

**CS 295:**  
**Optimal Control and**  
**Reinforcement Learning**  
**Winter 2020**

**Lecture 7: Actor–Critic Methods**

Roy Fox

Department of Computer Science

Bren School of Information and Computer Sciences

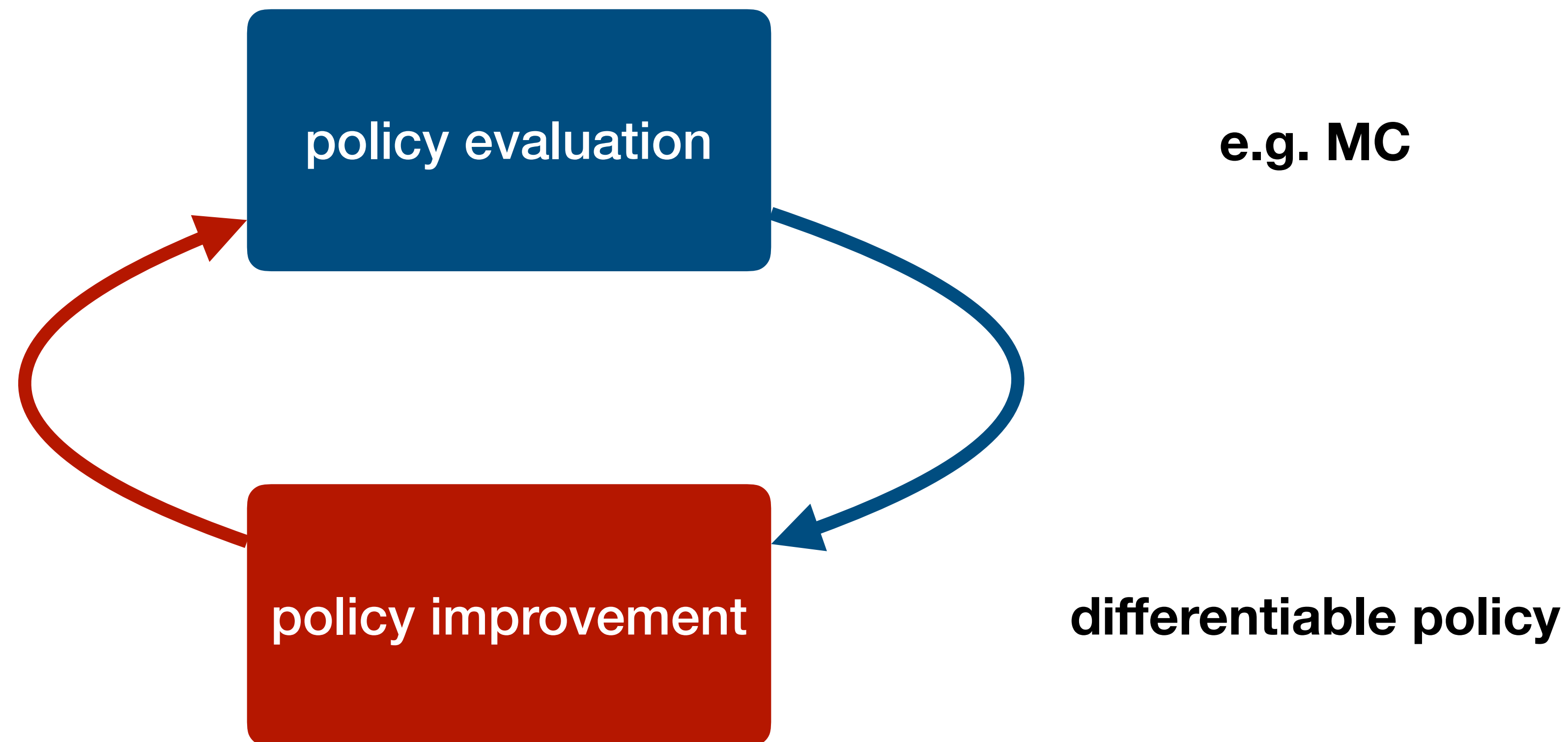
University of California, Irvine

# Today's lecture

---

- Policy Gradient with learned value function
- Actor–Critic methods
- Advantage estimation
- n-step TD and TD( $\lambda$ )

