



Robotics

Spring 2023

Abhishek Gupta

TAs: Yi Li, Srivatsa GS

Recap: Course Overview

Filtering/Smoothing

Localization

Mapping

SLAM

Search

Motion Planning

TrajOpt

Stability/Certification

MDPs and RL

Imitation Learning

Off-Policy/MBRL

Lecture Outline

Recap + Policy Gradient

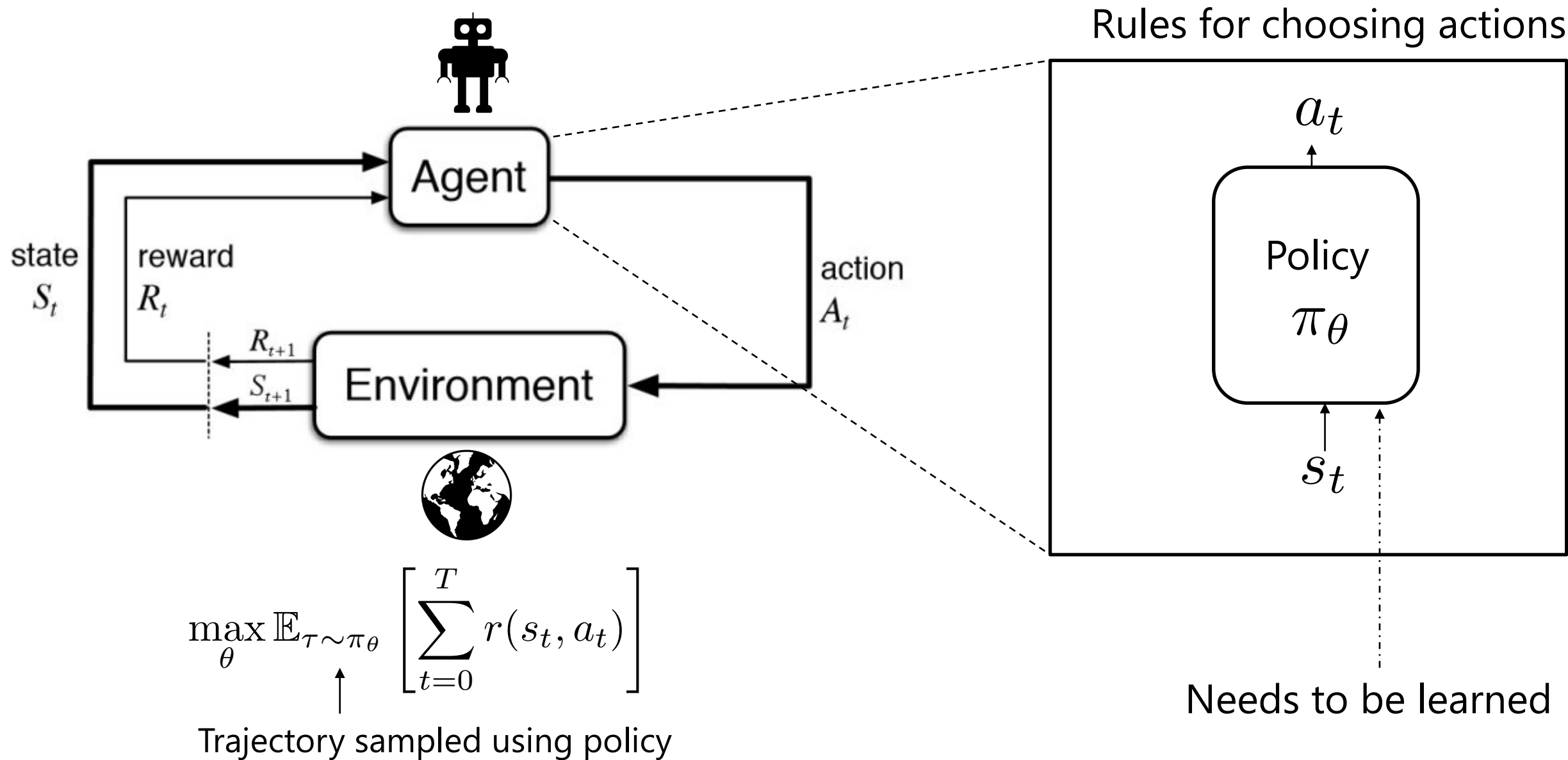


Basic Actor Critic Algorithms



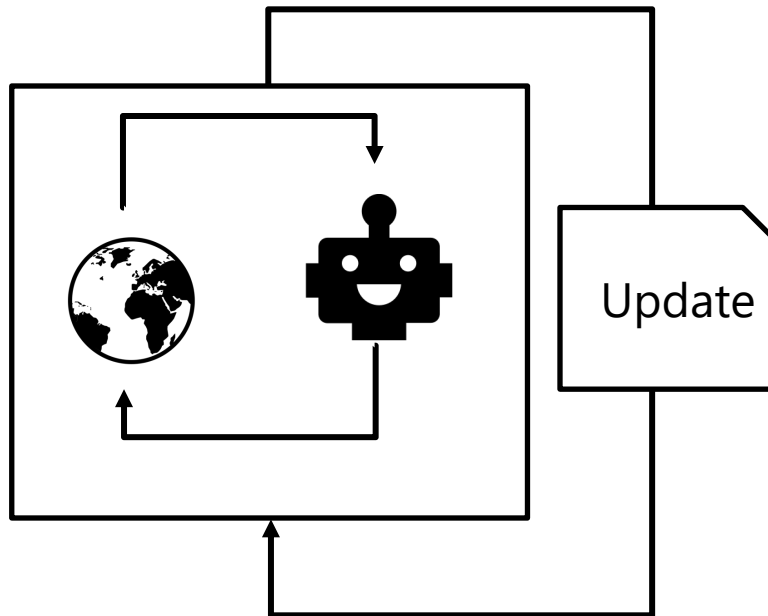
Making Actor-Critic Practical

Reinforcement Learning Formalism



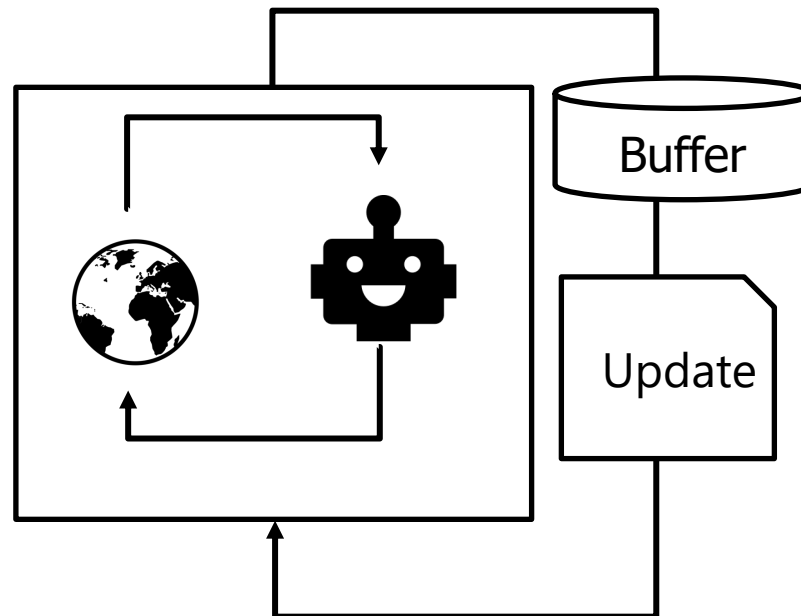
Learning Algorithms for Robotics

On-Policy Algorithms



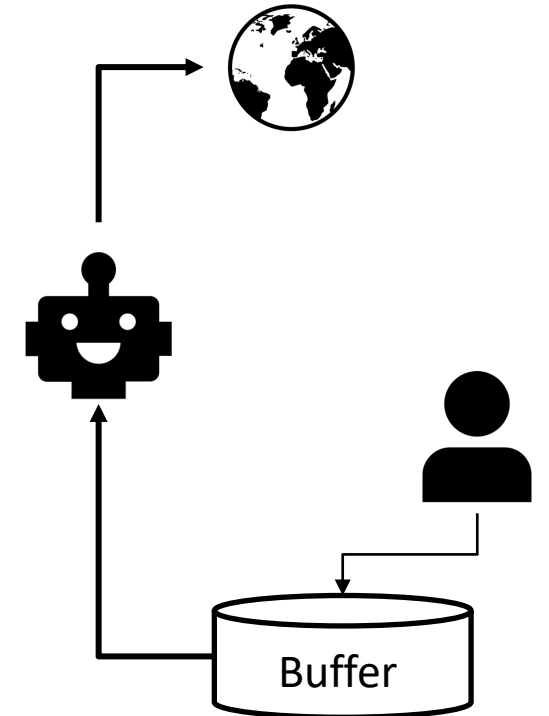
Simple, performant,
Data inefficient

Off-Policy Algorithms



Data-efficient,
sometimes unstable

Imitation Learning

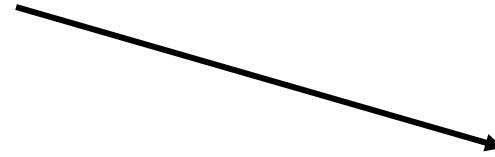


Performant, efficient, but
compounding error and
expensive data collection

What if we just performed gradient ascent?

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T r(s_t, a_t) \right]$$

$$= \int p_{\theta}(\tau) R(\tau) d\tau$$



Standard gradient descent (supervised learning)

$$\nabla_{\theta} \mathbb{E}_{x \sim g(x)} [f_{\theta}(x)]$$

REINFORCE gradient descent (RL)

$$\nabla_{\theta} \mathbb{E}_{x \sim p_{\theta}(x)} [f(x)]$$

Gradient wrt expectation variable, not of integrand!

