

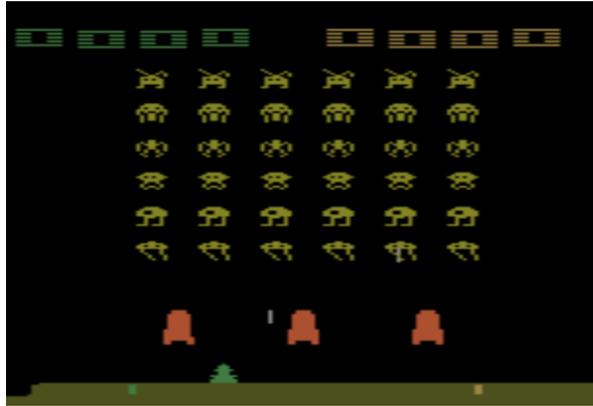
Deep Reinforcement Learning

Example: Playing Video Game

Start with
observation s_1

Observation s_2

Observation s_3



Obtain reward
 $r_1 = 0$

Action a_1 : "right"



Obtain reward
 $r_2 = 5$

Action a_2 : "fire"

(kill an alien)

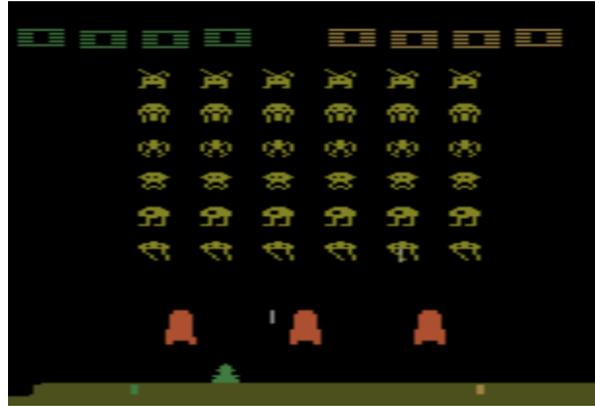
Usually there is some randomness in the environment

Example: Playing Video Game

Start with
observation s_1



Observation s_2



Observation s_3



After many turns



Obtain reward r_T

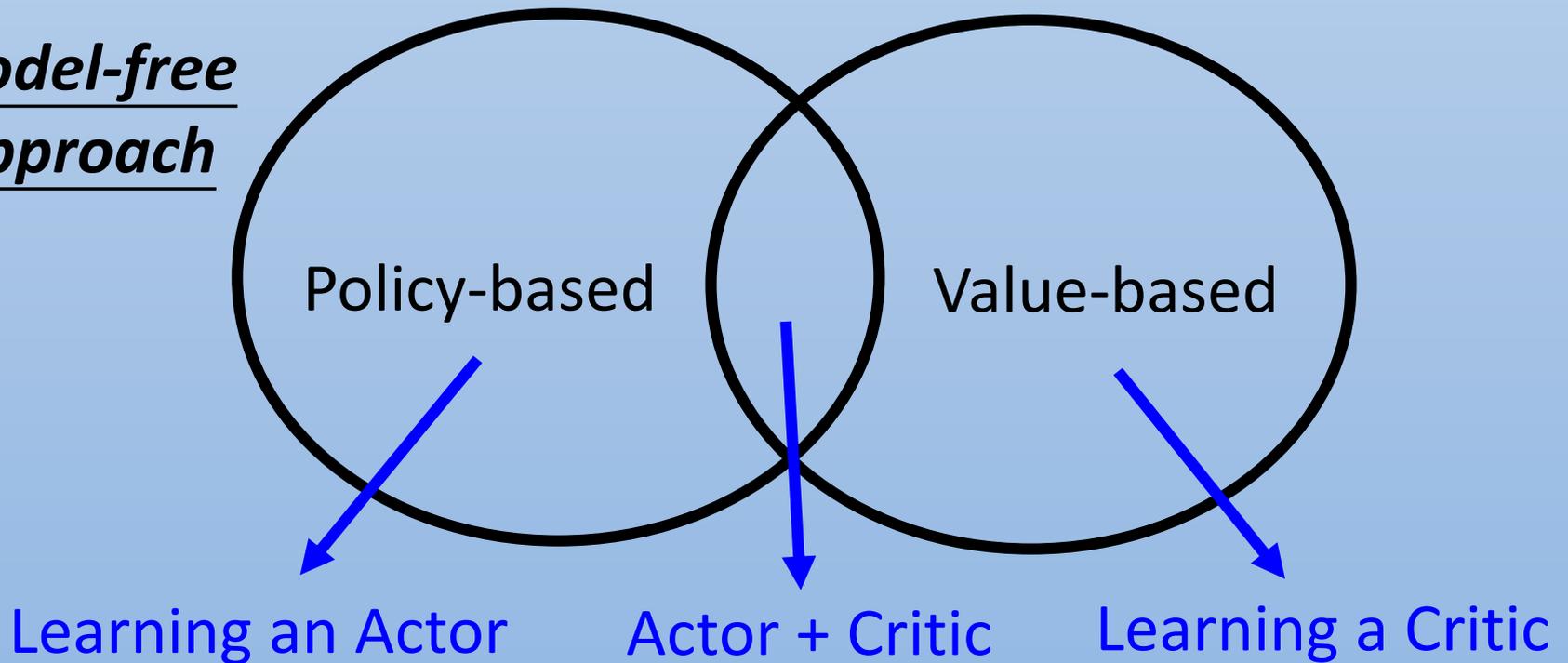
Action a_T

This is an episode.

Learn to maximize the
expected cumulative
reward per episode

Approaches

Model-free Approach



Model-based Approach

On-policy v.s. Off-policy

- On-policy: The agent learned and the agent interacting with the environment is the same.
- Off-policy: The agent learned and the agent interacting with the environment is different.



阿光下棋



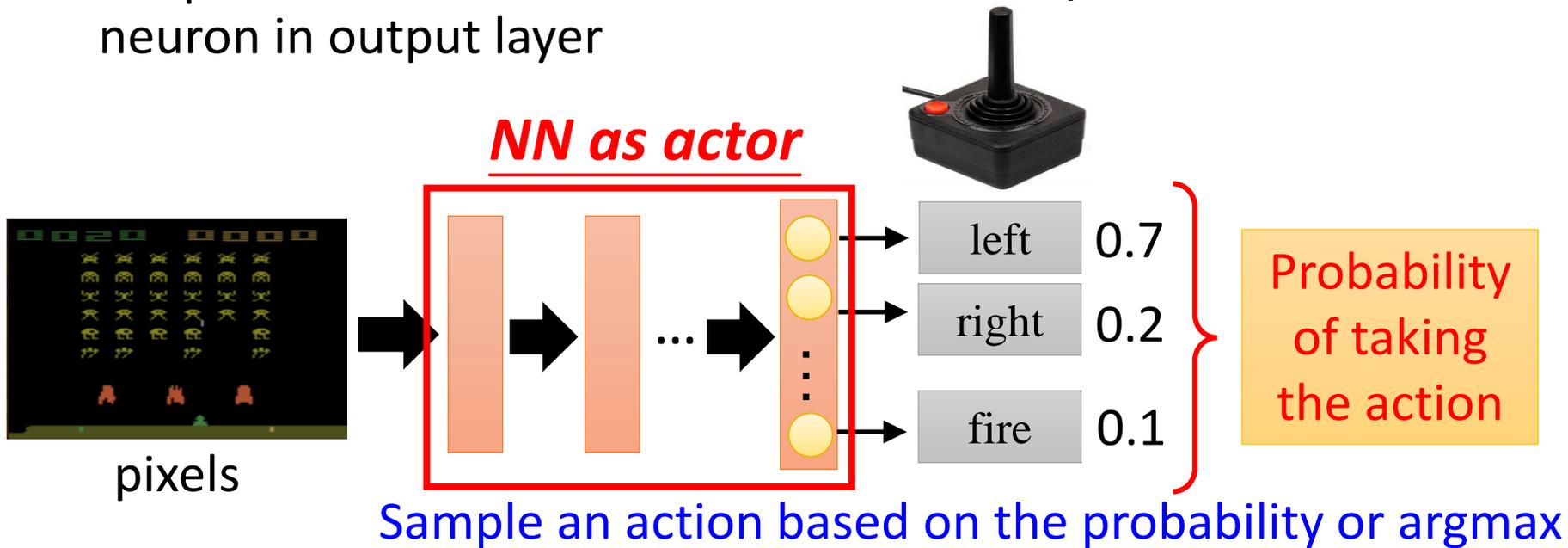
佐為下棋、阿光在旁邊看

Asynchronous Advantage Actor-Critic (A3C)

Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, Koray Kavukcuoglu, "Asynchronous Methods for Deep Reinforcement Learning", ICML, 2016

Actor is a Neural network

- Input of neural network: the observation of machine represented as a vector or a matrix
- Output neural network : each action corresponds to a neuron in output layer



Actor can also have continuous action.

