

This lecture will be recorded!!!

Welcome to

DS595

Reinforcement Learning

Prof. Yanhua Li

Time: 6:00pm –8:50pm W
Zoom Lecture
Spring 2022

Quiz 5 Today

- ❖ 30 minutes on Policy Gradient (PG)

Project 4

- ❖ <https://github.com/yingxue-zhang/DS595-RL-Projects/tree/master/Project4>
- ❖ Important Dates:
- ❖ Progressive report: Wed. April 13, 2022 (23:59)
- ❖ Final Project:
 - Mon April 25, 2022 team project report is due
 - Wed. April 27, 2022 Virtual Poster Session

	Reinforcement Learning	Inverse Reinforcement Learning
Single Agent	Tabular representation of reward <i>Model-based control</i> <i>Model-free control</i> <i>(MC, SARSA, Q-Learning)</i>	Linear reward function learning <i>Imitation learning</i> <i>Apprenticeship learning</i> <i>Inverse reinforcement learning</i> MaxEnt IRL MaxCausalEnt IRL MaxRelEnt IRL
	Function representation of reward <i>1. Linear value function approx (MC, SARSA, Q-Learning)</i> <i>2. Value function approximation (Deep Q-Learning, Double DQN, prioritized DQN, Dueling DQN)</i> <i>3. Policy function approximation (Policy gradient, PPO, TRPO)</i> <i>4. Actor-Critic methods (A2C, A3C, Pathwise Derivative PG)</i>	
	Review of Deep Learning <i>As bases for non-linear function approximation (used in 2-4).</i>	Non-linear reward function learning Generative adversarial imitation learning (GAIL) Adversarial inverse reinforcement learning (AIRL) Review of Generative Adversarial nets As bases for non-linear IRL
Multiple Agents	Multi-Agent Reinforcement Learning Multi-agent Actor-Critic etc.	Multi-Agent Inverse Reinforcement Learning MA-GAIL MA-AIRL AMA-GAIL



This Lecture

❖ Policy Gradient

- Intro and Stochastic Policy
- Basic Policy Gradient Algorithm
- REINFORCE and Vanilla Policy Gradient
- PPO, TRPO, PPO2

❖ Actor-Critic methods

- A2C
- A3C
- Pathwise Derivative Policy Gradient

This Lecture

- ❖ Policy Gradient (Review Quickly)
 - Intro and Stochastic Policy
 - Basic Policy Gradient Algorithm
 - REINFORCE and Vanilla Policy Gradient
 - PPO, TRPO, PPO2
- ❖ Actor-Critic methods
 - A2C
 - A3C
 - Pathwise Derivative Policy Gradient

This Lecture

- ❖ Policy Gradient
 - Intro and Stochastic Policy
 - Basic Policy Gradient Algorithm
 - REINFORCE and Vanilla Policy Gradient
 - PPO, TRPO, PPO2
- ❖ Actor-Critic methods
 - A2C
 - A3C
 - Pathwise Derivative Policy Gradient
- ❖ Generative Adversarial Networks (GAN)
- ❖ Deep Inverse Reinforcement Learning

Review – Policy Gradient

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n : \text{obtained via interaction}} - \overset{\text{baseline}}{\underline{b}} \right) \nabla \log \pi_\theta(a_t^n | s_t^n)$$

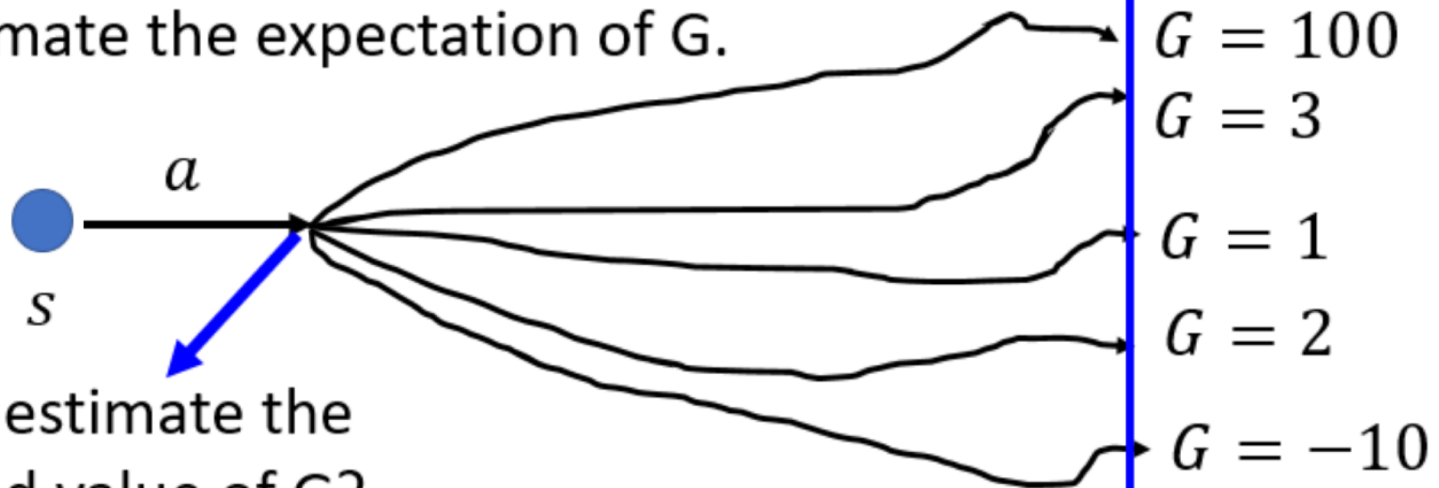
Review – Policy Gradient

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n} - \underline{b} \right) \nabla \log \pi_\theta(a_t^n | s_t^n)$$

G_t^n : obtained via interaction

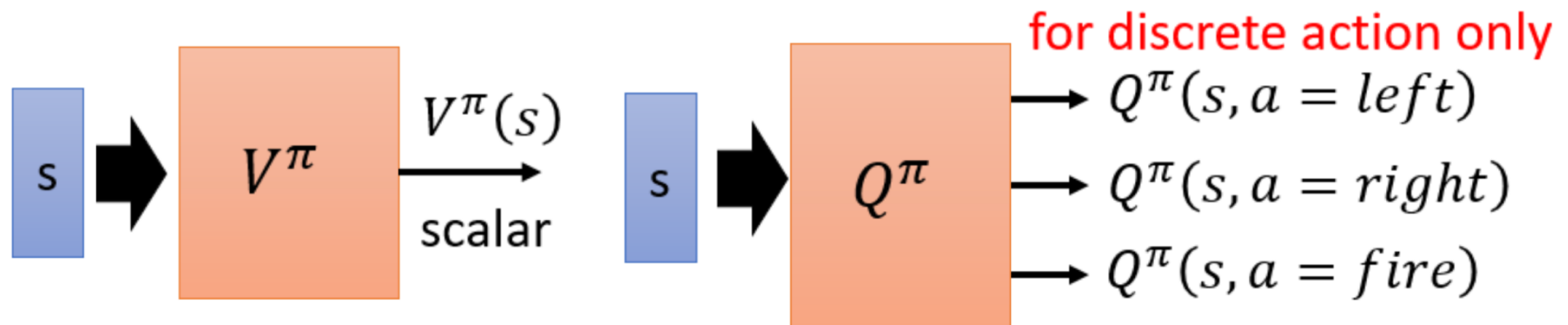
Very unstable

With sufficient samples,
approximate the expectation of G .