# Policy-based methods



# Policy Gradient (PG) with reward-to-go

$$\begin{aligned}\nabla_{\theta} \mathcal{J}_{\theta} &= \mathbb{E}_{\xi \sim p_{\theta}} [\nabla_{\theta} \log p_{\theta}(\xi) R] \\ &= \mathbb{E}_{\xi \sim p_{\theta}} \left[ \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R \right] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta}} [\mathbb{E}_{a_t | s_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R]] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta}} [\nabla_{\theta} \mathbb{E}_{a_t | s_t \sim \pi_{\theta}} [R]] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta}} \left[ \nabla_{\theta} \mathbb{E}_{a_t | s_t \sim \pi_{\theta}} \left[ \sum_{t'} r_{t'} \right] \right] \\ &= \sum_t \mathbb{E}_{s_t \sim p_{\theta}} \left[ \nabla_{\theta} \mathbb{E}_{a_t | s_t \sim \pi_{\theta}} \left[ \sum_{t' \geq t} r_{t'} \right] \right]\end{aligned}$$

- With finite horizon:

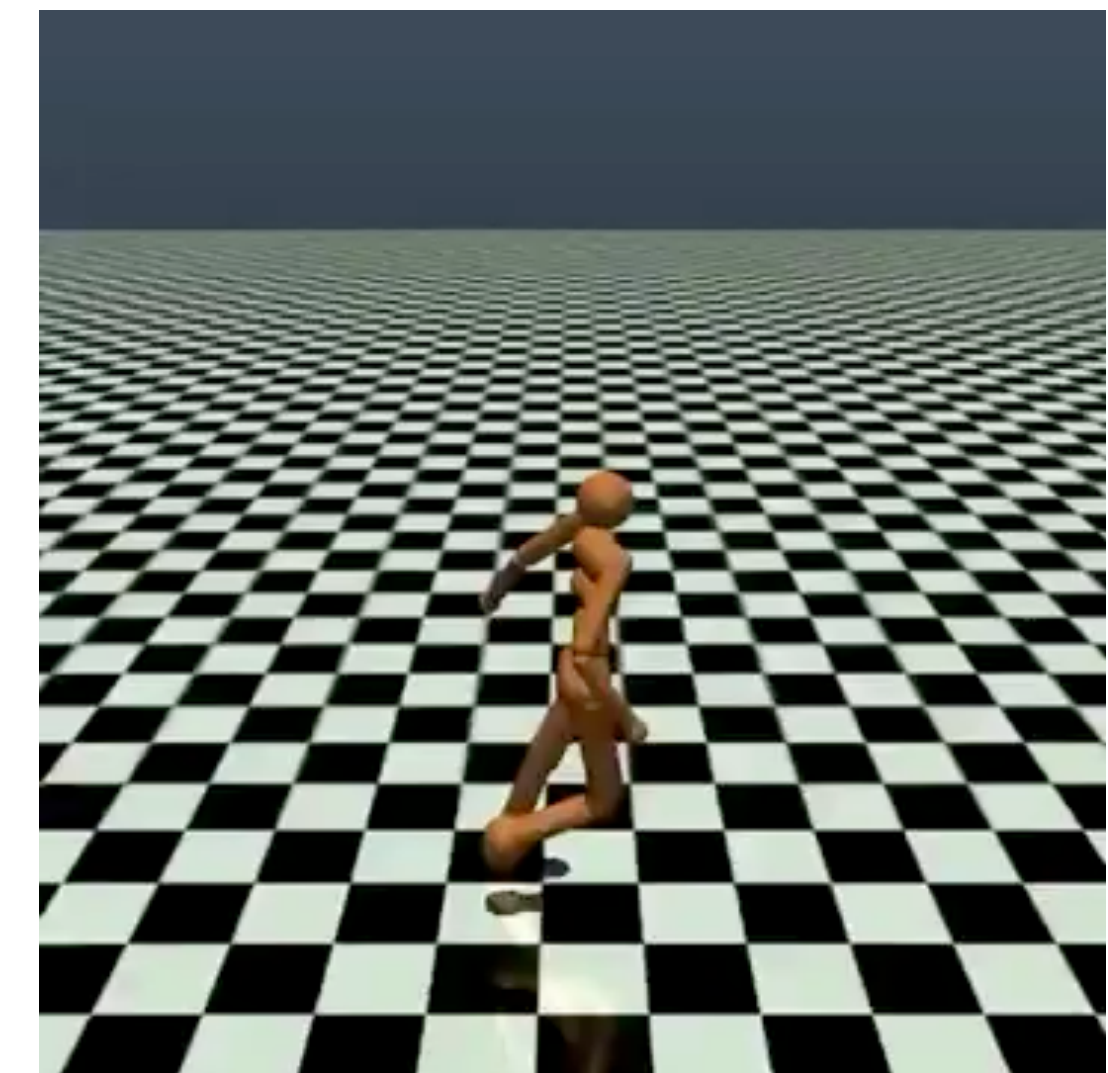
- Gradient estimator

independent of past rewards

