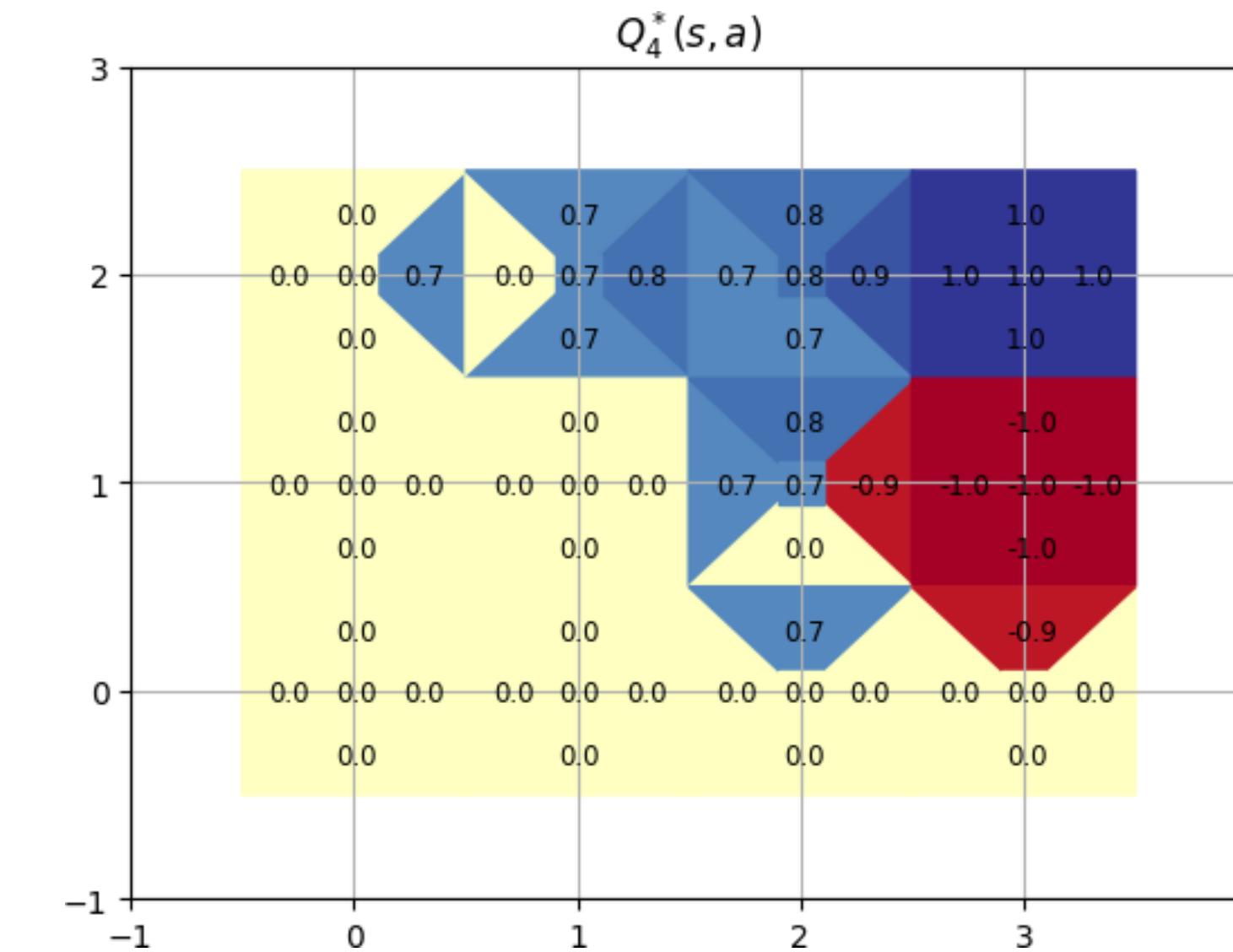
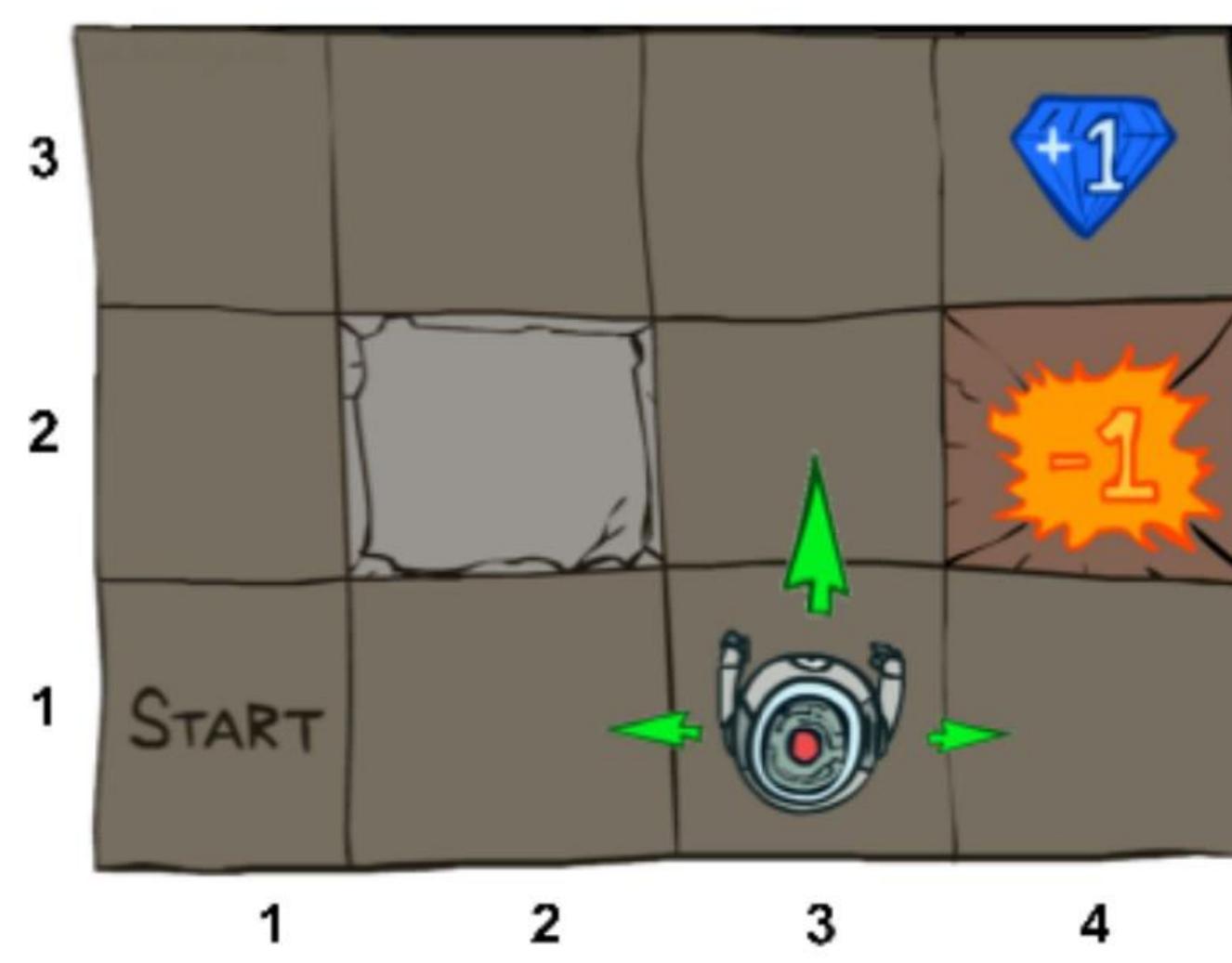
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



# Fitted Q Iteration with the MDP

Discount = 0.9

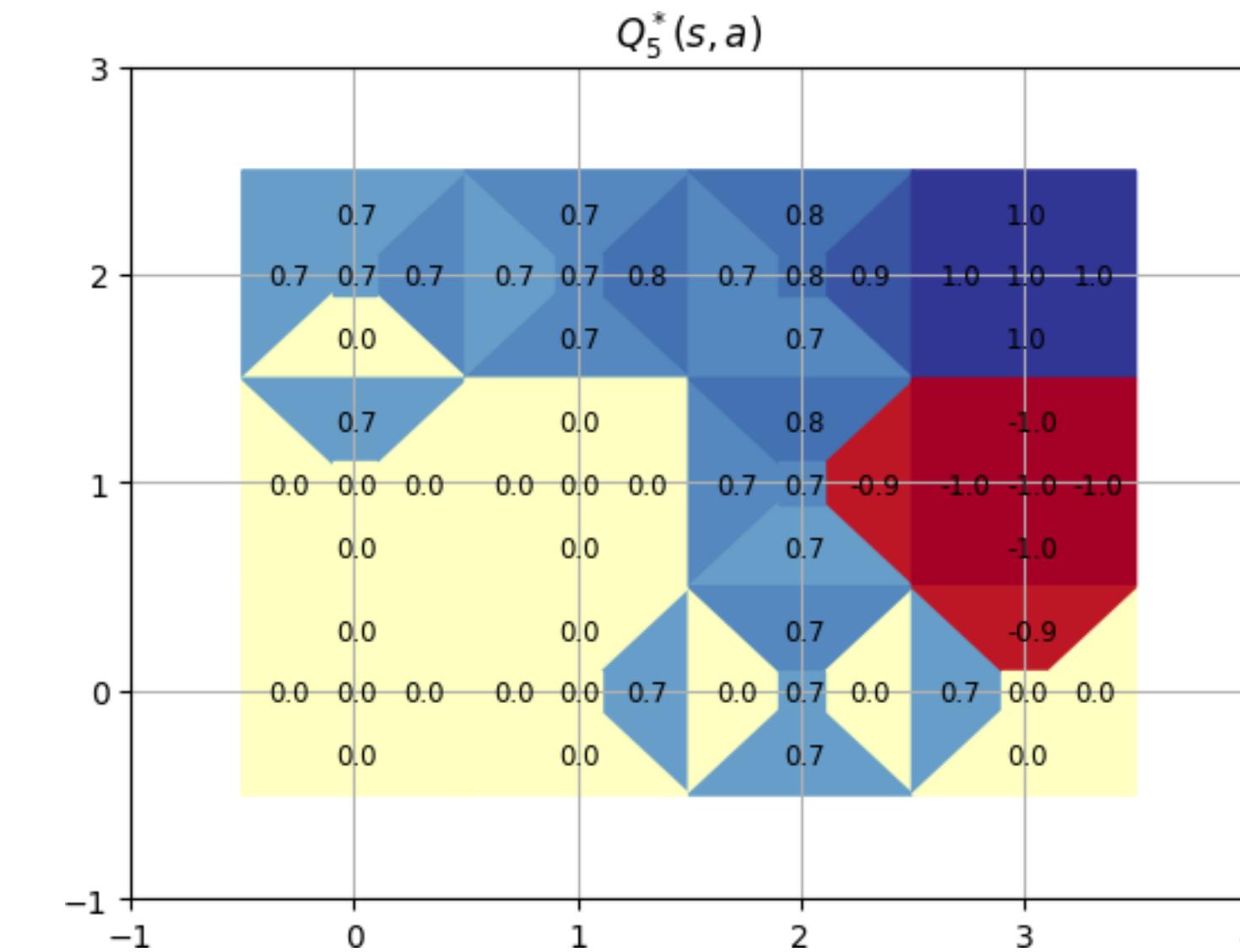
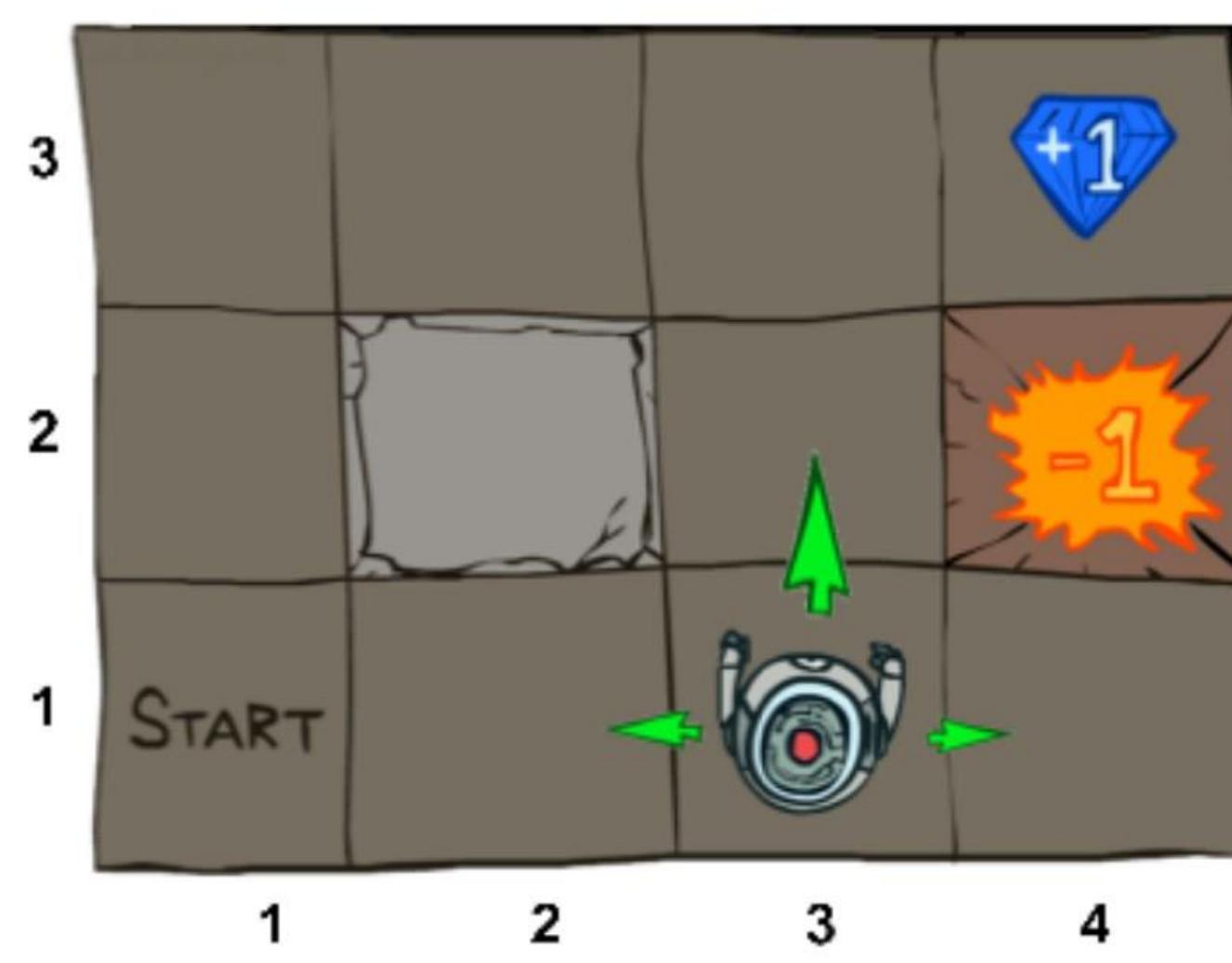
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