Review – Q-Learning

- State value function $V^\pi(s)$
 - When using actor π , the *cumulated* reward expects to be obtained after visiting state s
- State-action value function $Q^\pi(s, a)$
 - When using actor π , the *cumulated* reward expects to be obtained after taking a at state s




Estimated by TD or MC

Actor-Critic

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n : \text{obtained via interaction}} - \overset{\text{baseline}}{\underline{b}} \right) \nabla \log \pi_\theta(a_t^n | s_t^n)$$

Actor-Critic

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n : \text{obtained via interaction}} - \overset{\text{baseline}}{\underline{b}} \right) \nabla \log \pi_\theta(a_t^n | s_t^n)$$



$$E[G_t^n] = Q^{\pi_\theta}(s_t^n, a_t^n)$$

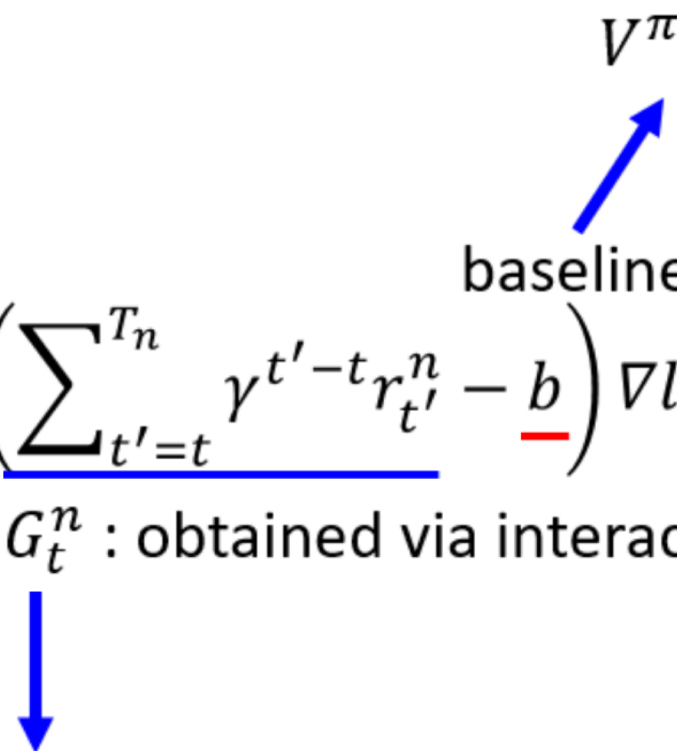
Actor-Critic

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n : \text{obtained via interaction}} - \underbrace{b}_{\text{baseline}} \right) \nabla \log \pi_\theta(a_t^n | s_t^n)$$

$G_t^n : \text{obtained via interaction}$

$V^{\pi_\theta}(s_t^n)$

$E[G_t^n] = Q^{\pi_\theta}(s_t^n, a_t^n)$



Actor-Critic

$$Q^{\pi_{\theta}}(s_t^n, a_t^n) - V^{\pi_{\theta}}(s_t^n)$$

$$V^{\pi_{\theta}}(s_t^n)$$

$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\underbrace{\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n}_{G_t^n : \text{obtained via interaction}} - \underbrace{b}_{\text{baseline}} \right) \nabla \log \pi_{\theta}(a_t^n | s_t^n)$$

$$E[G_t^n] = Q^{\pi_{\theta}}(s_t^n, a_t^n)$$

Advantage Actor-Critic

$$Q^{\pi}(s_t^n, a_t^n) - V^{\pi}(s_t^n)$$

Estimate two networks? We can only estimate one.

Advantage Actor-Critic

$$Q^{\pi}(s_t^n, a_t^n) - V^{\pi}(s_t^n)$$



$$r_t^n + V^{\pi}(s_{t+1}^n) - V^{\pi}(s_t^n)$$

Estimate two networks? We can only estimate one.

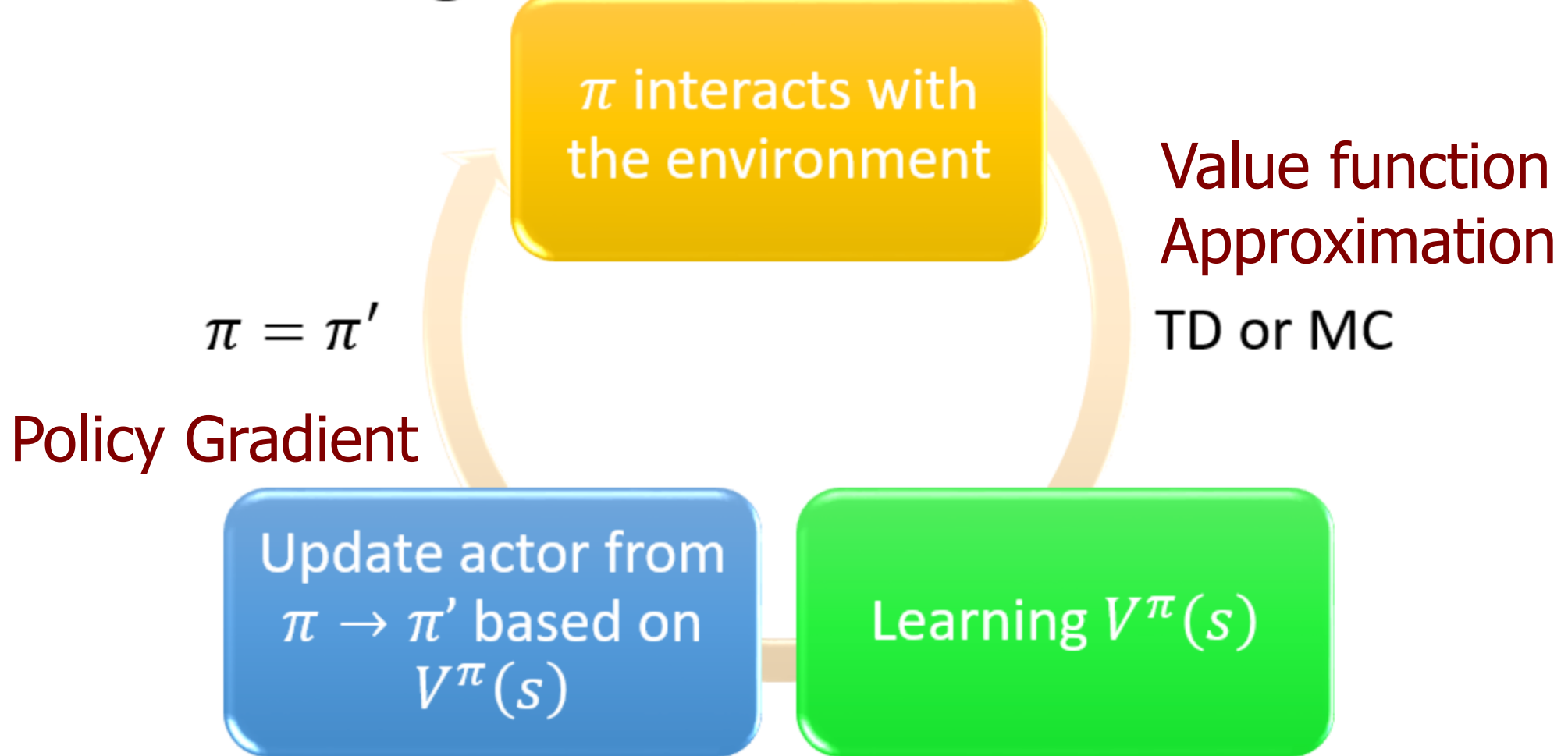
Only estimate state value
A little bit variance