# Discounted case

$$\begin{aligned}\nabla_{\theta} \mathcal{J}_{\theta} &= \sum_t \mathbb{E}_{s_t \sim p_{\theta}} \mathbb{E}_{a_t | s_t \sim \pi_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t' \geq t} \gamma^{t'} r_{t'} \right] \\ &= \sum_t \gamma^t \mathbb{E}_{s_t \sim p_{\theta}} \mathbb{E}_{a_t | s_t \sim \pi_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t' \geq t} \gamma^{t'-t} r_{t'} \right]\end{aligned}$$

- Does this make sense?
- Discounting is effectively finite-horizon
- If any state could be initial state
  - Discounting is just a computational / statistical trick
  - Don't discount the contribution of  $s_t$  to the loss
  - That's what most algorithms do



# Policy-Gradient Theorem

$$\begin{aligned}\nabla_{\theta} \mathcal{J}_{\theta} &= \sum_t \gamma^t \mathbb{E}_{s_t \sim p_{\theta}} [\mathbb{E}_{a_t | s_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R_{\geq t}]] \\ &\stackrel{?!}{=} \sum_t \gamma^t \mathbb{E}_{s_t \sim p_{\theta}} [\mathbb{E}_{a_t | s_t \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q_{\pi_{\theta}}(s_t, a_t)]]\end{aligned}$$