# Fitted Q Iteration with the MDP

Discount = 0.9

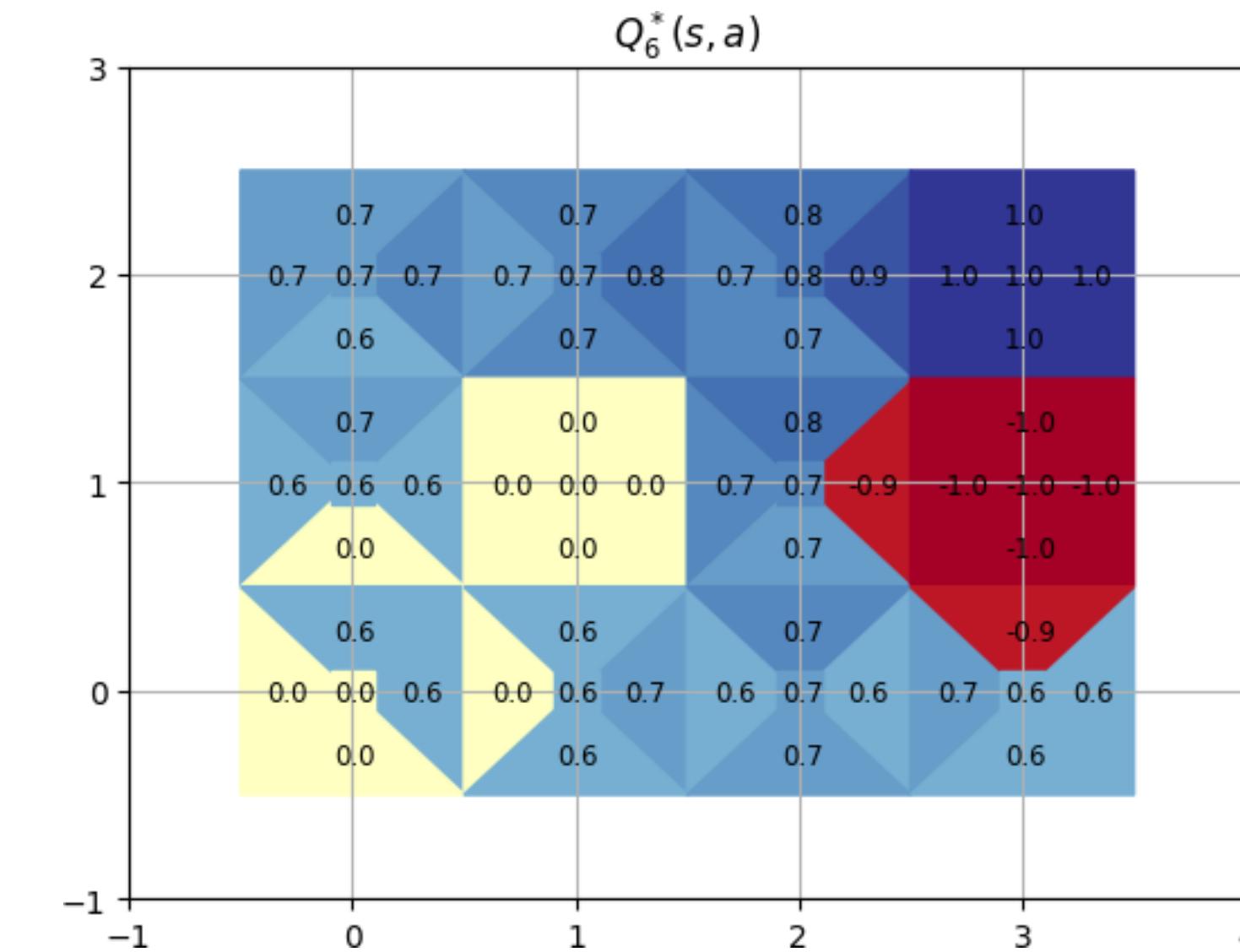
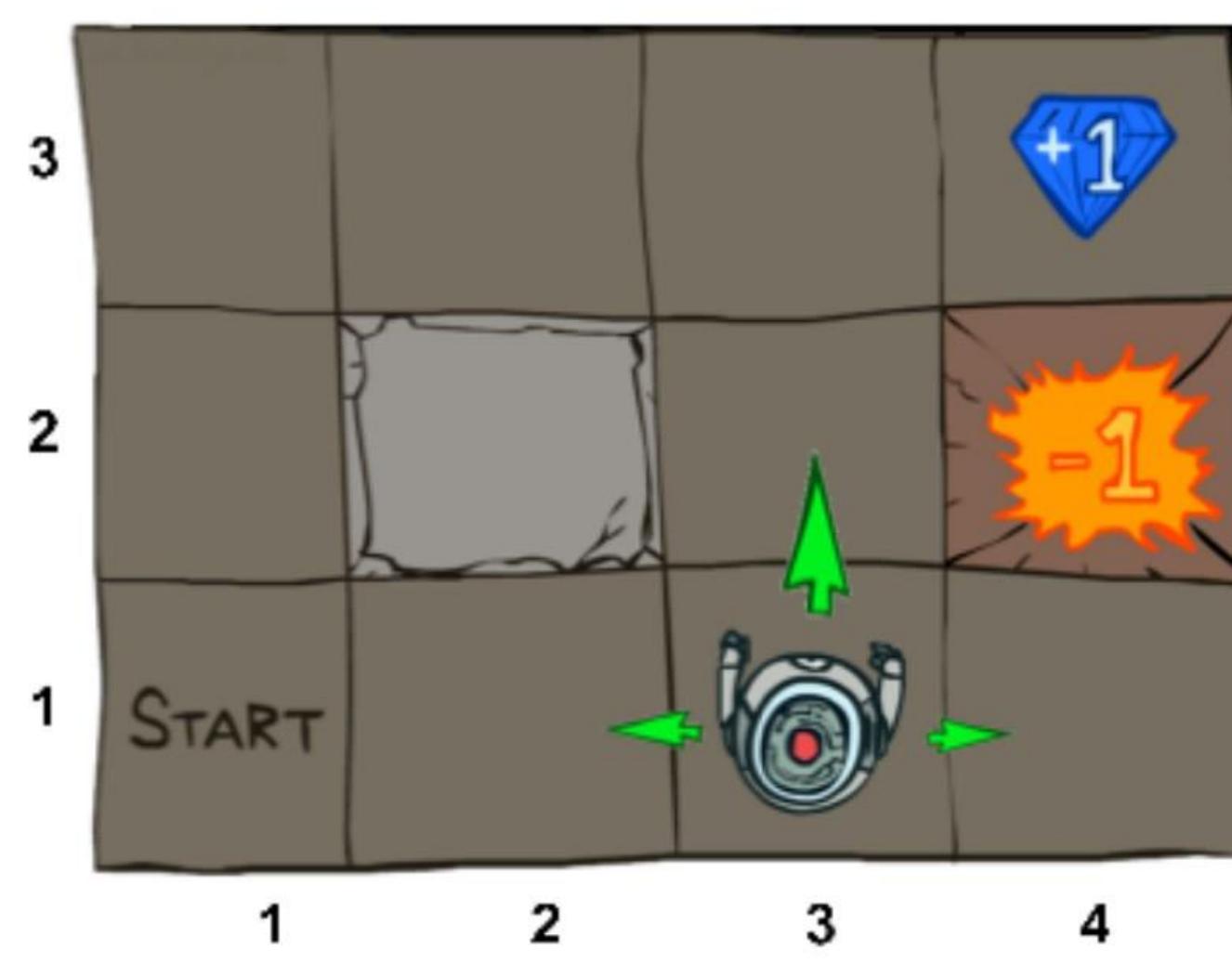
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