Actor – Goodness of an Actor

- Given an actor $\pi(s)$ with network parameter θ^π
- Use the actor $\pi(s)$ to play the video game
 - Start with observation s_1
 - Machine decides to take a_1
 - Machine obtains reward r_1
 - Machine sees observation s_2
 - Machine decides to take a_2
 - Machine obtains reward r_2
 - Machine sees observation s_3
 -
 - Machine decides to take a_T
 - Machine obtains reward r_T

END

Total reward: $R = \sum_{t=1}^T r_t$

Even with the same actor,
 R is different each time

Randomness in the actor
and the game

We define \bar{R}_{θ^π} as the
expected total reward

\bar{R}_{θ^π} evaluates the goodness of an actor $\pi(s)$

Actor – Policy Gradient

$\theta^{\pi'} \leftarrow \theta^{\pi} + \eta \nabla \bar{R}_{\theta^{\pi}}$ Using θ^{π} to obtain $\{\tau^1, \tau^2, \dots, \tau^N\}$

$$\begin{aligned} \nabla \bar{R}_{\theta^{\pi}} &\approx \frac{1}{N} \sum_{n=1}^N R(\tau^n) \nabla \log P(\tau^n | \theta^{\pi}) = \frac{1}{N} \sum_{n=1}^N R(\tau^n) \sum_{t=1}^{T_n} \nabla \log p(a_t^n | s_t^n, \theta^{\pi}) \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta^{\pi}) \end{aligned}$$

What if we replace $R(\tau^n)$ with r_t^n

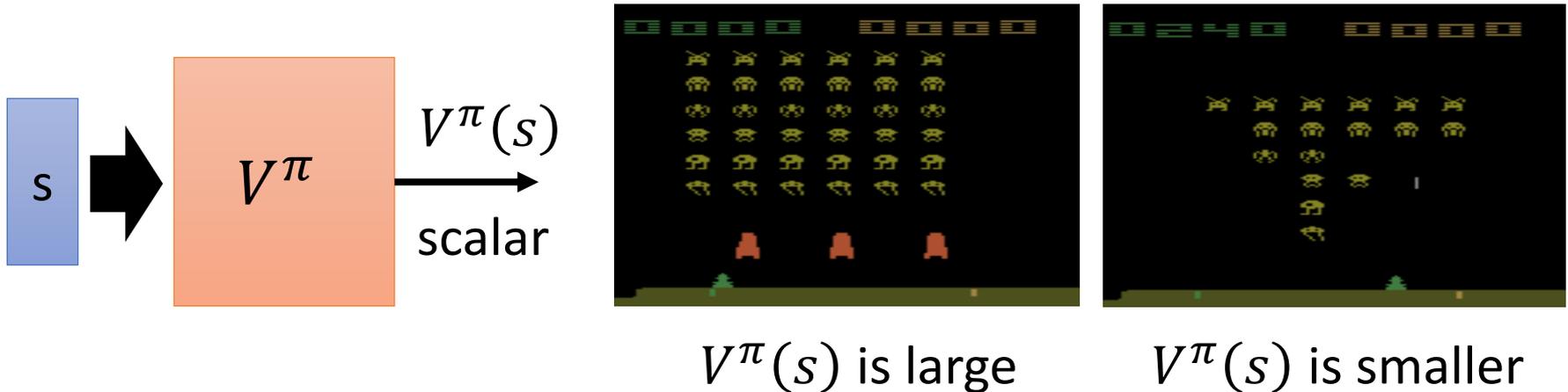
If in τ^n machine takes a_t^n when seeing s_t^n

$R(\tau^n)$ is positive  Tuning θ to increase $p(a_t^n | s_t^n)$
 $R(\tau^n)$ is negative  Tuning θ to decrease $p(a_t^n | s_t^n)$

It is very important to consider the cumulative reward $R(\tau^n)$ of the whole trajectory τ^n instead of immediate reward r_t^n

Critic

- A critic does not determine the action.
- Given an actor π , it evaluates the how good the actor is
- State value function $V^\pi(s)$
 - When using actor π , the *cumulated* reward expects to be obtained after seeing observation (state) s



