**Basic Q-Learning**

Apr 16th, 2019

# Applied Deep Learning

YUN-NUNG (VIVIAN) CHEN   HTTP://ADL.MIULAB.TW

National Taiwan University

1

# Reinforcement Learning Approach

Value-based RL
- Estimate the optimal value function $Q^*(s, a)$

$Q^*(s, a)$ is maximum value achievable under any policy
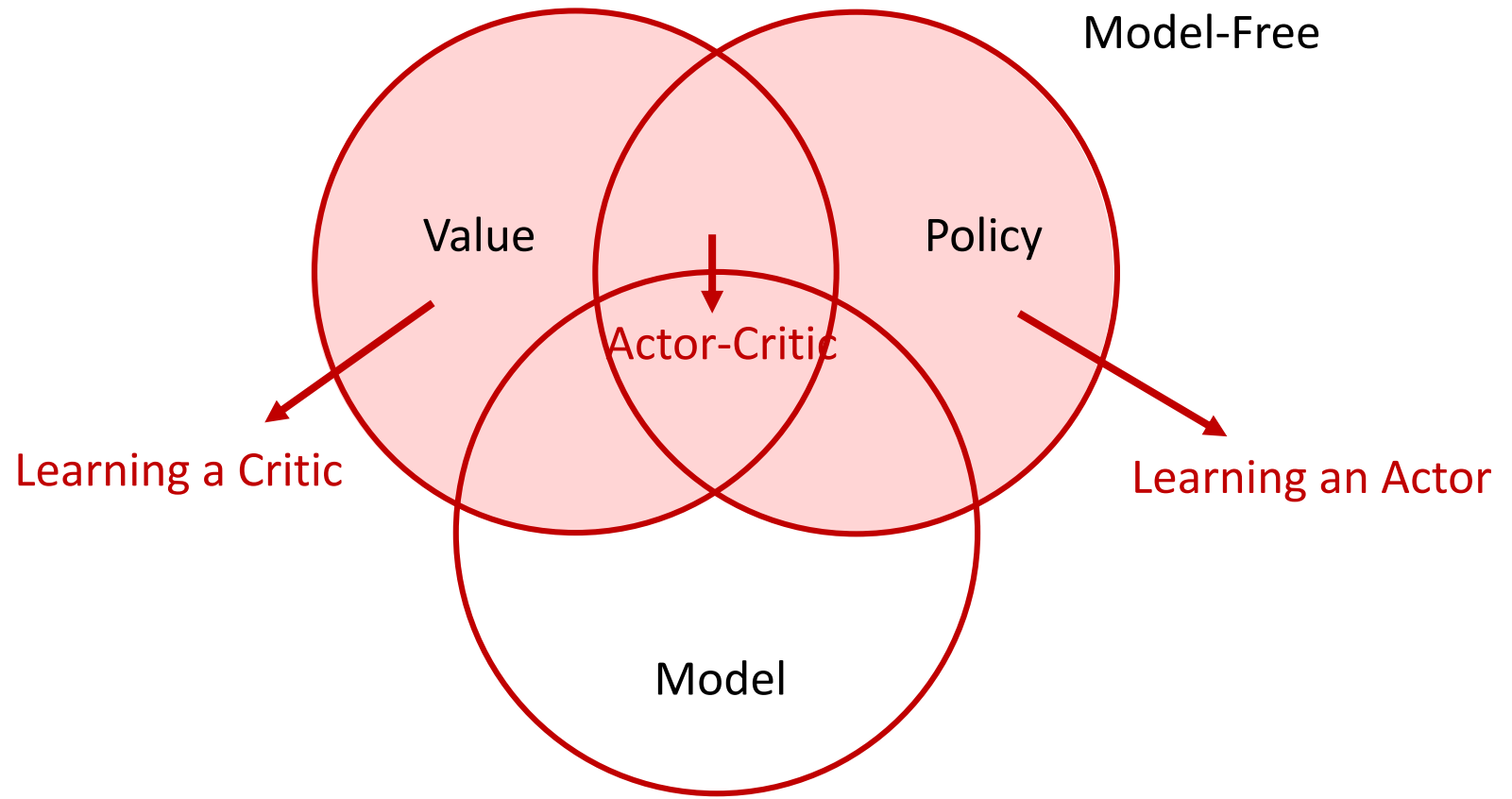
Policy-based RL
- Search directly for optimal policy $\pi^*$

$\pi^*$ is the policy achieving maximum future reward

Model-based RL
- Build a model of the environment
- Plan (e.g. by lookahead) using model

# RL Agent Taxonomy

# Value-Based Approach

LEARNING A CRITIC

# Value Function

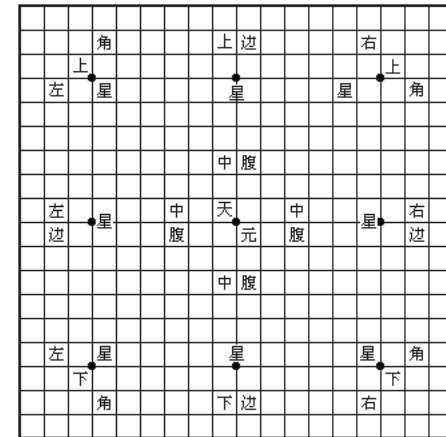A value function is a prediction of future reward (with action *a* in state *s*)