Taking the gradient of sum of rewards

$$J(\theta) = \int p_{\theta}(\tau) R(\tau) d(\tau)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \int p_{\theta}(\tau) R(\tau) d(\tau)$$

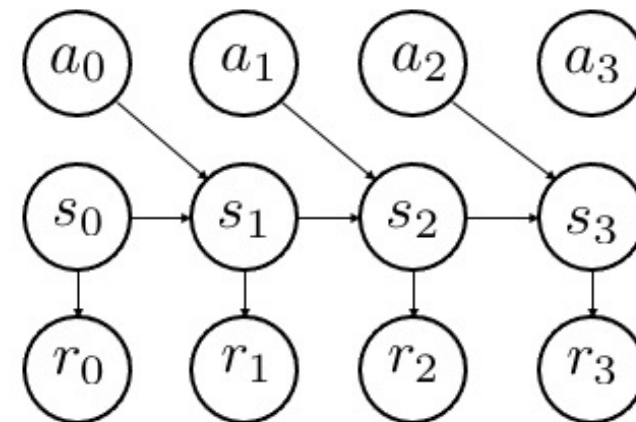
$$= \int \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau) = \int \frac{p_{\theta}(\tau)}{p_{\theta}(\tau)} \nabla_{\theta} p_{\theta}(\tau) R(\tau) d(\tau)$$

$$= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) R(\tau) d(\tau) = \mathbb{E}_{p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) R(\tau)]$$

REINFORCE trick

Taking the gradient of return

Initial State	Dynamics	Policy
$p_{\theta}(\tau) = p(s_0) \prod_{t=0}^{T-1} p(s_{t+1} s_t, a_t) \pi(a_t s_t)$		



$$\log p_{\theta}(\tau) = \log p(s_0) + \sum_{t=0}^{T-1} \log p(s_{t+1} | s_t, a_t) + \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \cancel{\nabla_{\theta} \log p(s_0)} + \sum_{t=0}^{T-1} \cancel{\nabla_{\theta} \log p(s_{t+1} | s_t, a_t)} + \nabla_{\theta} \log \pi(a_t | s_t)$$

$$\nabla_{\theta} \log p_{\theta}(\tau) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t)$$

Model Free!!

Taking the gradient of return

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\nabla_{\theta} \log p_{\theta}(\tau) \sum_{t=0}^T r(s_t, a_t) \right]$$

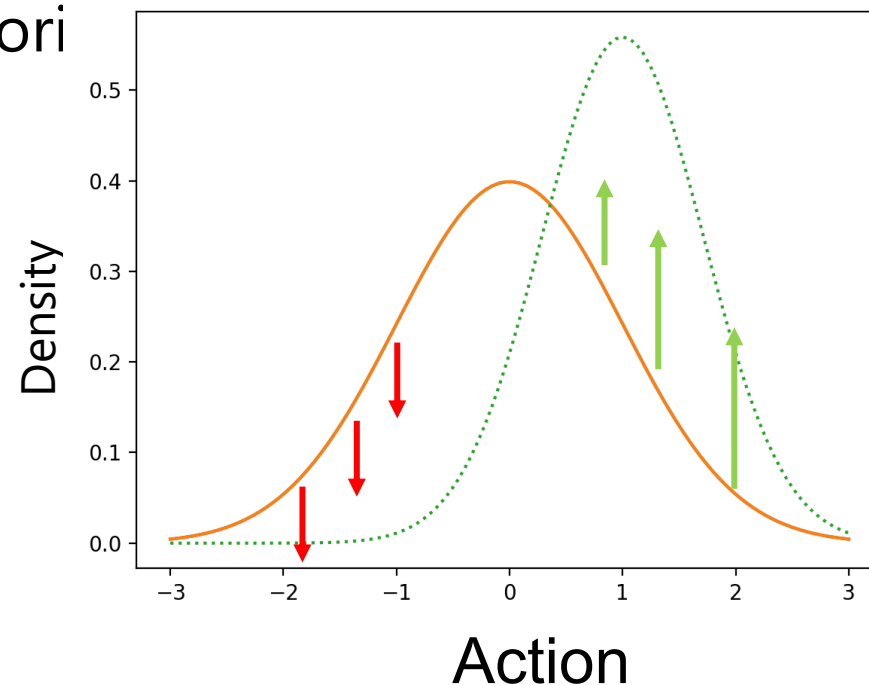
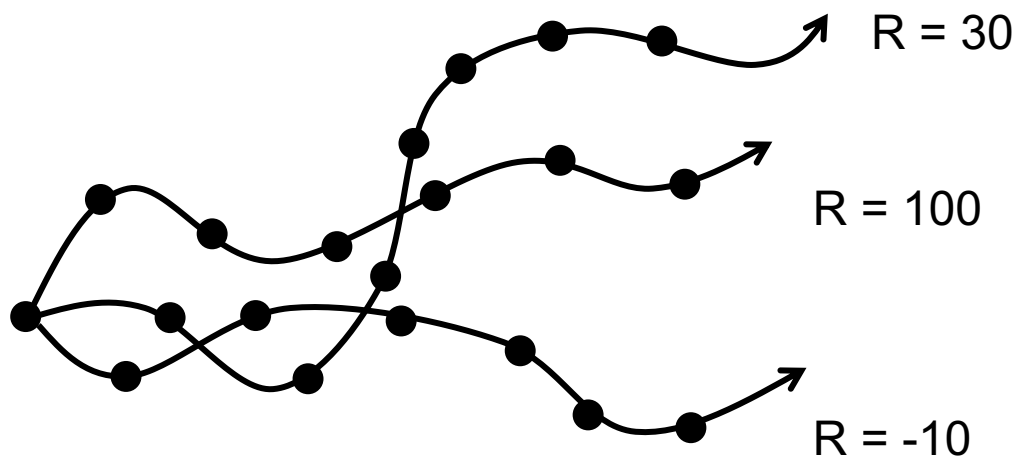
$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t) \\ a_t \sim \pi(a_t | s_t)}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=0}^T r(s_t, a_t) \right]$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i) \quad (\text{approximating using samples})$$

What does this mean?

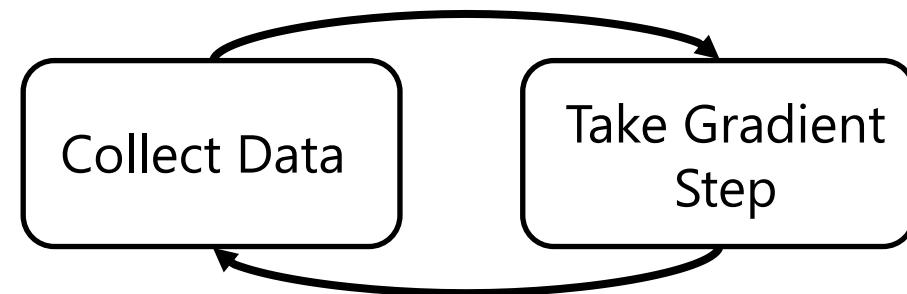
$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Increase the likelihood of actions in high return trajectories



Resulting Algorithm (REINFORCE)

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



REINFORCE algorithm:

On-policy



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

REINFORCE Pseudocode

Algorithm 1 Vanilla Policy Gradient Algorithm

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Estimate policy gradient as

→ Sum up reward to go

→ Simply subtract baseline

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 7: Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

or via another gradient ascent algorithm like Adam.

- 8: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

→ Regression to learn baseline

typically via some gradient descent algorithm.

- 9: **end for**
-

How is this related to supervised learning?

Reinforcement Learning

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

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

Supervised Learning

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log p_{\theta}(y|x)]$$

$$\approx \frac{1}{N} \sum_i \nabla_{\theta} \log p_{\theta}(y^i | x^i)$$

PG = select good data + increase likelihood of selected data

What makes policy gradient challenging?

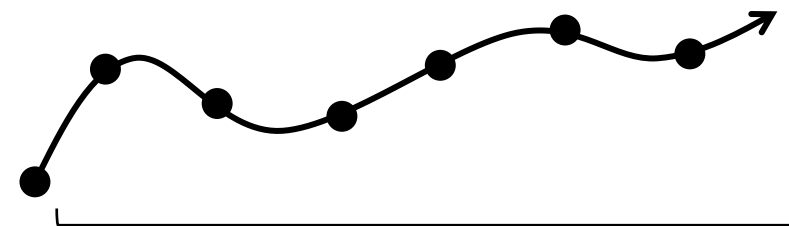
$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)$$

High variance estimator!!

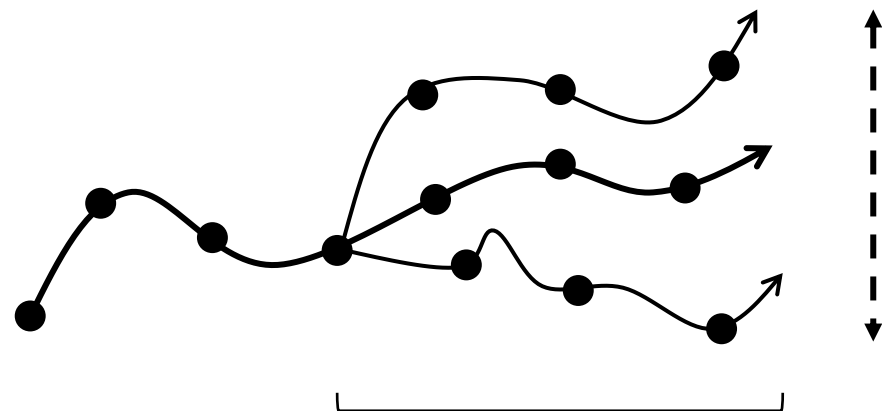
Hard to tell what matters without many samples

What we do



Single sample estimate

What we actually want



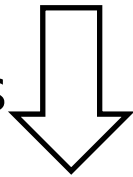
Averaged return estimate

Variance Reduction with Causality

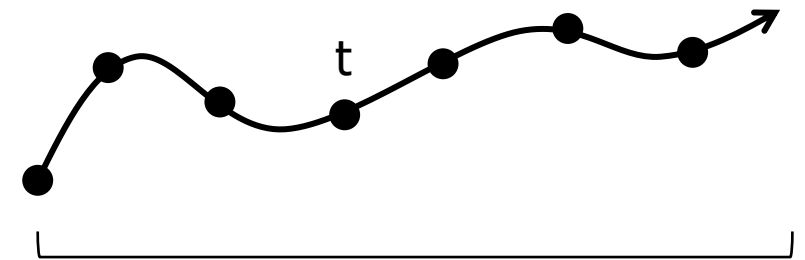
Idea: Trajectory returns depend on past and future, but we only care about the future, since actions cannot affect the past. Instead, consider “return-to-go”

$$\approx \frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \underbrace{\sum_{t'=0}^T r(s_{t'}^i, a_{t'}^i)}_{\text{Includes } t' < t}$$