$$Q^{\pi}(s_t^n, a_t^n) = E[r_t^n + V^{\pi}(s_{t+1}^n)]$$

$$Q^{\pi}(s_t^n, a_t^n) = r_t^n + V^{\pi}(s_{t+1}^n)$$

Advantage Actor-Critic

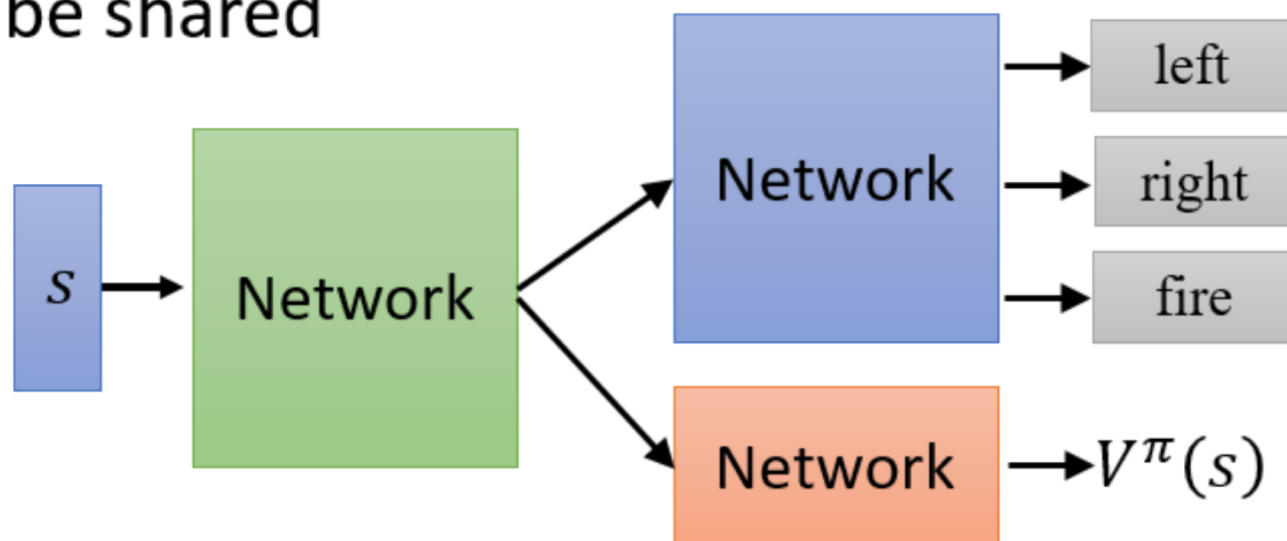
(A2C algorithm)



$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (r_t^n + V^\pi(s_{t+1}^n) - V^\pi(s_t^n)) \nabla \log \pi_\theta(a_t^n | s_t^n)$$

Advantage Actor-Critic

- Tips
 - The parameters of actor $\pi(s)$ and critic $V^\pi(s)$ can be shared



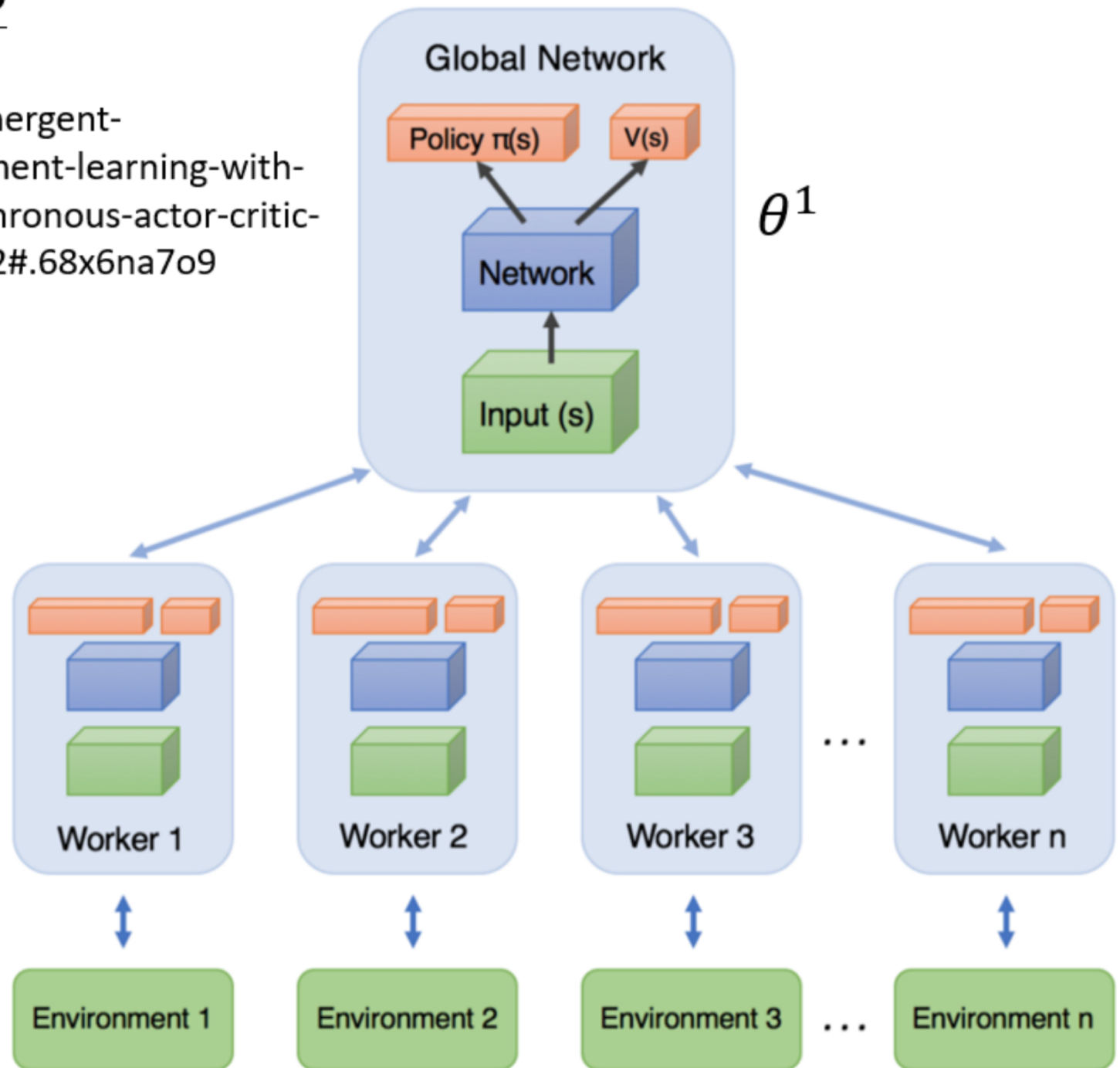
Asynchronous Advantage Actor-Critic (A3C)