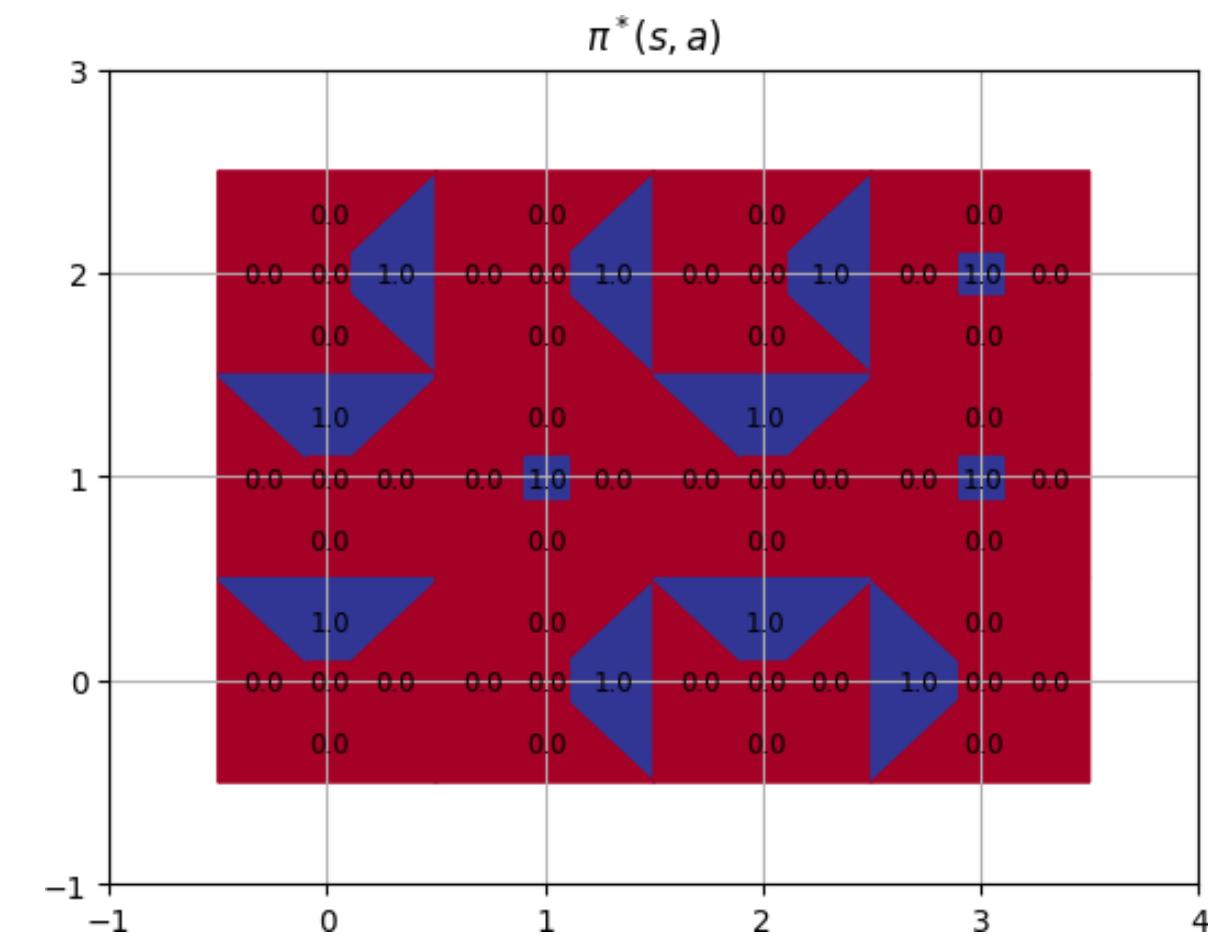
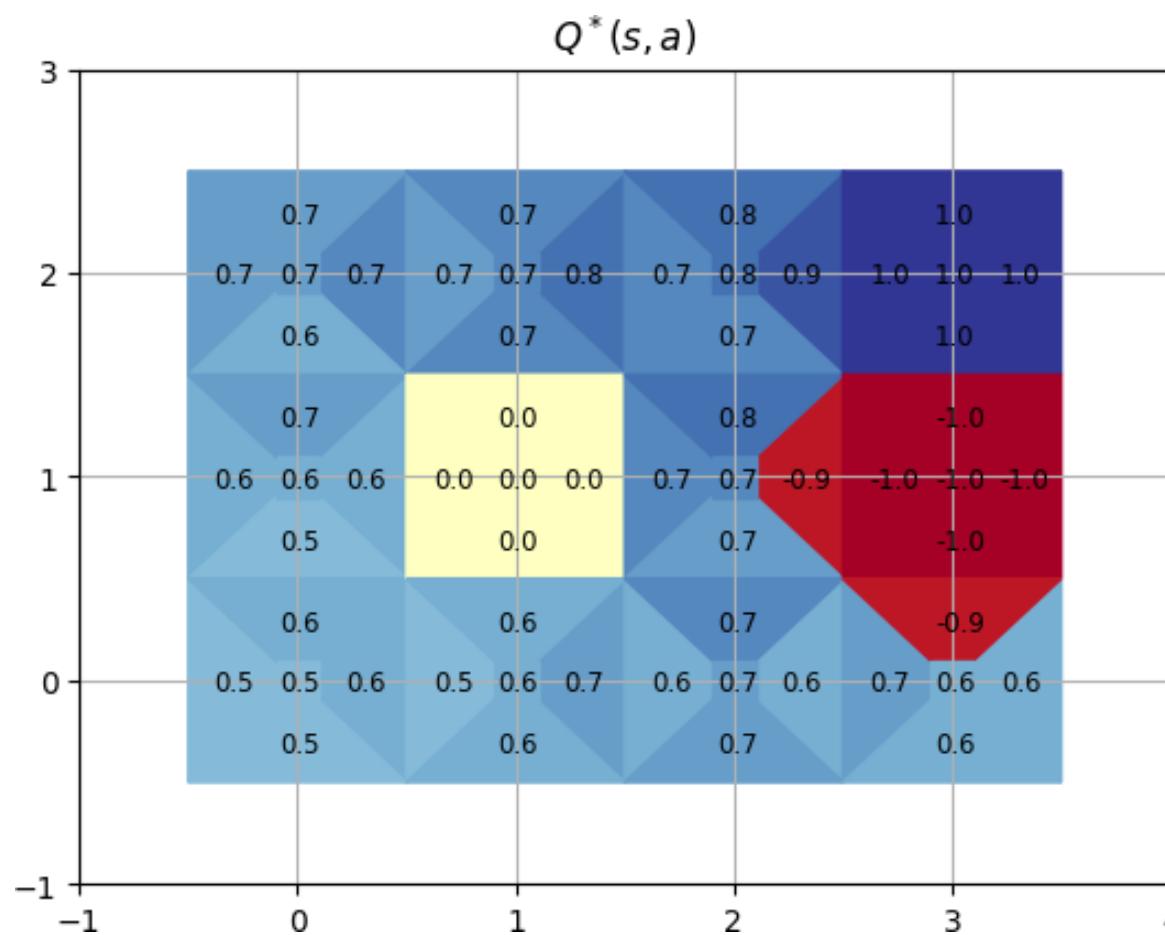
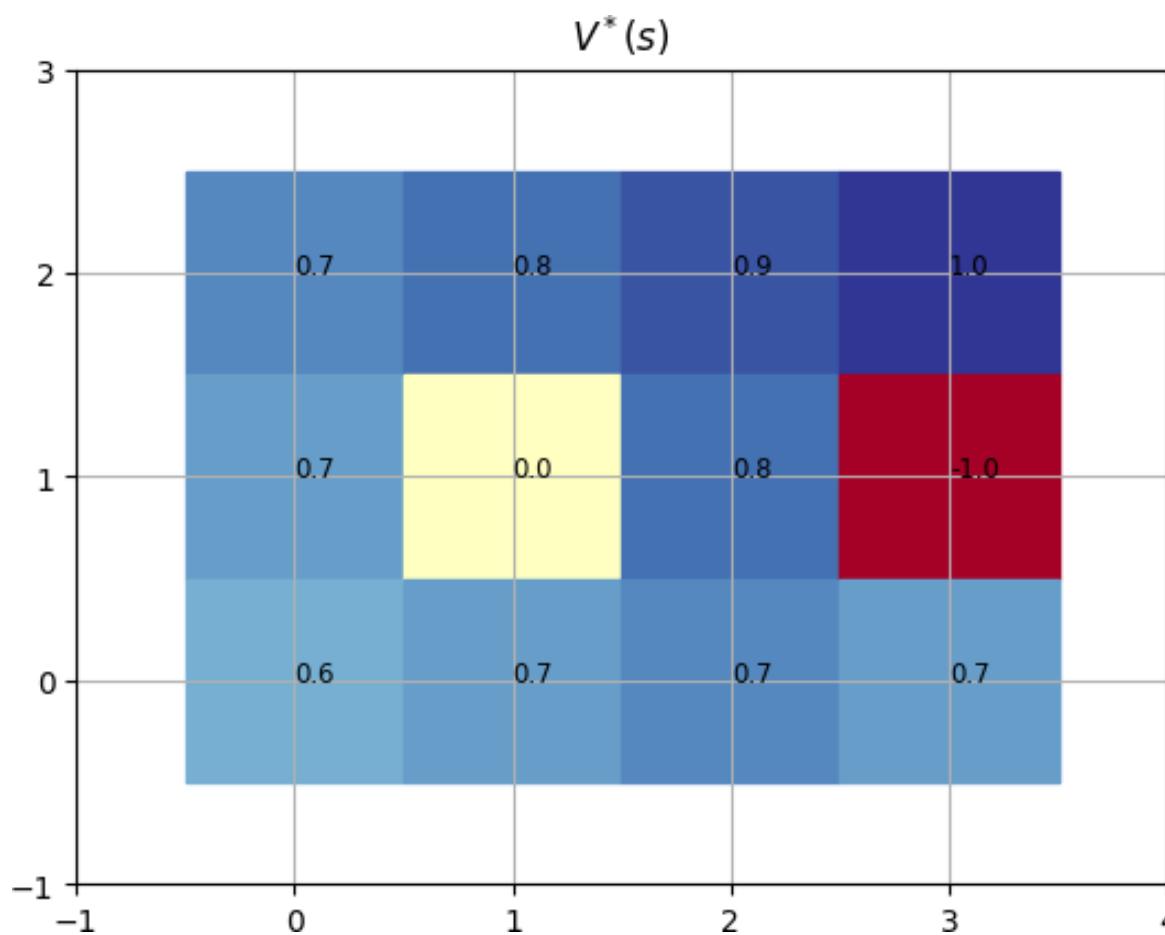
# Fitted Q Iteration with the MDP

Discount = 0.9

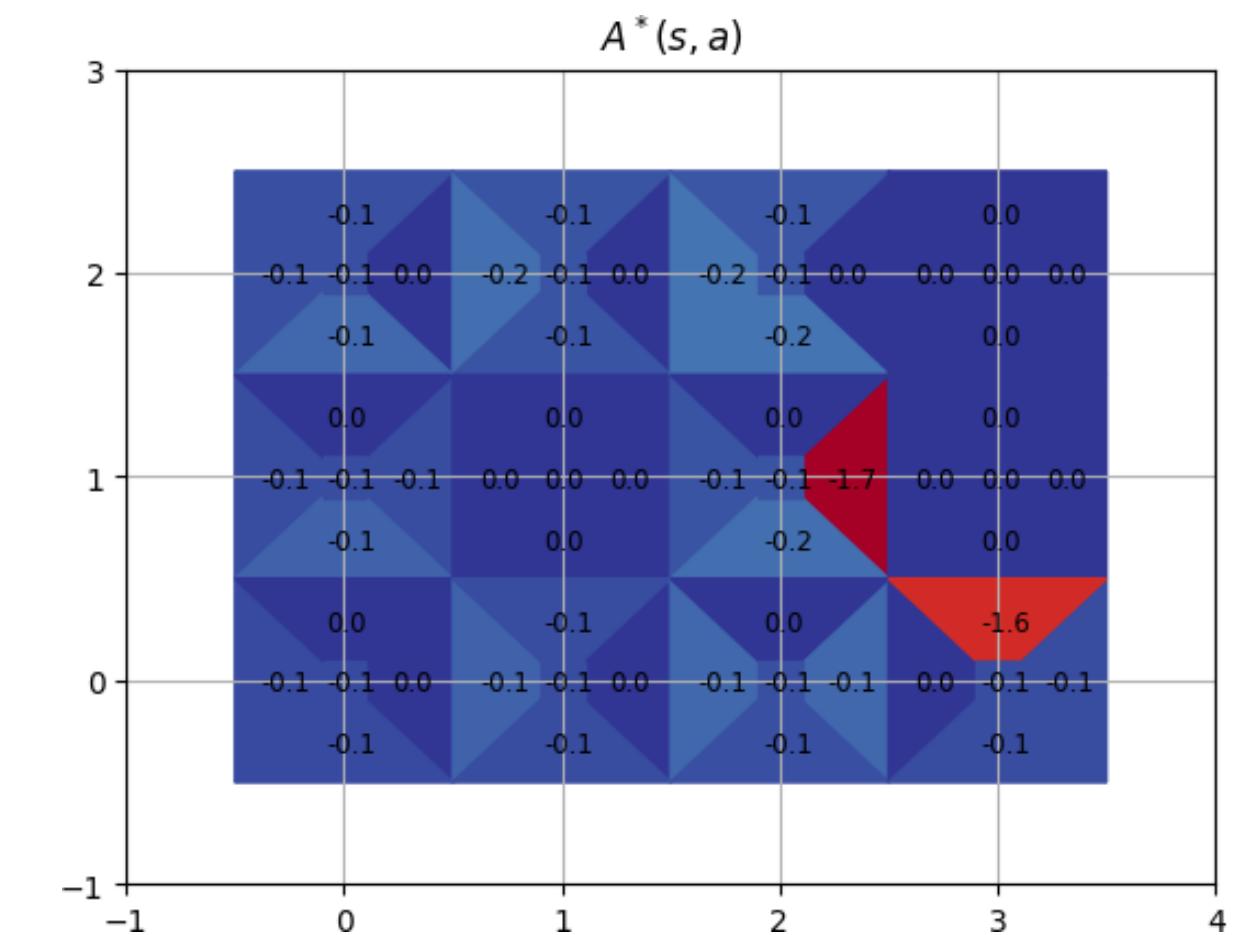
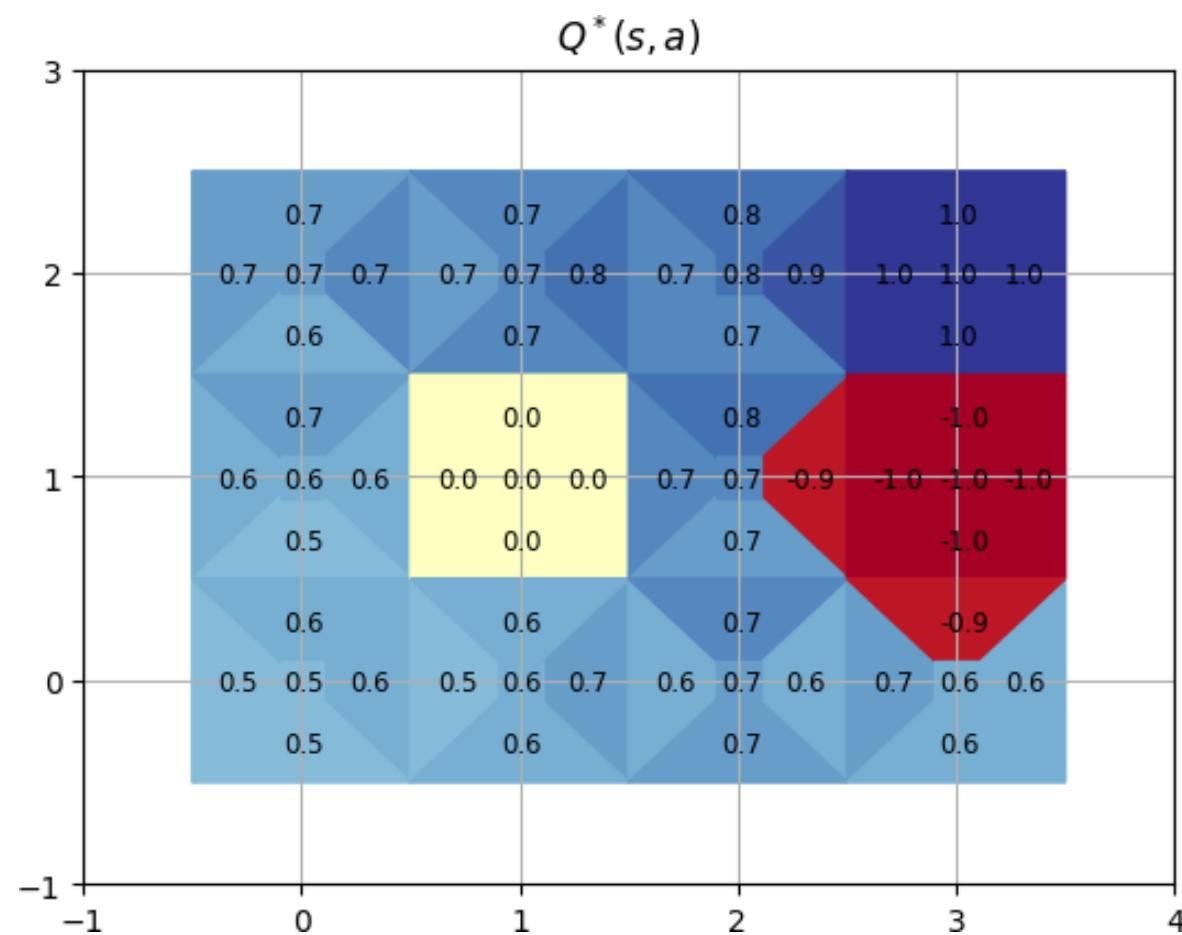
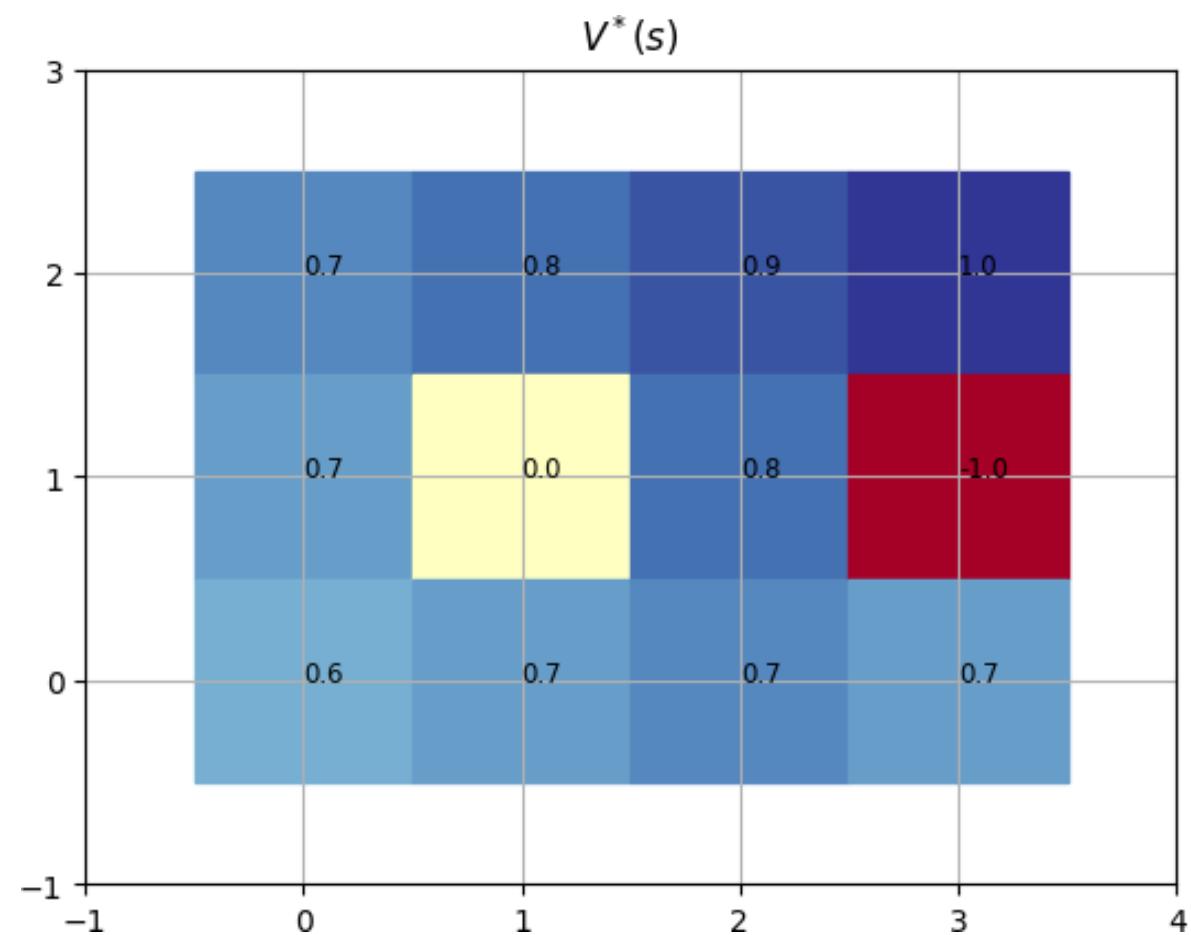
$$Q_{k+1}^*(s, a) = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \left( r(s_t, a_t) + \gamma \max_{a_t} Q_k^*(s_{t+1}, a_{t+1}) \right)$$