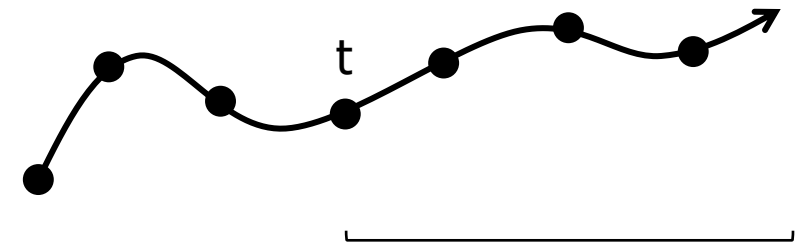
Ignore past terms



$$\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)$$

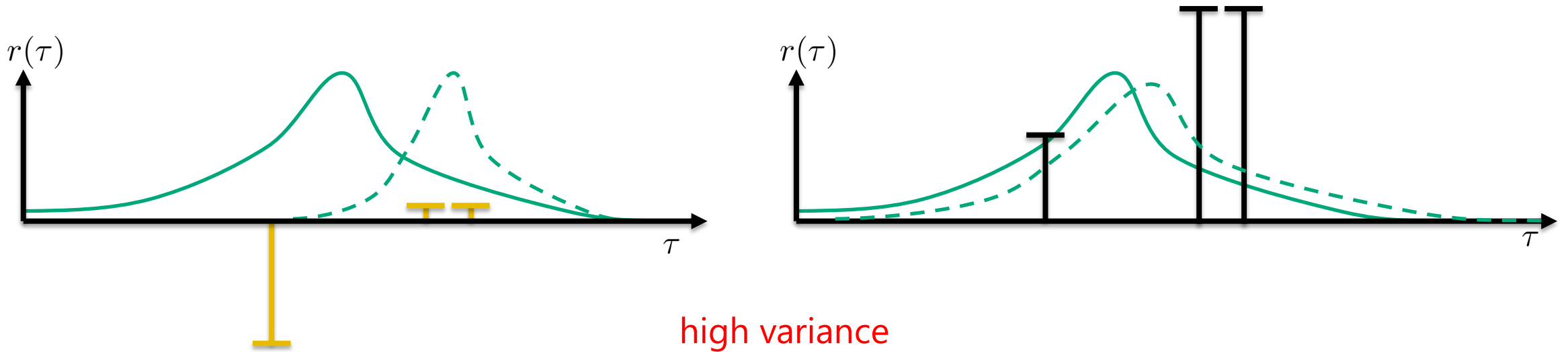


Full trajectory return



Return to go

Can we reduce variance further?




Arbitrarily uncentered, scaled returns can lead to huge variance:

- Imagine all rewards were positive, every action would be pushed up, some more than others
- What if instead, we pushed down some actions and pushed up some others (even if rewards are positive)

Variance Reduction with a Baseline

Idea: We can reduce variance by subtracting a current state dependent function from the policy gradient return

$$\frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \left[\sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i) - b(s_t) \right]$$


Baseline: Centers the returns, reduces variance

But does this increase bias??

Variance Reduction with a Baseline

$$\int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) - b(s_t) \right] ds_t da_t$$

$$\int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) \right] ds_t da_t - \int_{\mathcal{S}} \int_{\mathcal{A}} p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t) ds_t da_t$$

Let us show this is 0!

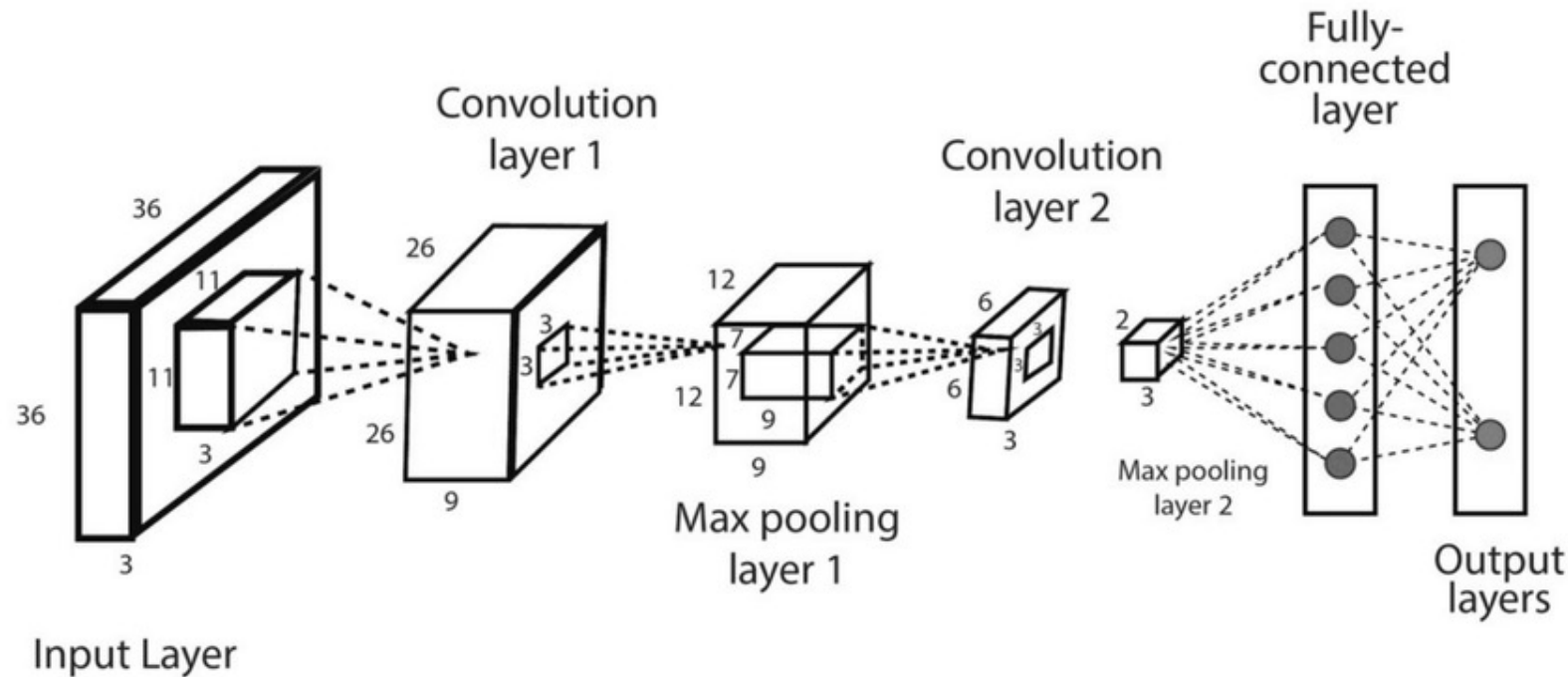
Variance Reduction with a Baseline

$$\begin{aligned}\int \int p(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [b(s_t)] ds_t da_t &= \int \int p(s_t) \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [b(s_t)] ds_t da_t \\&= \int p(s_t) b(s_t) \int \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) da_t ds_t \\&= \int p(s_t) b(s_t) \int \nabla_{\theta} \pi_{\theta}(a_t | s_t) da_t ds_t \\&= \int p(s_t) b(s_t) \nabla_{\theta} \int \pi_{\theta}(a_t | s_t) da_t ds_t = \int p(s_t) b(s_t) \nabla_{\theta} (1) ds_t = 0\end{aligned}$$

Unbiased!

Learning Baselines

Baselines are typically learned as deep neural nets from $R^s \rightarrow R^1$



$$\frac{1}{N} \sum_{j=1}^N \|\hat{V}(s_t^j, a_t^j) - \sum_{t=1}^H r(s_t^j, a_t^j)\|$$

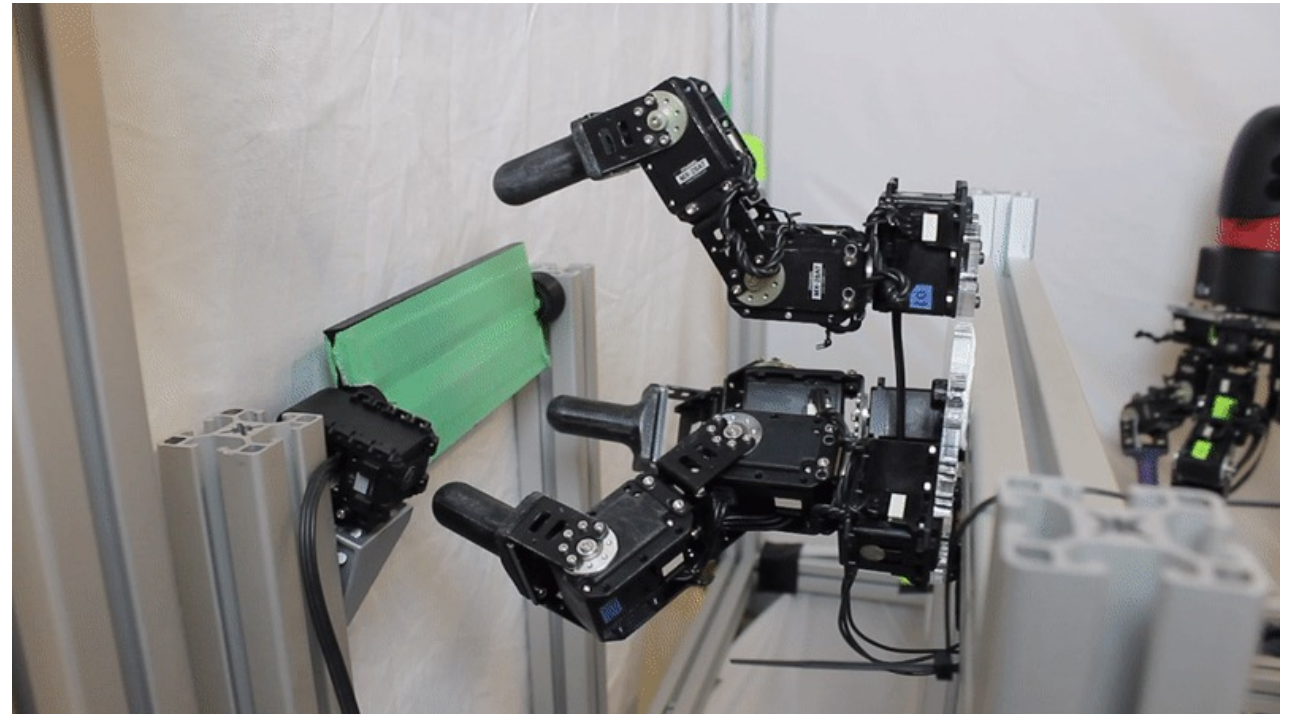
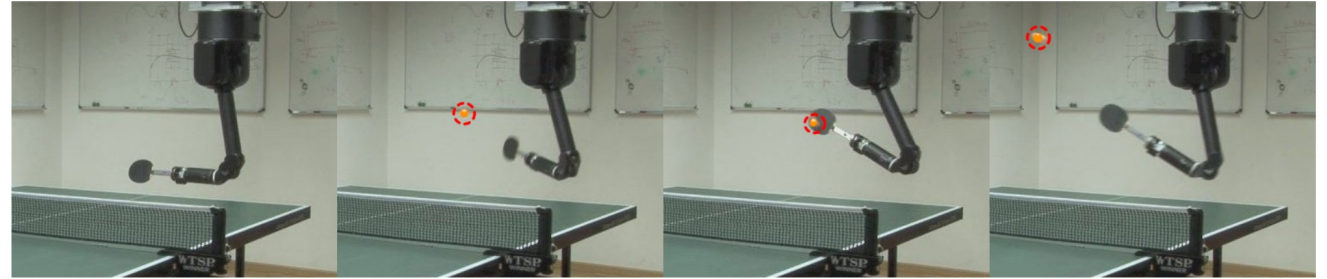
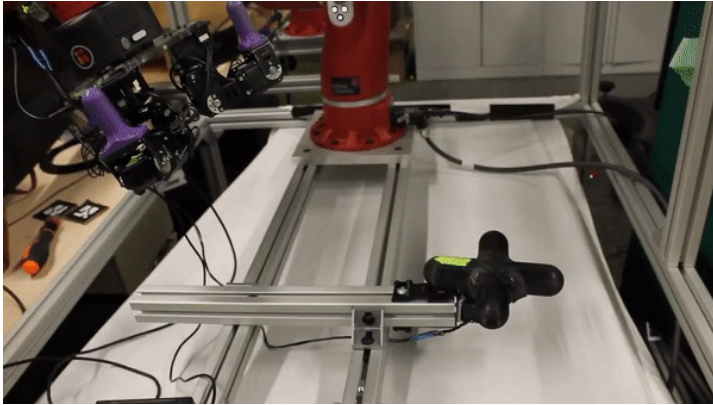
Minimize with Monte-carlo regression at every iteration, club with policy loss