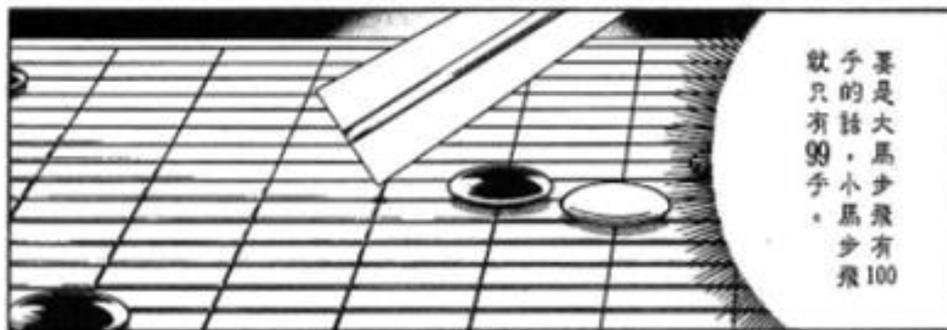
Critic

v 以前的阿光(大馬步飛) = bad

v 變強的阿光(大馬步飛) = good



* 小馬步飛：按兩棋一樣，兩棋子放在同一格；大馬步飛則是放在對好格。

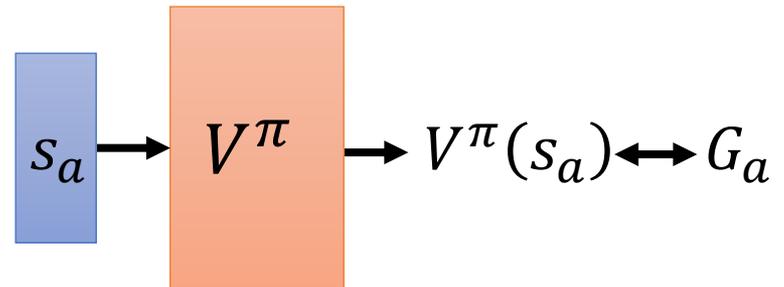


How to estimate $V^\pi(s)$

- Monte-Carlo based approach
 - The critic watches π playing the game

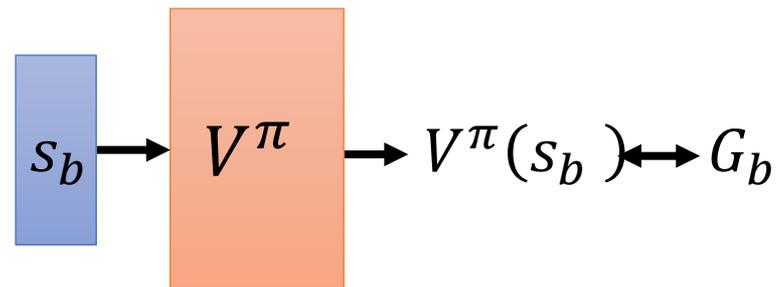
After seeing s_a ,

Until the end of the episode,
the cumulated reward is G_a



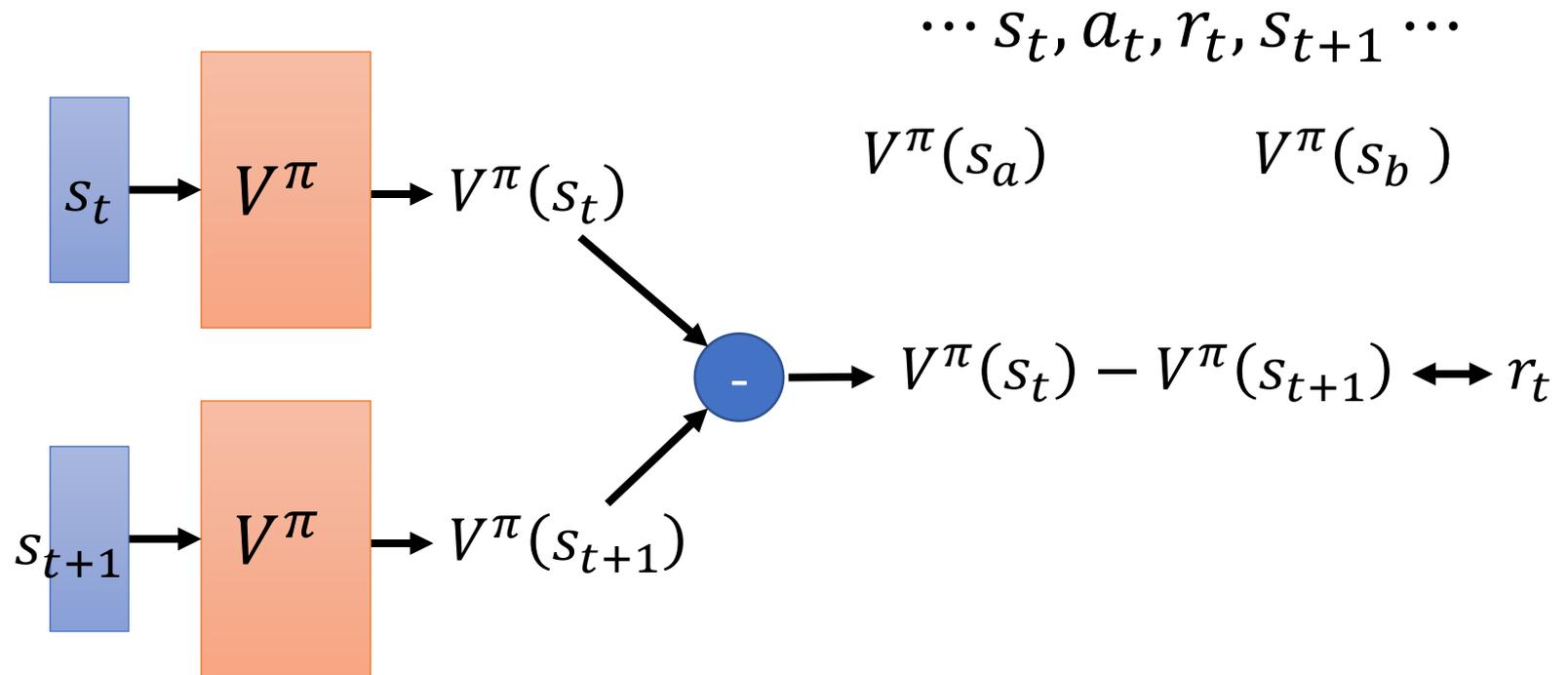
After seeing s_b ,

Until the end of the episode,
the cumulated reward is G_b



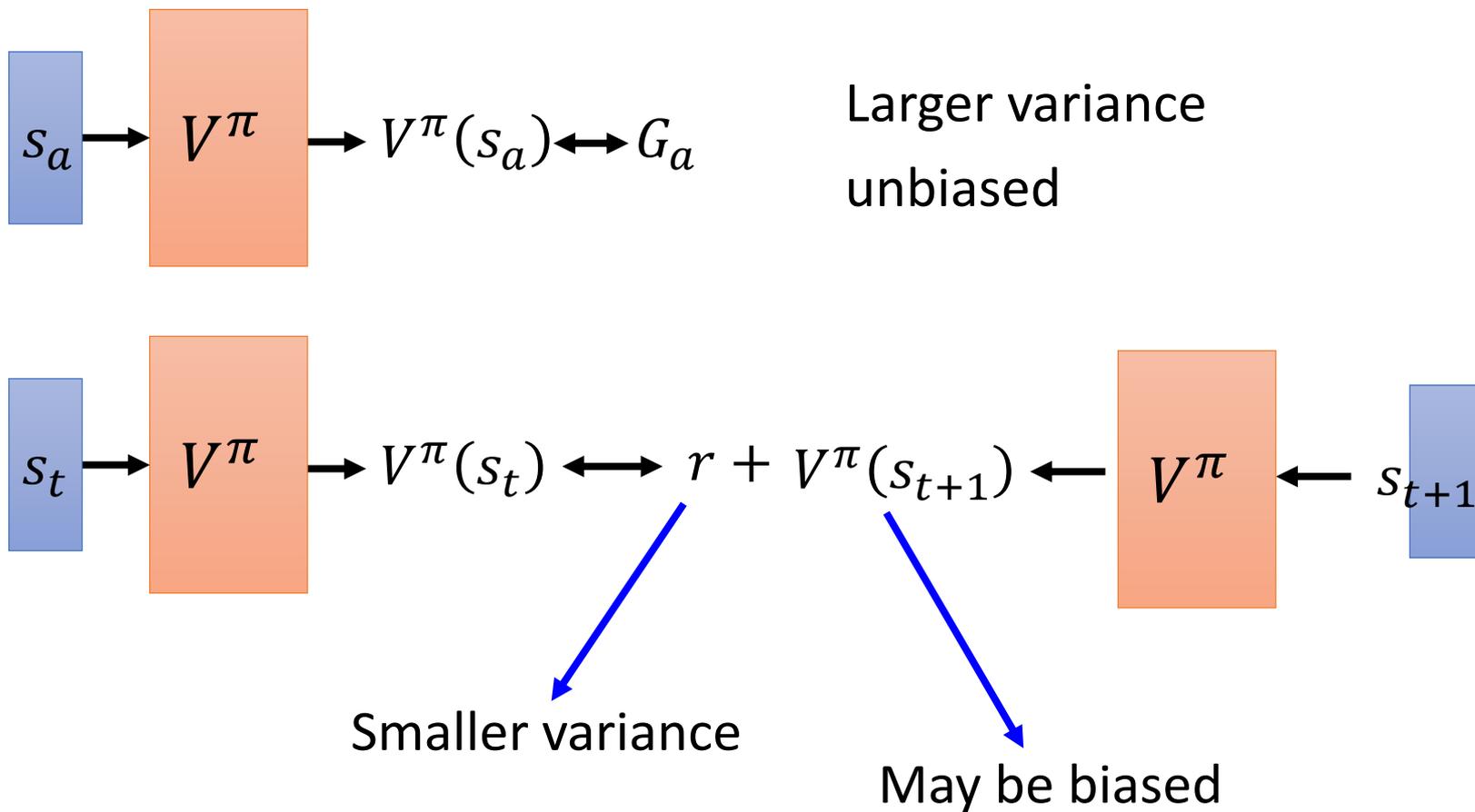
How to estimate $V^\pi(s)$

- Temporal-difference approach



Some applications have very long episodes, so that delaying all learning until an episode's end is too slow.

MC v.s. TD



MC v.s. TD

[Sutton, v2,
Example 6.4]

- The critic has the following 8 episodes

- $s_a, r = 0, s_b, r = 0, \text{END}$

- $s_b, r = 1, \text{END}$

$$V^\pi(s_b) = 3/4$$

- $s_b, r = 1, \text{END}$

- $s_b, r = 1, \text{END}$

$$V^\pi(s_a) = ? \quad 0? \quad 3/4?$$

- $s_b, r = 1, \text{END}$

- $s_b, r = 1, \text{END}$

Monte-Carlo: $V^\pi(s_a) = 0$

- $s_b, r = 1, \text{END}$

- $s_b, r = 0, \text{END}$

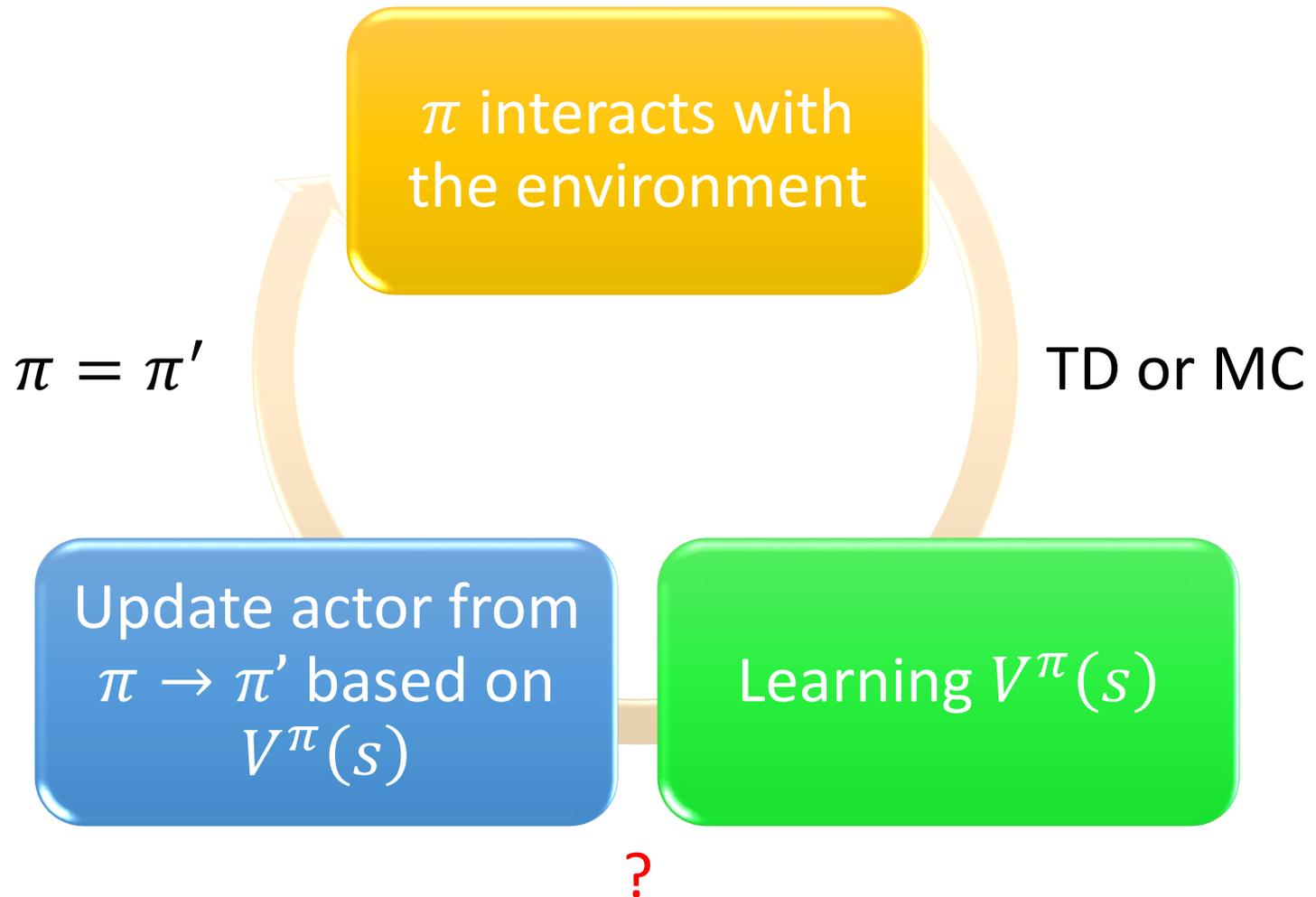
Temporal-difference:

$$V^\pi(s_a) + r = V^\pi(s_b)$$

$$3/4 \quad 0 \quad 3/4$$

(The actions are ignored here.)

