# 强化学习
# (Reinforcement Learning)

梁毅雄

yxliang@csu.edu.cn

Some materials from Hong-yi Lee, Andrew Ng and others
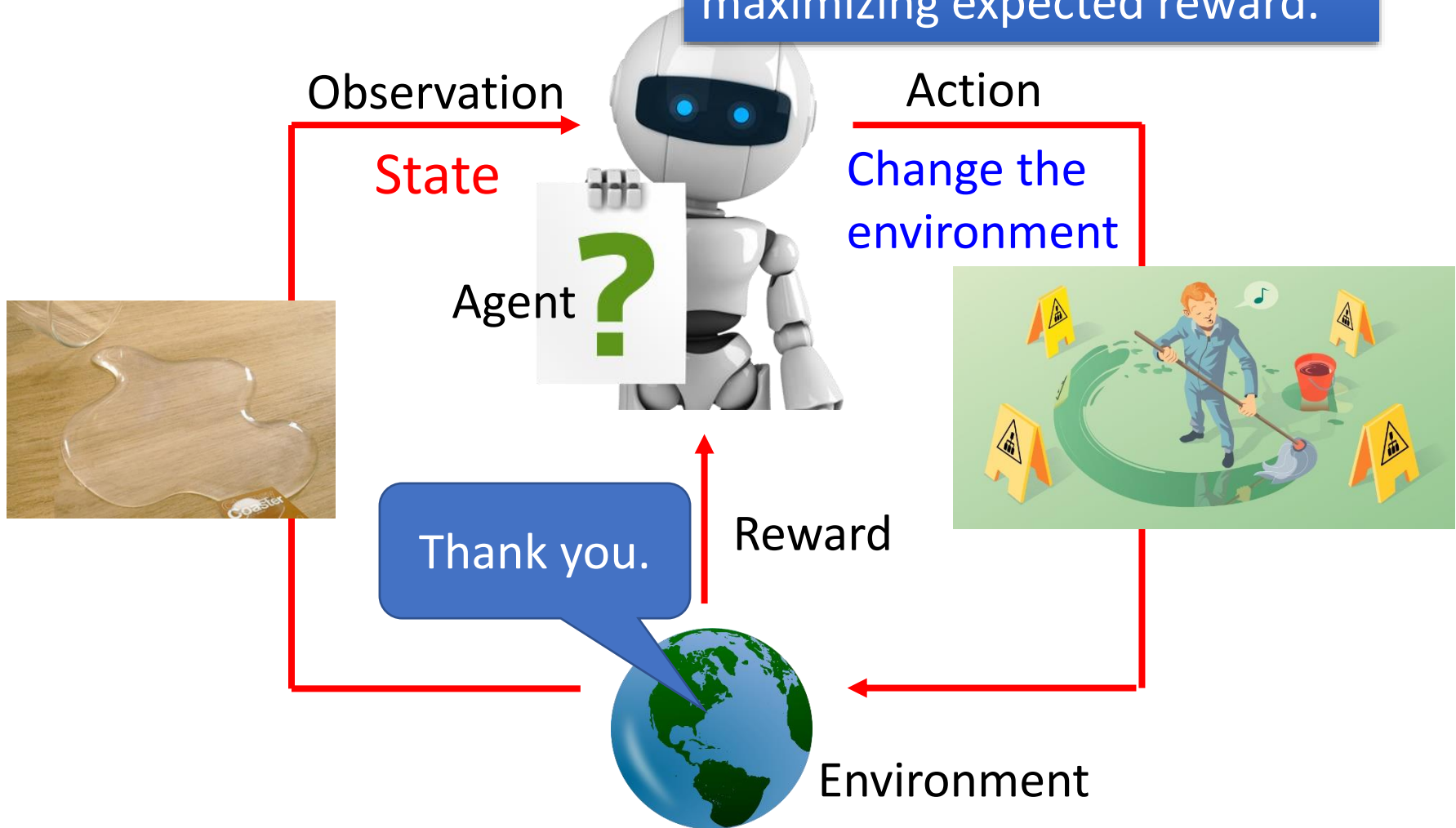
**Machine Learning**

# 强化学习

- 监督学习中对任意的$x$均给定了正确的标签$y$, 但在许多实际的应用中无法即时给出标准的正确答案，而需要反复去尝试后方可得到好的策略(实践出真知)

- 如何种西瓜: 选种，浇水、施肥、除草、杀虫等，我们其实不知道何时进行哪种操作能给出最正确的指导，因此也无法即时提供正确的监督信息去进行学习Action

- 我们并没有直接告诉Agent采取何种Action, 而是让Agent来自己发现最佳的Action。方案：取消监督信息，用reward奖励函数代替，表示对应action产生反馈

- 强化（reinforcement）学习是指从环境状态到行为映射的学习，以使系统行为从环境中获得的累积奖励值最大。在强化学习中，我们设计算法来把外界环境转化为最大化奖励量的方式的动作。
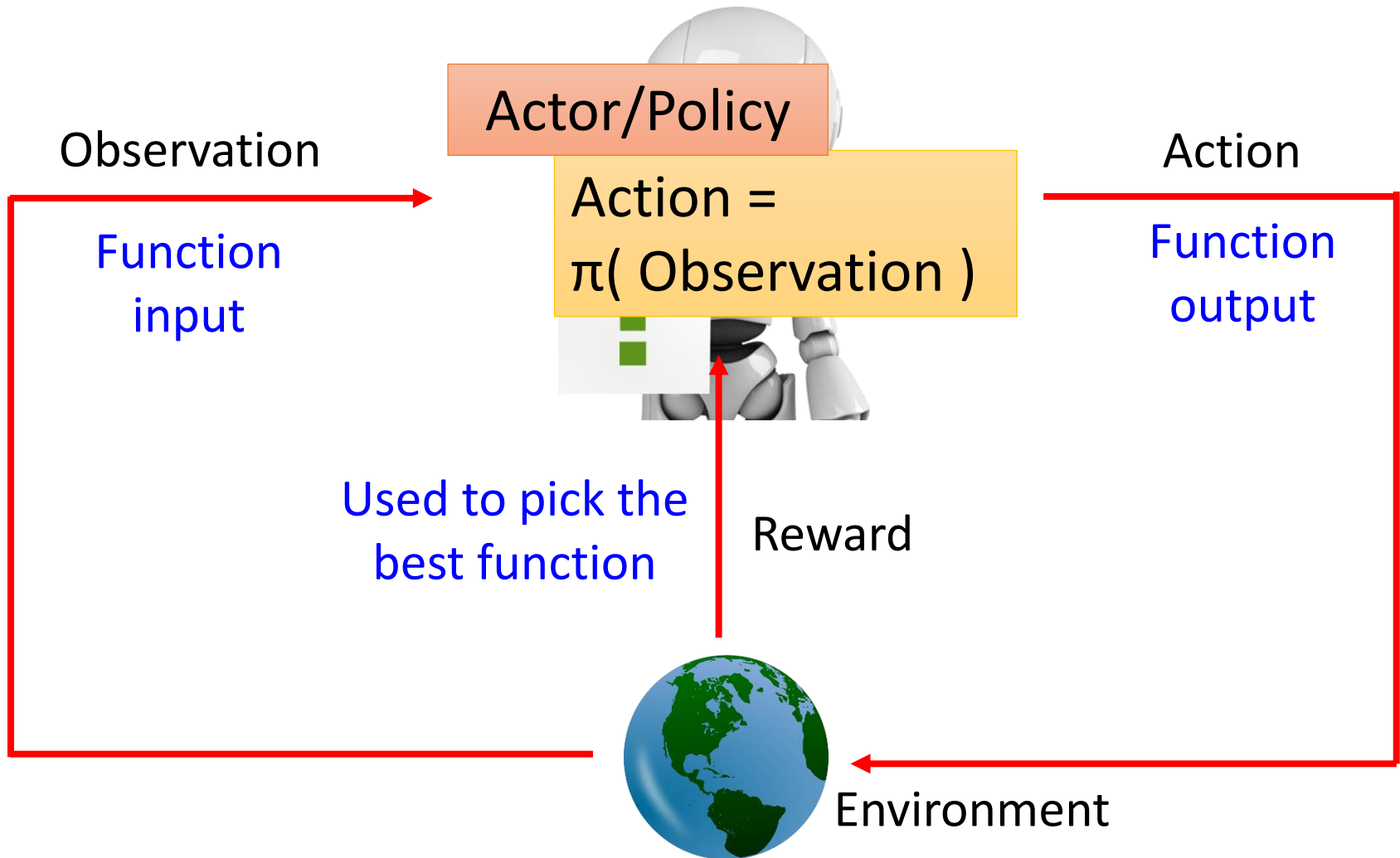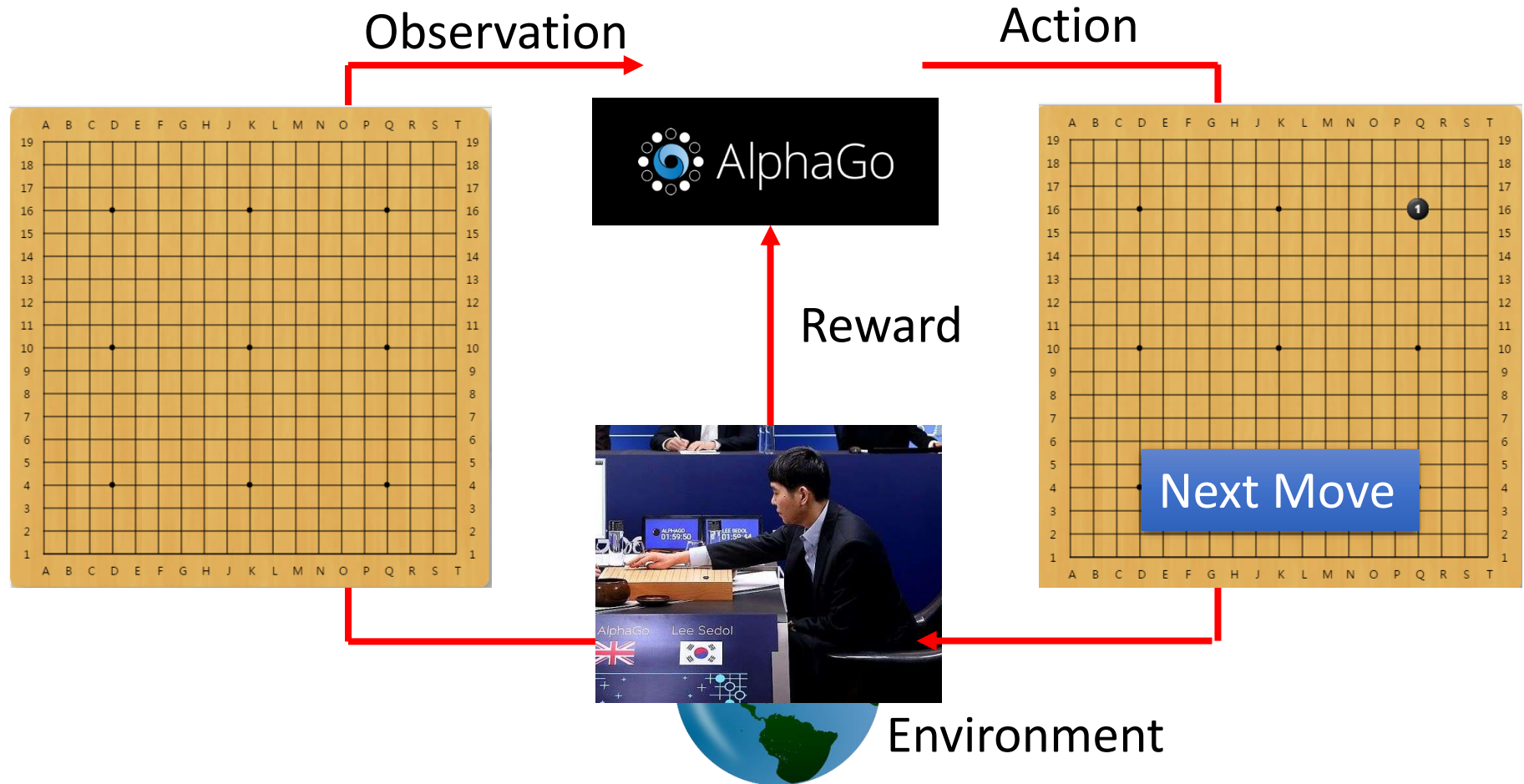
# 强化学习