$$A(s_t, a_t) = \sum_{t'=t}^T r(s_{t'}, a_{t'}) - V(s_t)$$

Allows us to define advantages

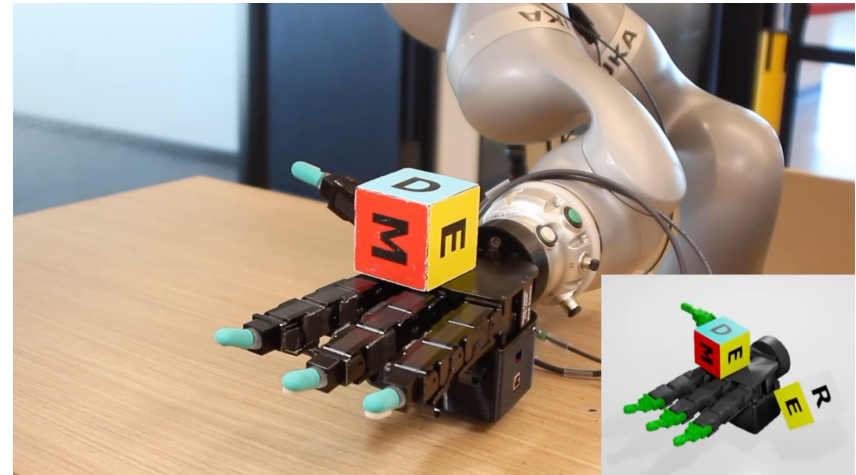
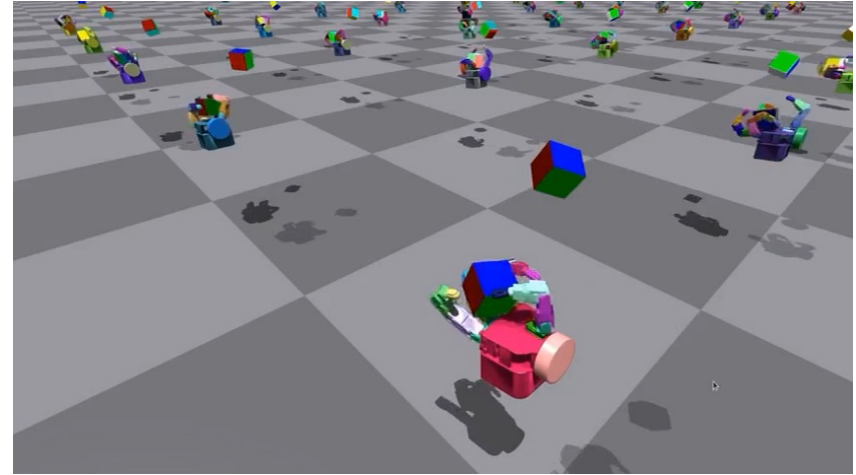
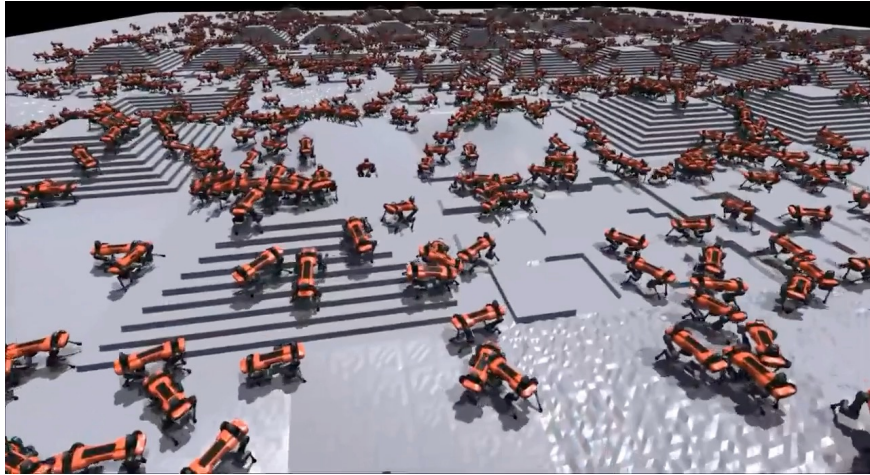
How is this useful for robotics?

Can be used to train robots in the real world but only in limited settings



How is this useful for robotics?

Largely useful for pretraining in simulation

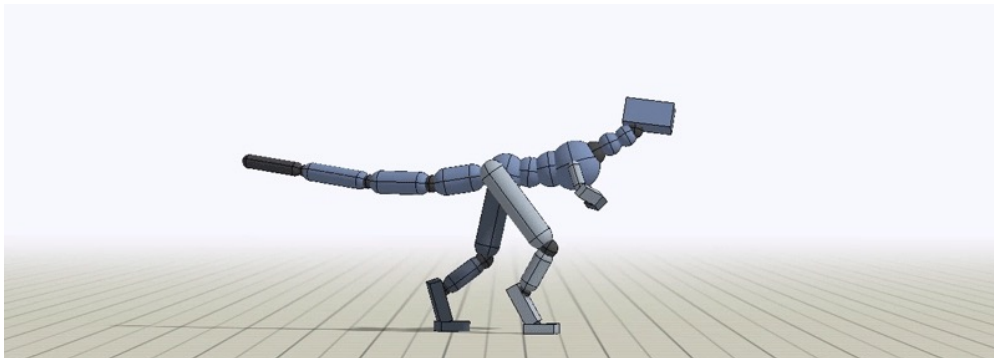
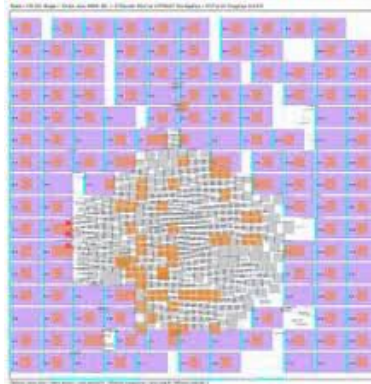


More in the sim2real lecture!

Pros/Cons of Policy Gradient Methods

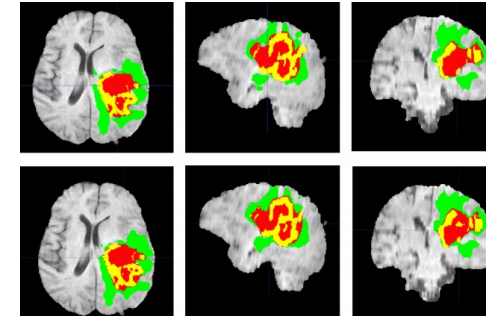
Pros

- Conceptually simple, easy to implement
- Stable, good asymptotic performance
- Compatible with deep models
- Require minimal modeling



Cons

- Sample inefficient
- Unable to reuse prior data effectively → on-policy
- Blackbox, can be hard to debug



Lecture Outline

Recap + Policy Gradient




Basic Actor Critic Algorithms



Making Actor-Critic Practical

Why is Policy Gradient sample inefficient?

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)\end{aligned}$$



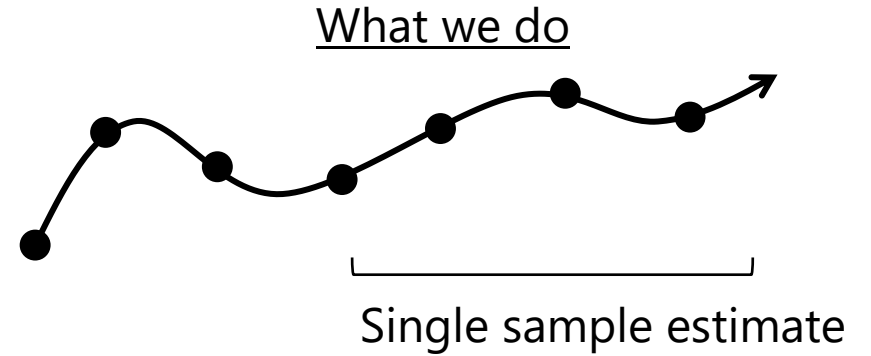
On-policy, unable to
effectively use past data

High Variance Estimator

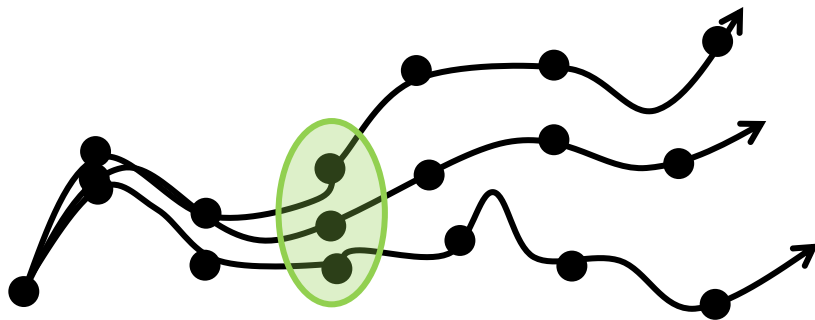
Can we develop a **low variance off-policy** RL algorithm that can bootstrap from prior data?