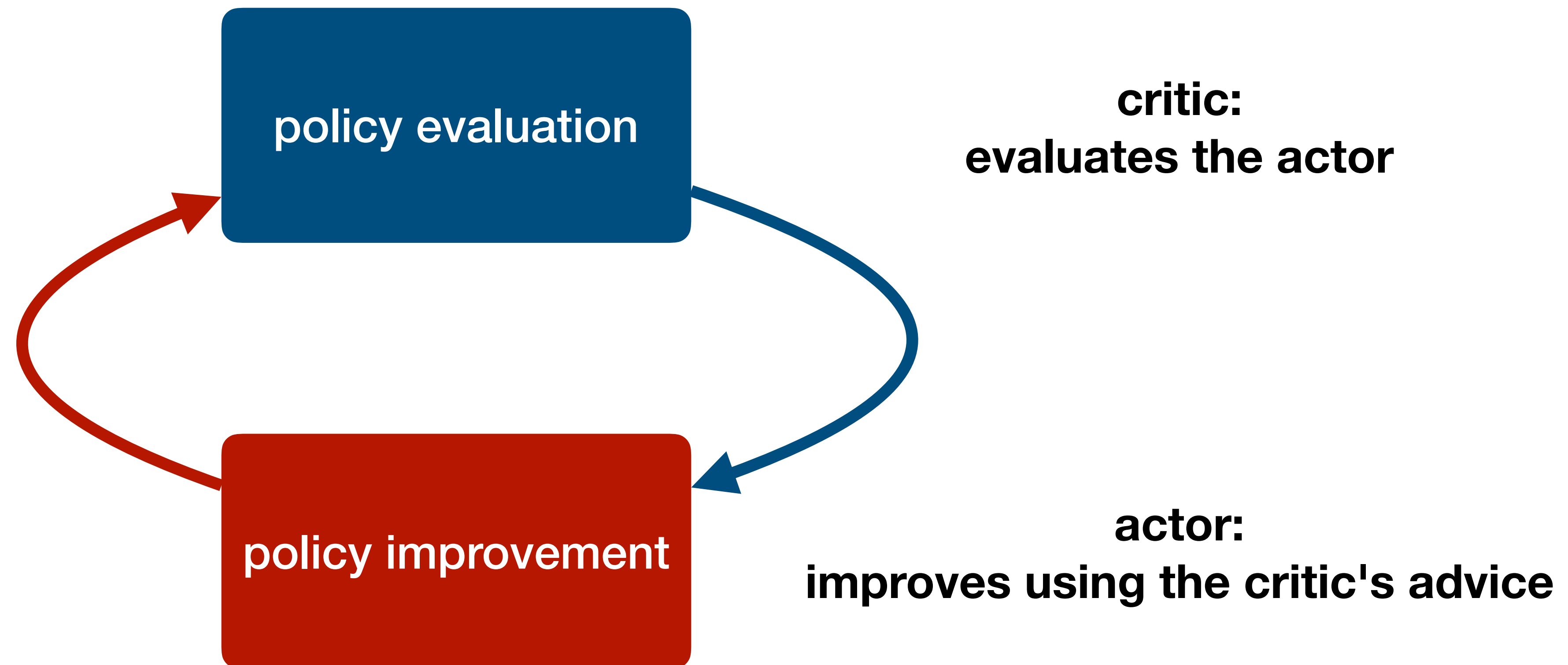
$$\begin{aligned}\nabla_{\theta} V_{\pi_{\theta}}(s) &= \nabla_{\theta} \mathbb{E}_{a | s \sim \pi_{\theta}} [Q_{\pi_{\theta}}(s, a)] \\ &= \sum_a (\nabla_{\theta} \pi_{\theta}(a | s) Q_{\pi_{\theta}}(s, a) + \pi_{\theta}(a | s) \nabla_{\theta} Q_{\pi_{\theta}}(s, a)) \\ &= \mathbb{E}_{a | s \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\pi_{\theta}}(s, a) + \nabla_{\theta} (r(s, a) + \gamma \mathbb{E}_{s' | s, a} [V_{\pi_{\theta}}(s')])] \\ &= \mathbb{E}_{a | s \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\pi_{\theta}}(s, a) + \gamma \mathbb{E}_{s' | s, a} [\nabla_{\theta} V_{\pi_{\theta}}(s')]]\end{aligned}$$

- Back-prop ~ Bellman recursion

# Actor–Critic methods

$$\mathcal{L}_\theta(s, a) = -\nabla_\theta \log \pi_\theta(a|s) Q_\phi(s, a)$$

$$\mathcal{L}_\phi(s, a, r, s') = (r + \gamma \mathbb{E}_{a'|s' \sim \pi_\theta} [Q_{\bar{\phi}}(s', a')] - Q_\phi(s, a))^2$$



# Baselines

---

$$\nabla_{\theta} \mathcal{J}_{\theta} = \mathbb{E}[\nabla_{\theta} \log \pi_{\theta}(a|s) (Q_{\pi_{\theta}}(s, a) - b) | s, a]$$

- $b$  can be any variable independent of  $a$  given  $s$ 
  - Can depend on the past, but not the future
- Previously, we used  $b = \frac{1}{N} \sum_i R_i$
- This suggests using  $b = V_{\pi_{\theta}}(s)$



# Advantage estimation

---

$$\nabla_{\theta} \mathcal{J}_{\theta} = \mathbb{E}[\nabla_{\theta} \log \pi_{\theta}(a|s) A_{\pi_{\theta}}(s, a) | s, a]$$