Q-value function gives expected total reward
◦ from state $s$ and action $a$
◦ under policy $\pi$
◦ with discount factor $\gamma$

$$Q^{\pi}(s, a) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

Value functions decompose into a Bellman equation

$$Q^{\pi}(s, a) = \mathbb{E}_{s', a'}[r + \gamma Q^{\pi}(s', a') \mid s, a]$$

# Optimal Value Function

An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_\pi Q^\pi(s, a) = Q^{\pi^*}(s, a)$$

The optimal value function allows us act optimally

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

The optimal value informally maximizes over all decisions

$$Q^*(s, a) = r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots$$
$$= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$$

Optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'}[r + \gamma \max_{a'} Q^*(s', a') \mid s, a]$$
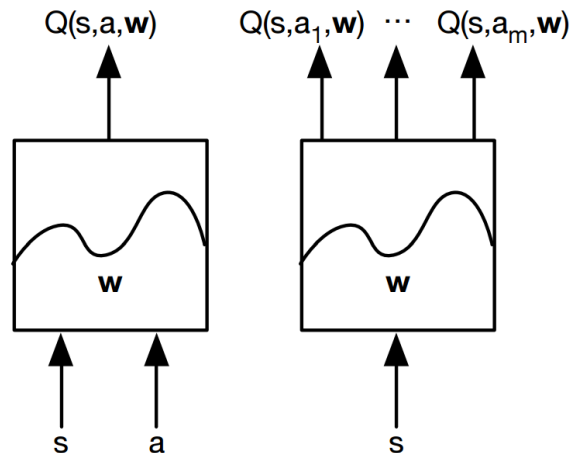
# Value Function Approximation

Value functions are represented by a *lookup table*

$$Q(s, a) \quad \forall s, a$$

◦ too many states and/or actions to store

◦ too slow to learn the value of each entry individually

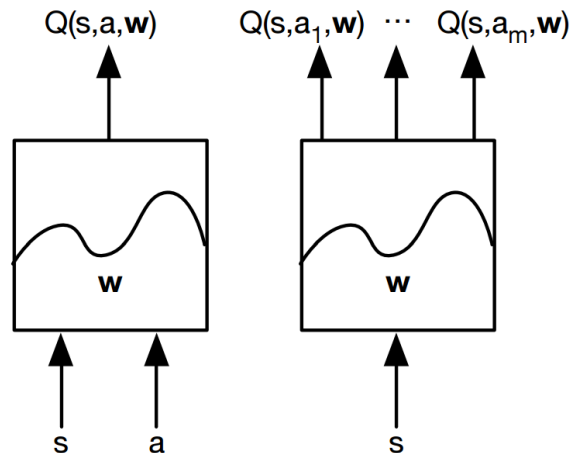Values can be estimated with *function approximation*

# Q-Networks

Q-networks represent value functions with weights $w$

$$Q(s, a, w) \approx Q^*(s, a)$$

- ◦ generalize from seen states to unseen states
- ◦ update parameter $w$ for function approximation

# Q-Learning

Goal: estimate optimal Q-values
  ◦ Optimal Q-values obey a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'}\left[\boxed{r + \gamma \max_{a'} Q^*(s', a')} \mid s, a\right]$$

learning target

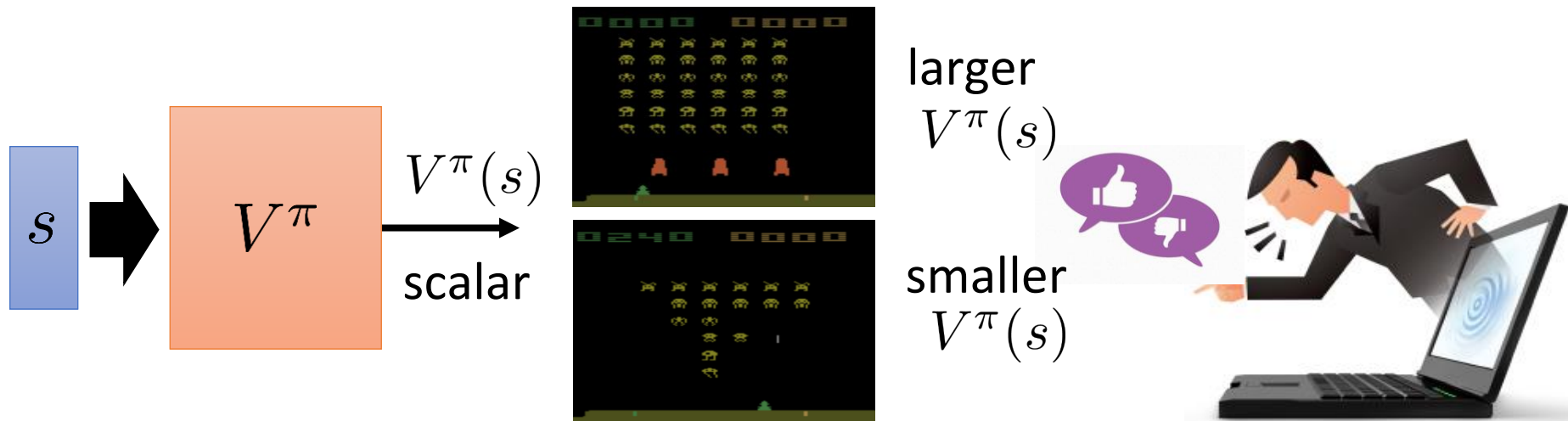  ◦ *Value iteration* algorithms solve the Bellman equation

$$Q_{i+1}(s, a) = \mathbb{E}_{s'}\left[r + \gamma \max_{a'} Q_i(s', a') \mid s, a\right]$$