# 强化学习



Agent learns to take actions maximizing expected reward.

Observation — State

Action — Change the environment

Agent

Thank you. — Reward

Environment

# Machine Learning
# ≈ Looking for a Function



Observation

Action

**Actor/Policy**

Action =
π( Observation )

**Function input**

**Function output**

**Used to pick the best function**

Reward

Environment

# Learning to play Go

Observation

Action

AlphaGo

Reward

Next Move

Environment

# Learning to play Go

Agent learns to take actions maximizing expected reward.

Observation

Action

AlphaGo

Reward

reward = 0 in most cases

If win, reward = 1

If loss, reward = -1

Environment

# Learning to play Go

- Supervised:  <span style="background-color:green; color:white">Learning from teacher</span>



Next move: "5-5"



Next move: "3-3"

- Reinforcement Learning  <span style="background-color:orange; color:white">Learning from experience</span>

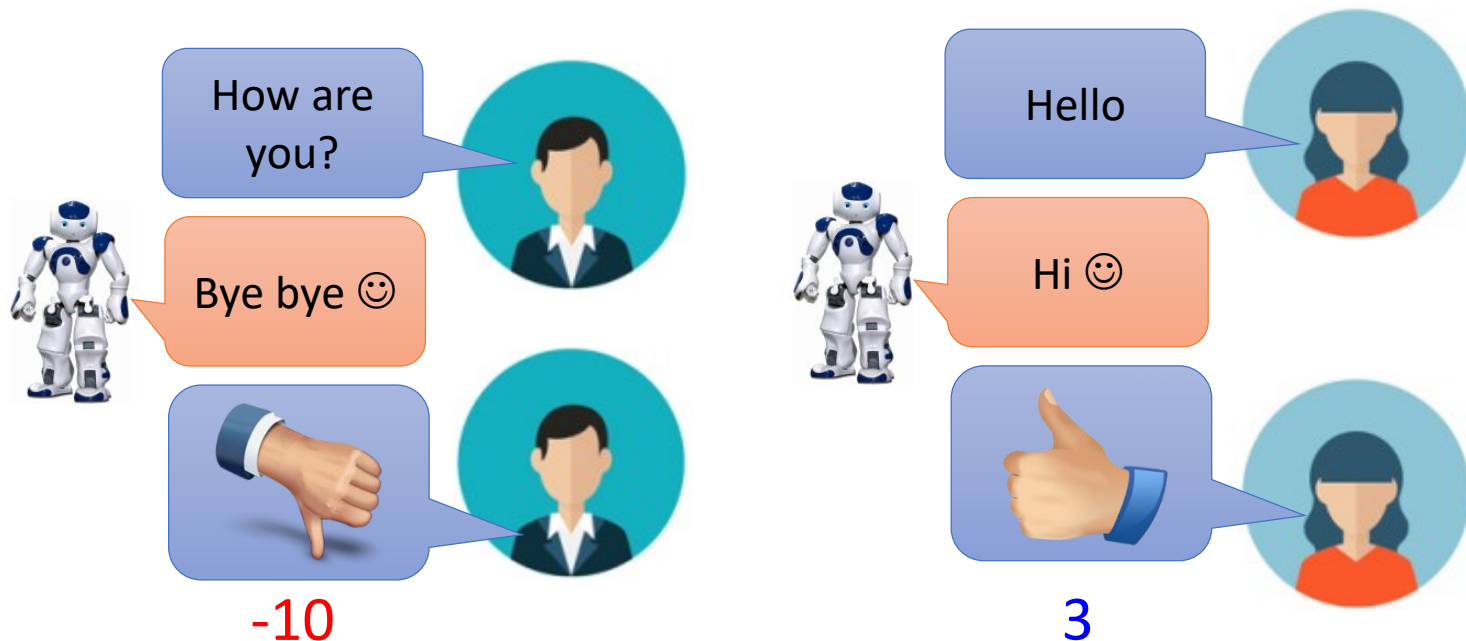First move ➡ …… many moves …… ➡ Win!

(Two agents play with each other.)

<span style="background-color:blue; color:white">Alpha Go is supervised learning + reinforcement learning.</span>

# Learning a chat-bot

- Machine obtains feedback from user



- Chat-bot learns to maximize the *expected reward*

# Learning a chat-bot

- Let two agents talk to each other (sometimes generate good dialogue, sometimes bad)

How old are you?
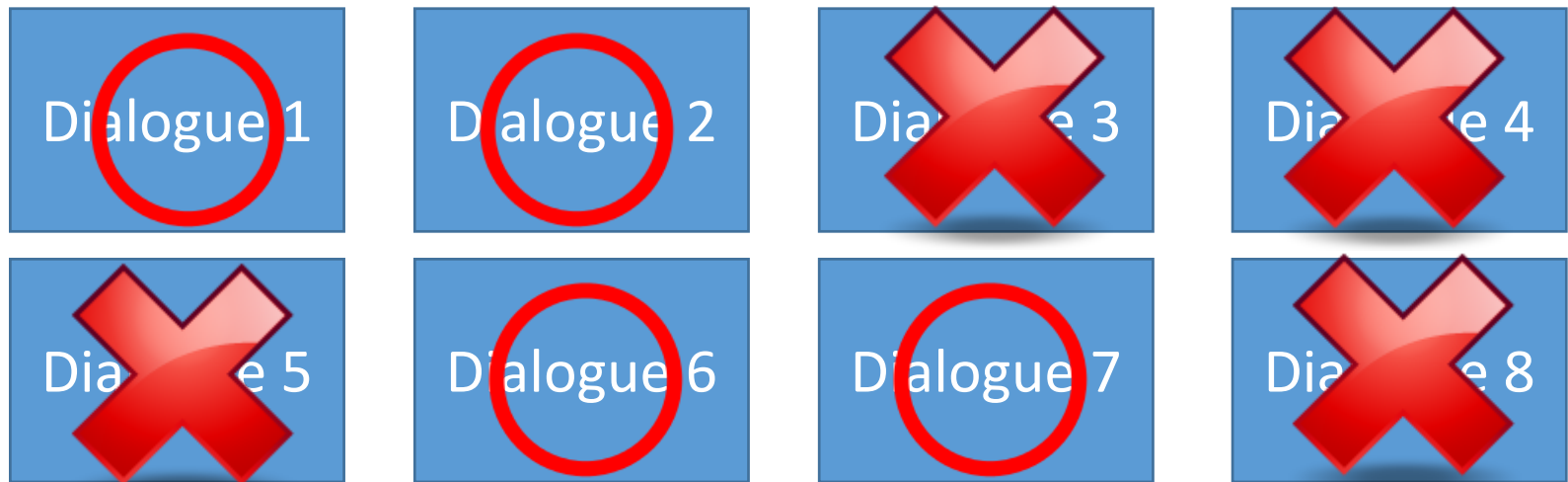
See you.

How old are you?

I am 16.

See you.

I though you were 12.

See you.

What make you think so?
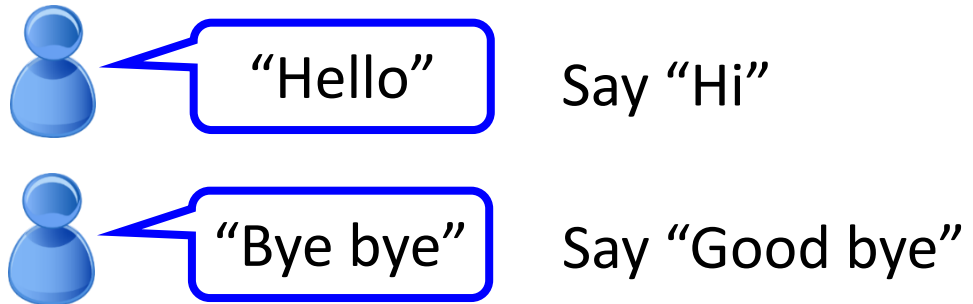
# Learning a chat-bot