Asynchronous

Source of image:

<https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-8-asynchronous-actor-critic-agents-a3c-c88f72a5e9f2#.68x6na7o9>

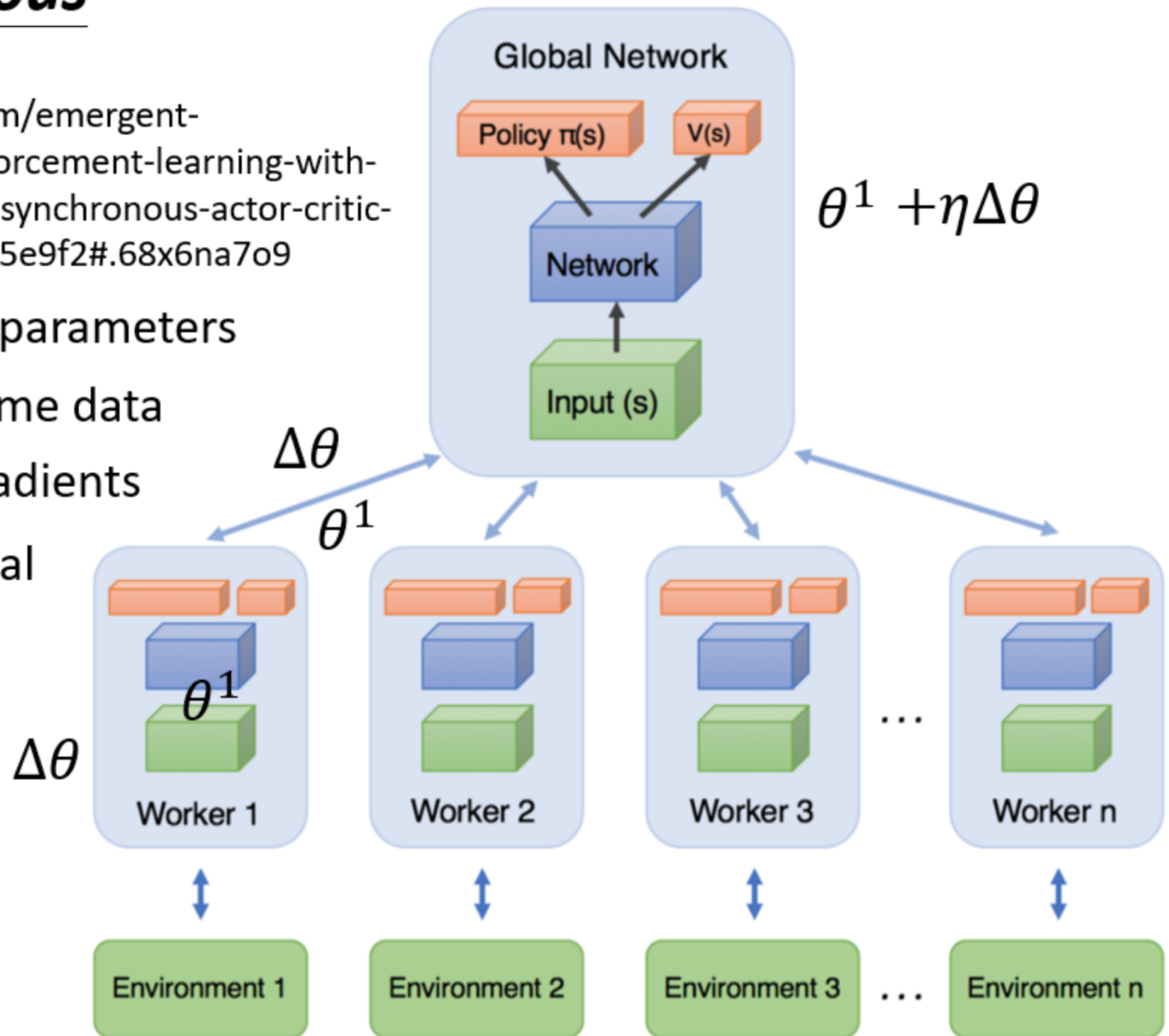


Asynchronous

Source of image:

<https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-8-asynchronous-actor-critic-agents-a3c-c88f72a5e9f2#.68x6na7o9>