- How to estimate  $A_{\pi_{\theta}}(s, a)$ ?

$$A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s) = r(s, a) + \gamma \mathbb{E}_{s'|s, a \sim p}[V_{\pi}(s')] - V(s)$$

$$\hat{A}(s, a) \approx r + \gamma V_{\phi}(s') - V_{\phi}(s)$$

# An Actor–Critic algorithm

---

---

## Algorithm 1 Actor–Critic

---

get on-policy sample  $(s, a, r, s')$

take gradient step on  $\mathcal{L}_\phi = (r + \gamma V_\phi(s') - V_\phi(s))^2$

compute  $\hat{A}(s, a) = r + \gamma V_\phi(s') - V_\phi(s)$

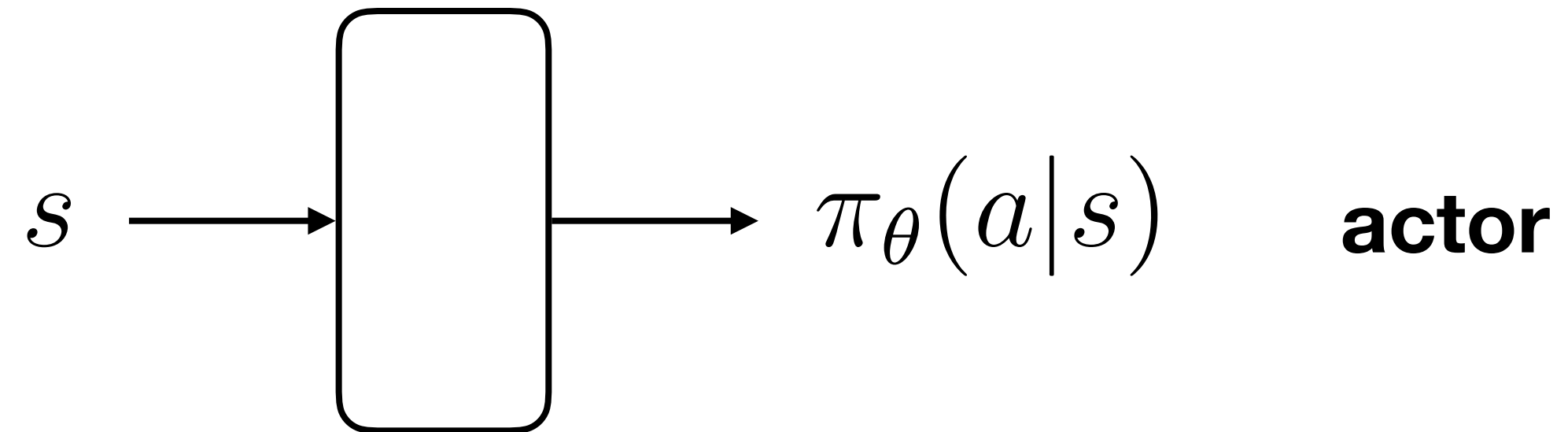
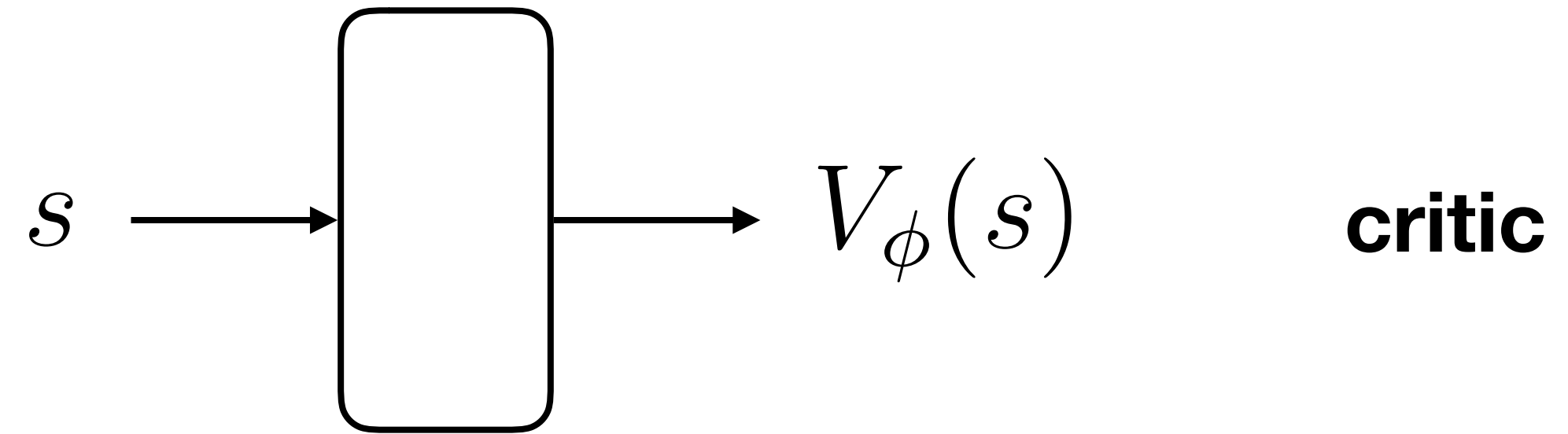
take gradient step  $\nabla_\theta \log \pi_\theta(a|s) \hat{A}(s, a)$