- By this approach, we can generate a lot of dialogues.
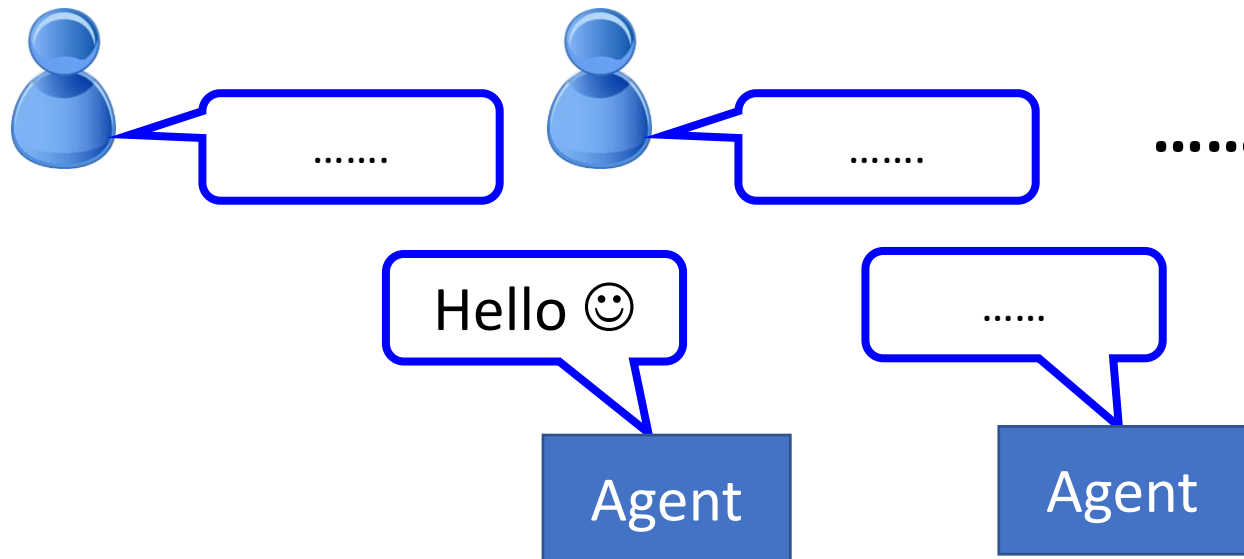- Use some pre-defined rules to evaluate the goodness of a dialogue



Machine learns from the evaluation

Deep Reinforcement Learning for Dialogue Generation
https://arxiv.org/pdf/1606.01541v3.pdf

# Learning a chat-bot

- Supervised

"Hello" — Say "Hi"

"Bye bye" — Say "Good bye"

- Reinforcement

……    ……    ……

Hello ☺    ……

Agent    Agent

Bad

# More applications

- Flying Helicopter
  - https://www.youtube.com/watch?v=0JL04JJjocc
- Driving
  - https://www.youtube.com/watch?v=0xo1Ldx3L5Q
- Robot
  - https://www.youtube.com/watch?v=370cT-OAzzM
- Google Cuts Its Giant Electricity Bill With DeepMind-Powered AI
  - http://www.bloomberg.com/news/articles/2016-07-19/google-cuts-its-giant-electricity-bill-with-deepmind-powered-ai
- Text generation
  - https://www.youtube.com/watch?v=pbQ4qe8EwLo

# Example: Playing Video Game

- Widely studies:
  - Gym: https://gym.openai.com/
  - Universe: https://openai.com/blog/universe/

  Machine learns to play video games as human players
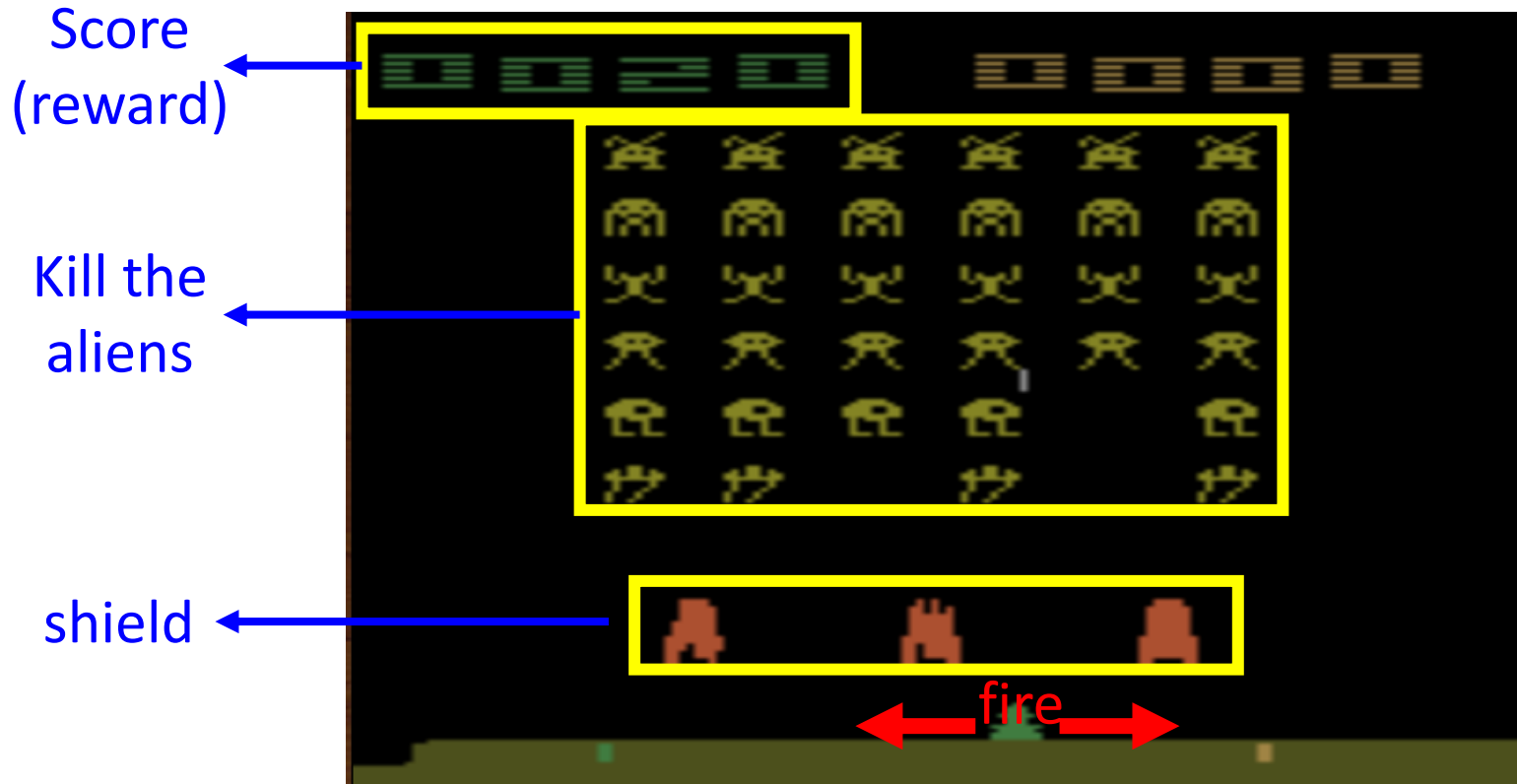
  ➢ What machine observes is pixels

  ➢ Machine learns to take proper action itself

# Example: Playing Video Game

- Space invader

Termination: all the aliens are killed, or your spaceship is destroyed.

Score (reward)

Kill the aliens

shield

fire

# Example: Playing Video Game

- Space invader
  - Play yourself: http://www.2600online.com/spaceinvaders.html
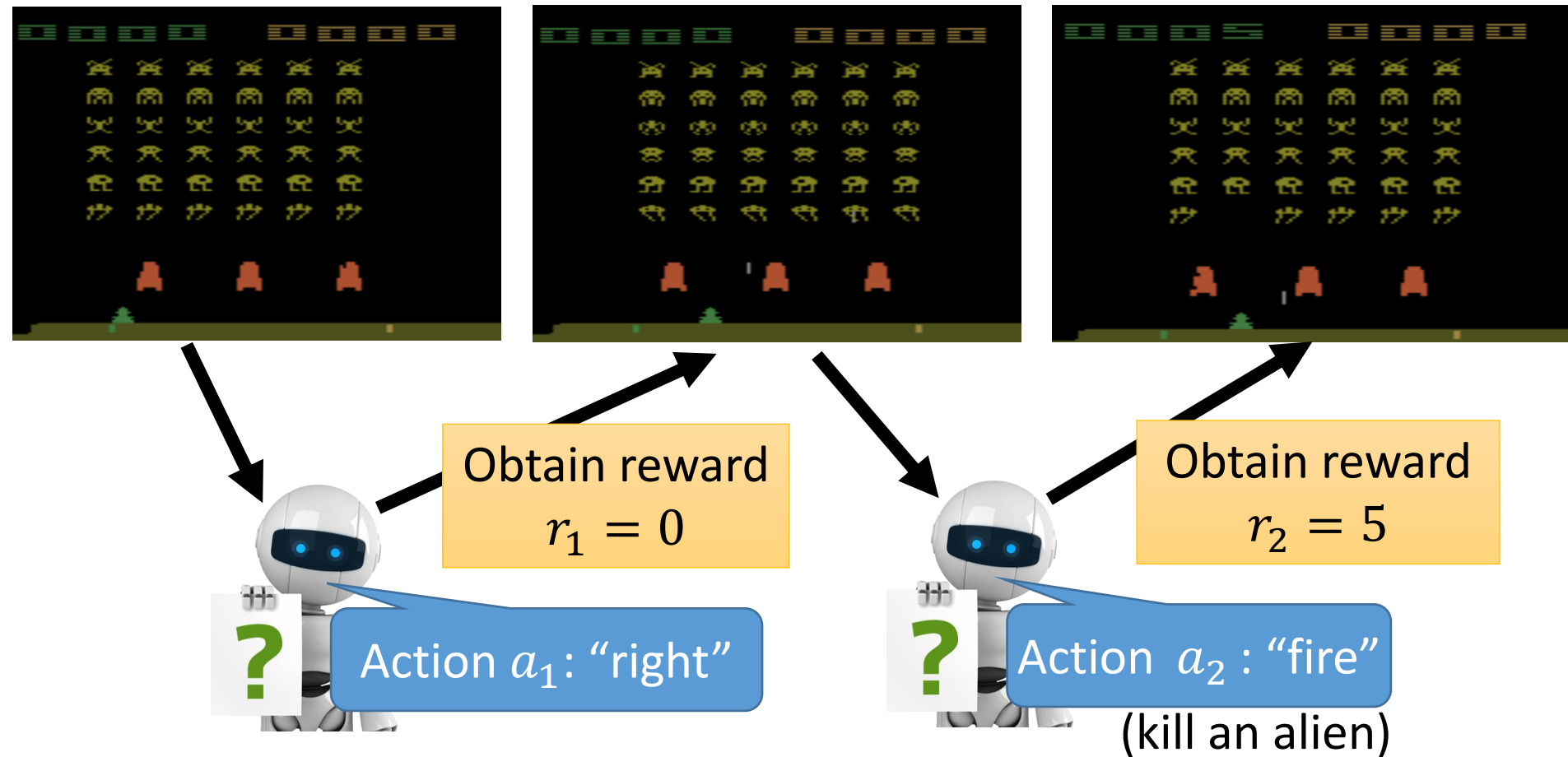  - How about machine: https://gym.openai.com/evaluations/eval_Eduozx4HRyqgTCVk9ltw

# 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

Game Over (spaceship destroyed)

Obtain reward $r_T$

Action $a_T$

This is an **episode**.