What can we do to lower variance?

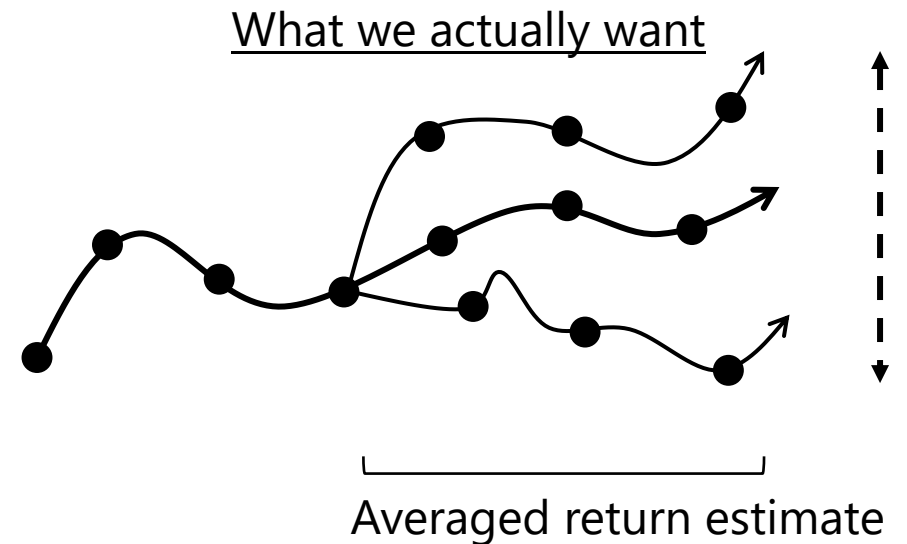
$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \underbrace{\sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)}\end{aligned}$$



Idea: bundle this across many (s, a) with a function approximator



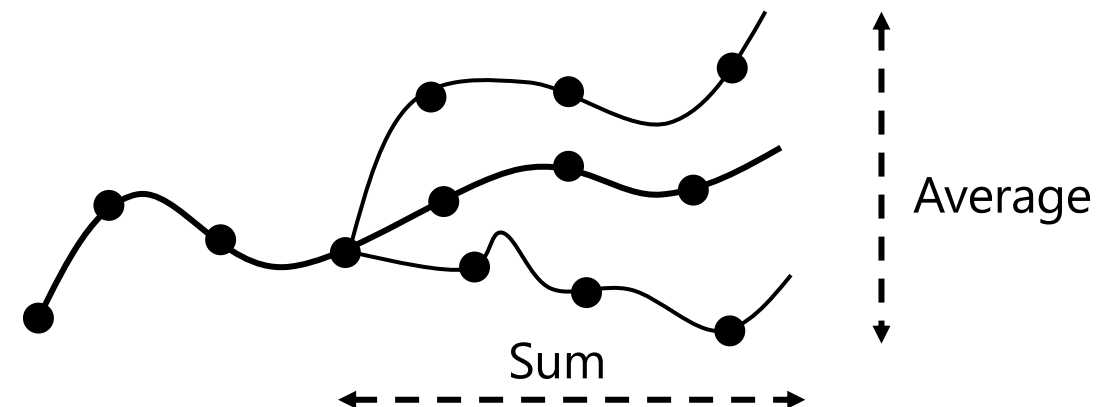
Function approximator bundles return estimates across states



Notation: Q functions

2
3

$$\frac{1}{N} \sum_{i=1}^N \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)$$



Expected sum of rewards in the future, starting from (s, a) on first step, then π

$$Q^{\pi}(s_t, a_t) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t'=t}^T r(s_{t'}', a_{t'}') | s_t, a_t \right] \quad \text{Bundles estimates across } (s, a)$$

Use the magic of (deep) function approximation

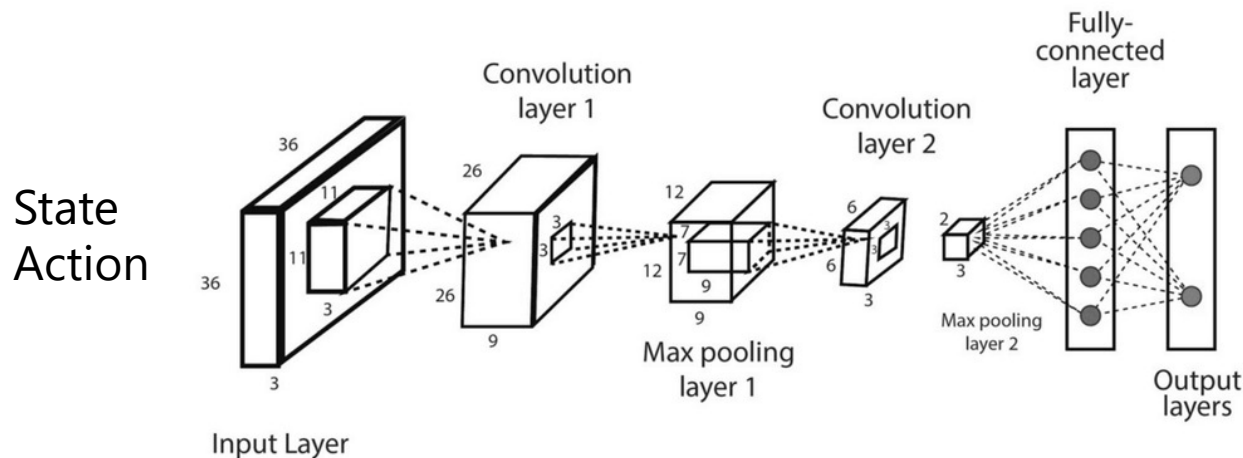
Estimation of Q-Functions

2
4

$$\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) Q^{\pi}(s_{t'}^i, a_{t'}^i)$$

$$Q^{\pi}(s_t, a_t) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) | s_t, a_t \right] \leftarrow \text{Monte-carlo approximation}$$

Idea: Regression from (s, a) to Monte-Carlo estimate

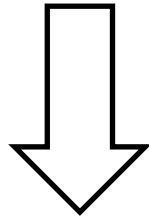


Unbiased, but high variance!

Can we do better?

2
4

$$\frac{1}{N} \sum_{i=1}^N \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i)$$



Much lower variance if estimated well

$$\frac{1}{N} \sum_{i=0}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) Q^{\pi}(s_{t'}^i, a_{t'}^i)$$

Can be learned off-policy!

Has special structure we can exploit!!

Recursive structure in Q functions

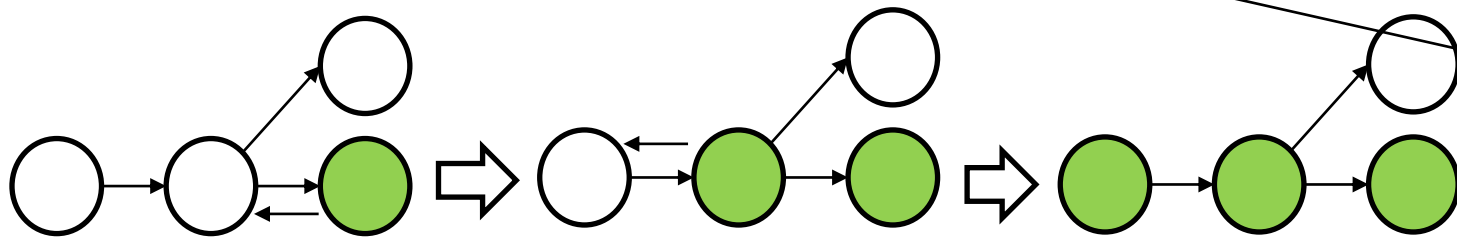
Q functions have special recursive structure!

$$Q^\pi(s_t, a_t) = \mathbb{E}_{\pi_\theta} \left[\sum_{t'=t}^T r(s'_t, a'_t) | s_t, a_t \right]$$

$$= r(s_t, a_t) + \mathbb{E}_\pi \left[\sum_{t'=t+1}^T r(s_{t'}, a_{t'}) | s_{t+1}, a_{t+1} \sim \pi(\cdot | s_{t+1}) \right]$$

$$Q^\pi(s_t, a_t) = r(s_t, a_t) + \mathbb{E}_{\substack{s_{t+1} \sim p(\cdot | s_t, a_t) \\ a_{t+1} \sim \pi_\theta(\cdot | s_{t+1})}} [Q^\pi(s_{t+1}, a_{t+1})]$$

Bellman equation



Decompose temporally via dynamic programming

Can be from different policies

Learning Q-functions via Dynamic Programming

Policy Evaluation: Try to minimize Bellman Error
(almost)

Bellman
equation

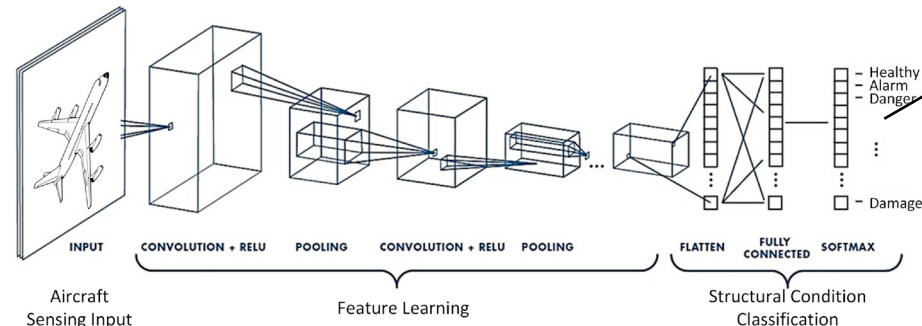
$$Q^\pi(s_t, a_t) = r(s_t, a_t) + \mathbb{E}_{\substack{s_{t+1} \sim p(\cdot | s_t, a_t) \\ a_{t+1} \sim \pi_\theta(\cdot | s_{t+1})}} [Q^\pi(s_{t+1}, a_{t+1})]$$

Same function approximator

How can we convert this recursion into a learning objective?

$$\min_{\phi} \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left(Q_{\phi}^{\pi}(s_t, a_t) - \left(r(s_t, a_t) + \mathbb{E}_{a_{t+1} \sim \pi_{\theta}(a_{t+1} | s_{t+1})} \left[Q_{\hat{\phi}}^{\pi}(s_{t+1}, a_{t+1}) \right] \right) \right)^2$$

Off-
policy



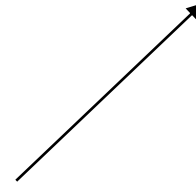
Can train via
GD!