# Optimal Value, Q-Value, Policy



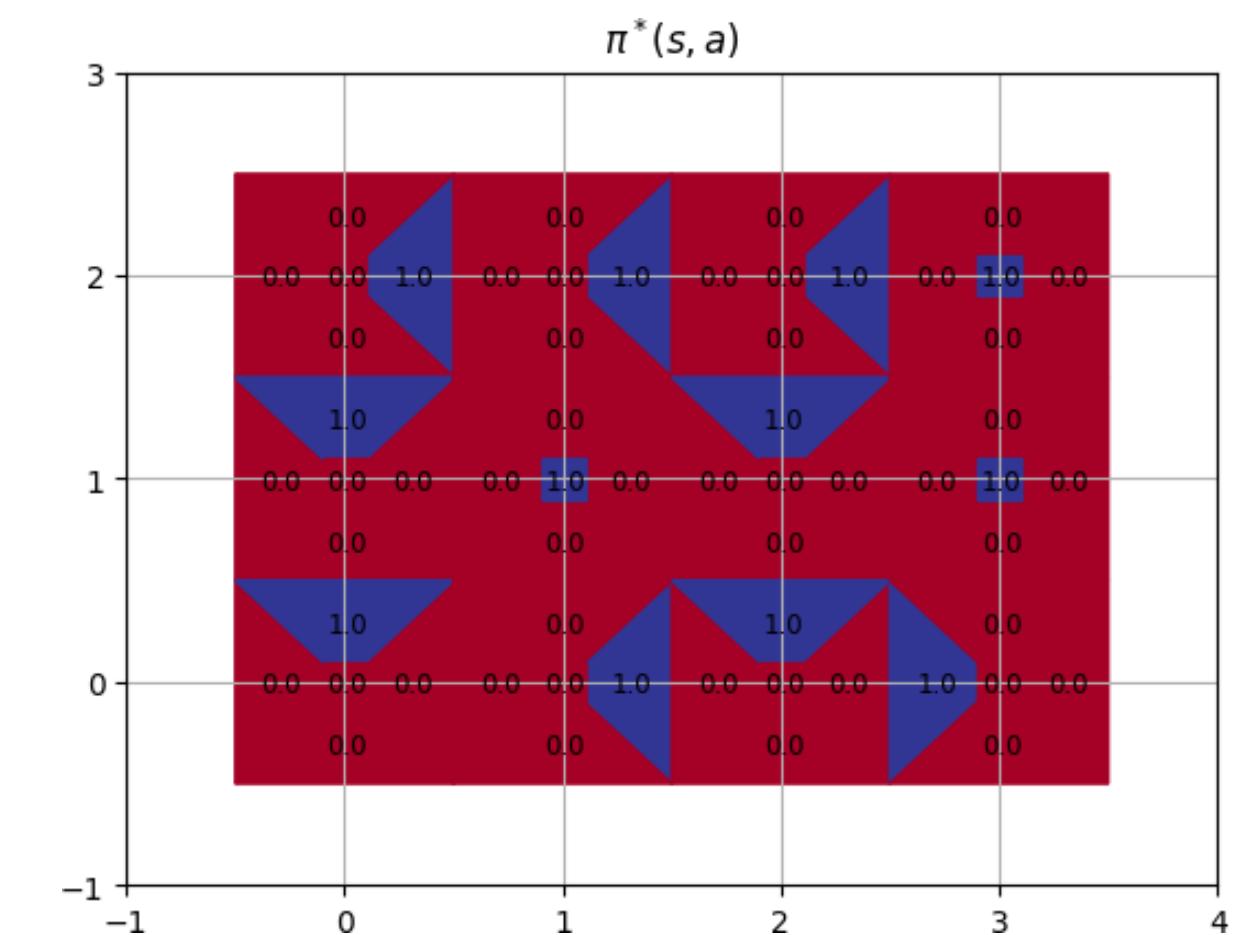
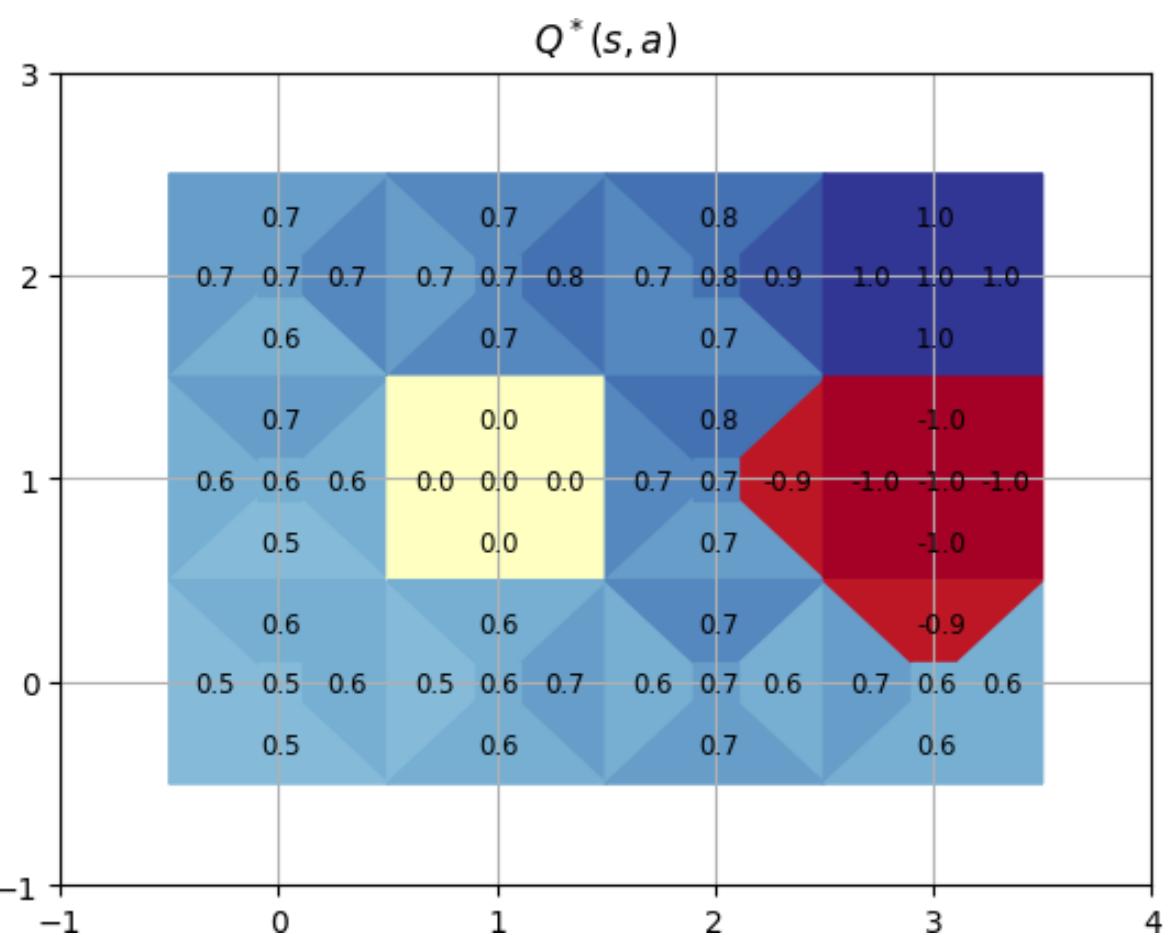
# Advantage (Optimal Policy)