Learn to maximize the expected cumulative reward per episode
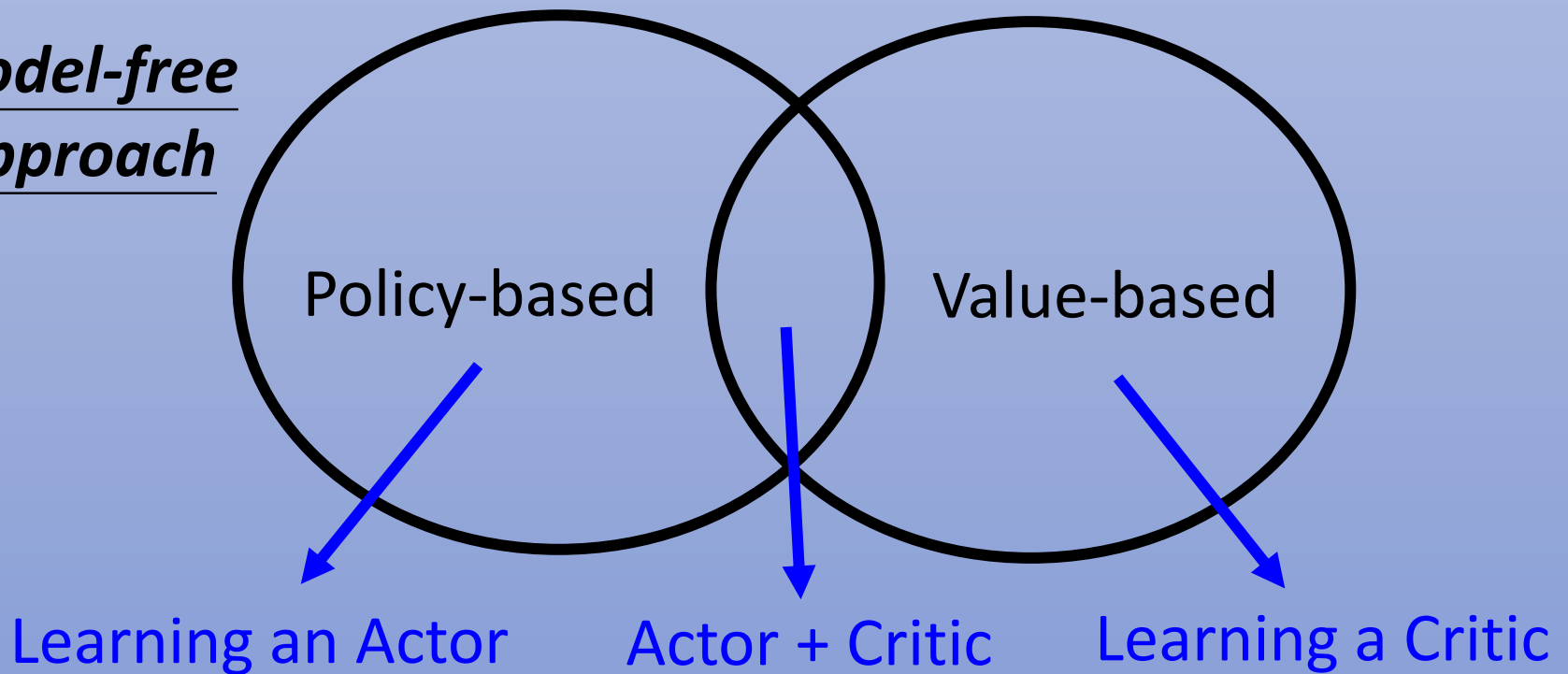
# Properties of Reinforcement Learning

- **Reward delay**
  - In space invader, only "fire" obtains reward
    - Although the moving before "fire" is important
  - In Go playing, it may be better to sacrifice immediate reward to gain more long-term reward
- Agent's actions **affect the subsequent data it receives**
  - E.g. Exploration

# Outline

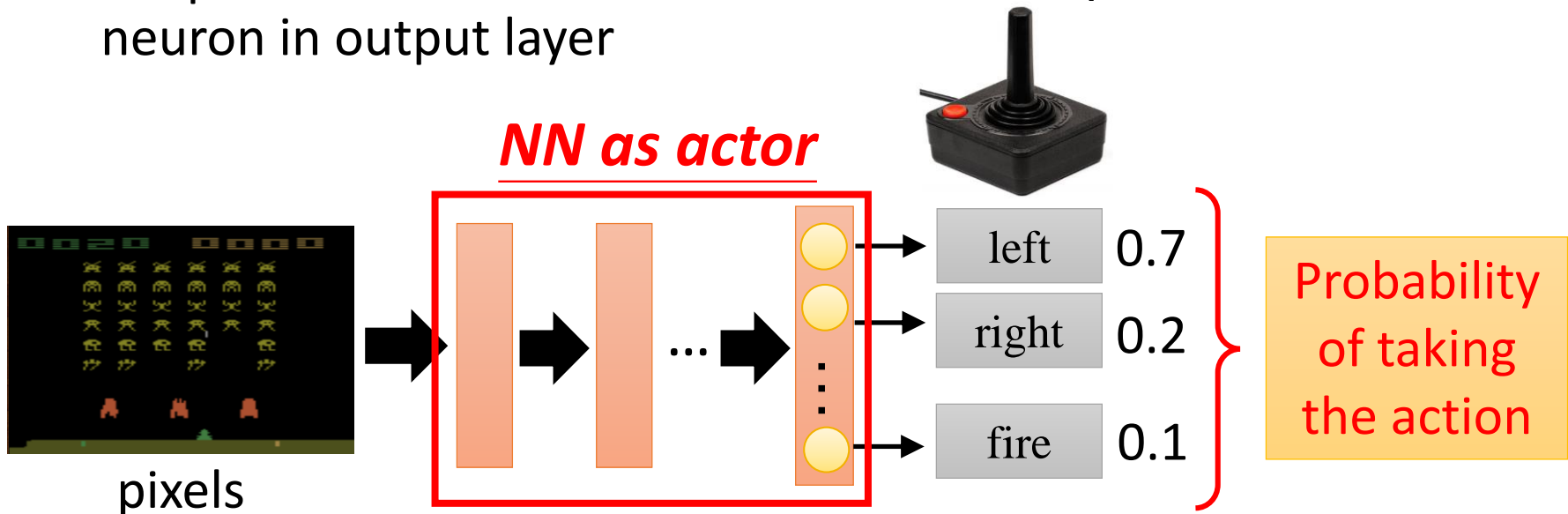Alpha Go: policy-based + value-based + model-based

***Model-free Approach***

Policy-based

Value-based

Learning an Actor

Actor + Critic

Learning a Critic

***Model-based Approach***

# Policy-based Approach

## Learning an Actor

# Neural network as Actor

- 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

**NN as actor**



| | | |
|---|---|---|
| left | 0.7 | |
| right | 0.2 | |
| fire | 0.1 | |

pixels

Probability of taking the action

What is the benefit of using network instead of lookup table?

generalization

# Goodness of Actor

- Review: Supervised learning

Total Loss:

$$L = \sum_{n=1}^{N} l_n$$

Training Example



"1"

$x_1$
$x_2$
$\vdots$
$x_{256}$

Softmax

$y_1$
$y_2$
$\vdots$
$y_{10}$

Given a set of parameters $\theta$

As close as possible

Loss $l$

1
0
$\vdots$
0

target

# Goodness of Actor

- Given an actor $\pi_\theta(s)$ with network parameter $\theta$
- Use the actor $\pi_\theta(s)$ 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: $R_\theta = \sum_{t=1}^{T} r_t$

Even with the same actor, $R_\theta$ is different each time

Randomness in the actor and the game

We define $\bar{R}_\theta$ as the *expected value* of $R_\theta$

$\bar{R}_\theta$ evaluates the goodness of an actor $\pi_\theta(s)$

# Goodness of Actor

- $\tau = \{s_1, a_1, r_1, s_2, a_2, r_2, \cdots, s_T, a_T, r_T\}$

$$P(\tau|\theta) =$$

$$p(s_1)p(a_1|s_1,\theta)p(r_1,s_2|s_1,a_1)p(a_2|s_2,\theta)p(r_2,s_3|s_2,a_2)\cdots$$

$$= p(s_1)\prod_{t=1}^{T} p(a_t|s_t,\theta)p(r_t,s_{t+1}|s_t,a_t)$$

$p(a_t = "fire"|s_t,\theta) = 0.7$

not related to your actor

Control by your actor $\pi_\theta$

$s_t$ → Actor $\pi_\theta$

left → 0.1

right → 0.2

fire → 0.7

# Goodness of Actor

- An episode is considered as a trajectory $\tau$
  - $\tau = \{s_1, a_1, r_1, s_2, a_2, r_2, \cdots, s_T, a_T, r_T\}$
  - $R(\tau) = \sum_{t=1}^{T} r_t$
  - If you use an actor to play the game, each $\tau$ has a probability to be sampled
    - The probability depends on actor parameter $\theta$: $P(\tau|\theta)$

$$\bar{R}_\theta = \boxed{\sum_\tau} R(\tau) \boxed{P(\tau|\theta)} \approx \boxed{\frac{1}{N} \sum_{n=1}^{N}} R(\tau^n)$$

Sum over all possible trajectory

Use $\pi_\theta$ to play the game N times, obtain $\{\tau^1, \tau^2, \cdots, \tau^N\}$