1. Copy global parameters
2. Sampling some data
3. Compute gradients
4. Update global models

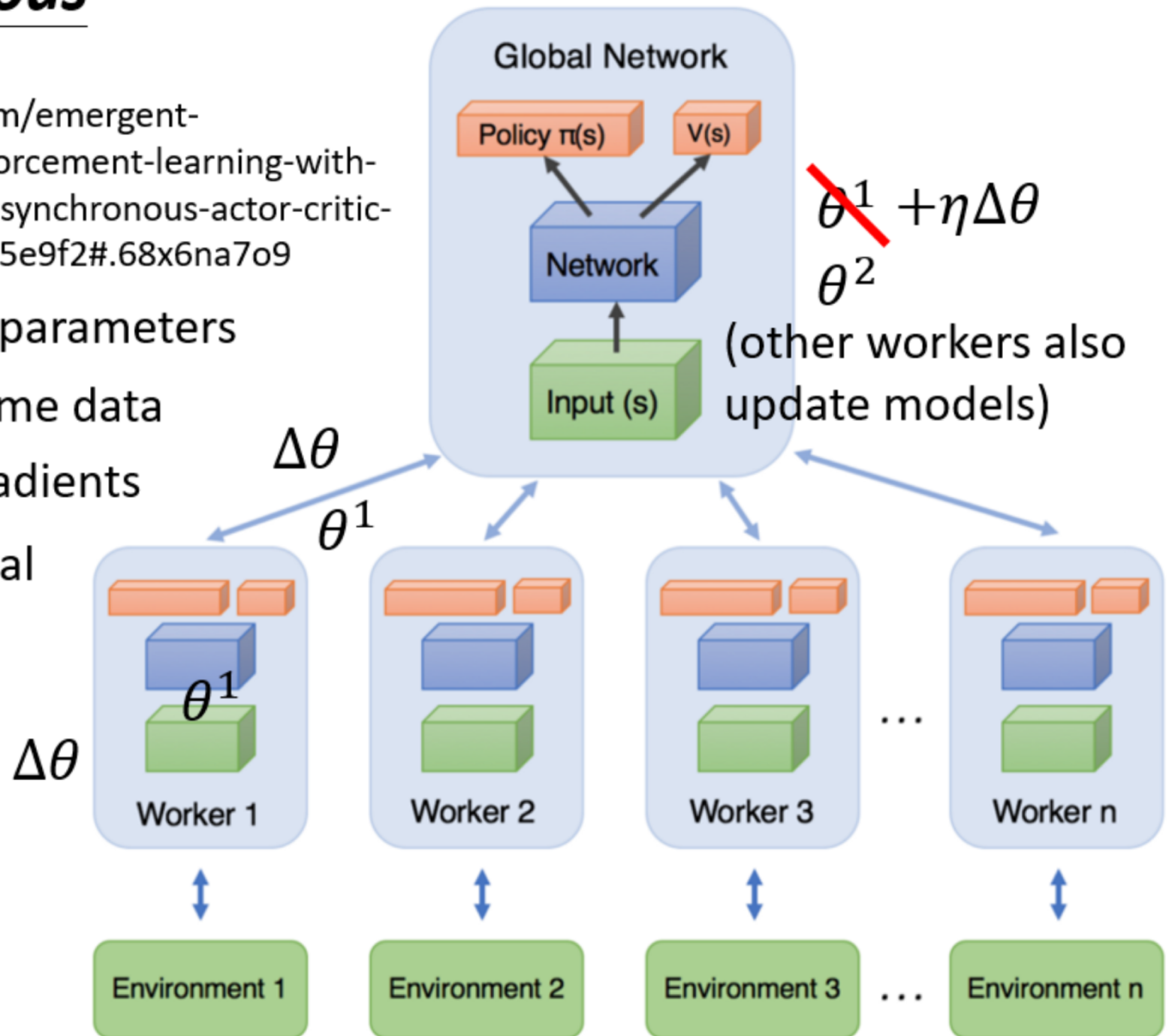


Asynchronous

Source of image:

<https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-8-asynchronous-actor-critic-agents-a3c-c88f72a5e9f2#.68x6na7o9>

1. Copy global parameters
2. Sampling some data
3. Compute gradients
4. Update global models



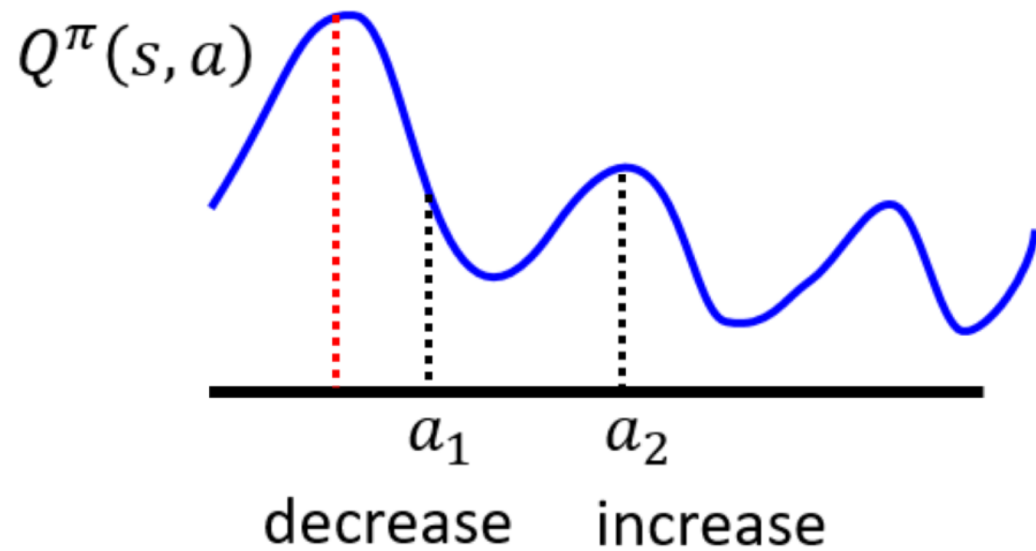
Pathwise Derivative Policy Gradient

David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, Martin Riedmiller,
“Deterministic Policy Gradient Algorithms”, ICML, 2014

Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess,
Tom Erez, Yuval Tassa, David Silver, Daan Wierstra, “CONTINUOUS CONTROL WITH DEEP
REINFORCEMENT LEARNING”, ICLR, 2016