# Critic = Value Function

Idea: how good the actor is

State value function: when using actor $\pi$, the *expected total reward* after seeing observation (state) s

$$V^\pi(s) \; \forall s \quad = \mathbb{E}[G_t \mid s_t = s]$$



larger $V^\pi(s)$

smaller $V^\pi(s)$

$s$ → $V^\pi$ → $V^\pi(s)$

scalar

A critic does not determine the action
An actor can be found from a critic
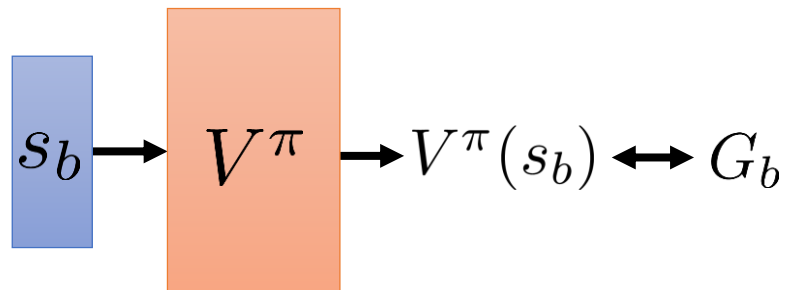
## Monte-Carlo (MC)

◦ The critic watches $\pi$ playing the game

◦ MC learns directly from *complete* episodes: no bootstrapping

Idea: value = *empirical mean* return

After seeing $s_a$,
until the end of the episode,
the cumulated reward is $G_a$

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

After seeing $s_b$,
until the end of the episode,
the cumulated reward is $G_b$

$$s_b \rightarrow V^\pi \rightarrow V^\pi(s_b) \leftrightarrow G_b$$
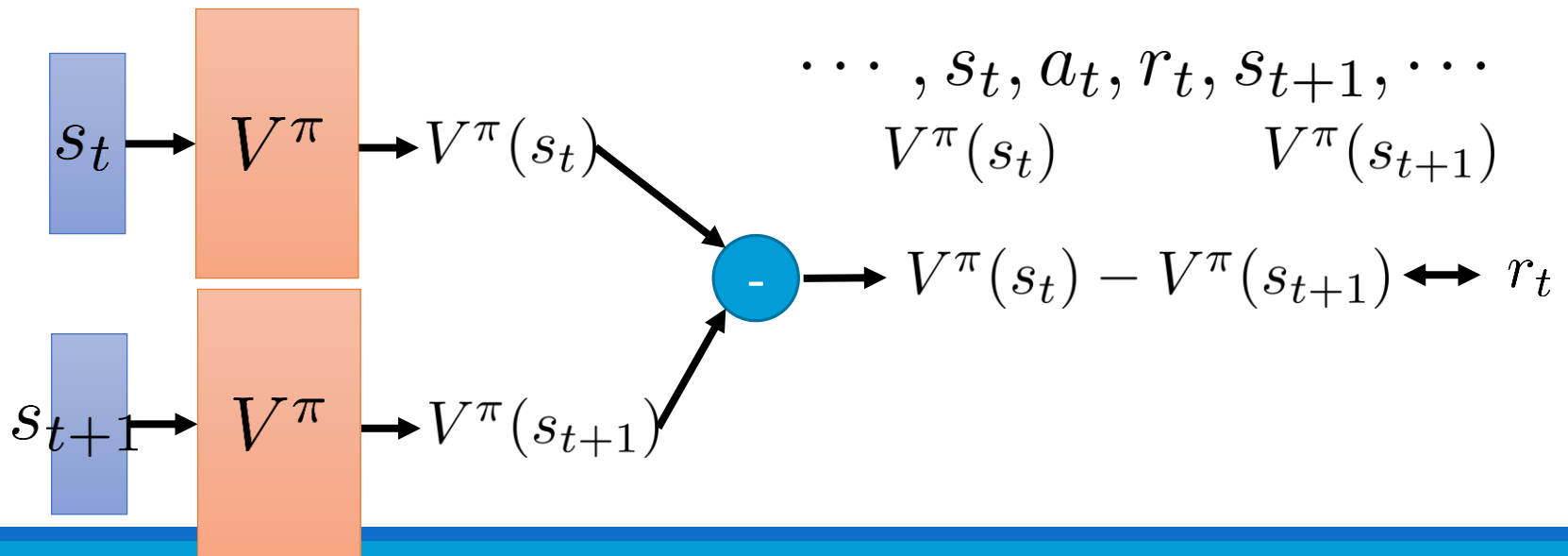
Issue: long episodes delay learning