Sampling $\tau$ from $P(\tau|\theta)$ N times

# Gradient Ascent

- Problem statement

$$\theta^* = arg \max_{\theta} \bar{R}_{\theta}$$

- Gradient ascent
  - Start with $\theta^0$
  - $\theta^1 \leftarrow \theta^0 + \eta \nabla \bar{R}_{\theta^0}$
  - $\theta^2 \leftarrow \theta^1 + \eta \nabla \bar{R}_{\theta^1}$
  - ……

$$\theta = \{w_1, w_2, \cdots, b_1, \cdots\}$$

$$\nabla \bar{R}_{\theta} = \begin{bmatrix} \partial \bar{R}_{\theta} / \partial w_1 \\ \partial \bar{R}_{\theta} / \partial w_2 \\ \vdots \\ \partial \bar{R}_{\theta} / \partial b_1 \\ \vdots \end{bmatrix}$$

# Policy Gradient

$$\bar{R}_\theta = \sum_\tau R(\tau)P(\tau|\theta) \quad \nabla\bar{R}_\theta = ?$$

$$\nabla\bar{R}_\theta = \sum_\tau R(\tau)\nabla P(\tau|\theta) = \sum_\tau R(\tau)P(\tau|\theta)\frac{\nabla P(\tau|\theta)}{P(\tau|\theta)}$$

$R(\tau)$ do not have to be differentiable

It can even be a black box.

$$= \sum_\tau R(\tau)P(\tau|\theta)\nabla log P(\tau|\theta) \qquad \boxed{\frac{dlog(f(x))}{dx} = \frac{1}{f(x)}\frac{df(x)}{dx}}$$

$$\approx \frac{1}{N}\sum_{n=1}^{N} R(\tau^n)\nabla log P(\tau^n|\theta)$$

Use $\pi_\theta$ to play the game N times,
Obtain $\{\tau^1, \tau^2, \cdots, \tau^N\}$

# Policy Gradient

$\nabla log P(\tau | \theta) = ?$

- $\tau = \{s_1, a_1, r_1, s_2, a_2, r_2, \cdots, s_T, a_T, r_T\}$

$$P(\tau | \theta) = p(s_1) \prod_{t=1}^{T} p(a_t | s_t, \theta) p(r_t, s_{t+1} | s_t, a_t)$$

$$log P(\tau | \theta)$$

$$= log p(s_1) + \sum_{t=1}^{T} log p(a_t | s_t, \theta) + log p(r_t, s_{t+1} | s_t, a_t)$$

$$\nabla log P(\tau | \theta) = \sum_{t=1}^{T} \nabla log p(a_t | s_t, \theta)$$

Ignore the terms not related to $\theta$

# Policy Gradient

$$\boxed{\begin{aligned}&\nabla log P(\tau|\theta)\\ &= \sum_{t=1}^{T} \nabla log p(a_t|s_t,\theta)\end{aligned}}$$

$$\theta^{new} \leftarrow \theta^{old} + \eta \nabla \bar{R}_{\theta^{old}}$$

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

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

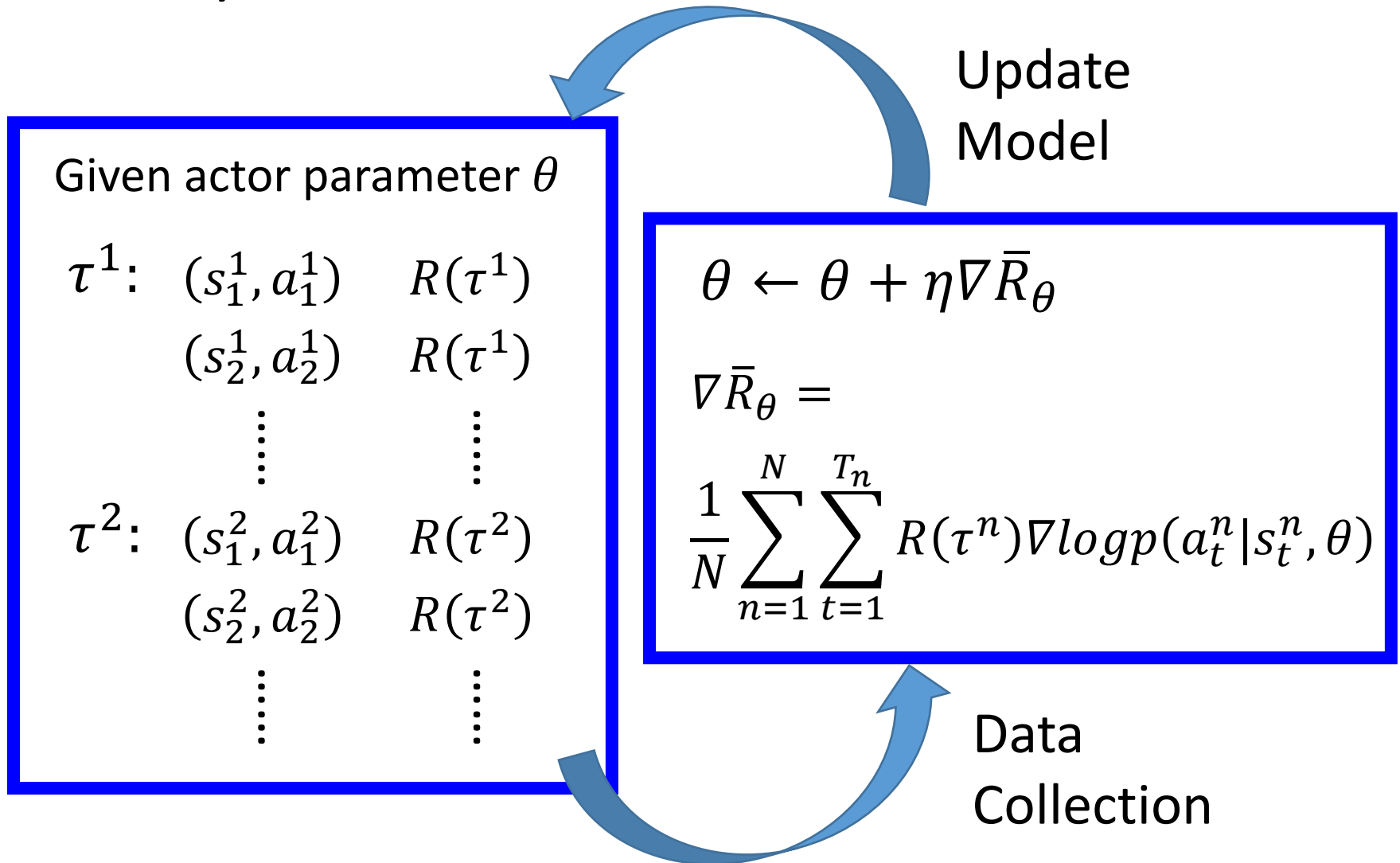What if we replace $R(\tau^n)$ with $r_t^n$ ......

If in $\tau^n$ machine takes $a_t^n$ when seeing $s_t^n$ in

$R(\tau^n)$ is positive ➡ Tuning $\theta$ to increase $p(a_t^n|s_t^n)$

$R(\tau^n)$ is negative ➡ Tuning $\theta$ 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 $\tau^n$ instead of immediate reward $r_t^n$
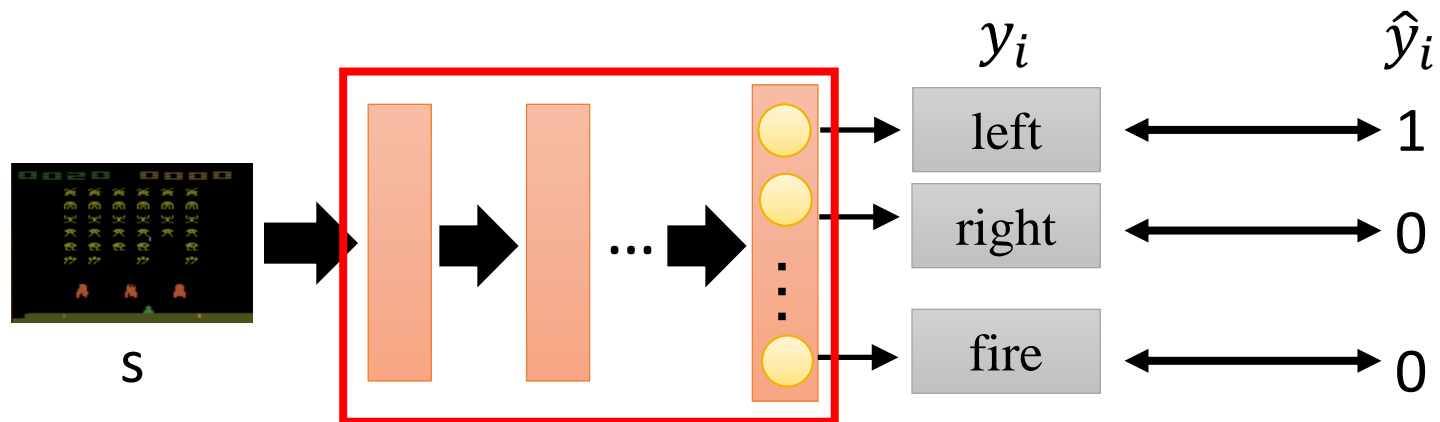
# Policy Gradient

Given actor parameter $\theta$

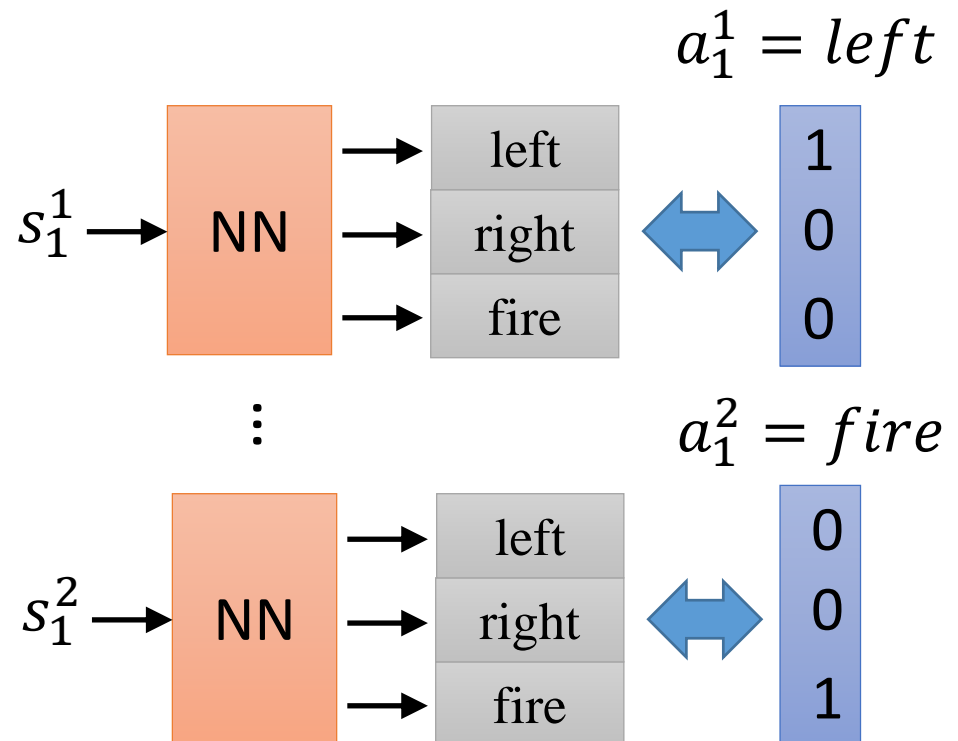$\tau^1$: $(s_1^1, a_1^1)$     $R(\tau^1)$