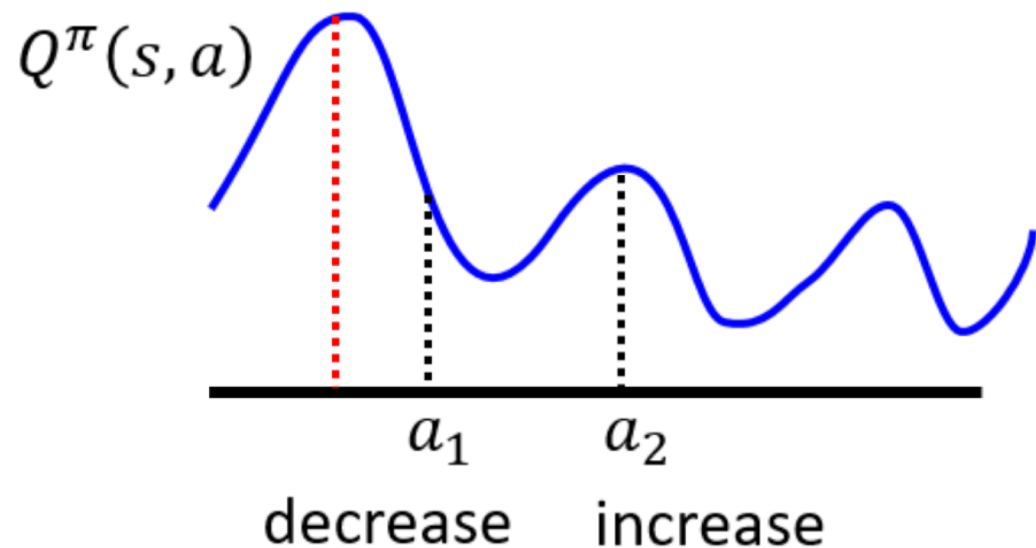
Another Way to use Critic

Original Actor-critic



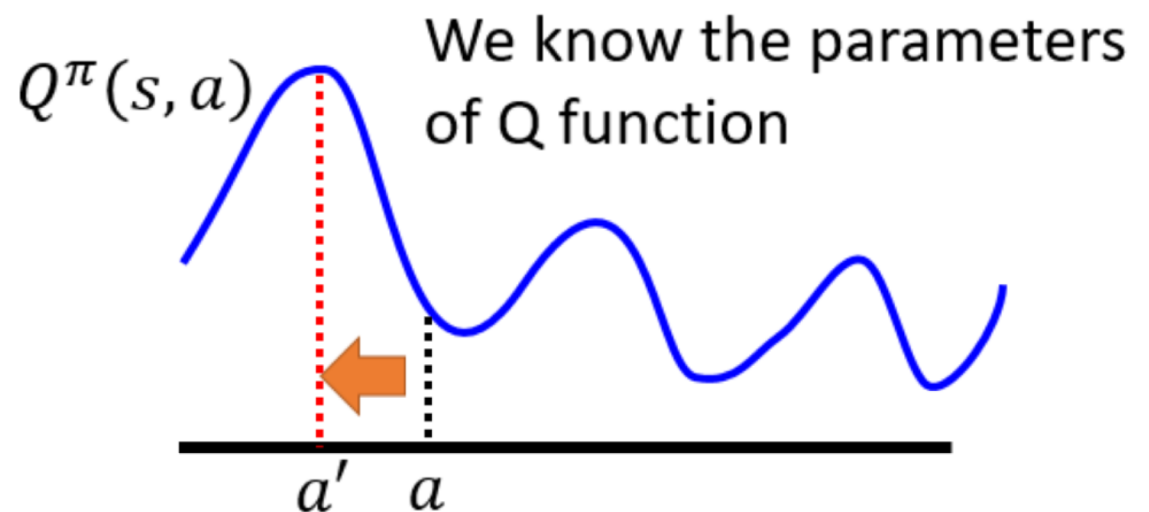
Another Way to use Critic

Original Actor-critic



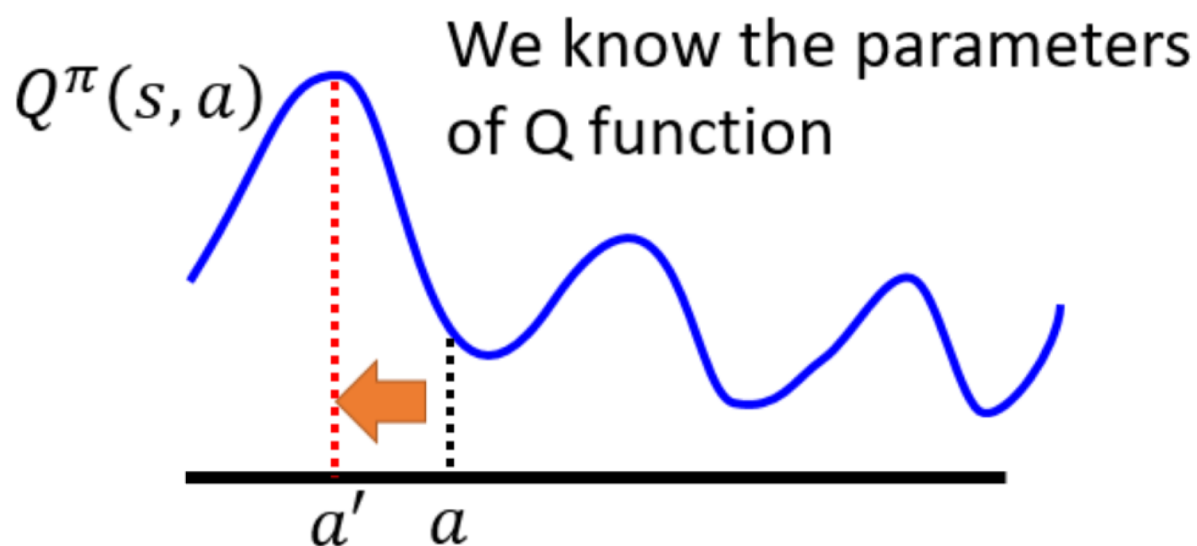
Pathwise derivative policy gradient

From Q function we know that taking a' at state s is better than a



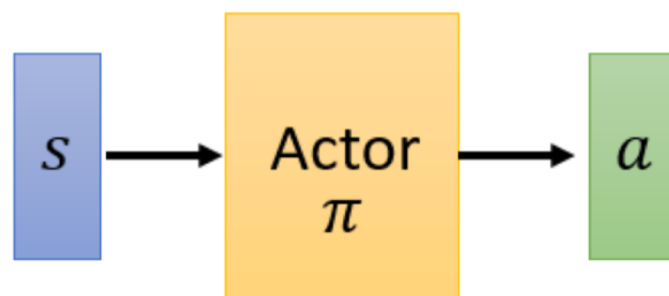
Pathwise derivative policy gradient