Actor-Critic



Advantage Actor-Critic

$$\theta^{\pi'} \leftarrow \theta^{\pi} + \eta \nabla \bar{R}_{\theta^{\pi}}$$

$$\nabla \bar{R}_{\theta^{\pi}} = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p(a_t^n | s_t^n, \theta^{\pi})$$

Evaluated by critic

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

Baseline is added

The reward r_t^n we truly obtain when taking action a_t^n

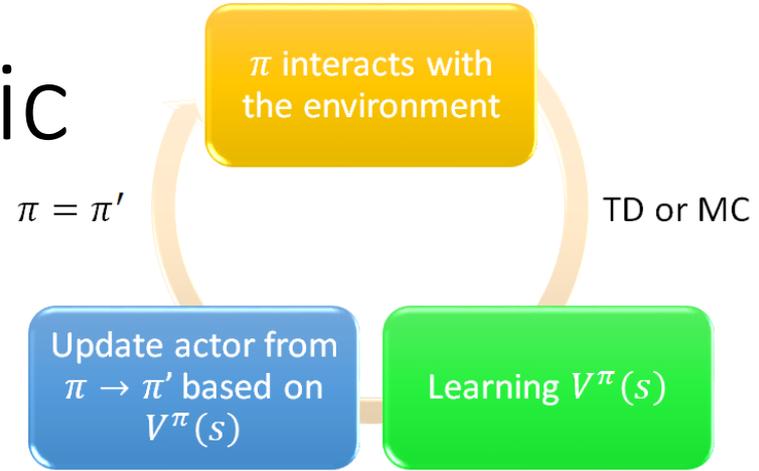
Expected reward r_t^n we obtain if we use actor π

Positive advantage function

Increasing the prob. of action a_t^n

Negative advantage function

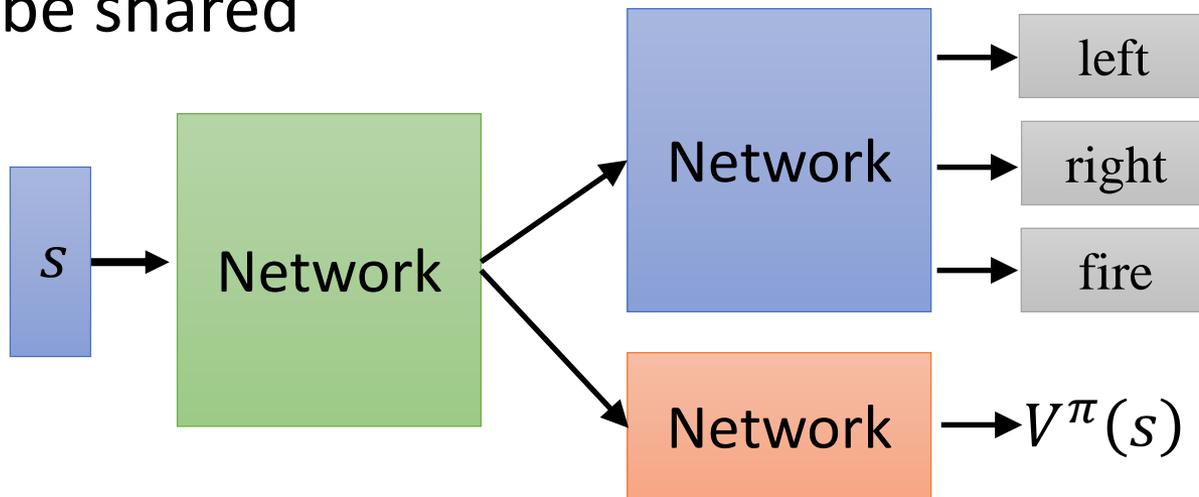
decreasing the prob. of action a_t^n



Advantage Actor-Critic

- Tips

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



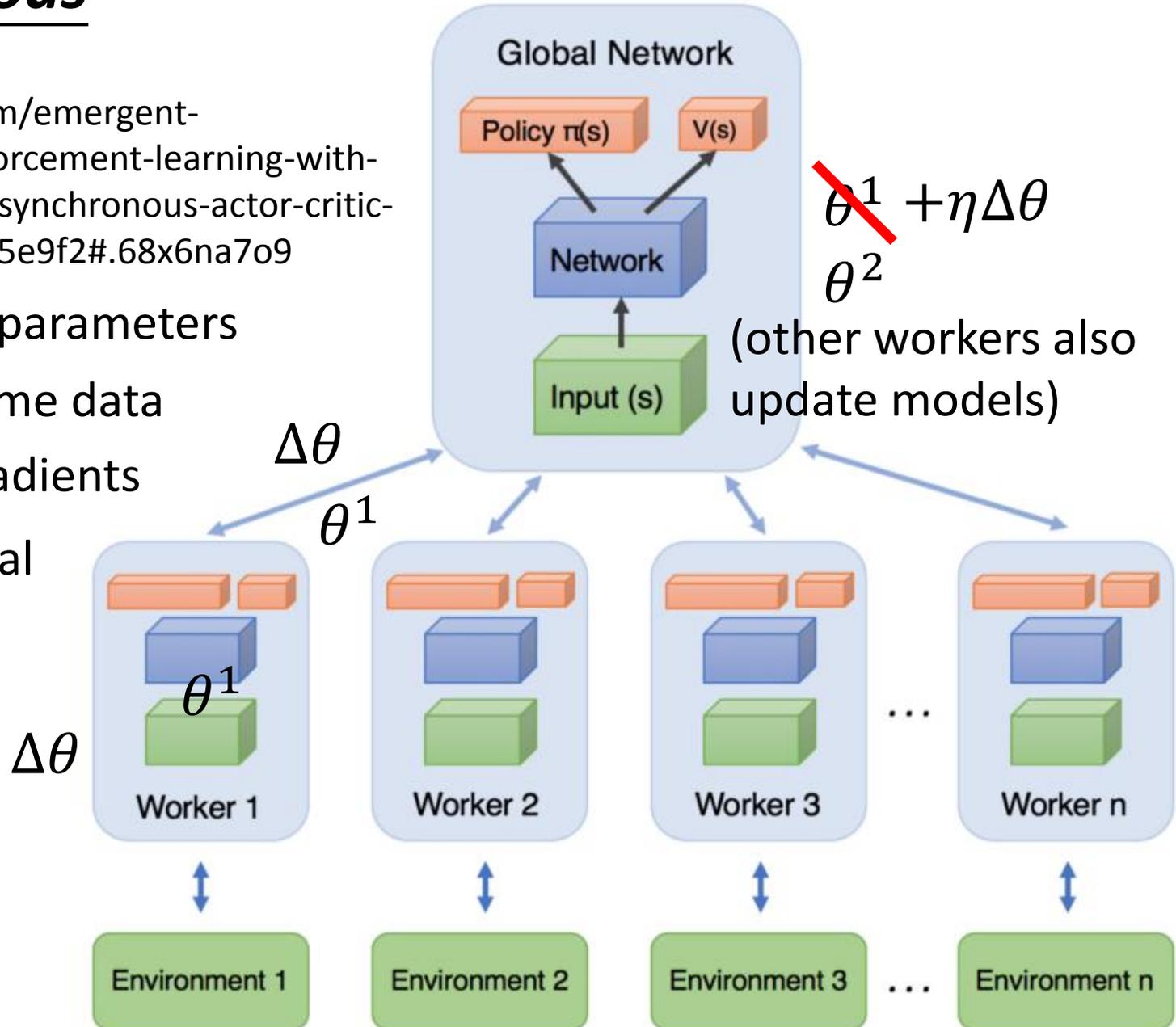
- Use output entropy as regularization for $\pi(s)$
 - Larger entropy is preferred \rightarrow exploration