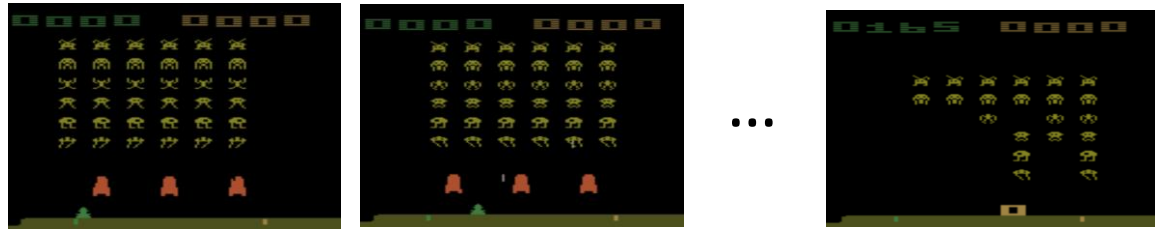# Temporal-Difference for Estimating $V^\pi(s)$

## Temporal-difference (TD)

◦ The critic watches $\pi$ playing the game

◦ TD learns directly from *incomplete* episodes by *bootstrapping*

◦ TD updates a guess towards a guess

Idea: update value toward *estimated* return

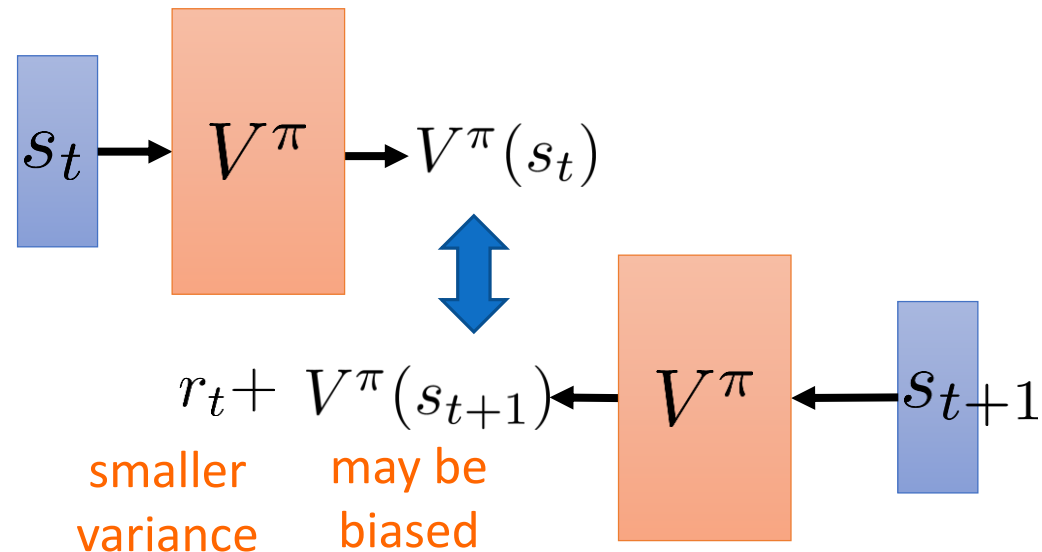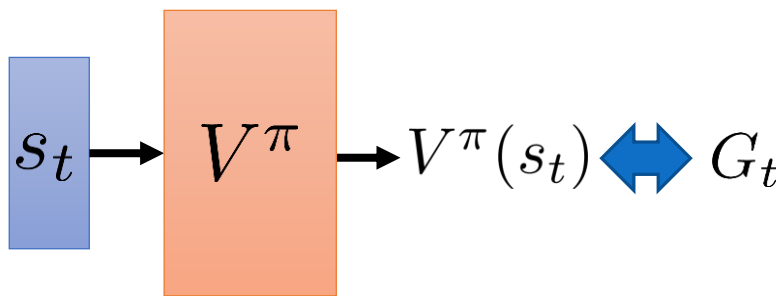$$\cdots, s_t, a_t, r_t, s_{t+1}, \cdots$$