       $(s_2^1, a_2^1)$     $R(\tau^1)$

       $\vdots$         $\vdots$

$\tau^2$: $(s_1^2, a_1^2)$     $R(\tau^2)$

       $(s_2^2, a_2^2)$     $R(\tau^2)$

       $\vdots$         $\vdots$

Update Model

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

$$\nabla \bar{R}_\theta =$$

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

Data Collection

# Policy Gradient

***Considered as***
***Classification Problem***

Minimize: $-\displaystyle\sum_{i=1}^{3} \hat{y}_i log y_i$

$y_i$                $\hat{y}_i$



s

left $\longleftrightarrow$ 1

right $\longleftrightarrow$ 0

fire $\longleftrightarrow$ 0

Maximize:   $log y_i =$

$log P("left"|s)$

$\theta \leftarrow \theta + \eta \nabla log P("left"|s)$

# Policy Gradient

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

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^{N} \sum_{t=1}^{T_n} \boxed{\phantom{xx}} \nabla log p(a_t^n | s_t^n, \theta)$$
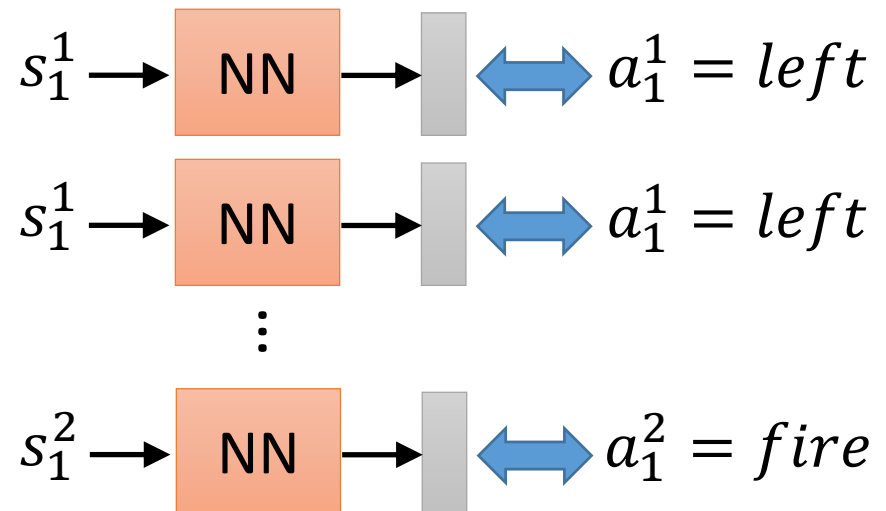
Given actor parameter $\theta$

$\tau^1:$  $(s_1^1, a_1^1)$  $R(\tau^1)$

$\quad\quad (s_2^1, a_2^1)$  $R(\tau^1)$

$\quad\quad \vdots \quad\quad\quad \vdots$

$\tau^2:$  $(s_1^2, a_1^2)$  $R(\tau^2)$

$\quad\quad (s_2^2, a_2^2)$  $R(\tau^2)$

$\quad\quad \vdots \quad\quad\quad \vdots$

$a_1^1 = left$

$s_1^1 \rightarrow$ NN $\rightarrow$ left / right / fire $\Longleftrightarrow$ 1 0 0

$a_1^2 = fire$

$s_1^2 \rightarrow$ NN $\rightarrow$ left / right / fire $\Longleftrightarrow$ 0 0 1

# Policy Gradient

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

$$\nabla \bar{R}_\theta =$$

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

Given actor parameter $\theta$

$\tau^1:$ $(s_1^1, a_1^1)$ $R(\tau^1)$ **2**

$(s_2^1, a_2^1)$ $R(\tau^1)$ **2**

$\vdots$ $\vdots$

$\tau^2:$ $(s_1^2, a_1^2)$ $R(\tau^2)$ **1**

$(s_2^2, a_2^2)$ $R(\tau^2)$ **1**

$\vdots$ $\vdots$

Each training data is weighted by $R(\tau^n)$



$s_1^1 \rightarrow$ NN $\rightarrow$ $\longleftrightarrow$ $a_1^1 = left$

$s_1^1 \rightarrow$ NN $\rightarrow$ $\longleftrightarrow$ $a_1^1 = left$

$\vdots$

$s_1^2 \rightarrow$ NN $\rightarrow$ $\longleftrightarrow$ $a_1^2 = fire$

# Add a Baseline

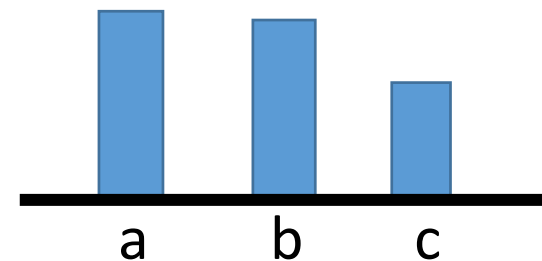It is possible that $R(\tau^n)$ is always positive.

$$\theta^{new} \leftarrow \theta^{old} + \eta \nabla \bar{R}_{\theta^{old}}$$

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

Ideal case
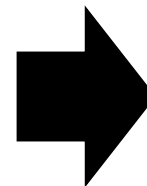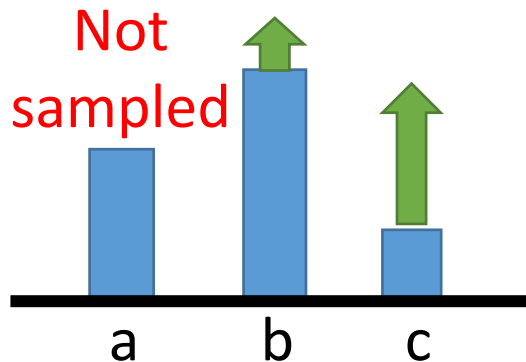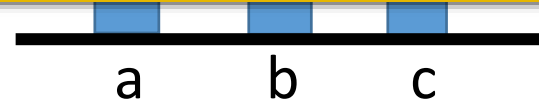
It is probability …

a    b    c

a    b    c

Sampling ……

Not sampled

The probability of the actions not sampled will decrease.

a    b    c

a    b    c

# Value-based Approach

## Learning a Critic

# Critic

- A critic does not determine the action.
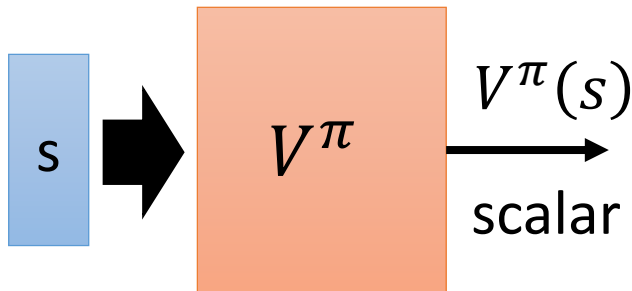- Given an actor π, it evaluates the how good the actor is

An actor can be
found from a critic.

e.g. Q-learning

http://combiboilersleeds.com/picaso/critics/critics-4.html

# Critic

- State value function $V^\pi(s)$
  - When using actor $\pi$, the *cumulated* reward expects to be obtained after seeing observation (state) s
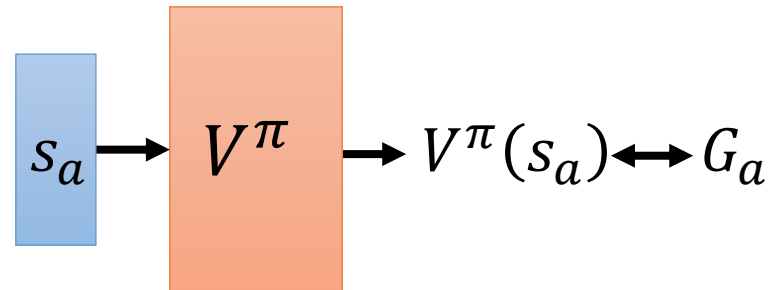


$$V^\pi(s)$$ is large   $$V^\pi(s)$$ is smaller

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

- Monte-Carlo based approach
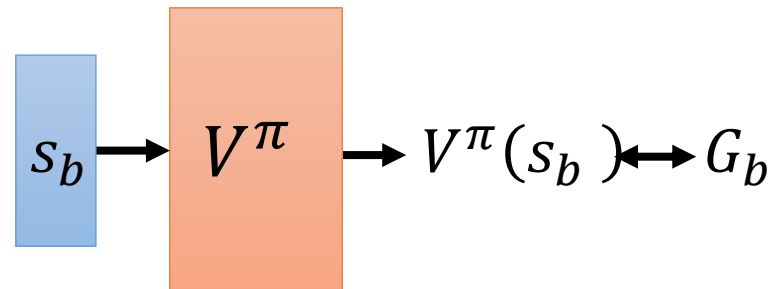  - The critic watches $\pi$ playing the game