Note: this may look like gradient descent on Bellman error, it is not!

Improving Policies with Learned Q-functions

Policy Improvement: Improve policy with **policy gradient**

$$\max_{\theta} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\theta}(a|s)} [Q^{\pi_{\theta}}(s, a)]$$



Replace Monte-Carlo sum of rewards with learned Q function

Lowers variance compared to policy gradient!

Policy Updates – REINFORCE or Reparameterization

Let's look a little deeper into the policy update

$$\max_{\theta} J(\theta) = \max_{\theta} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [Q^{\pi}(s, a)]$$

Likelihood Ratio/Score Function

Pathwise derivative/Reparameterization

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a)]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{z \sim p(z)} [\nabla_a Q^{\pi}(s, a)|_{a=\mu_{\theta}+z\sigma_{\theta}} \nabla_{\theta}(\mu_{\theta} + z\sigma_{\theta})]$$

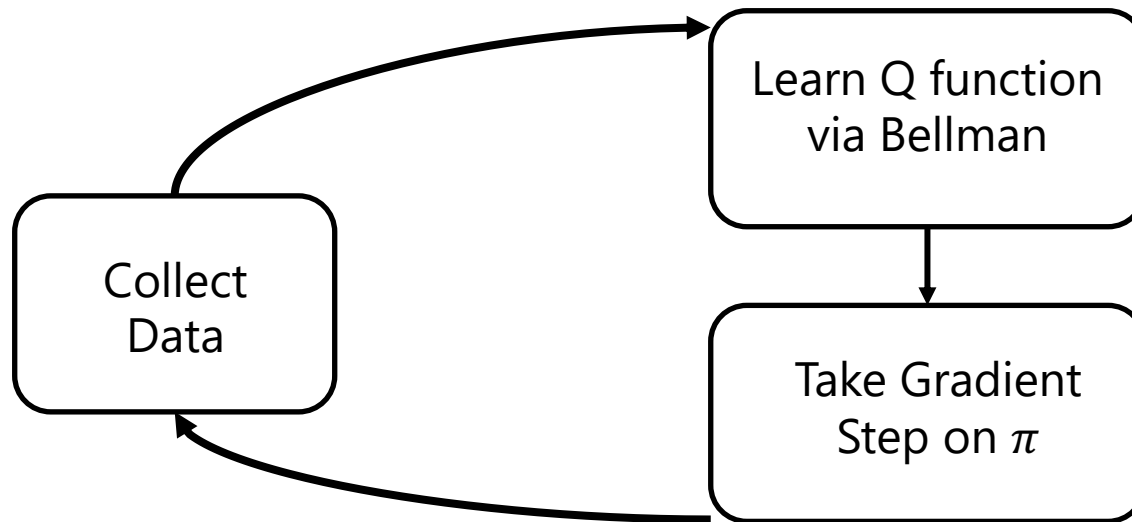
Easier to Apply to Broad Policy Class

Lower variance

Actor-Critic: Policy Gradient in terms of Q functions

Critic: learned via the Bellman update (Policy Evaluation)

$$\min_{\phi} \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left[\left(Q_{\phi}^{\pi}(s_t, a_t) - \left(r(s_t, a_t) + \mathbb{E}_{a_{t+1} \sim \pi(\cdot | s_{t+1})} [Q_{\phi}(s_{t+1}, a_{t+1})] \right) \right)^2 \right]$$

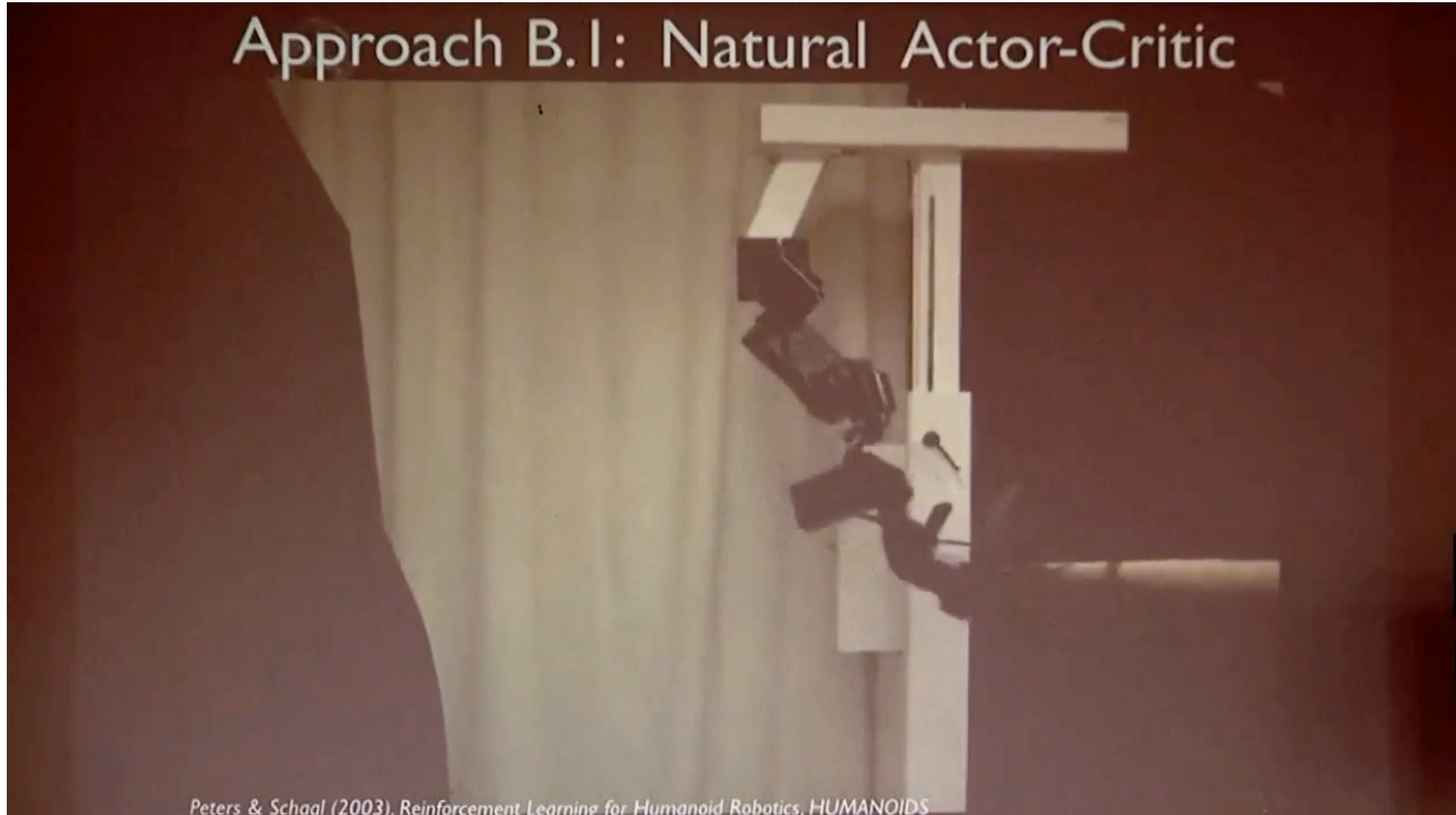


Lowers variance and is off-policy!

Actor: updated using learned critic (Policy Improvement)

$$\max_{\pi} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot | s)} [Q^{\pi}(s, a)]$$

Actor-Critic in Action



Lecture Outline

Recap + Policy Gradient



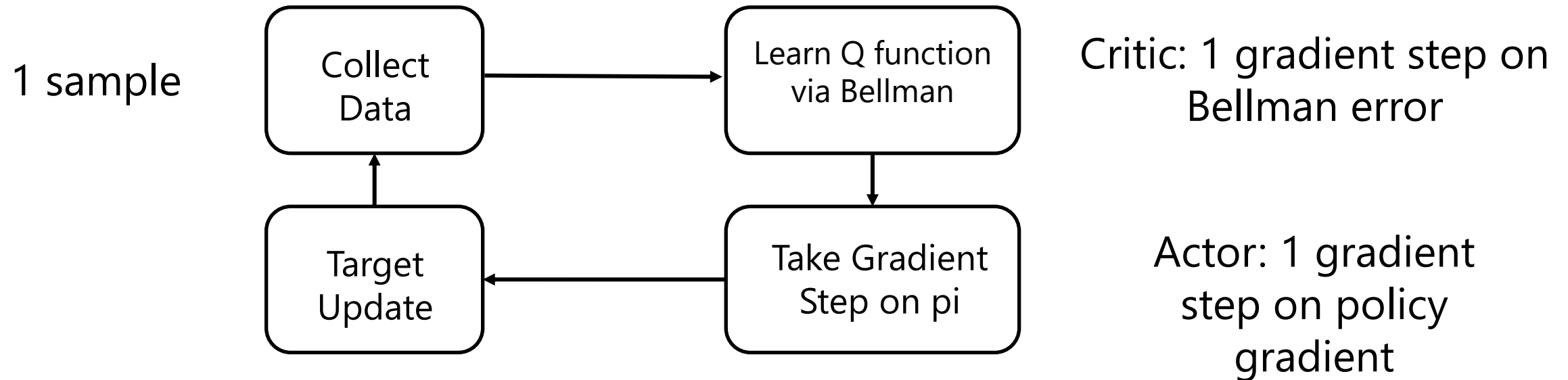
Basic Actor Critic Algorithms



Making Actor-Critic Practical

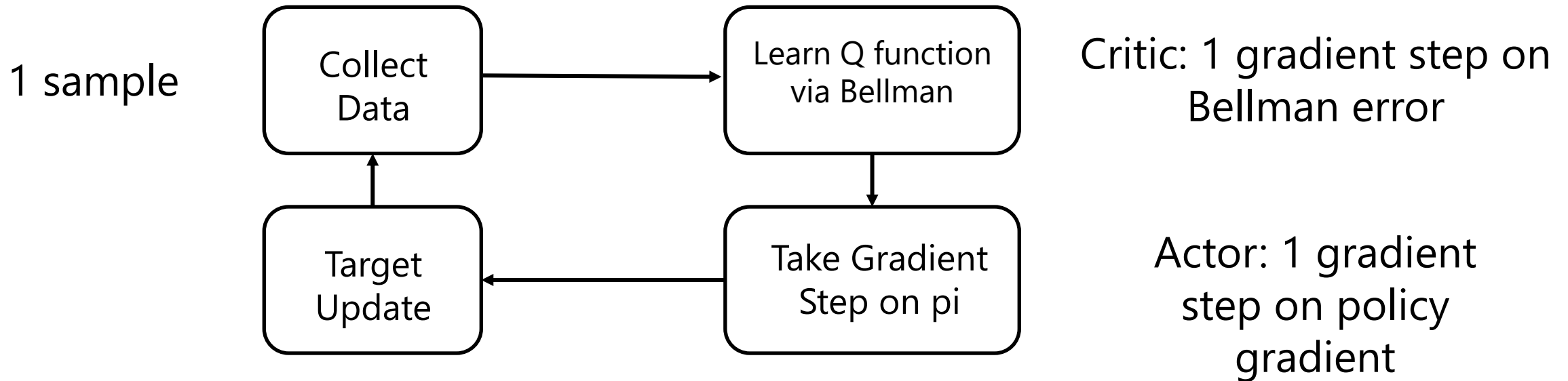
Going from Batch Updates to Online Updates