From Q function we know that taking a' at state s is better than a



Action a is a *continuous vector*

$$a = \arg \max_a Q(s, a)$$



Actor as the solver of this optimization problem

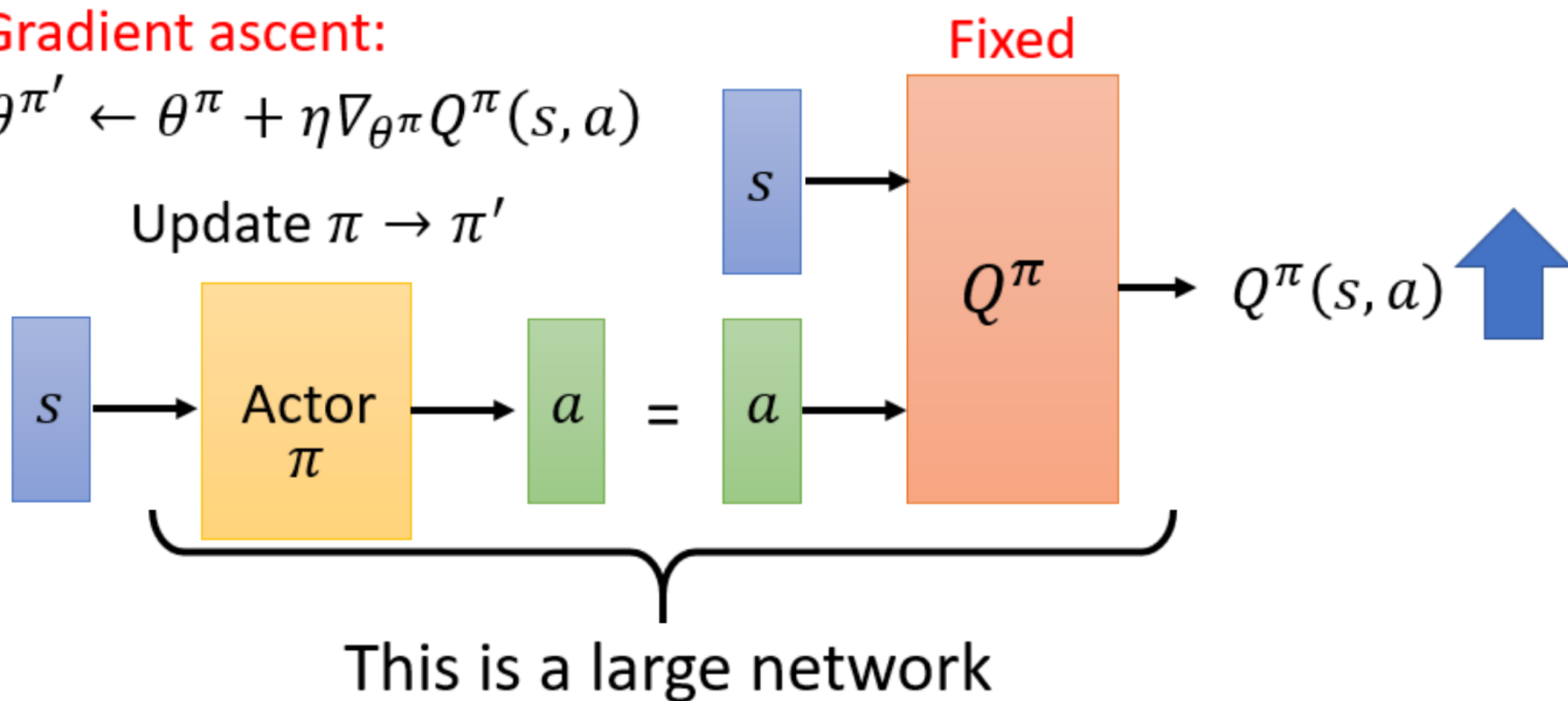
Pathwise Derivative Policy Gradient

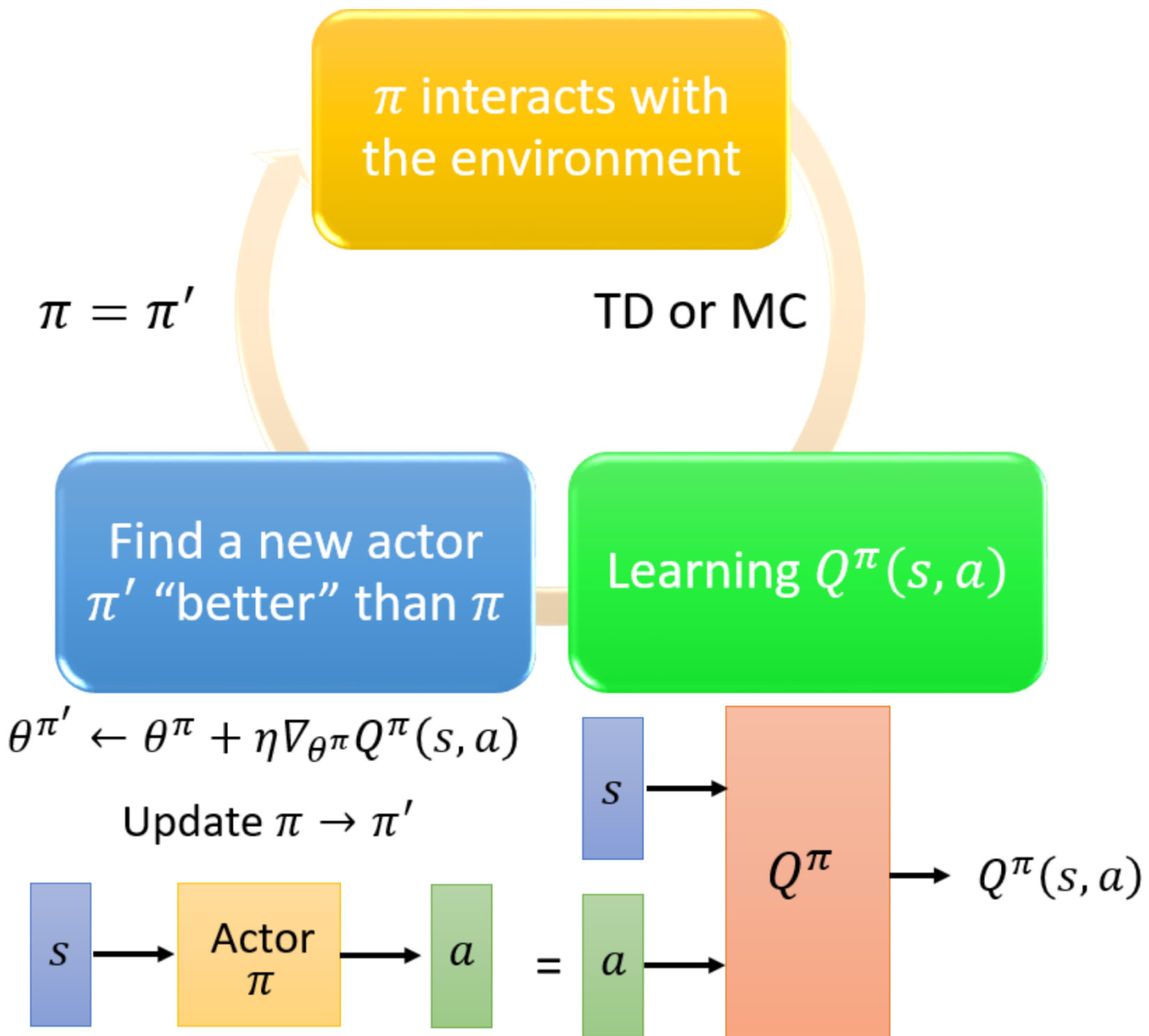
$$\pi'(s) = \arg \max_a Q^\pi(s, a) \quad \leftarrow a \text{ is the output of an actor}$$