After seeing $s_a$,

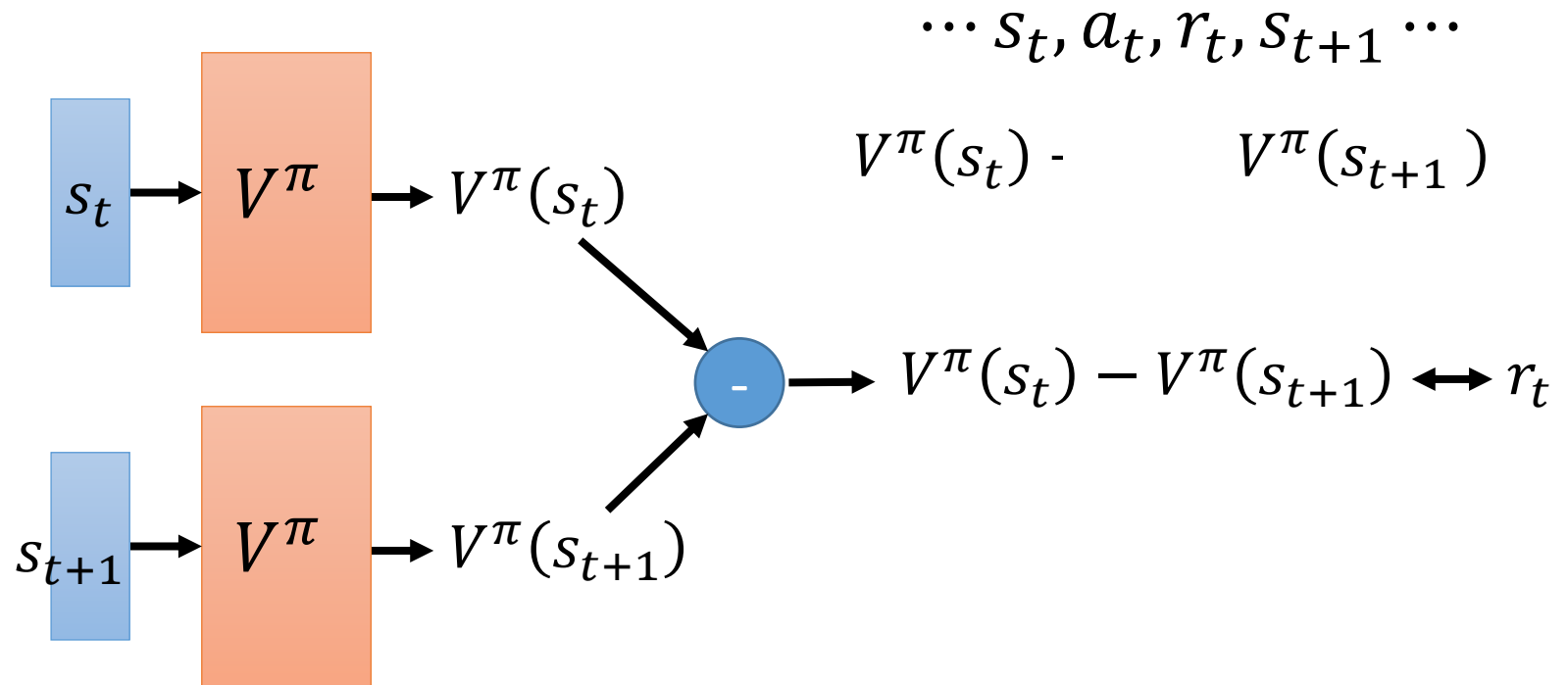Until the end of the episode, the cumulated reward is $G_a$

$$s_a \rightarrow V^\pi \rightarrow V^\pi(s_a) \longleftrightarrow G_a$$

After seeing $s_b$,

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

$$s_b \rightarrow V^\pi \rightarrow V^\pi(s_b) \longleftrightarrow G_b$$

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

- Temporal-difference approach

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

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

$s_t \rightarrow V^\pi \rightarrow V^\pi(s_t)$

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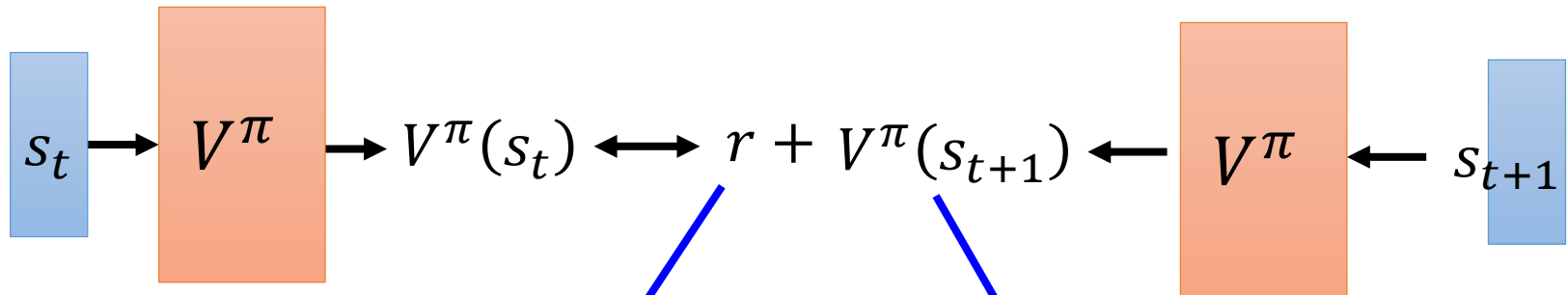
$$V^\pi(s_t) - V^\pi(s_{t+1}) \longleftrightarrow r_t$$

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

# MC v.s. TD



$$s_a \rightarrow V^\pi \rightarrow V^\pi(s_a) \longleftrightarrow G_a$$

Larger variance
unbiased

$$s_t \rightarrow V^\pi \rightarrow V^\pi(s_t) \longleftrightarrow r + V^\pi(s_{t+1}) \leftarrow V^\pi \leftarrow s_{t+1}$$

Smaller variance

May be biased

# MC v.s. TD

- The critic has the following 8 episodes
  - $s_a, r = 0, s_b, r = 0$, END
  - $s_b, r = 1$, END
  - $s_b, r = 1$, END
  - $s_b, r = 1$, END
  - $s_b, r = 1$, END
  - $s_b, r = 1$, END
  - $s_b, r = 1$, END
  - $s_b, r = 0$, END

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

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

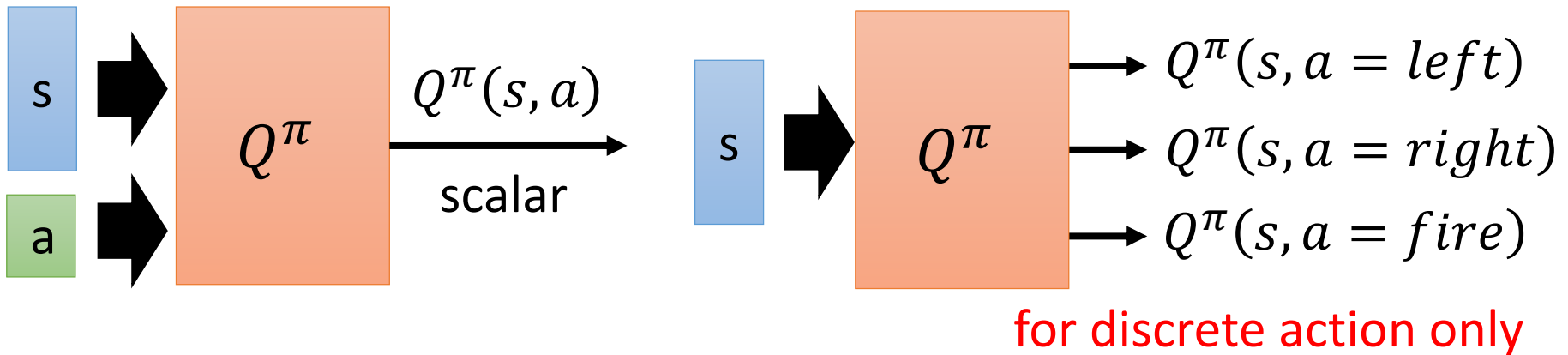Monte-Carlo: $\quad V^\pi(s_a) = 0$

Temporal-difference:

$$V^\pi(s_b) + r = V^\pi(s_a)$$
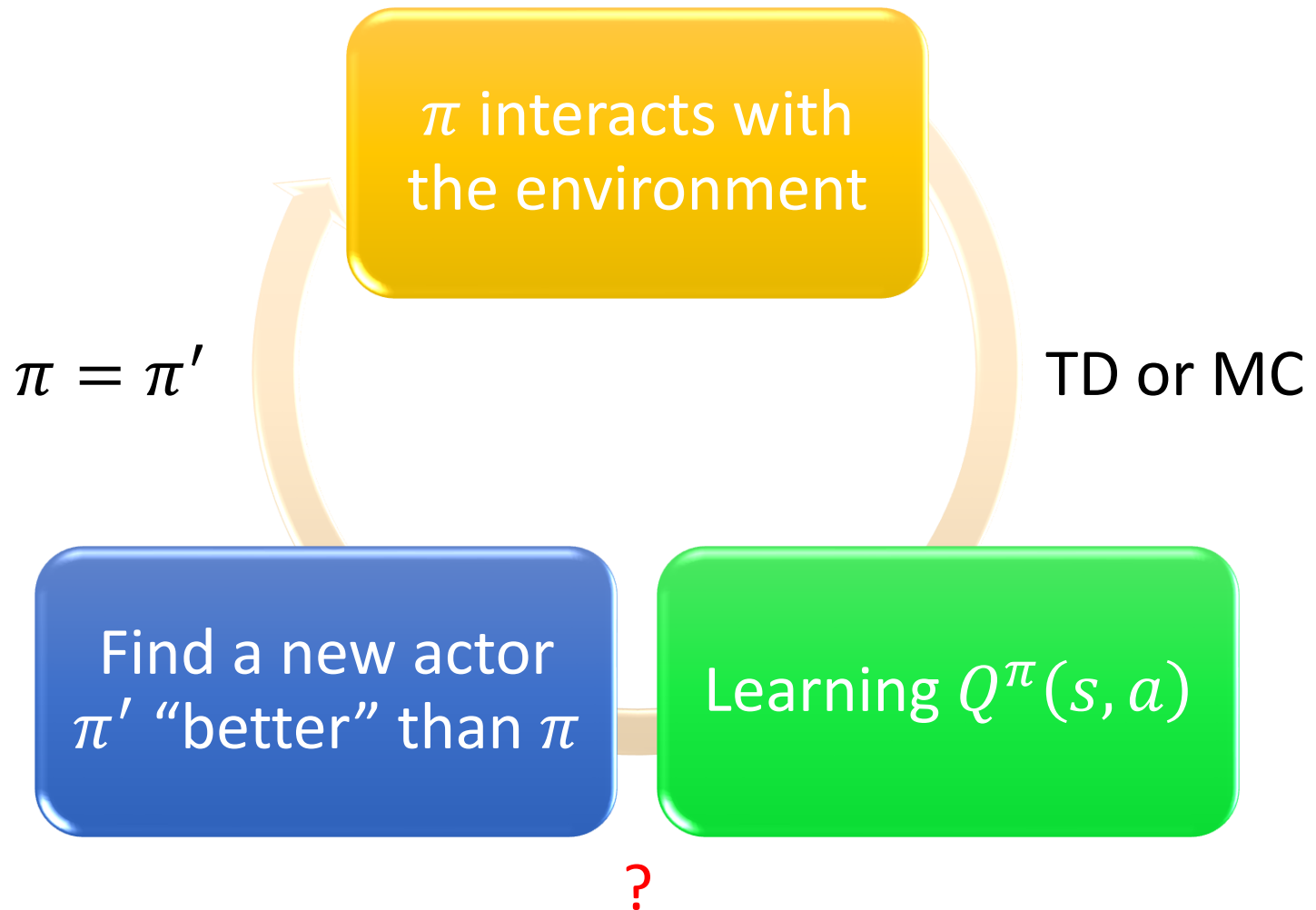
$$\text{3/4} \qquad \text{0} \qquad \text{3/4}$$
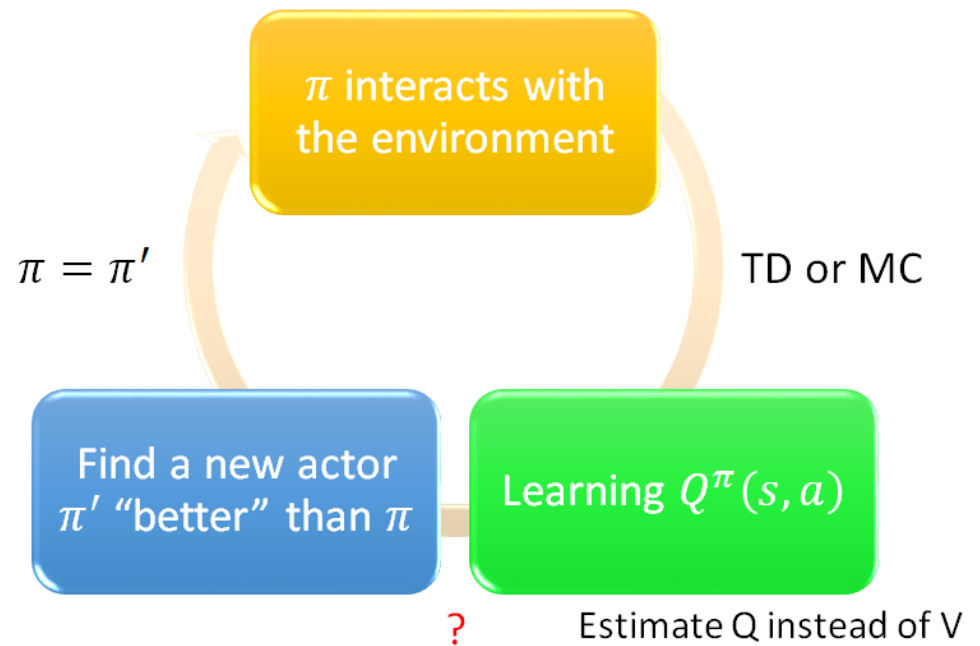
(The actions are ignored here.)

# Another Critic

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



for discrete action only

# Q-Learning



$\pi$ interacts with the environment

TD or MC

Learning $Q^\pi(s, a)$

?

Find a new actor $\pi'$ "better" than $\pi$

$\pi = \pi'$

# Q-Learning



$\pi$ interacts with the environment

$\pi = \pi'$

TD or MC

Find a new actor $\pi'$ "better" than $\pi$

Learning $Q^\pi(s,a)$

? Estimate Q instead of V

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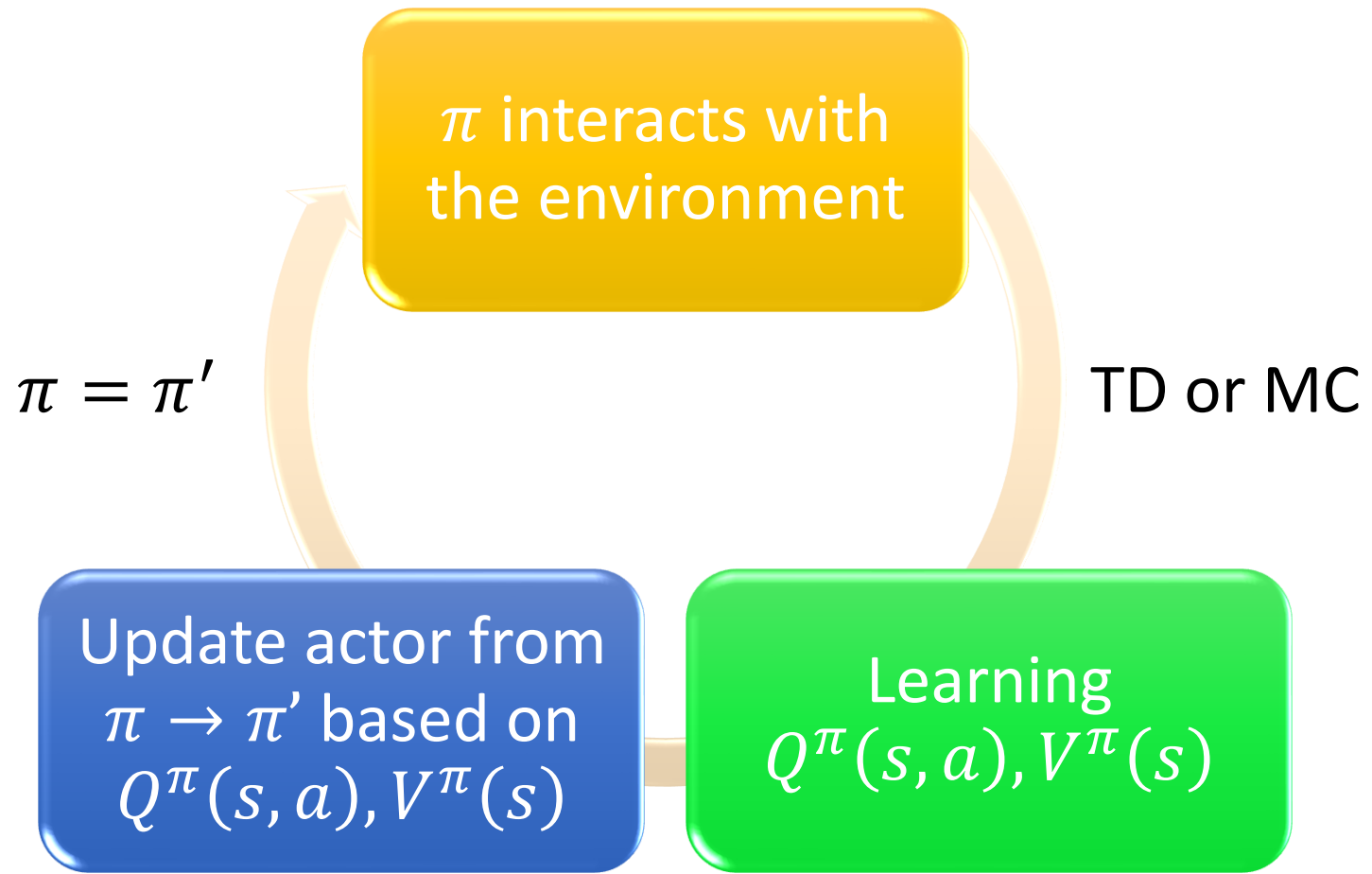  - "Better": $V^{\pi'}(s) \geq V^\pi(s)$, for all state s

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

➢ $\pi'$ does not have extra parameters. It depends on Q

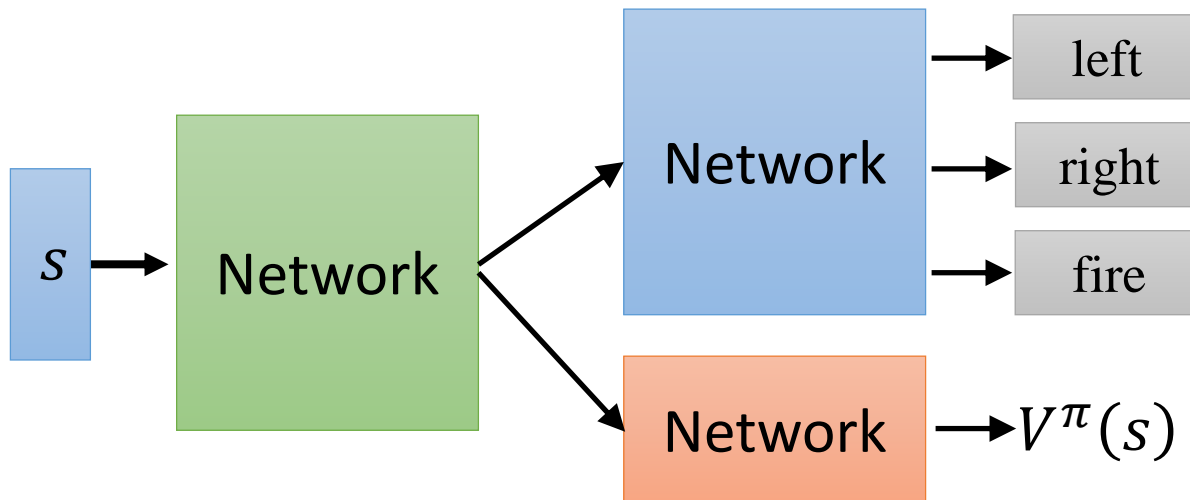➢ Not suitable for continuous action a

# Deep Reinforcement Learning

Actor-Critic

# Actor-Critic



$\pi$ interacts with the environment

TD or MC

Learning $Q^\pi(s, a), V^\pi(s)$

Update actor from $\pi \rightarrow \pi'$ based on $Q^\pi(s, a), V^\pi(s)$

$\pi = \pi'$

# Actor-Critic

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

# *Asynchronous*

1. Copy global parameters

2. Sampling some data

3. Compute gradients

4. Update global models



Global Network

Policy π(s)  V(s)

Network

Input (s)

$$\theta^1 + \eta \Delta\theta$$

$$\theta^2$$

(other workers also update models)

$\Delta\theta$

$\theta^1$

$\theta^1$

$\Delta\theta$

Worker 1   Worker 2   Worker 3   ...   Worker n

Environment 1   Environment 2   Environment 3 ...   Environment n

# Concluding Remarks

# Thanks!

Any questions?