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

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

$$s_{t+1} \rightarrow V^\pi \rightarrow V^\pi(s_{t+1})$$

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

# MC v.s. TD

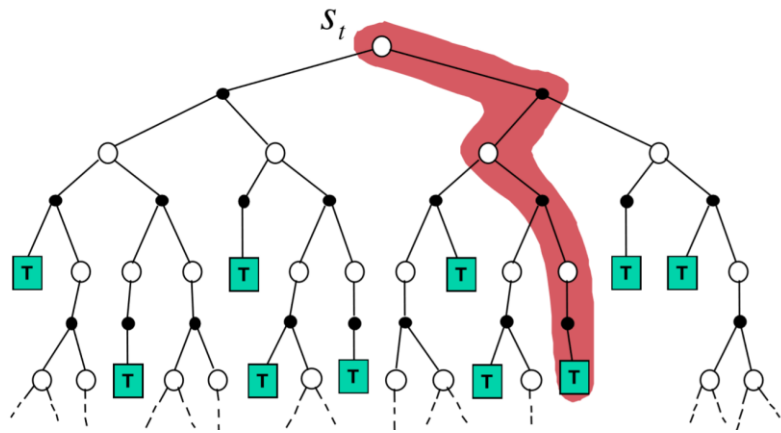## Monte-Carlo (MC)
◦ Large variance

◦ Unbiased

◦ No Markov property

## Temporal-Difference (TD)
◦ Small variance

◦ Biased

◦ Markov property

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

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

$$r_t + V^\pi(s_{t+1}) \leftarrow V^\pi \leftarrow s_{t+1}$$

smaller variance    may be biased

$$V'^{\pi}(s_t)$$

$$= V^{\pi}(s_t) + \alpha(G_t - V^{\pi}(s_t))$$
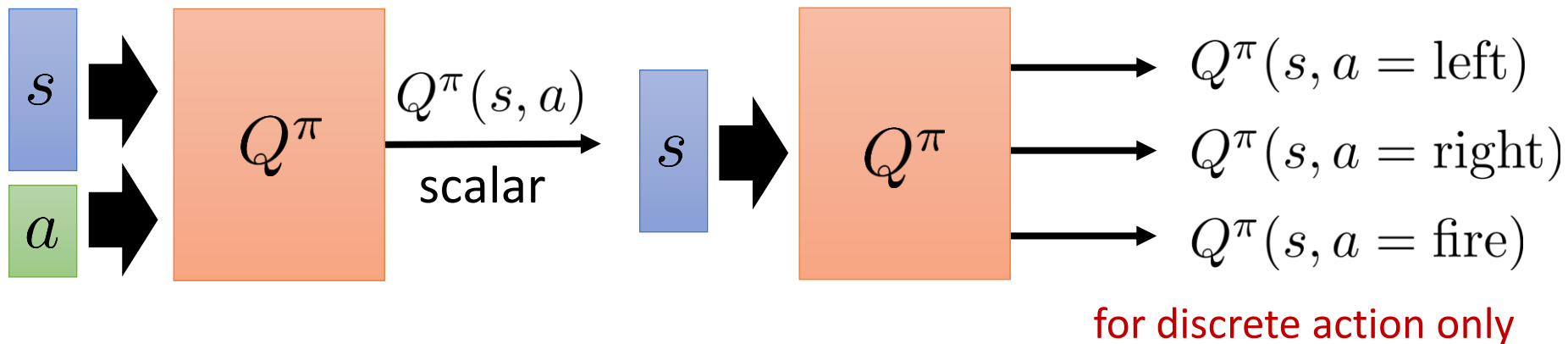
$$V'^{\pi}(s_t)$$

$$= V^{\pi}(s_t) + \alpha(r_{t+1} + \gamma V^{\pi}(s_{t+1}) - V^{\pi}(s_t))$$

# MC v.s. TD

# Critic = Value Function

State-action value function: when using actor $\pi$, the *expected total reward* after seeing observation (state) $s$ and taking action $a$

$$Q^\pi(s, a) \; \forall s, a = \mathbb{E}[G_t \mid s_t = s, a_t = a]$$



for discrete action only

# Q-Learning

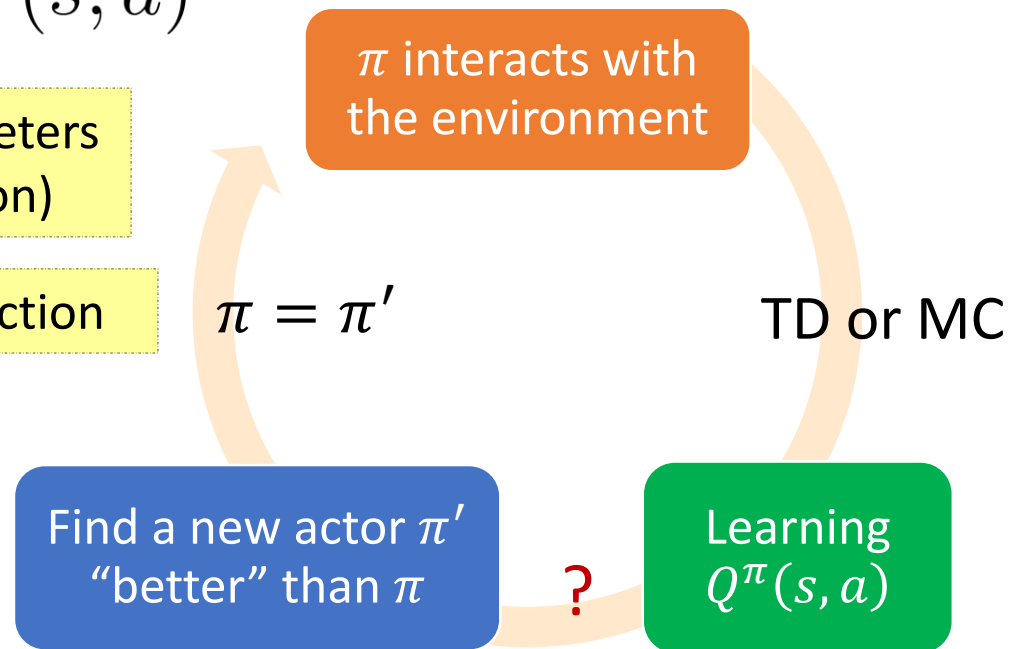Given $Q^{\pi}(s, a)$, find a new actor $\pi'$ "better" than $\pi$

$$V^{\pi'}(s) \geq V^{\pi}(s) \ \forall s$$

$$\pi'(s) = \arg \max_{a} Q^{\pi}(s, a)$$

$\pi$ interacts with
the environment

$\pi'$ does not have extra parameters
(depending on value function)