This algorithm can go from full batch mode to fully online updates



Allows for much more immediate updates

Challenges of doing online updates

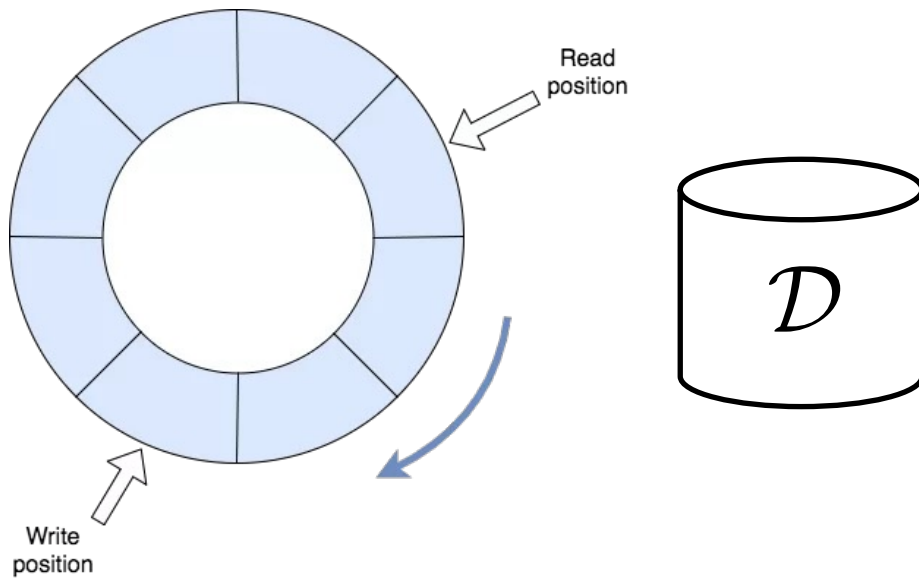


- When updates are performed online, two issues persist:
1. Correlated updates since samples are correlated
 2. Optimization objective changes constantly, unstable

Decorrelating updates with replay buffers

2
7

Updates can be decorrelated by storing and shuffling data in a replay buffer



Instead of doing updates in order,
sample batches from replay buffer

How?

Sampled from replay buffer

$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[Q_{\phi}^{\pi}(s_t, a_t) - (r(s_t, a_t) + \mathbb{E}_{a_{t+1} \sim \pi(\cdot | s_{t+1})} [Q_{\bar{\phi}}(s_{t+1}, a_{t+1})]) \right]^2$$
$$\max_{\pi} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot | s)} [Q^{\pi}(s, a)]$$

1. Sample uniformly
2. Prioritize by TD-error
3. Prioritize by target error
4. ... open area of research!

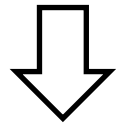
Slowing moving targets with target networks

2
7

Continuous updates can be unstable since there is a churn of projection and backup

$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[Q_{\phi}^{\pi}(s_t, a_t) - (r(s_t, a_t) + \mathbb{E}_{a_{t+1} \sim \pi(\cdot | s_{t+1})} [Q_{\bar{\phi}}(s_{t+1}, a_{t+1})]) \right]^2$$

If we set $\bar{\phi}$ to ϕ every update, the update becomes very unstable



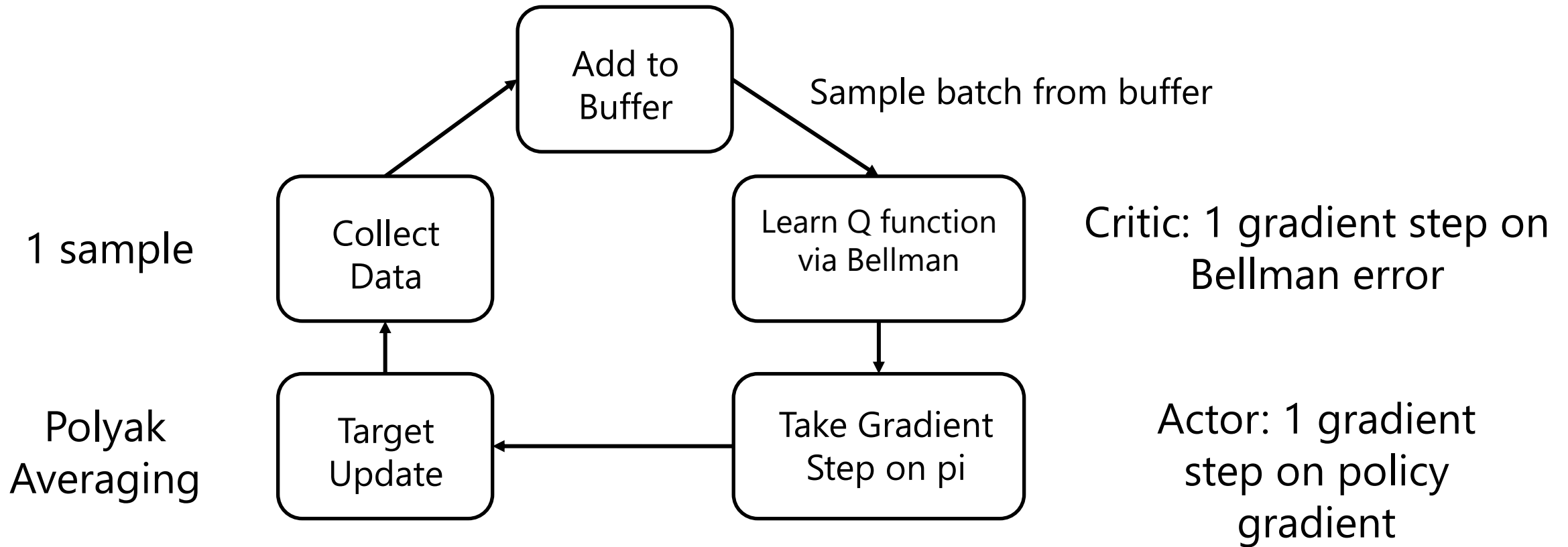
Move $\bar{\phi}$ to ϕ slowly!

$$\bar{\phi} = (1 - \tau)\phi + \tau\bar{\phi}$$

Polyak averaging

A Practical Off-Policy RL Algorithm

2
7



Simplify -- Can we get rid of a parametric actor?

Critic Update

$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[Q_{\phi}^{\pi}(s_t, a_t) - (r(s_t, a_t) + \mathbb{E}_{a_{t+1} \sim \pi(\cdot | s_{t+1})} [Q_{\bar{\phi}}(s_{t+1}, a_{t+1})]) \right]^2$$

Actor Update

$$\max_{\pi} \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot | s)} [Q^{\pi}(s, a)]$$

What if we removed this explicit actor completely?

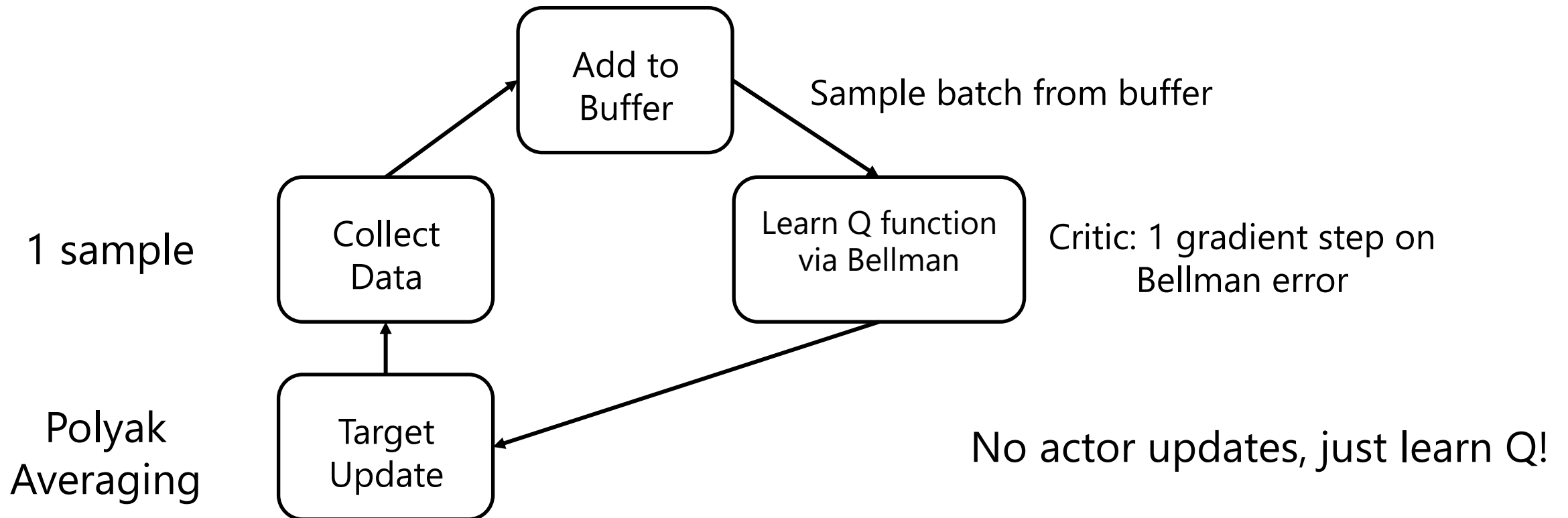


Directly Learning Q*


$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[\left[Q_{\phi}^{\pi}(s_t, a_t) - (r(s_t, a_t) + \max_{a_{t+1}} [Q_{\phi}(s_{t+1}, a_{t+1})]) \right]^2 \right]$$

$$\pi(a|s) = \max_a Q(s, a)$$

Directly do max in the Bellman update



How can we maximize w.r.t a ?

$$\pi(a|s) = \max_a Q(s, a)$$


Analytic maximization can be very difficult to perform in continuous action spaces
Much easier in discrete spaces! → just do categorical max!