repeat

---

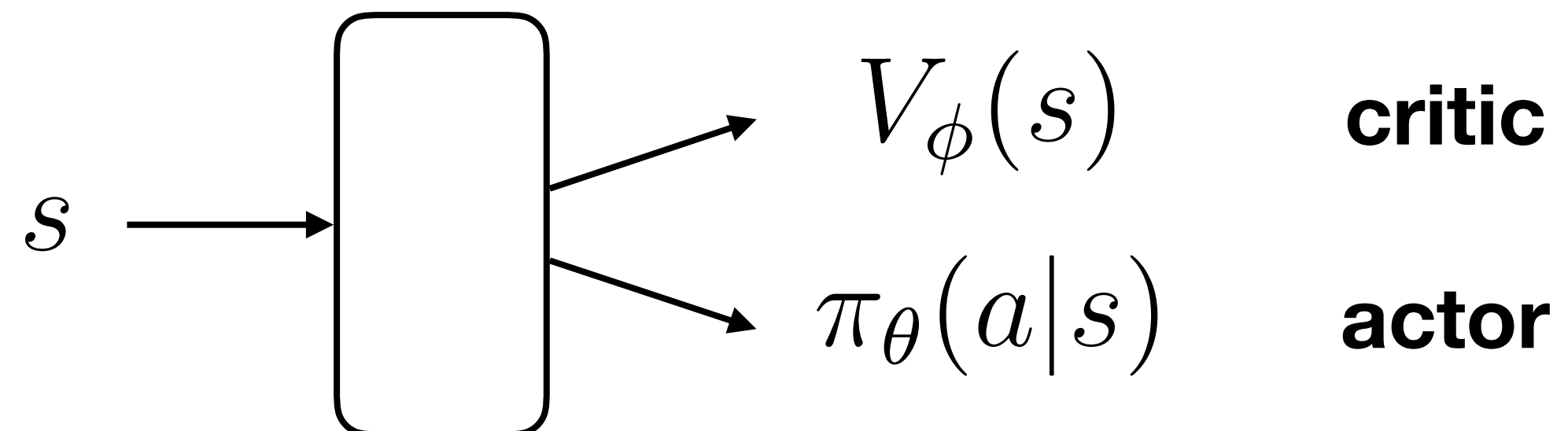
# Practical considerations: param sharing

- Separate parameters:



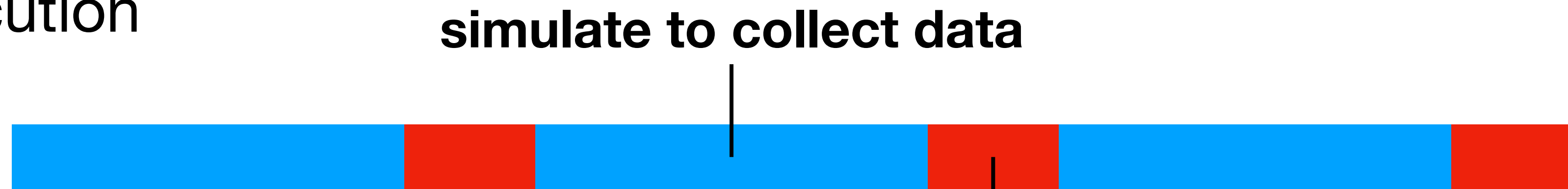
- Shared parameters:

- Can be more data efficient
- Can be less stable

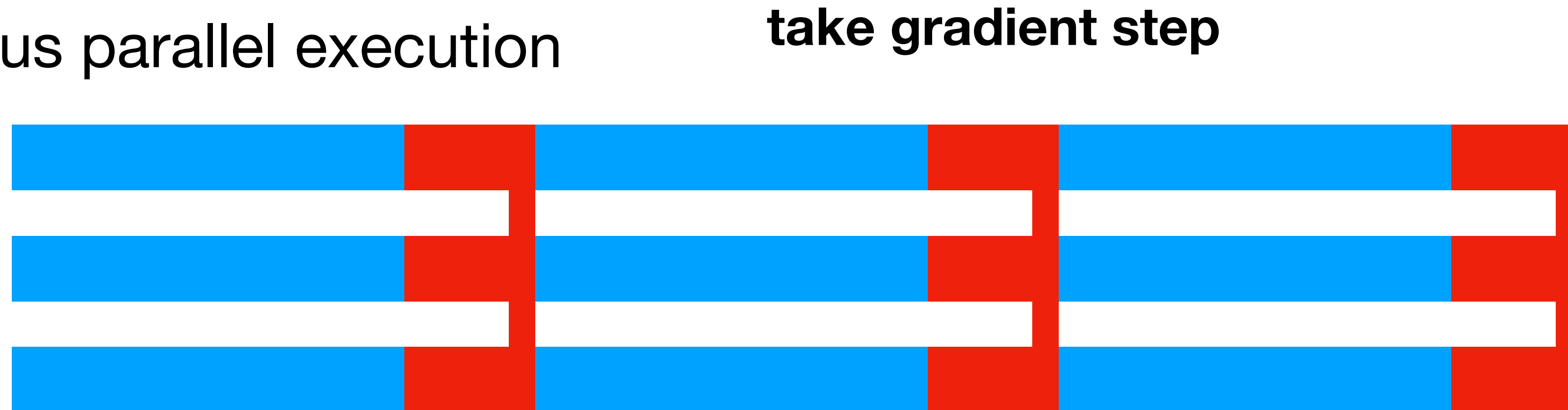


# Practical considerations: distributed comp.

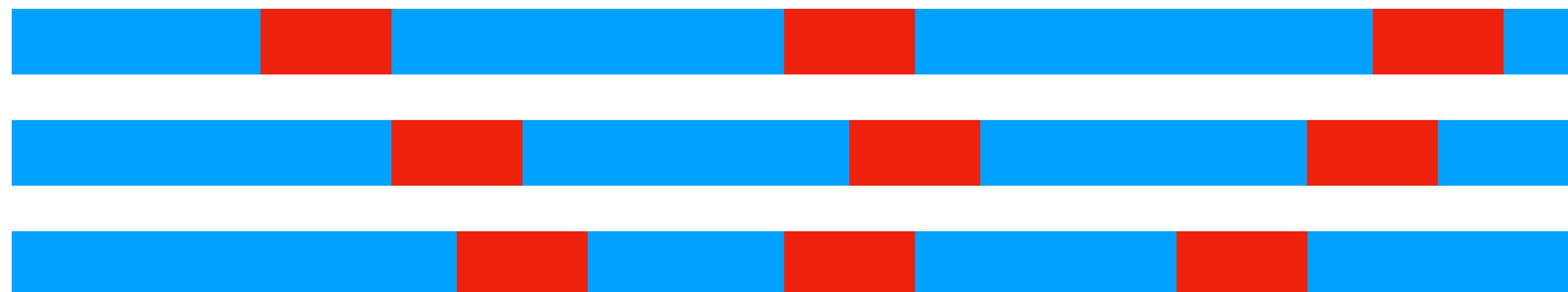
- Serial execution



- Synchronous parallel execution



- Asynchronous parallel execution



# More advantage estimation

---

$$\hat{A}(s_t, a_t) \approx \sum_{t' \geq t} \gamma^{t'-t} r_{t'} - V_{\phi}(s_t)$$

- Asynchronous Advantage Actor Critic (A3C):
  - MC advantage estimation + asynchronous parallel execution
- Advantage Actor Critic (A2C): same but serial

# Comparing advantage estimators

	<b>bias</b>	<b>variance</b>
$\nabla_{\theta} \mathcal{J}_{\theta} \approx \nabla_{\theta} \log \pi_{\theta}(a s) \left( \sum_{t' \geq t} \gamma^{t'-t} r_{t'} - b \right)$	<b>none</b>	<b>high</b> one grad per traj
$\nabla_{\theta} \mathcal{J}_{\theta} \approx \nabla_{\theta} \log \pi_{\theta}(a s) (r + \gamma V_{\phi}(s') - V_{\phi}(s))$	<b>some</b> approx value	<b>lower</b>
$\nabla_{\theta} \mathcal{J}_{\theta} \approx \nabla_{\theta} \log \pi_{\theta}(a s) \left( \sum_{t' \geq t} \gamma^{t'-t} r_{t'} - V_{\phi}(s) \right)$	<b>none</b>	<b>mid</b> state-dependent baseline

# n-step TD

- 1-step TD:  $\hat{A}_t^1 = r_t + \gamma V(s_{t+1}) - V(s_t)$
- 2-step TD:  $\hat{A}_t^2 = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) - V(s_t)$