not suitable for continuous action

$\pi = \pi'$

TD or MC

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

**?**

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

# Q-Learning

Goal: estimate optimal Q-values
- Optimal Q-values obey a Bellman equation

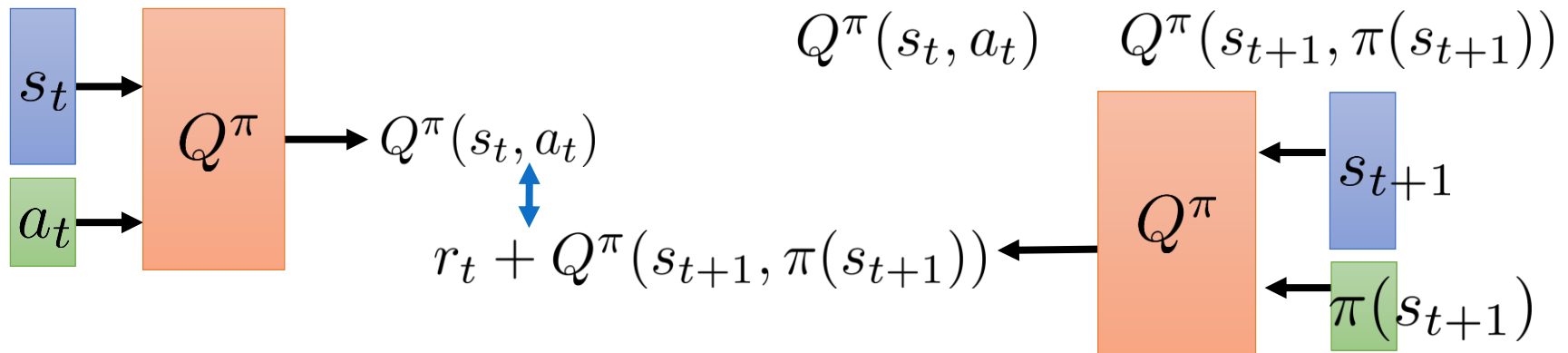$$Q^*(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

learning target

- *Value iteration* algorithms solve the Bellman equation

$$Q_{i+1}(s, a) = \mathbb{E}_{s'} \left[ r + \gamma \max_{a'} Q_i(s', a') \mid s, a \right]$$

# Deep Q-Networks (DQN)

Estimate value function by TD $\quad \cdots, s_t, a_t, r_t, s_{t+1}, \cdots$

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



$$Q^\pi(s_t, a_t)$$

$$r_t + Q^\pi(s_{t+1}, \pi(s_{t+1}))$$

Represent value function by deep Q-network with weights $w$

$$Q(s, a, w) \approx Q^*(s, a)$$

Objective is to minimize MSE loss by SGD

$$\mathcal{L}(w) = \mathbb{E}\left[ \left( r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right)^2 \right]$$

# Deep Q-Networks (DQN)

Objective is to minimize MSE loss by SGD

$$\mathcal{L}(w) = \mathbb{E}\left[ \left( r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right)^2 \right]$$

Leading to the following Q-learning gradient

$$\frac{\partial \mathcal{L}(w)}{\partial w} = \mathbb{E}\left[ \left( r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w) \right) \frac{\partial Q(s, a, w)}{\partial w} \right]$$

Issue: naïve Q-learning oscillates or diverges using NN due to:
1) correlations between samples 2) non-stationary targets

# Stability Issues with Deep RL

Naive Q-learning <span style="color:red">oscillates</span> or <span style="color:red">diverges</span> with neural nets

1. Data is sequential
   ◦ Successive samples are correlated, non-iid (independent and identically distributed)
2. Policy changes rapidly with slight changes to Q-values
   ◦ Policy may oscillate
   ◦ Distribution of data can swing from one extreme to another
3. Scale of rewards and Q-values is unknown
   ◦ Naive Q-learning gradients can be unstable when backpropagated
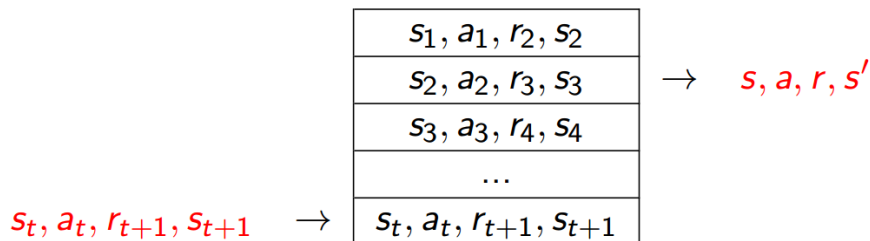
# Stable Solutions for DQN