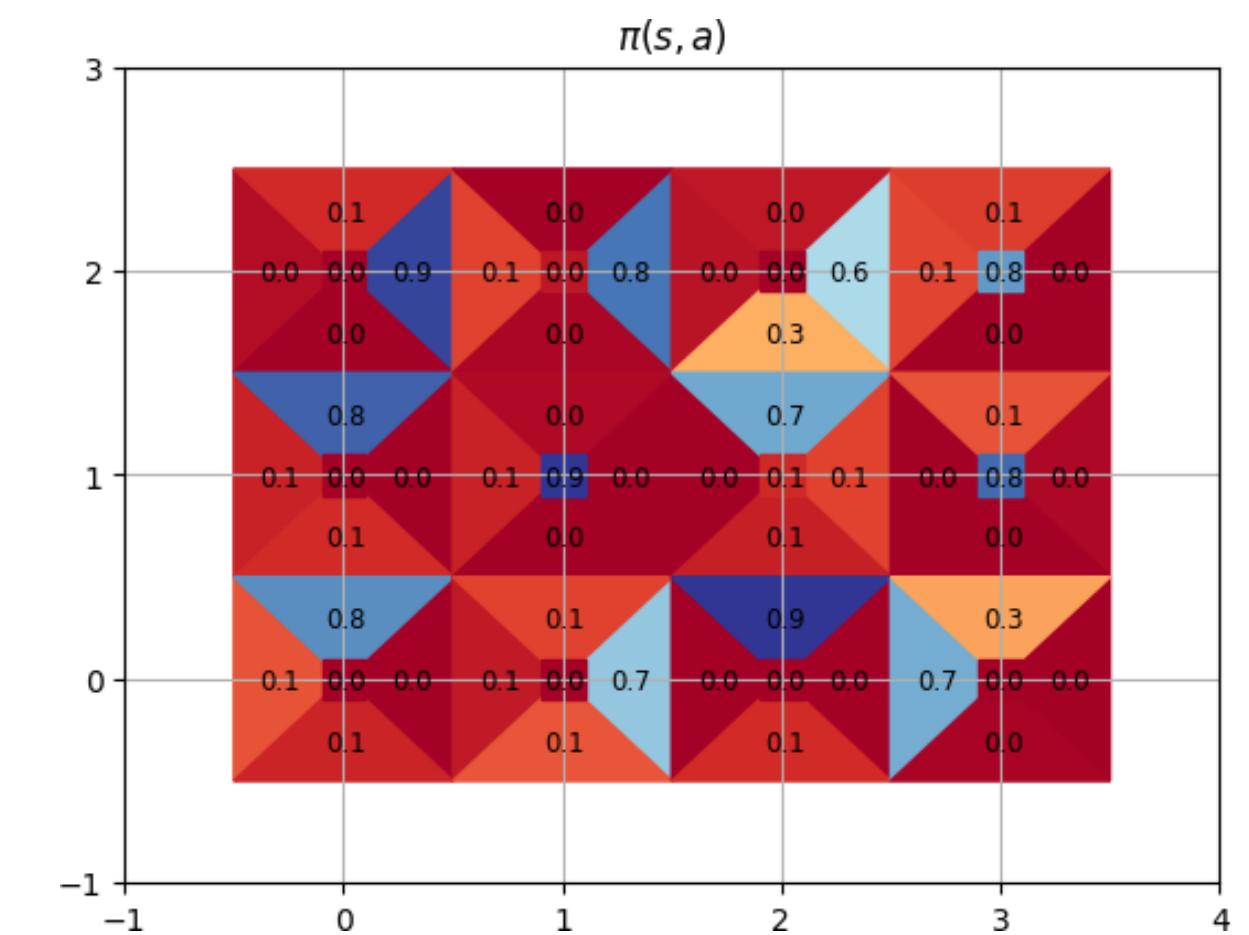
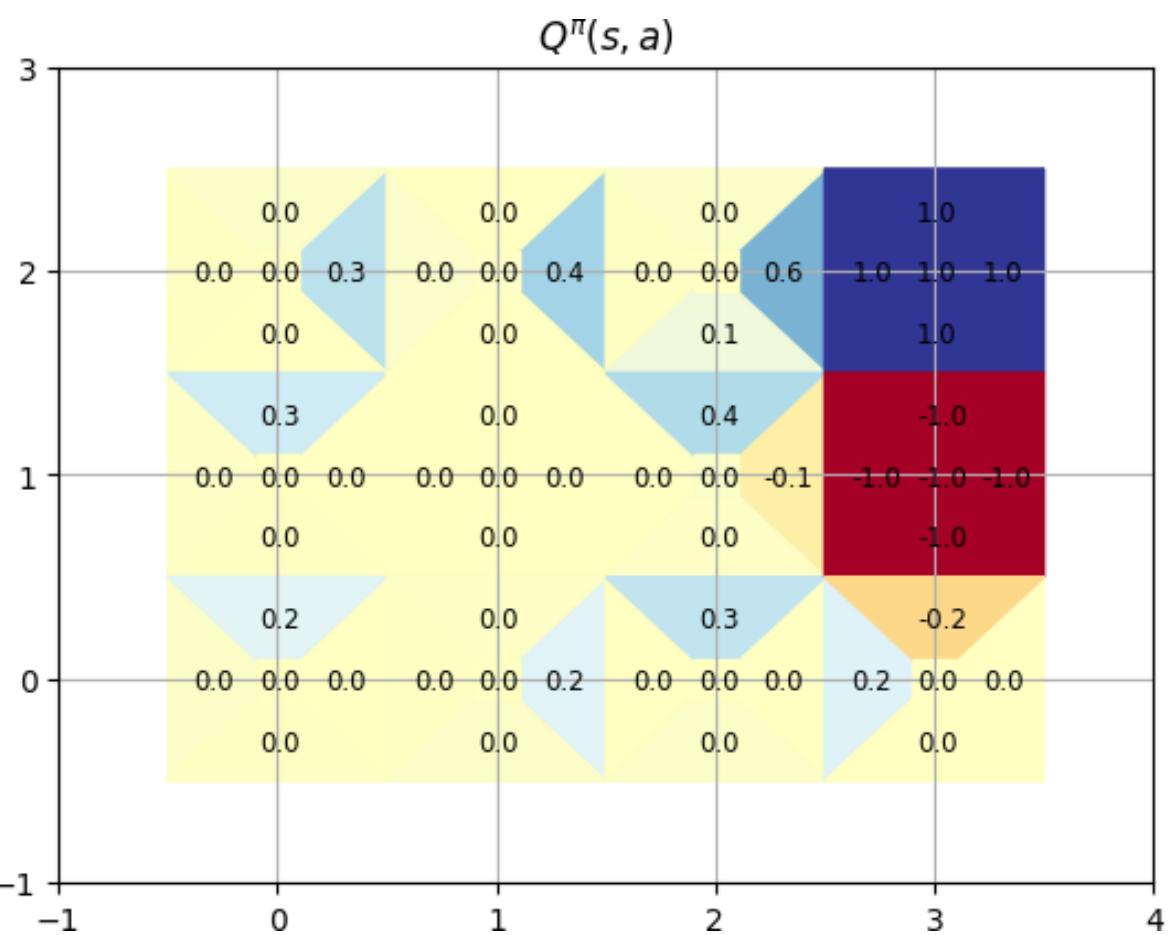
$A^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t)$ : how much better  $\mathbf{a}_t$  is

# $Q^*$ vs $Q^\pi$

Optimal Policy  
(Deterministic for MDP)



Suboptimal Policy



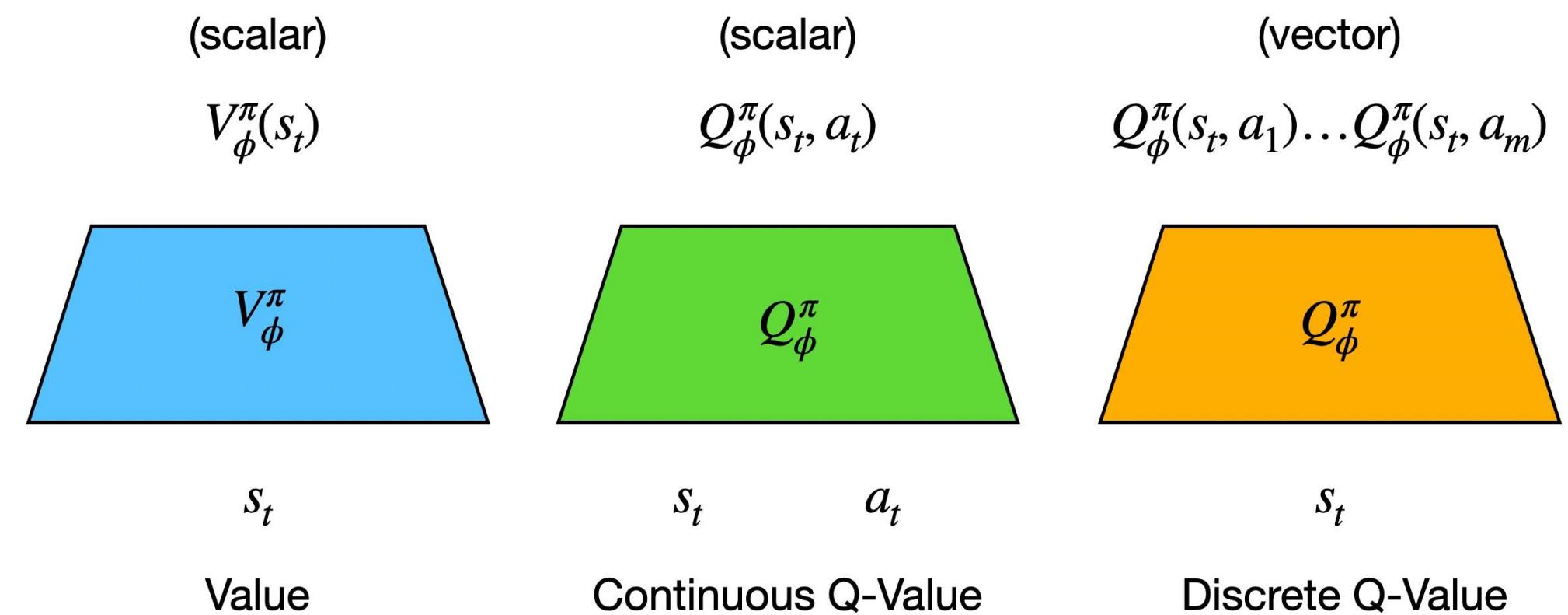
To play more with the Tabular Environment, checkout the [colab](#)

**Pause for Questions.**

# Limitations of Tabular Environment

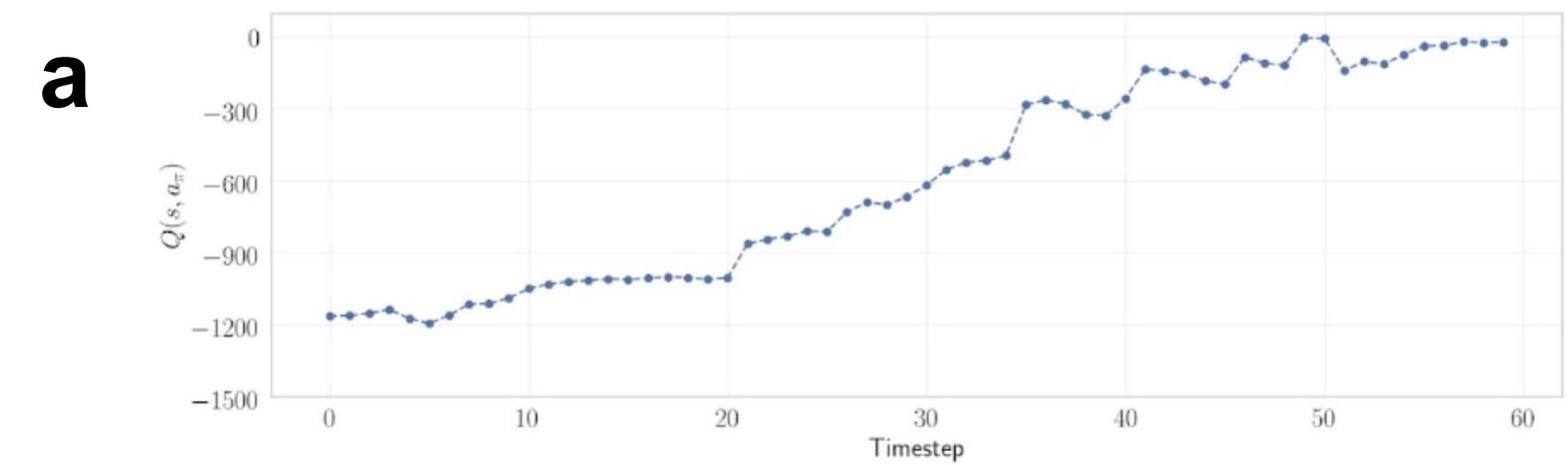
## Moving towards parametric value functions

- A tabular value function intractable as state/action space grows!
- Learning a parametric value function can also allow for generalization
  - Related states have similar q-values (smoothness)
- Can use regression to learn value functions with a neural network



# Discrete vs Continuous Q-values

- Different formulations for continuous and discrete action spaces:
  - In **continuous (a)**, typically condition on an action and output a scalar, where you sample actions from your policy for an expectation calculation
  - In **discrete (b)**, typically learn a value for each bin (output) and can enumerate over all actions for an expectation