- ...
- n-step TD:  $\hat{A}_t^n = r_t + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n}) - V(s_t)$
- In the limit: MC  $\hat{A}_t^\infty = r_t + \gamma r_{t+1} + \dots - V(s_t)$

# TD( $\lambda$ )

- How to choose  $n$ ?
- Any specific  $n$  is hard truncation of the window of evidence we consider
- Instead, use "exponential window"
  - Take  $n$ -step TD with weight proportional to  $\lambda^{n-1}$ , where  $0 \leq \lambda \leq 1$

- Now with  $\hat{A}_t^n = \sum_{t'=0}^{n-1} \gamma^{t'} \hat{A}_{t+t'}^1 = \sum_{t'=0}^{n-1} \gamma^{t'} (r_{t+t'} + \gamma V(s_{t+t'+1}) - V(s_{t+t'}))$   
$$\hat{A}_t^\lambda = (1 - \lambda) \sum_n \lambda^{n-1} \hat{A}_t^n = (1 - \lambda) \sum_n \lambda^{n-1} \sum_{t'=0}^{n-1} \gamma^{t'} \hat{A}_{t+t'}^1$$
$$= (1 - \lambda) \sum_{t'} \gamma^{t'} \hat{A}_{t+t'}^1 \sum_{n \geq t'+1} \lambda^{n-1} = \sum_{t'} (\lambda \gamma)^{t'} \hat{A}_{t+t'}^1$$



# Generalized Advantage Estimation (GAE( $\lambda$ ))

$$\nabla_{\theta} \mathcal{J}_{\theta} \approx \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'} (\lambda \gamma)^{t'} \hat{A}_{t+t'}^1$$

$$\hat{A}_t^1 = r_t + \gamma V(s_{t+1}) - V(s_t)$$

- GAE(0) = 1-step; GAE(1) = MC

# Recap

---

- Policy-Gradient Theorem
- State-dependent baselines
- Actor–Critic methods
  - Advantage estimation
  - Practical considerations
- n-step TD and TD( $\lambda$ )