DQN provides a stable solutions to deep value-based RL

1. Use experience replay
   ◦ Break correlations in data, bring us back to iid setting
   ◦ Learn from all past policies
2. Freeze target Q-network
   ◦ Avoid oscillation
   ◦ Break correlations between Q-network and target
3. Clip rewards or normalize network adaptively to sensible range
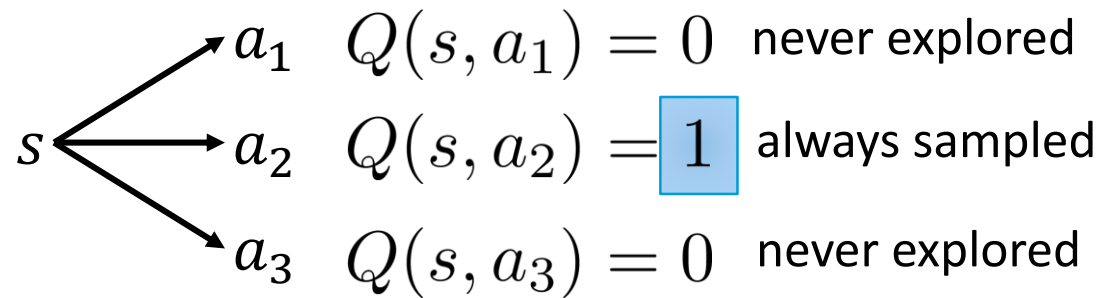   ◦ Robust gradients

# Stable Solution 1: Experience Replay

To remove correlations, build a dataset from agent's experience
- Take action at according to $\epsilon$-greedy policy    small prob for exploration
- Store transition $\left(s_t, a_t, r_t, s_{t+1}\right)$ in replay memory $D$
- Sample random mini-batch of transitions $\left(s, a, r, s'\right)$ from $D$

| $s_1, a_1, r_2, s_2$ |
| :---: |
| $s_2, a_2, r_3, s_3$ |
| $s_3, a_3, r_4, s_4$ |
| ... |
| $s_t, a_t, r_{t+1}, s_{t+1}$ |

$\rightarrow \quad s, a, r, s'$

$s_t, a_t, r_{t+1}, s_{t+1} \quad \rightarrow$

- Optimize MSE between Q-network and Q-learning targets

$$\mathcal{L}(w) = \mathbb{E}_{s,a,r,s' \sim D}\left[\left(r + \gamma \max_{a'} Q(s', a', w) - Q(s, a, w)\right)^2\right]$$

$a_1$  $Q(s, a_1) = 0$  never explored

$s$ $a_2$  $Q(s, a_2) = \boxed{1}$  always sampled

$a_3$  $Q(s, a_3) = 0$  never explored

# Exploration

The policy is based on Q-function

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

not good for data collection
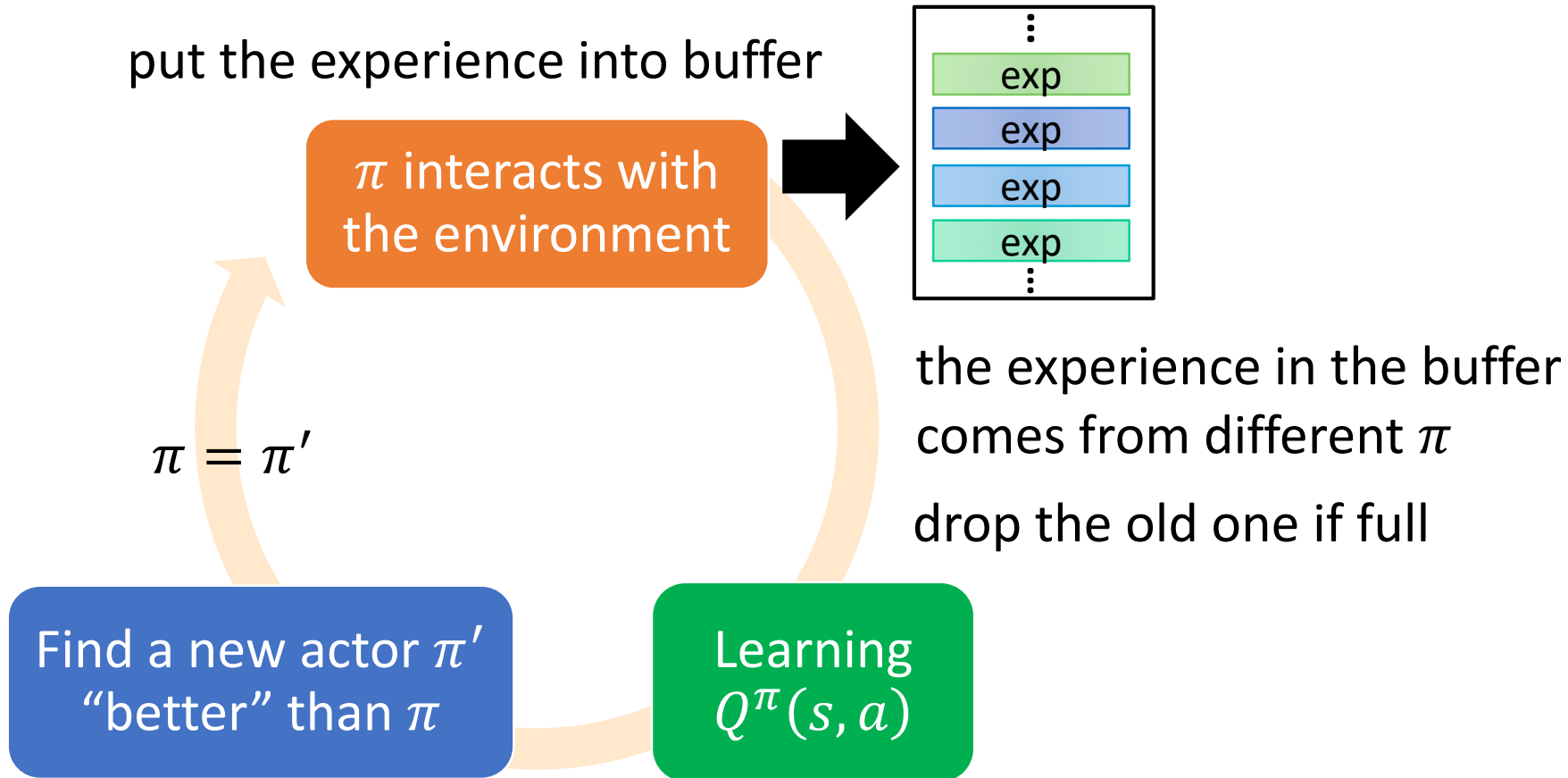→ inefficient learning