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

Actor

Critic

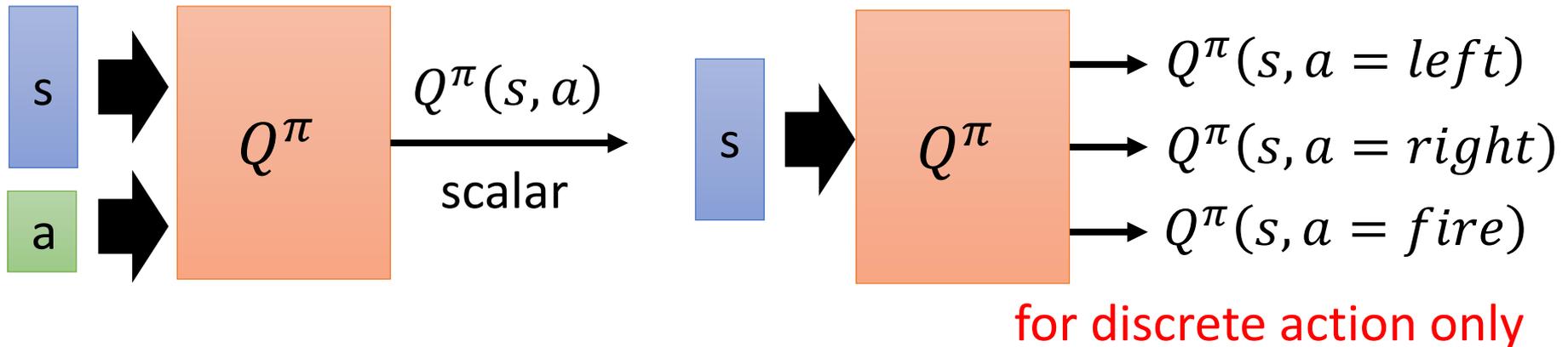


Pathwise derivative
policy gradient

Original Actor-critic

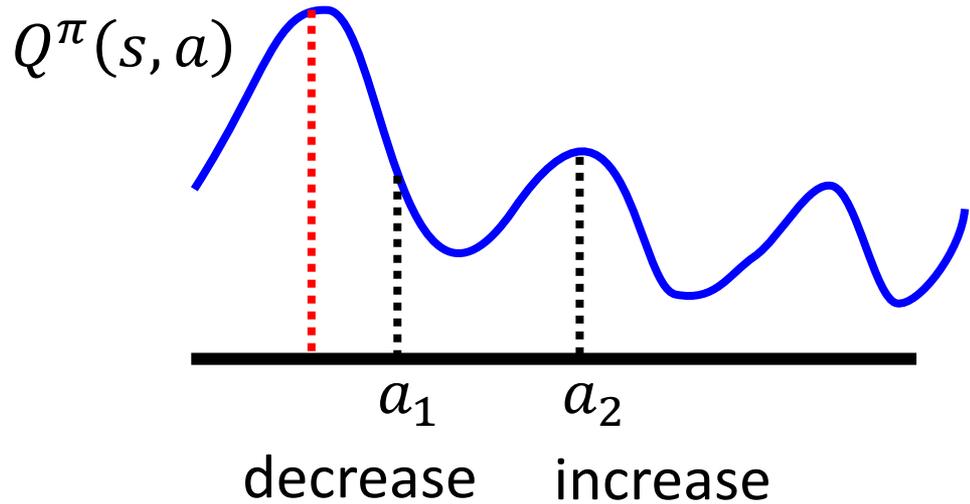
Another Critic

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



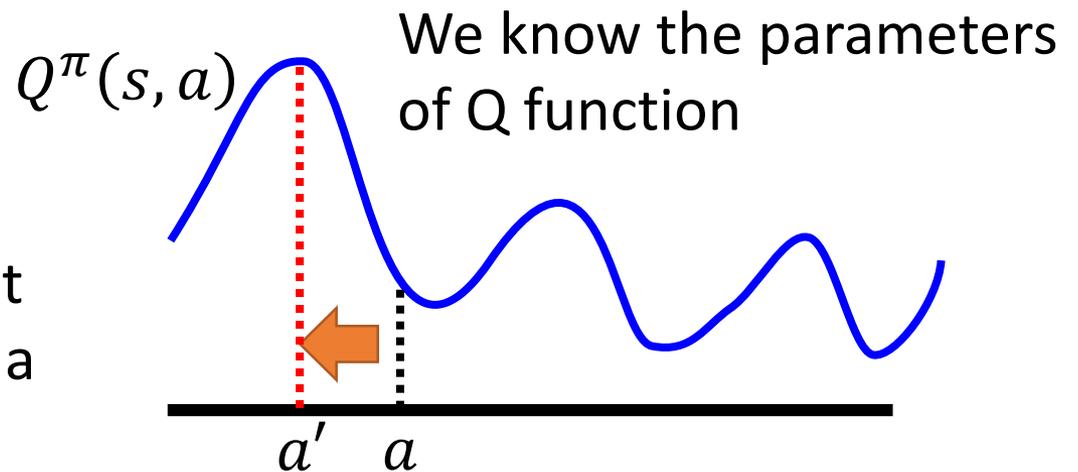
Another Way to use Critic

Original Actor-critic

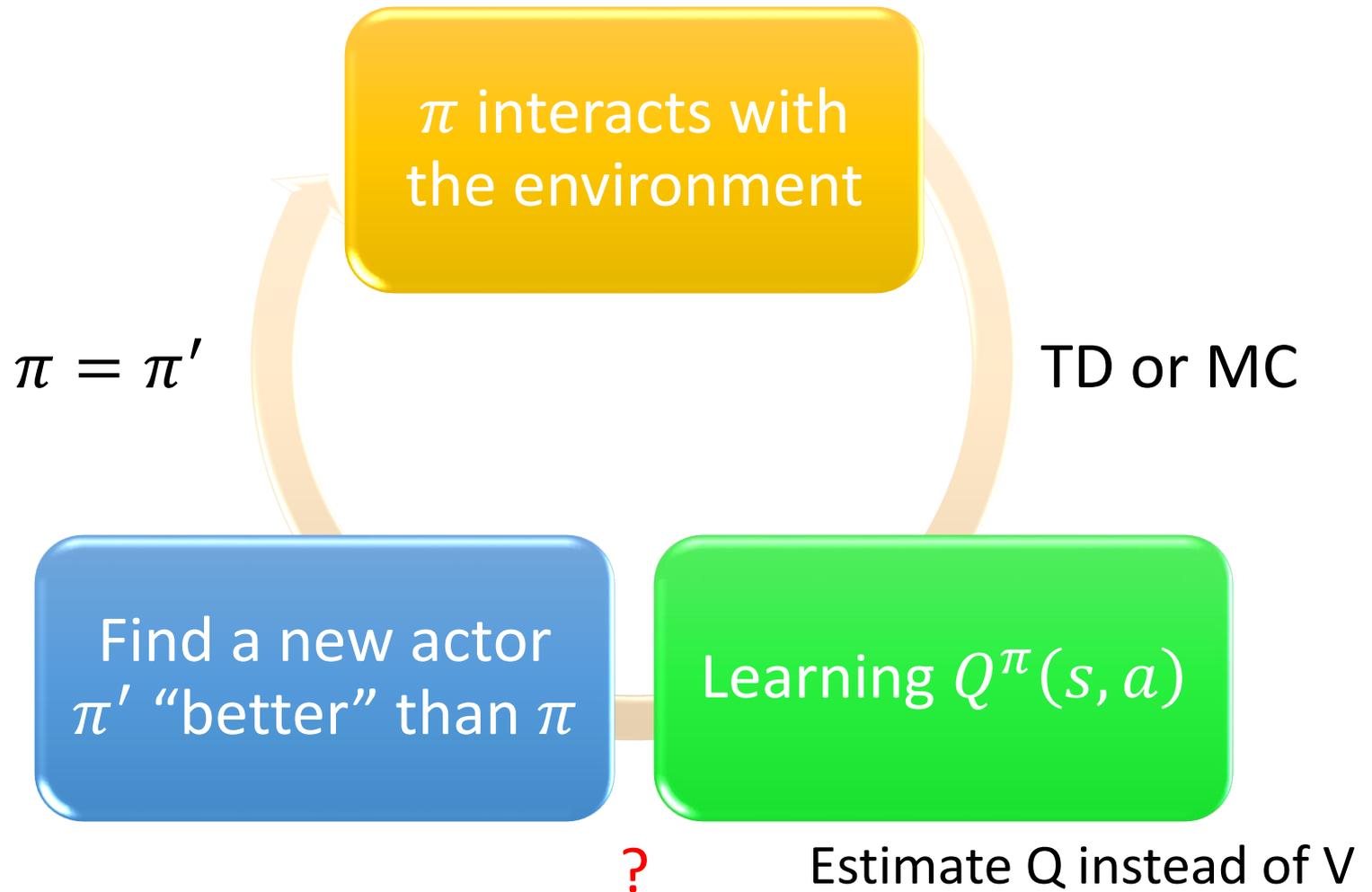


Q-Learning

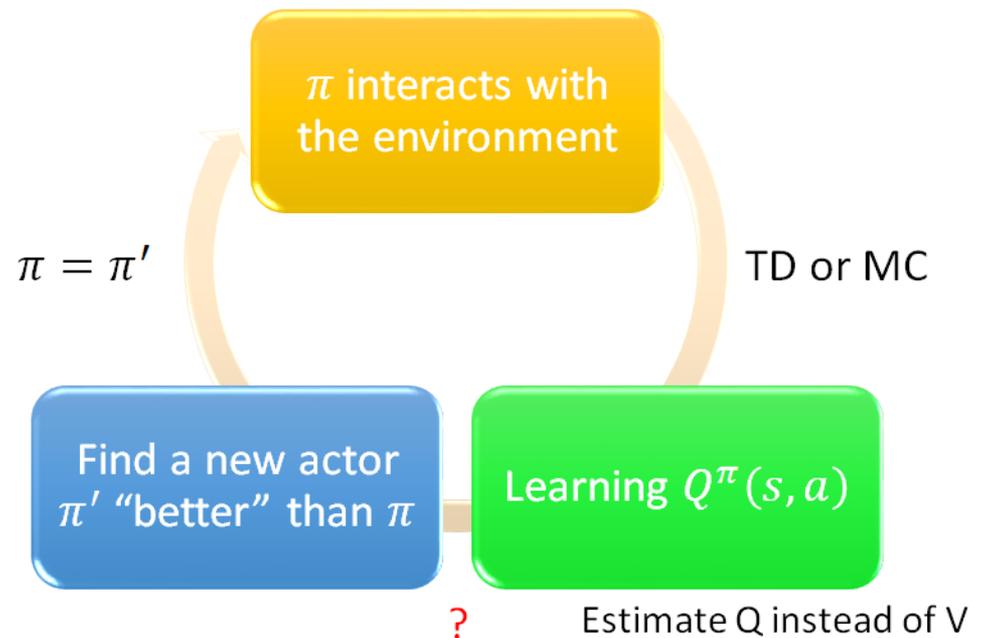
From Q function we know that taking a' at state s is better than a



Q-Learning



Q-Learning



- Given $Q^\pi(s, a)$, find a new actor π' "better" than π
 - "Better": $V^{\pi'}(s) \geq V^\pi(s)$, for all state s

$$\pi'(s) = \arg \max_a Q^\pi(s, a)$$

- π' does not have extra parameters. It depends on Q
- Not suitable for continuous action a (solve it later)

Q-Learning

$$\pi'(s) = \underset{a}{\operatorname{arg\,max}} Q^\pi(s, a)$$

$$V^{\pi'}(s) \geq V^\pi(s), \text{ for all state } s$$

$$V^\pi(s) = Q^\pi(s, \pi(s))$$