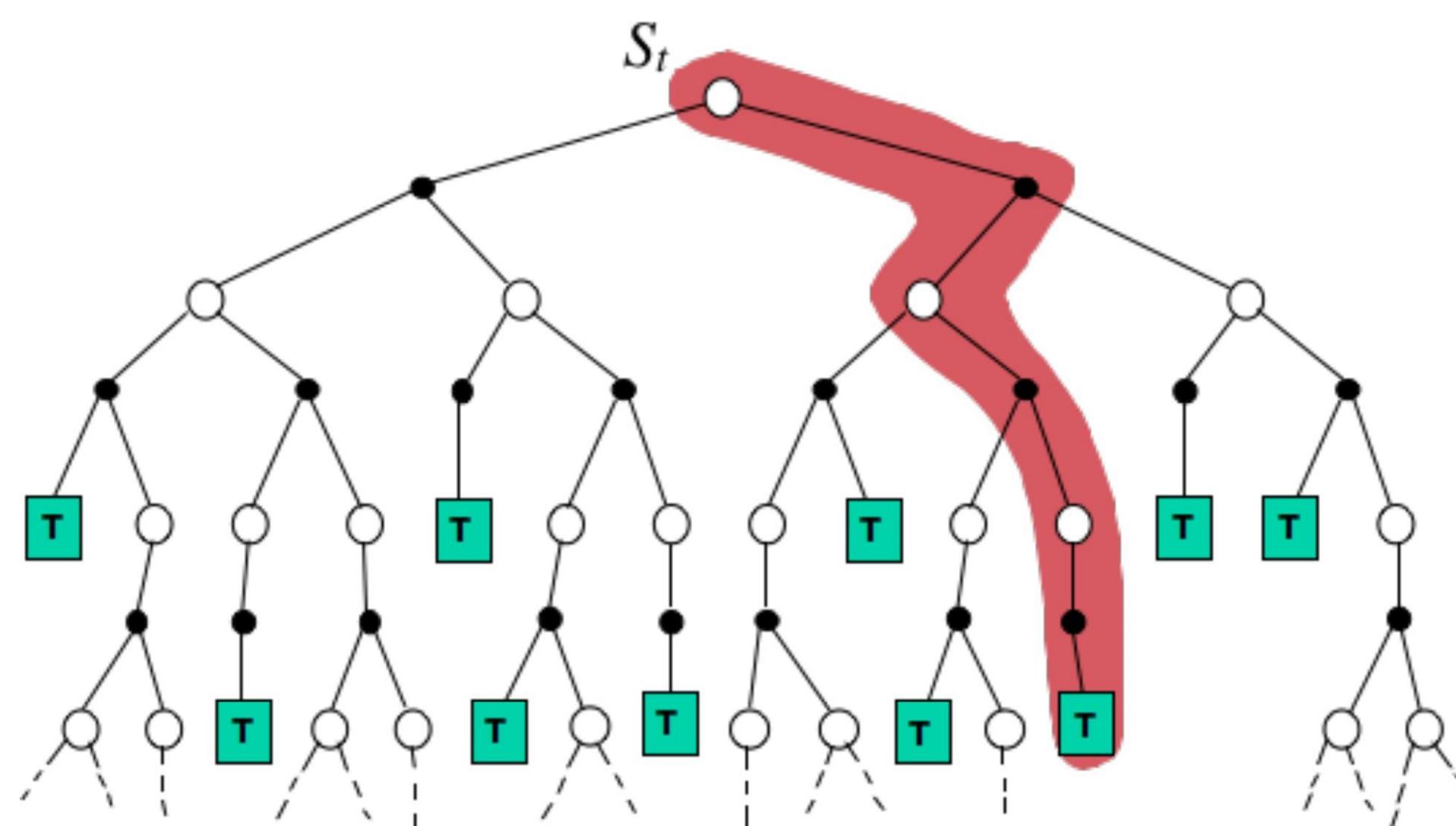
# Q-Learning as a Regression Problem

## Whiteboard

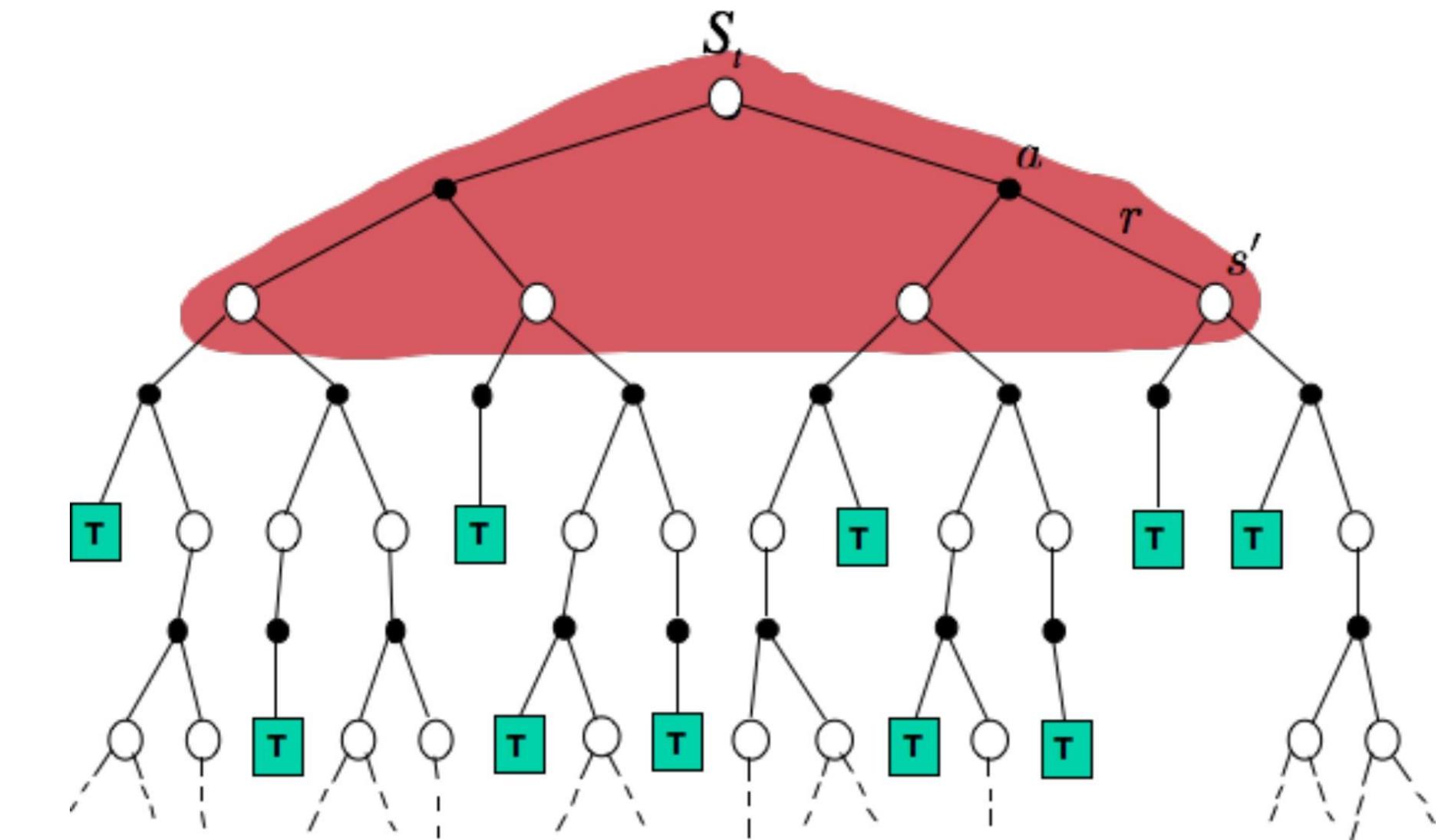
$$\begin{aligned} V^\pi(s) &= E_\pi [\mathcal{R}_t \mid S_t = s] && \text{(Definition)} \\ &= E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r(s_{t+k}, a_{t+k}) \mid S_t = s \right] && \text{(Monte-Carlo Rollout, MC)} \\ &= E_\pi \left[ r(s_t, a_t) + \gamma \sum_{k=0}^{\infty} \gamma^k r(s_{t+k+1}, a_{t+k+1}) \mid S_t = s \right] && \text{(Expansion)} \\ &= E_\pi [r(s_t, a_t) + \gamma V^\pi(s_{t+1}) \mid S_t = s] && \text{(Dynamic Programming, TD)} \\ &= E_\pi \left[ r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \dots + \gamma^{k+1} V^\pi(s_{t+k+1}) \mid S_t = s \right] && \text{(N-Step)} \end{aligned}$$

Bias vs Variance Tradeoff!

# MC vs TD Visualization



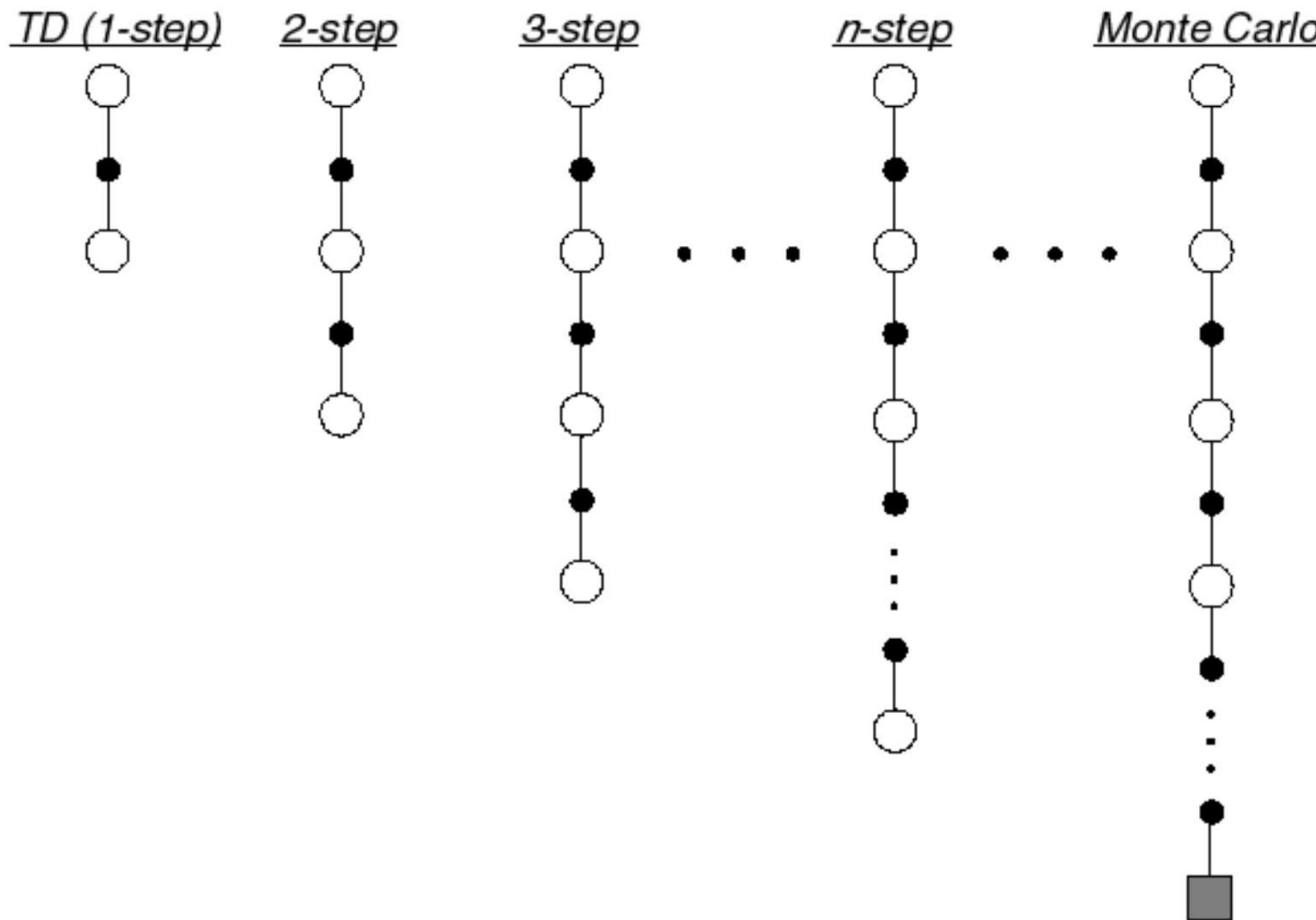
Monte-Carlo Rollouts  
(Regression)



TD Backup  
(Dynamic Programming)

# N-Step Returns

Let TD target look  $n$  steps into the future



# Bias vs Variance Tradeoff

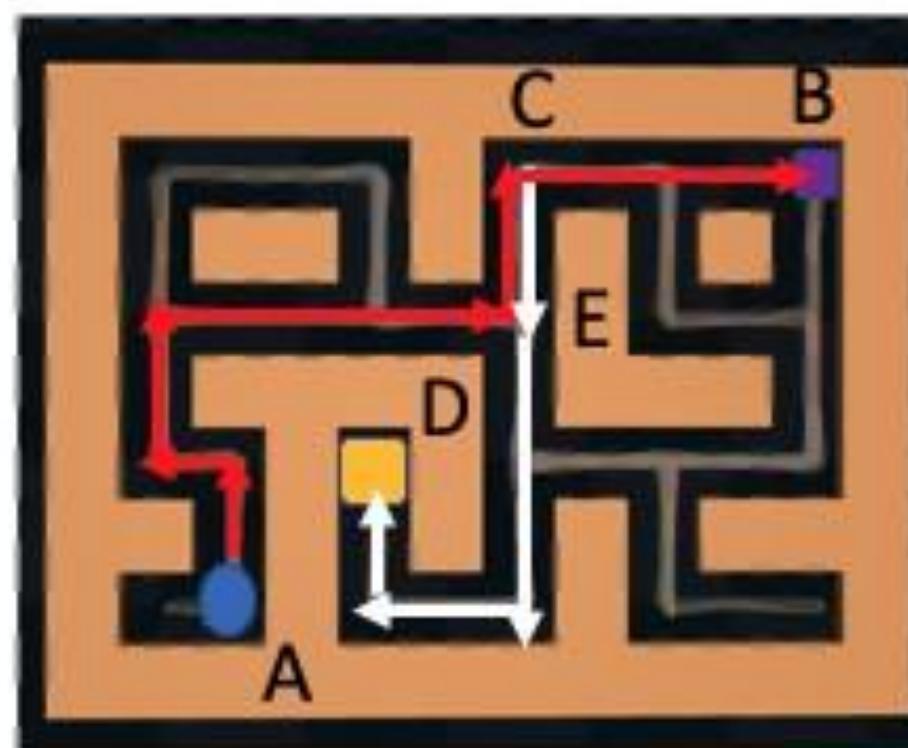
## TD, MC, N-Step Return

| Method      | Bootstrap Depth | Bias   | Variance |
|-------------|-----------------|--------|----------|
| TD          | 1 Steps         | High   | Low      |
| N-Step      | N Steps         | Medium | Medium   |
| Monte Carlo | Full Episode    | Low    | High     |