Some ideas to do maximization:

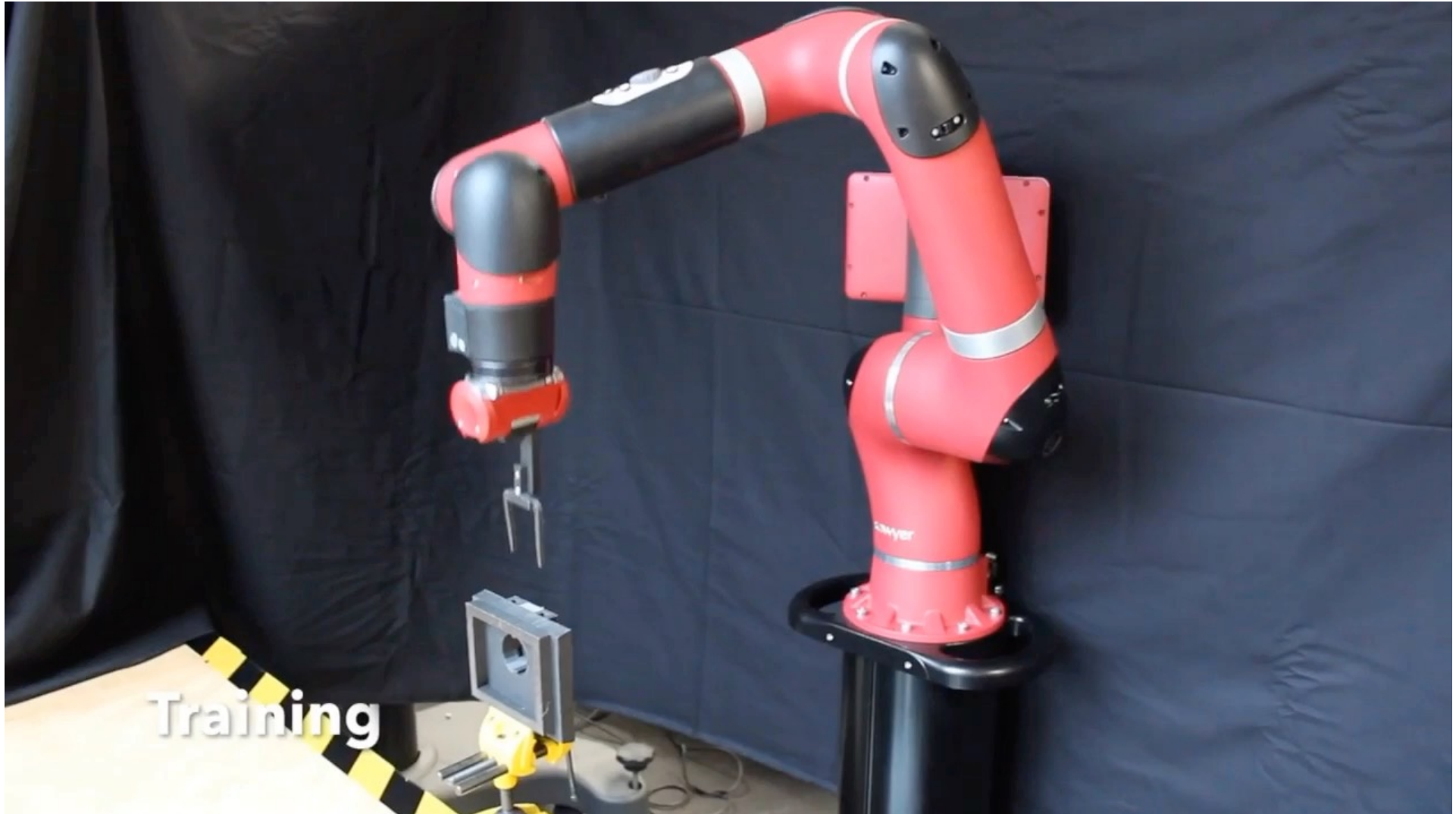
1. Sampling based (QT-opt (Kalashnikov et al))
2. Optimization based (NAF, Gu et al)

Practical Actor-Critic in Action



Trained using QT-Opt

Practical Actor-Critic in Action

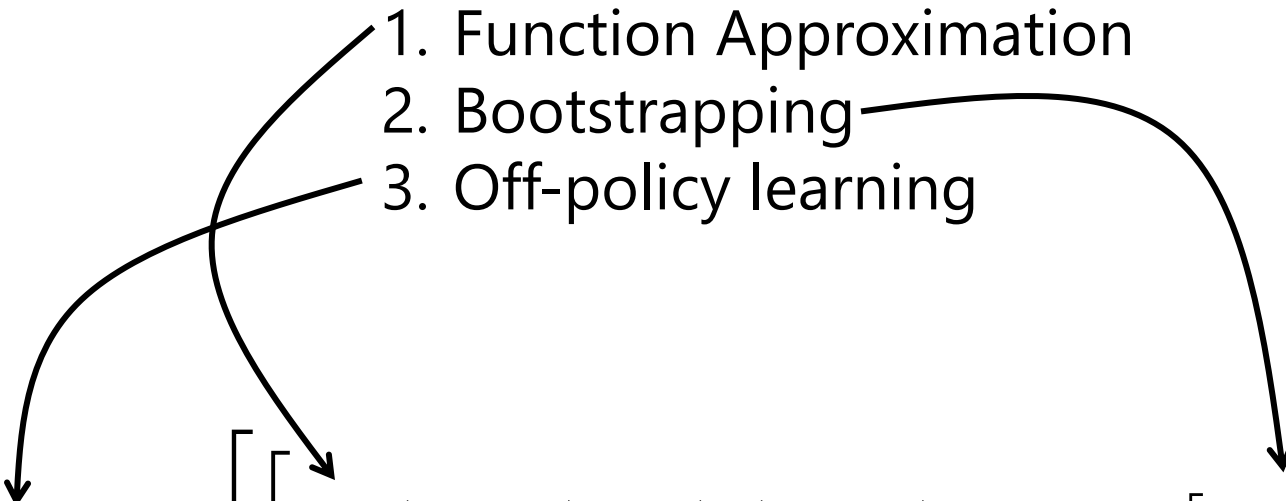


Trained using DDPG

What makes off-policy RL hard?

Deadly triad:

1. Function Approximation
2. Bootstrapping
3. Off-policy learning



The diagram illustrates the 'Deadly triad' of off-policy reinforcement learning. Three components are listed: 1. Function Approximation, 2. Bootstrapping, and 3. Off-policy learning. Arrows from each component point to a mathematical expression representing the off-policy RL loss function. The expression is the squared temporal difference error, averaged over the data distribution \mathcal{D} .

$$\min_{\phi} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[\left[Q_{\phi}^{\pi}(s_t, a_t) - (r(s_t, a_t) + \max_{a_{t+1}} [Q_{\bar{\phi}}(s_{t+1}, a_{t+1})]) \right]^2 \right]$$

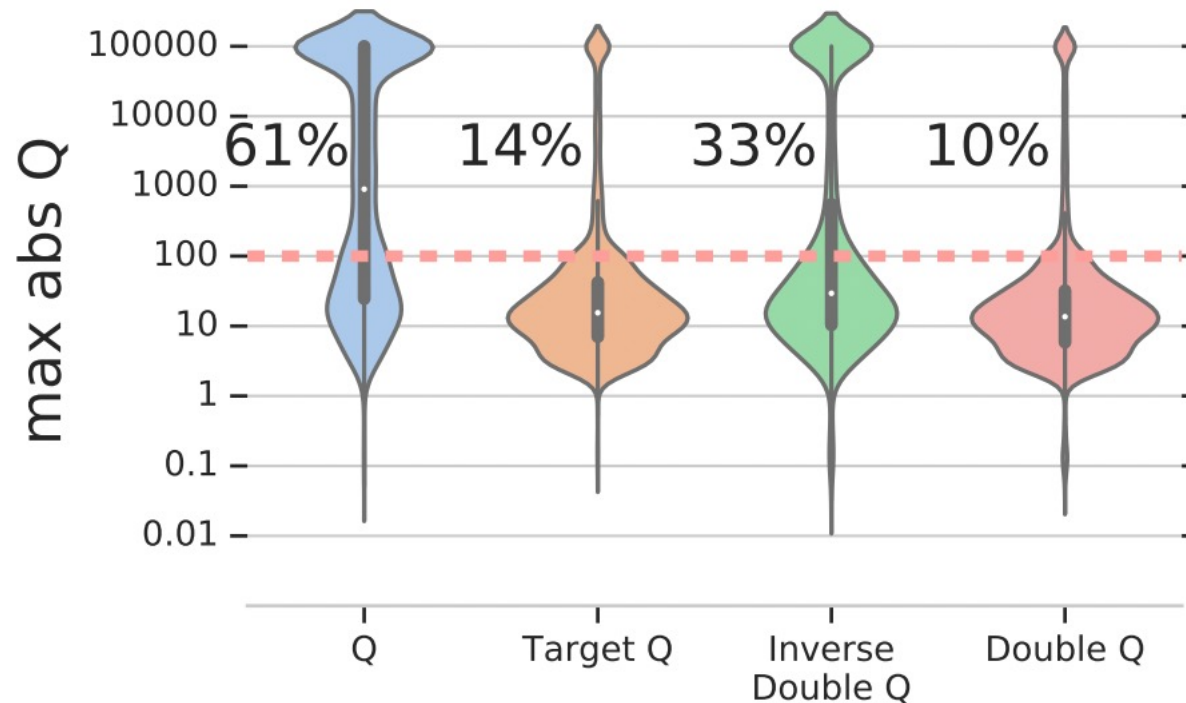
These in combination lead to many of the difficulties in stabilizing off-policy RL with function approximation

Zooming out – what makes off-policy RL hard?

Deadly triad:

1. Function Approximation
2. Bootstrapping
3. Off-policy learning

61% of runs show divergence of Q-values



Diverges even with linear function approximation, when off-policy + bootstrapping

Lecture Outline

Recap + Policy Gradient



Basic Actor Critic Algorithms



Making Actor-Critic Practical

Recap: Course Overview

Filtering/Smoothing

Localization

Mapping

SLAM

Search

Motion Planning

TrajOpt

Stability/Certification

MDPs and RL

Imitation Learning

Off-Policy/MBRL