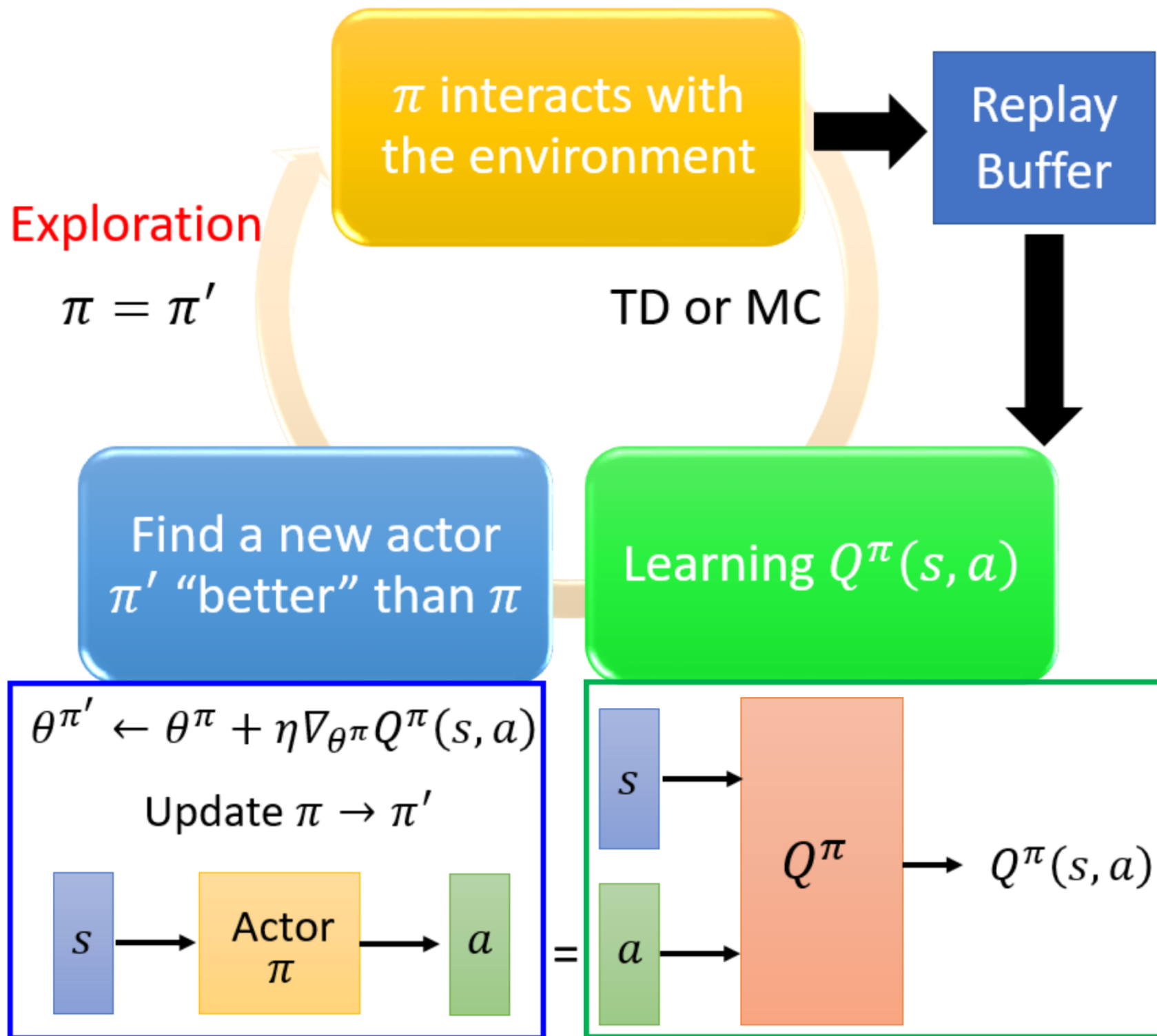
Gradient ascent:

$$\theta^{\pi'} \leftarrow \theta^\pi + \eta \nabla_{\theta^\pi} Q^\pi(s, a)$$

Update $\pi \rightarrow \pi'$







Q-Learning Algorithm

- Initialize Q-function Q , target Q-function $\hat{Q} = Q$
- In each episode
 - For each time step t
 - Given state s_t , take action a_t based on Q (exploration)
 - Obtain reward r_t , and reach new state s_{t+1}
 - Store (s_t, a_t, r_t, s_{t+1}) into buffer
 - Sample (s_i, a_i, r_i, s_{i+1}) from buffer (usually a batch)
 - Target $y = r_i + \max_a \hat{Q}(s_{i+1}, a)$
 - Update the parameters of Q to make $Q(s_i, a_i)$ close to y (regression)
 - Every C steps reset $\hat{Q} = Q$

Q-Learning Algorithm ➡ Pathwise Derivative Policy Gradient

- Initialize Q-function Q , target Q-function $\hat{Q} = Q$, actor π , target actor $\hat{\pi} = \pi$

Replaced ϵ -greedy policy with π network.

- In each episode

- For each time step t

- 1 • Given state s_t , take action a_t based on ~~Q~~ π (exploration)
 - Obtain reward r_t , and reach new state s_{t+1}
 - Store (s_t, a_t, r_t, s_{t+1}) into buffer
 - Sample (s_i, a_i, r_i, s_{i+1}) from buffer (usually a batch)
- 2 • Target $y = r_i + \max_a \hat{Q}(s_{i+1}, a) - \hat{Q}(s_i, \hat{\pi}(s_i)) + \hat{Q}(s_{i+1}, \hat{\pi}(s_{i+1}))$
 - Update the parameters of Q to make $Q(s_i, a_i)$ close to y (regression)
- 3 • Update the parameters of π to maximize $Q(s_i, \pi(s_i))$
 - Every C steps reset $\hat{Q} = Q$
- 4 • Every C steps reset $\hat{\pi} = \pi$

	Reinforcement Learning	Inverse Reinforcement Learning
Single Agent	Tabular representation of reward <i>Model-based control</i> <i>Model-free control</i> <i>(MC, SARSA, Q-Learning)</i>	Linear reward function learning <i>Imitation learning</i> <i>Apprenticeship learning</i> <i>Inverse reinforcement learning</i> MaxEnt IRL MaxCausalEnt IRL MaxRelEnt IRL
	Function representation of reward 1. <i>Linear value function approx (MC, SARSA, Q-Learning)</i> 2. <i>Value function approximation (Deep Q-Learning, Double DQN, prioritized DQN, Dueling DQN)</i> 3. <i>Policy function approximation (Policy gradient, PPO, TRPO)</i> 4. <i>Actor-Critic methods (A2C, A3C, Pathwise Derivative PG)</i>	
	Review of Deep Learning <i>As bases for non-linear function approximation (used in 2-4).</i>	Non-linear reward function learning Generative adversarial imitation learning (GAIL) Adversarial inverse reinforcement learning (AIRL)
Multiple Agents		Review of Generative Adversarial nets
	Multi-Agent Reinforcement Learning Multi-agent Actor-Critic etc.	Multi-Agent Inverse Reinforcement Learning MA-GAIL MA-AIRL AMA-GAIL



Questions?