$$\leq \underset{a}{\operatorname{max}} Q^\pi(s, a) = Q^\pi(s, \pi'(s)) = V^{\pi'}(s)$$

$$V^\pi(s) \leq Q^\pi(s, \pi'(s))$$

$$= E_{\pi'}[r_{t+1} + V^\pi(s_{t+1}) | s_t = s]$$

$$\leq E_{\pi'}[r_{t+1} + Q^\pi(s_{t+1}, \pi'(s_{t+1})) | s_t = s]$$

$$= E_{\pi'}[r_{t+1} + r_{t+2} + V^\pi(s_{t+2}) | s_t = s]$$

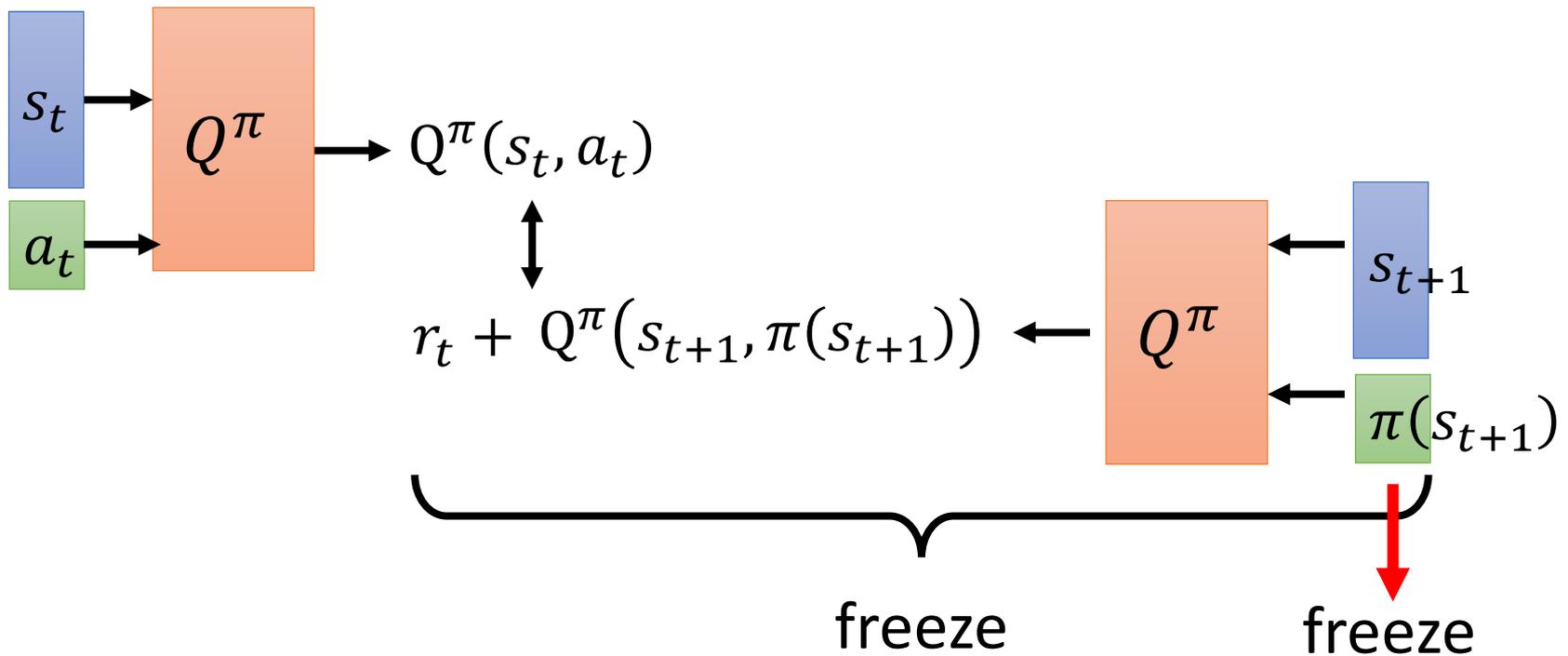
$$\leq E_{\pi'}[r_{t+1} + r_{t+2} + Q^\pi(s_{t+2}, \pi'(s_{t+2})) | s_t = s]$$

$$\dots \dots \leq V^{\pi'}(s)$$

Estimate $Q^\pi(s, a)$ by TD

$\dots s_t, a_t, r_t, s_{t+1} \dots$

$$Q^\pi(s_t, a_t) \quad Q^\pi(s_{t+1}, \pi(s_{t+1}))$$



Double DQN

- Q value is usually over estimate

$$Q(s_t, a_t) \longleftrightarrow r_t + \max_a Q(s_{t+1}, a)$$

Tend to select the action that is over-estimated



Double DQN

- Q value is usually over estimate

$$Q(s_t, a_t) \longleftrightarrow r_t + \max_a Q(s_{t+1}, a)$$

- Double DQN: two functions Q and Q'

$$Q(s_t, a_t) \longleftrightarrow r_t + Q' \left(s_{t+1}, \arg \max_a Q(s_{t+1}, a) \right)$$