# Why Dynamic Programming?

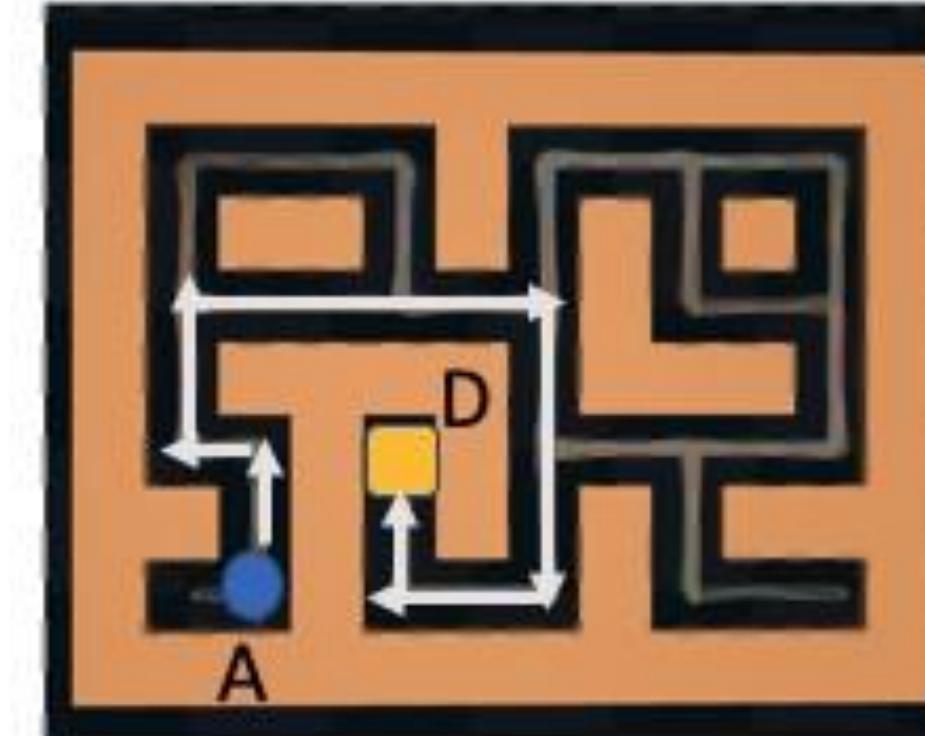
## Aside on Stitching

Our maze  
navigation task



- Start (A)    ■ Goal (D)
- trajectory from A → B
- trajectory from C → D

“Stitching” enables  
finding optimal behavior



Trajectory found by  
stitching segments of  
trajectories

**Pause for Questions.**

# Practical Implementation Detail

## Semi-Gradient and Target Network

Going back to our formulation:

$$Q(s, a) \leftarrow r(s, a) + \gamma \max_{\hat{a}} Q(s', \hat{a})$$

Semi-Gradient

$$Q_\theta(s, a) \leftarrow r(s, a) + \gamma \max_{\hat{a}} \text{stopgrad} (Q_\theta(s', \hat{a}))$$

Target Network

$$Q_\theta(s, a) \leftarrow r(s, a) + \gamma \max_{\hat{a}} \text{stopgrad} (Q_{\text{target}}(s', \hat{a}))$$

Update Types

$$w' \leftarrow w \text{ when } \text{mod}(n, N) = 0$$

Hard Target Update

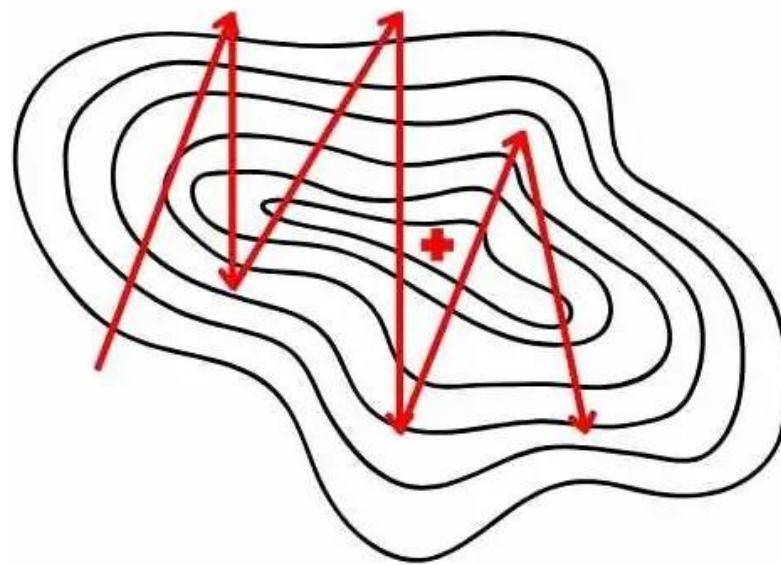
$$w' \leftarrow \tau \cdot w + (1 - \tau) \cdot w'$$

Soft Target Update (Polyak)

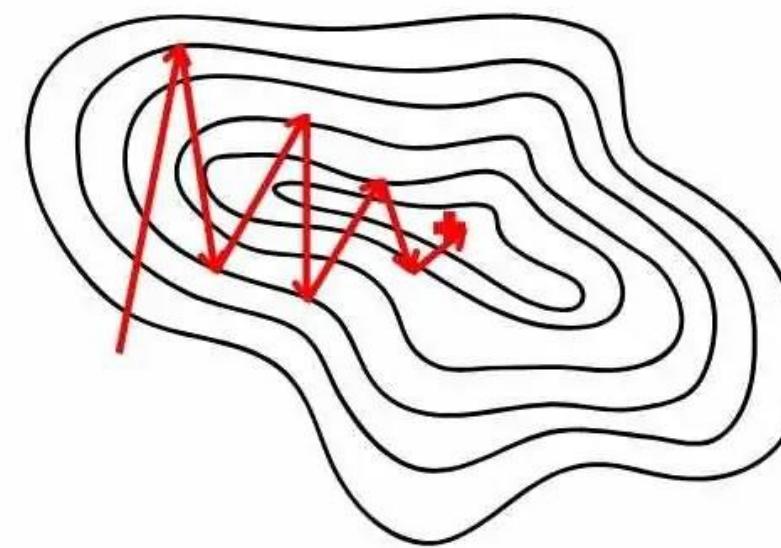
# Practical Implementation Details

## Dealing with TD Gradients

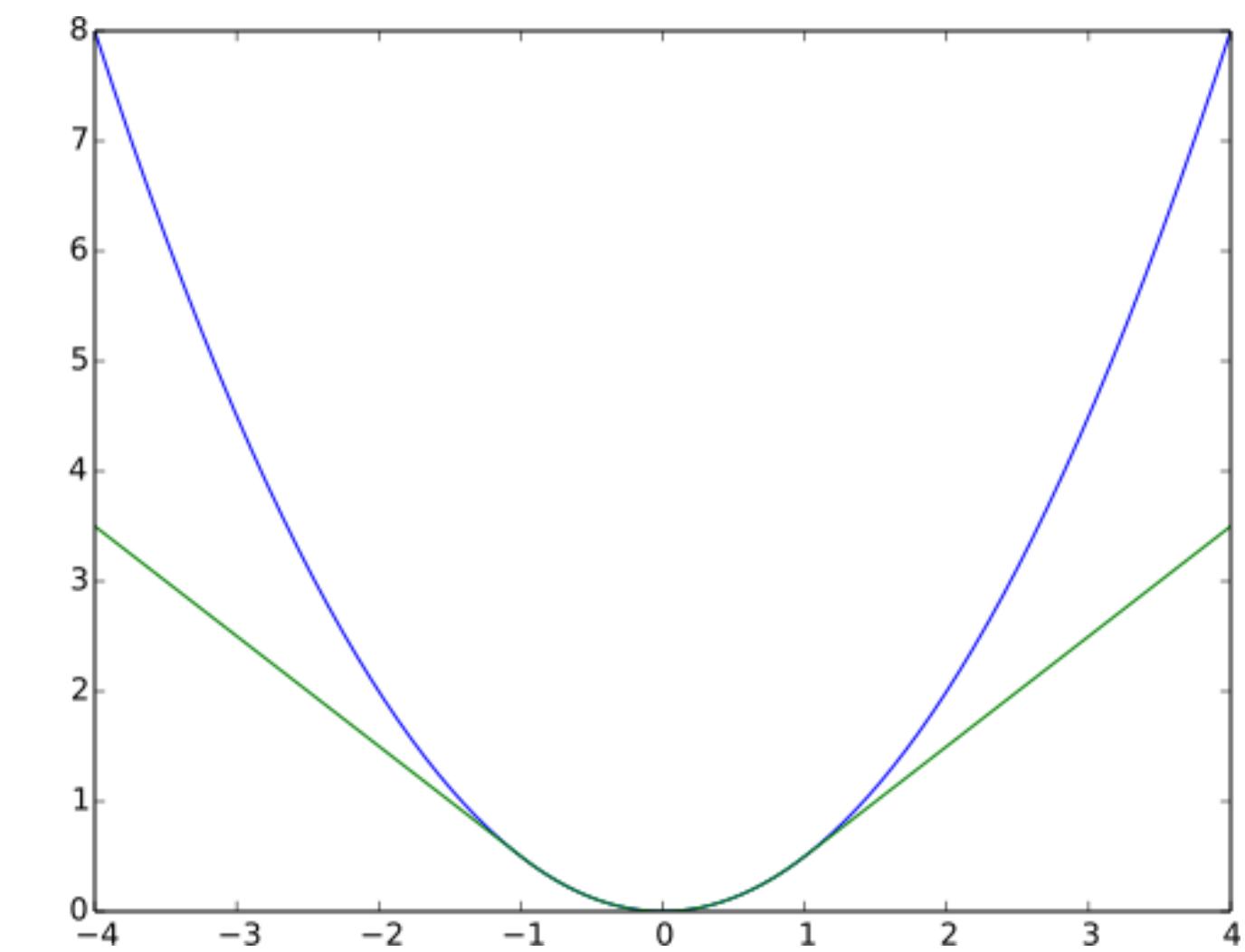
Without Gradient Clipping



With Gradient Clipping



Gradient Clipping

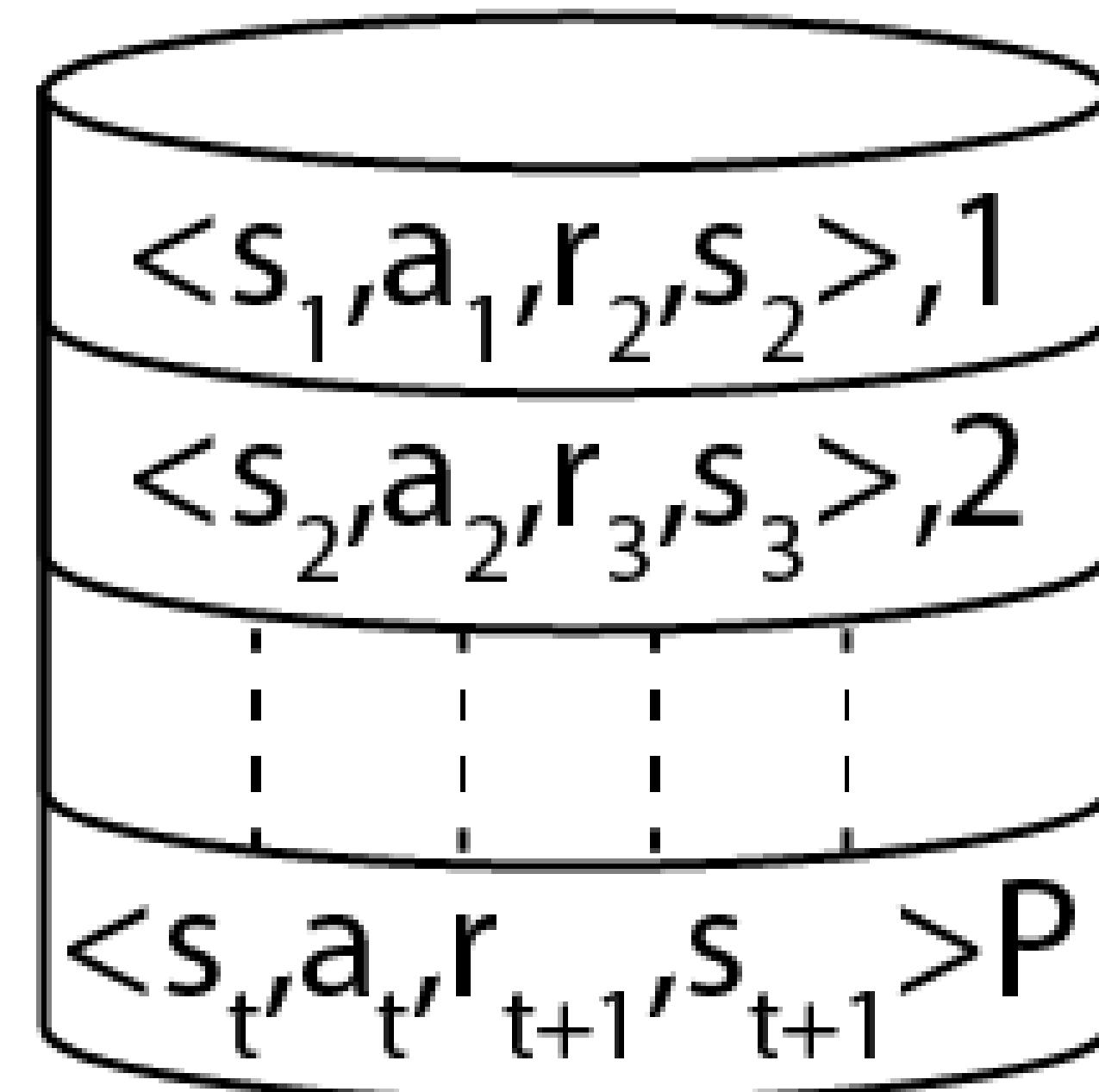


Huber Loss

# Practical Implementation Details

## Replay Buffers

- Replay Buffers are **memory structures** that store past experiences encountered by the agent, allowing the **experiences to be reused** instead of discarded after a single update
- They **break temporal correlations** between consecutive training samples and prevents recency bias
- **Improves sample efficiency** in Q-learning algorithms, allowing an agent to reuse its experience effectively