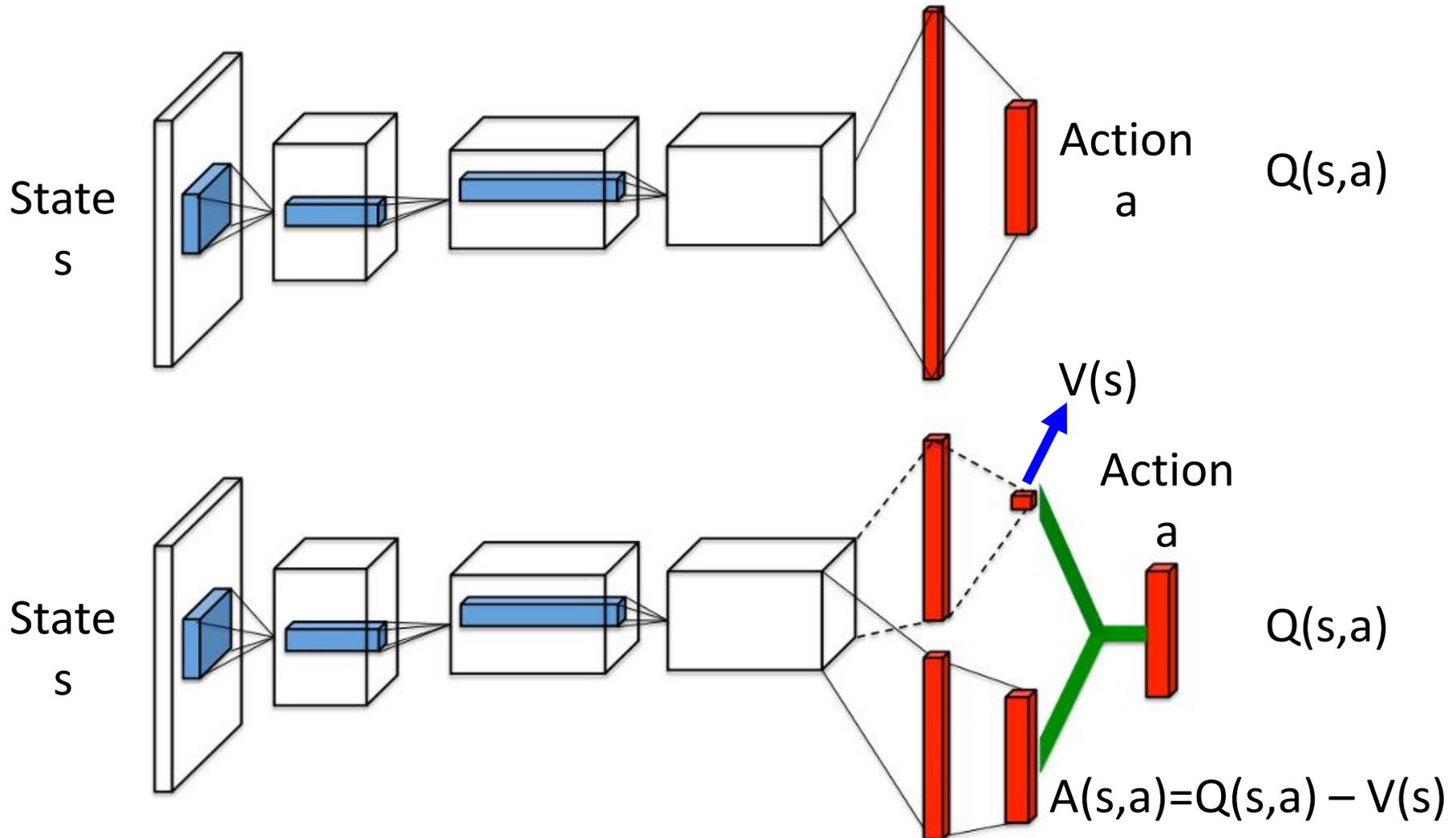
If Q over-estimate a, so it is selected. Q' would give it proper value.

How about Q' overestimate? The action will not select by Q.

Hado V. Hasselt, "Double Q-learning", NIPS 2010

Hado van Hasselt, Arthur Guez, David Silver, "Deep Reinforcement Learning with Double Q-learning", AAAI 2016

Dueling DQN



Dueling DQN - Visualization



(from the link in the original paper)

Dueling DQN - Visualization



(from the link in the original paper)

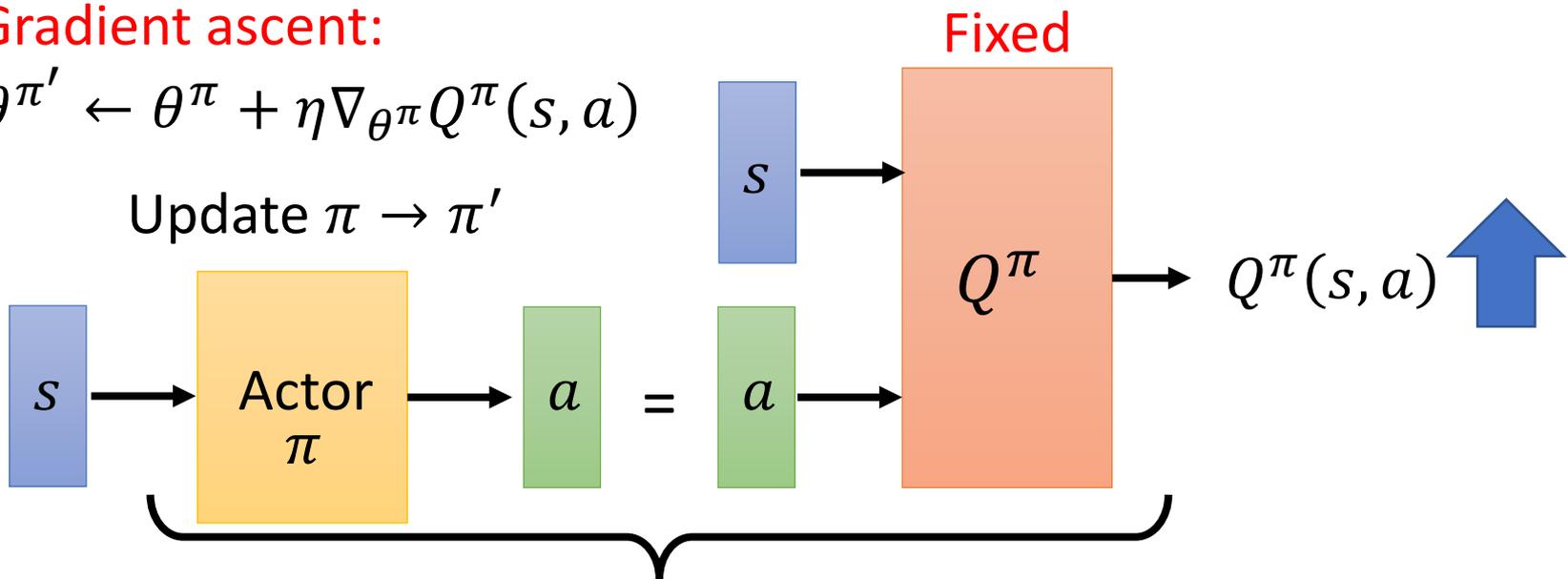
Pathwise Derivative Policy Gradient

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

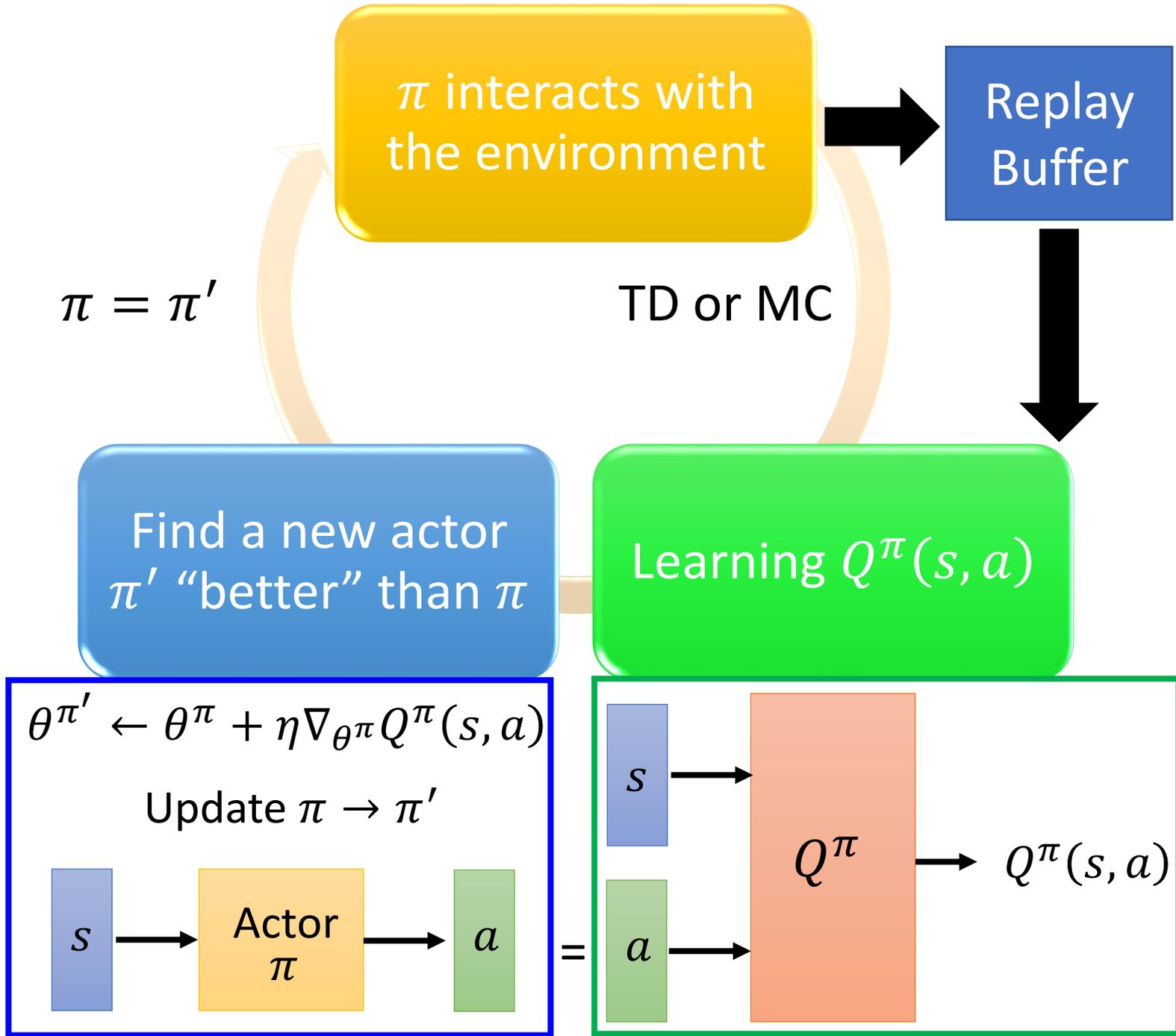
Gradient ascent:

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

Update $\pi \rightarrow \pi'$



This is a large network

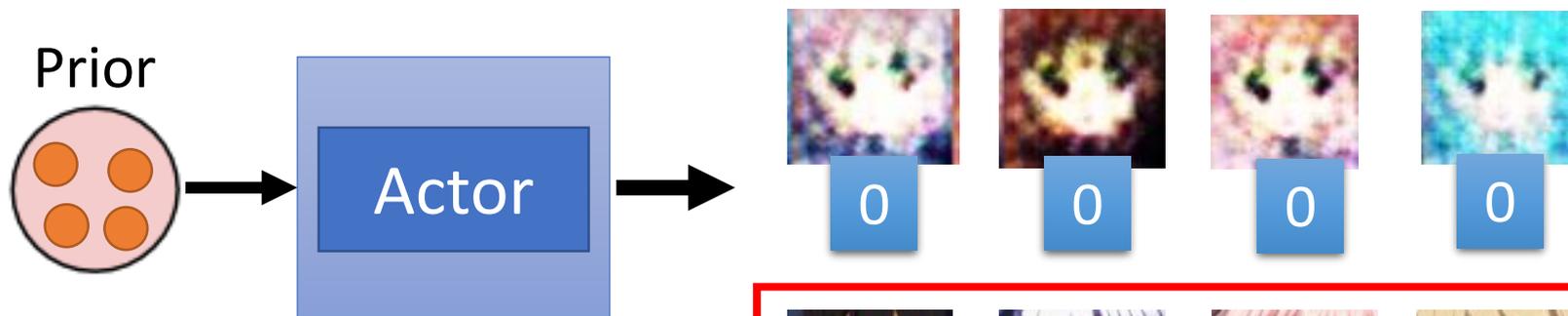


DDPG Algorithm

- Initialize critic network θ^Q and actor network θ^π
- Initialize target critic network $\theta^{Q'} = \theta^Q$ and target actor network $\theta^{\pi'} = \theta^\pi$
- Initialize replay buffer R
- In each iteration
 - Use $\pi(s) + noise$ to interact with the environment, collect a set of $\{s_t, a_t, r_t, s_{t+1}\}$, put them in R
 - Sample N examples $\{s_n, a_n, r_n, s_{n+1}\}$ from R
 - Update critic Q to minimize: $L = \sum_n (\hat{y}_n - Q(s_n, a_n))^2$
 - $\hat{y}_n = r_n + Q'(s_{n+1}, \pi'(s_{n+1}))$ **Using target networks**
 - Update actor π to maximize: $J = \sum_n Q(s_n, \pi(s_n))$
 - Update target networks:
The target networks update slower
 - $\theta^{\pi'} \leftarrow m\theta^\pi + (1 - m)\theta^{\pi'}$
 - $\theta^{Q'} \leftarrow m\theta^Q + (1 - m)\theta^{Q'}$

Connection with GAN

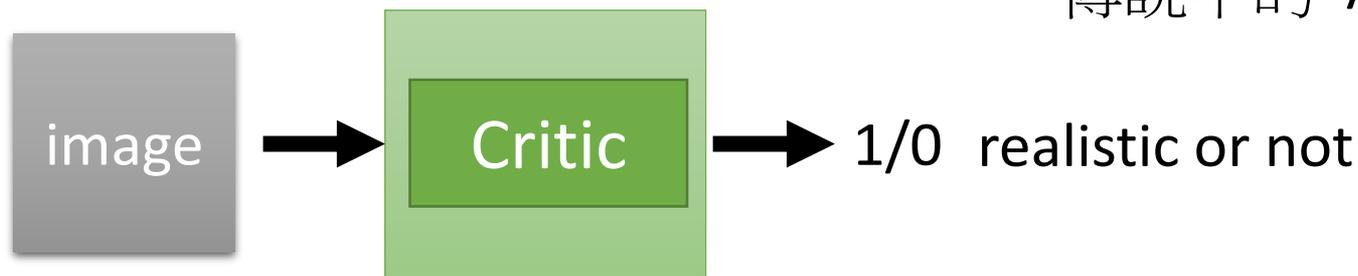
GAN as Actor-critic



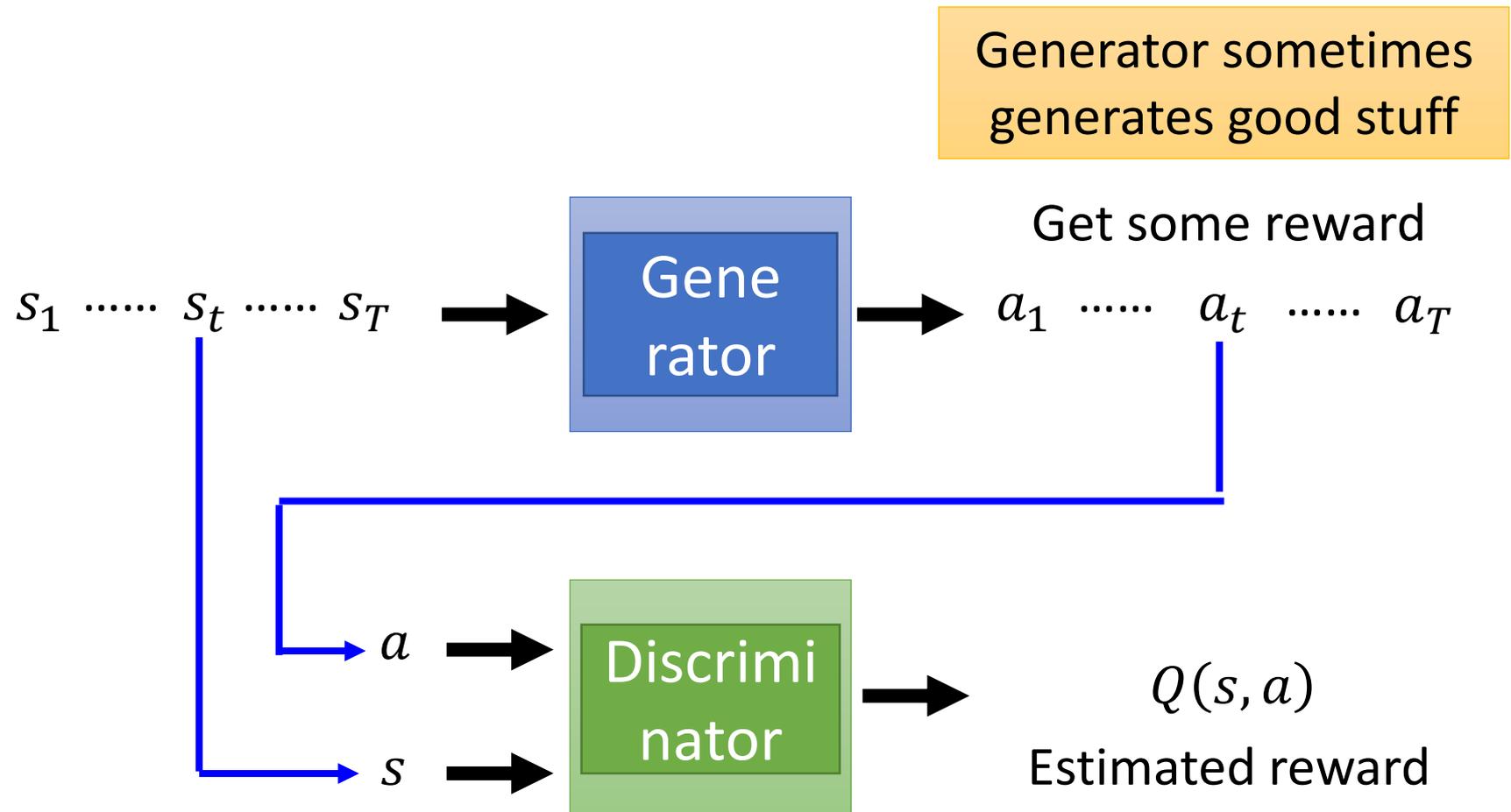
Action: generate an object



傳說中的 Action



Actor-critic as GAN



Method	GANs	AC
Freezing learning	yes	yes
Label smoothing	yes	no
Historical averaging	yes	no
Minibatch discrimination	yes	no
Batch normalization	yes	yes
Target networks	n/a	yes
Replay buffers	no	yes
Entropy regularization	no	yes
Compatibility	no	yes

David Pfau, Oriol Vinyals, "Connecting Generative Adversarial Networks and Actor-Critic Methods", arXiv preprint, 2016