# Q-Learning Overestimation

Q-Learning is typically formulated as follows:

$$Q(s, a) \leftarrow r(s, a) + \gamma \max_{\hat{a}} Q(s', \hat{a})$$

When parameterized by a function approximator like a NN, there is some inherent noise in the q-value estimates due to modeling errors

$$Q^{\text{approx}}(s', \hat{a}) = Q^{\text{target}}(s', \hat{a}) + Y_{s'}^{\hat{a}}$$

Due to the noise on the RHS ( $Y_{s'}^{\hat{a}}$ ), there results in error on the LHS, expressed as:

$$Z_s := \gamma \left( \max_{\hat{a}} Q^{\text{approx}}(s', \hat{a}) - \max_{\hat{a}} Q^{\text{target}}(s', \hat{a}) \right)$$

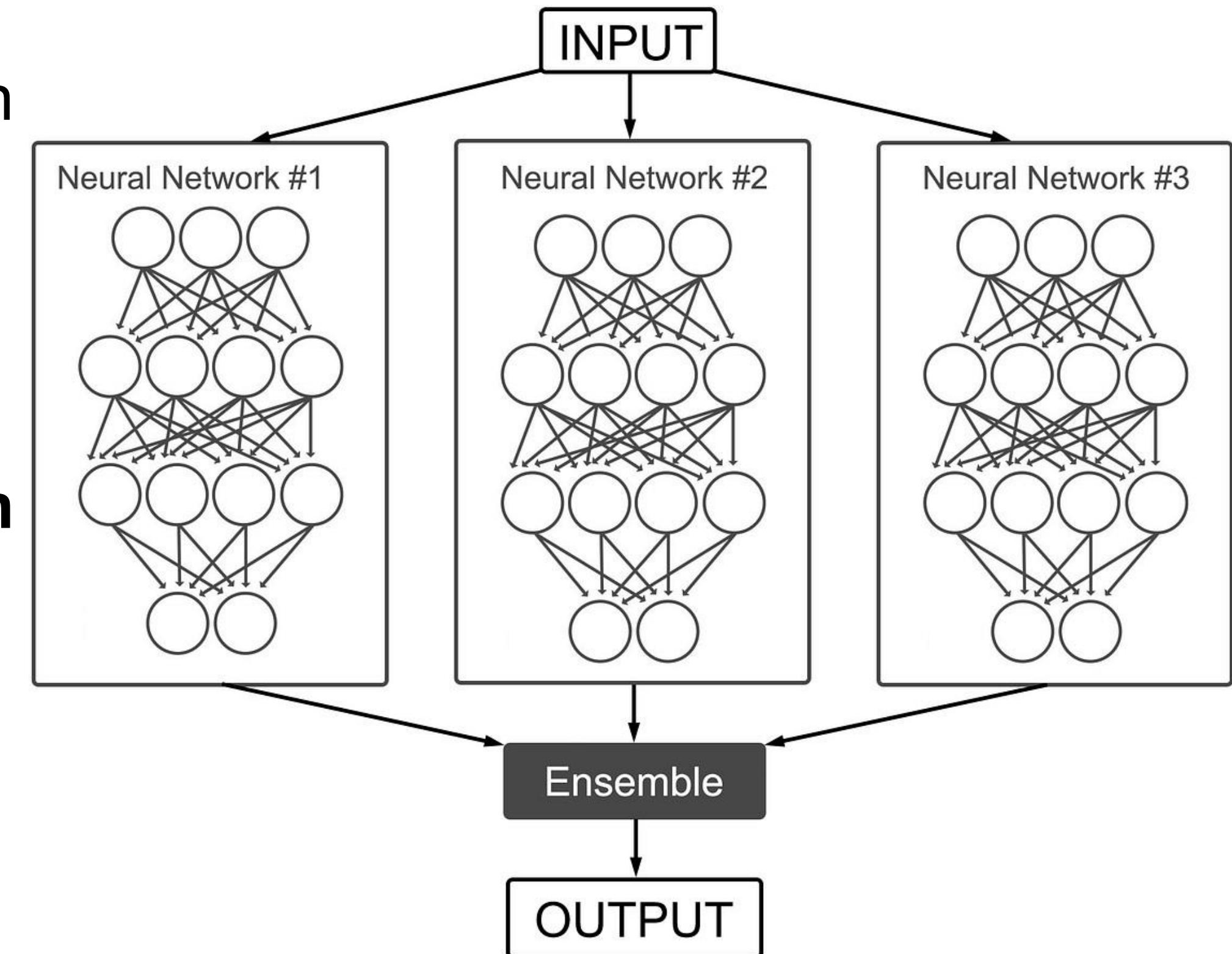
Even with zero mean ( $Y_{s'}^{\hat{a}}$ ),  $Z_s$  may have positive mean:

$$\mathbb{E}[Y_{s'}^{\hat{a}}] = 0 \forall \hat{a} \Rightarrow \mathbb{E}[Z_s] > 0 \text{ (often)}$$

# Practical Implementation Detail

## Critic Ensembling

- One approach to tackle overestimation is to measure **epistemic uncertainty** with an ensemble of networks
- Since different model initializations may have **different modeling errors**, this can lead to ensembling with a **min reduction** to be effective
- Has been even shown to be effective in **offline RL** where distribution shift is a larger concern (Sac-N and EDAC)



# Q-Learning vs Double Q-Learning

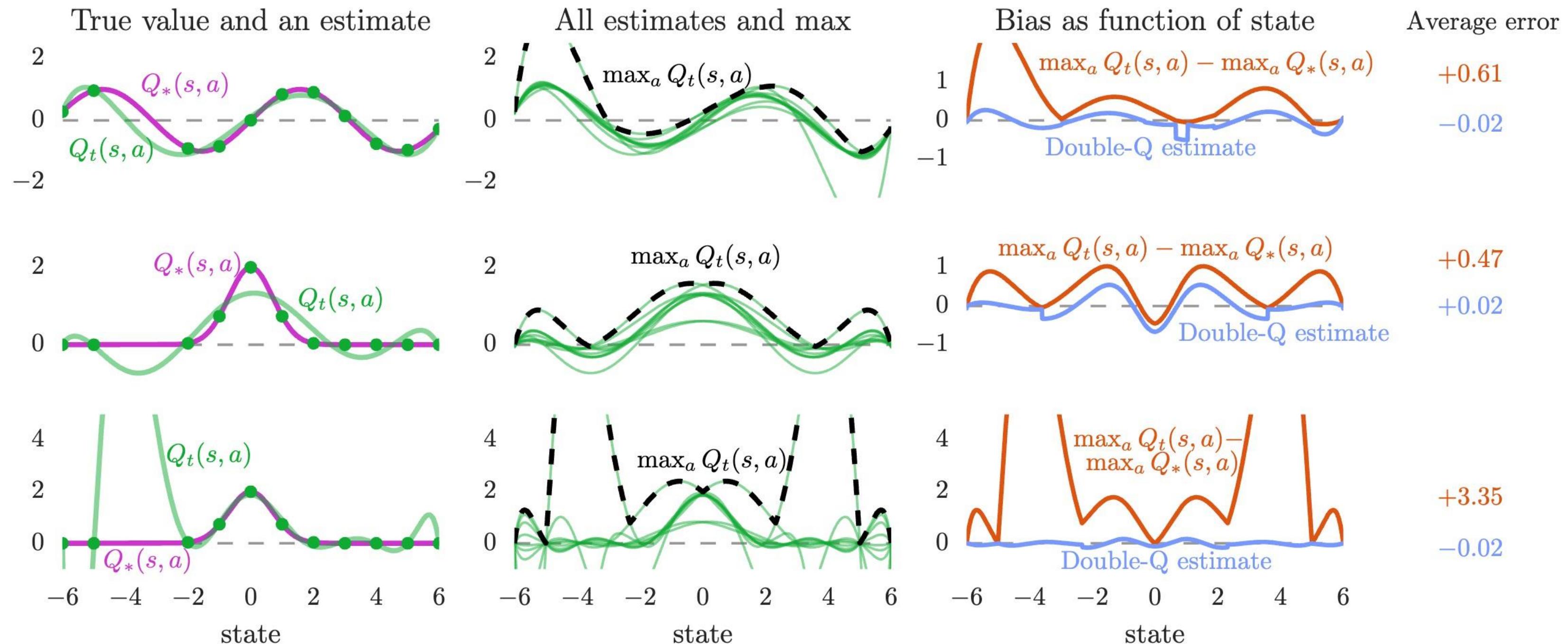


Figure 2: Illustration of overestimations during learning. In each state (x-axis), there are 10 actions. The **left column** shows the true values  $V_*(s)$  (purple line). All true action values are defined by  $Q_*(s, a) = V_*(s)$ . The green line shows estimated values  $Q(s, a)$  for one action as a function of state, fitted to the true value at several sampled states (green dots). The **middle column** plots show all the estimated values (green), and the maximum of these values (dashed black). The maximum is higher than the true value (purple, left plot) almost everywhere. The **right column** plots shows the difference in orange. The blue line in the right plots is the estimate used by Double Q-learning with a second set of samples for each state. The blue line is much closer to zero, indicating less bias. The three **rows** correspond to different true functions (left, purple) or capacities of the fitted function (left, green). (Details in the text)

# Understanding HW Algorithms

# Q-Learning : DQN

## Whiteboard

---

**Algorithm 1:** DQN

---

**Input:** Replay buffer  $\mathcal{D}$ , Q-network  $Q_\theta$ , target network  $Q_{\theta'}$  with  $\theta' \leftarrow \theta$ ,  
discount factor  $\gamma$ , batch size  $N$ , learning rate  $\alpha$ , target update  
frequency  $C$ , exploration schedule  $\epsilon$

**for** each interaction step  $t = 1, 2, \dots$  **do**

    Observe state  $s_t$

    Choose action  $a_t$  using  $\epsilon$ -greedy:

$$a_t = \begin{cases} \text{random action,} & \text{with probability } \epsilon, \\ \arg \max_a Q_\theta(s_t, a), & \text{otherwise} \end{cases}$$

    Execute  $a_t$ , observe  $(r_t, s_{t+1})$ , and store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$

    // Sample a mini-batch

    Sample  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N \sim \mathcal{D}$

    // Compute targets

**for**  $i = 1, \dots, N$  **do**

$y_i \leftarrow r_i + \gamma \max_{a'} Q_{\theta'}(s_{i+1}, a')$

    // Q-network update

$$\theta \leftarrow \theta - \alpha \nabla_\theta \frac{1}{N} \sum_{i=1}^N [Q_\theta(s_i, a_i) - y_i]^2$$

    // Periodic target update

**if**  $t \bmod C = 0$  **then**

$\theta' \leftarrow \theta$

# Actor Critic: SAC

## Whiteboard

---

**Algorithm 1:** Soft Actor-Critic (no entropy term)

---

**Input:** Replay buffer  $\mathcal{D}$ , critic nets  $Q_{\theta_1}, Q_{\theta_2}$ , actor net  $\pi_\phi$ , target critics  $\theta'_j \leftarrow \theta_j$ , discount  $\gamma$ , smoothing  $\tau$ , batch size  $N$

**for** each interaction step  $t = 1, 2, \dots$  **do**

- Observe state  $s_t$
- Sample action  $a_t \sim \pi_\phi(\cdot | s_t)$
- Execute  $a_t$ , observe  $(r_t, s_{t+1})$  and store in  $\mathcal{D}$
- // Sample a mini-batch
- Sample  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^N \sim \mathcal{D}$
- // Compute targets (no entropy term)
- for**  $i = 1, \dots, N$  **do**

  - $a'_{i+1} \leftarrow \pi_\phi(s_{i+1});$
  - $y_i \leftarrow r_i + \gamma \min_{j=1,2} Q_{\theta'_j}(s_{i+1}, a'_{i+1})$

- // Critic update
- for**  $j = 1, 2$  **do**

  - $\theta_j \leftarrow \theta_j - \lambda_Q \nabla_{\theta_j} \frac{1}{N} \sum_i [Q_{\theta_j}(s_i, a_i) - y_i]^2$

- // Actor update (maximize Q only)
- $\phi \leftarrow \phi + \lambda_\pi \nabla_\phi \frac{1}{N} \sum_i Q_{\theta_1}(s_i, \pi_\phi(s_i))$
- // Target networks soft-update
- for**  $j = 1, 2$  **do**

  - $\theta'_j \leftarrow \tau \theta_j + (1 - \tau) \theta'_j$