Exploration algorithms

ε would decay during learning

◦ Epsilon greedy

$$a = \begin{cases} \arg\max_a Q(s, a), & \text{with } p = (1 - \epsilon) \\ \text{random}, & \text{otherwise} \end{cases}$$
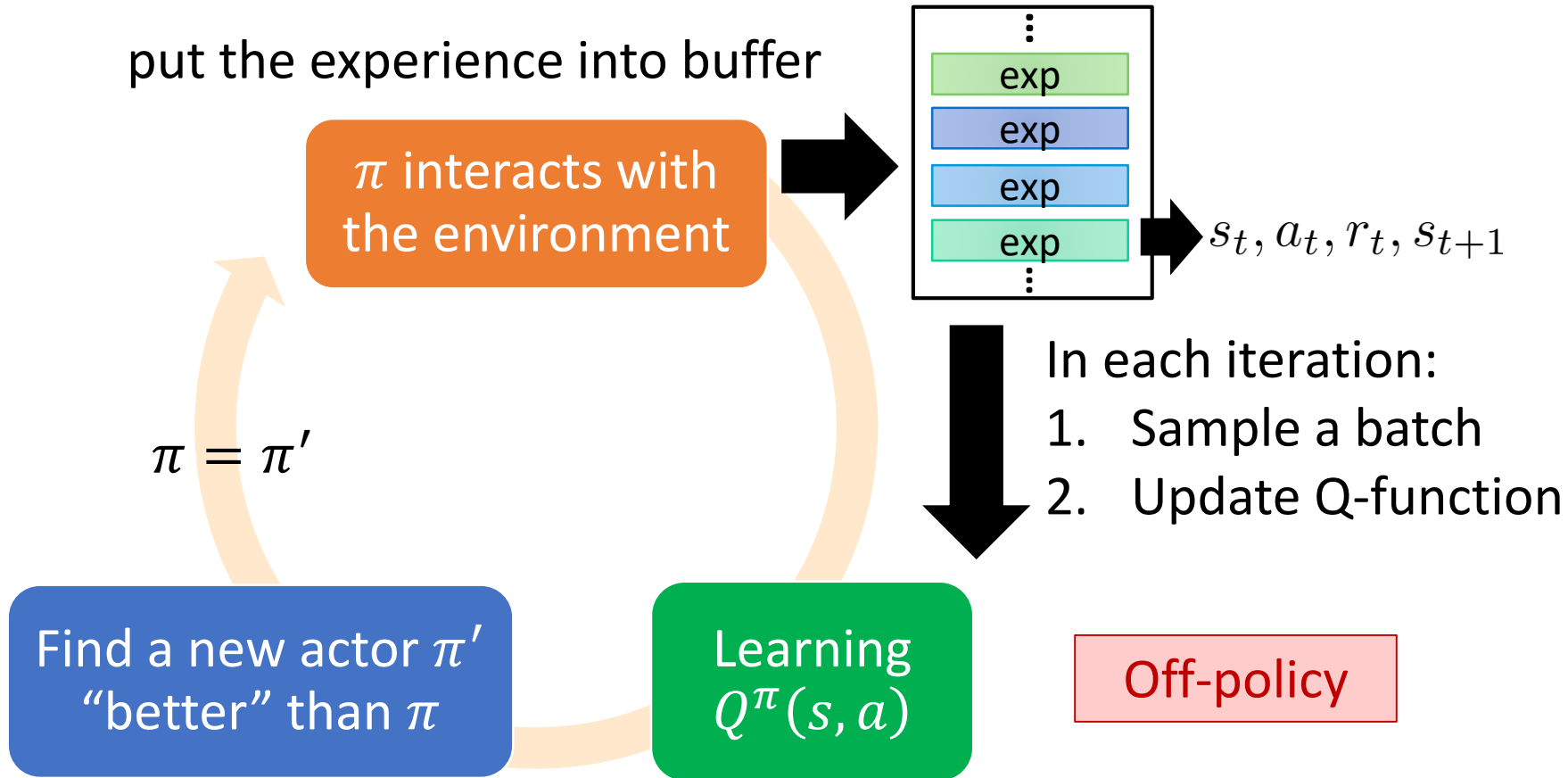
◦ Boltzmann sampling

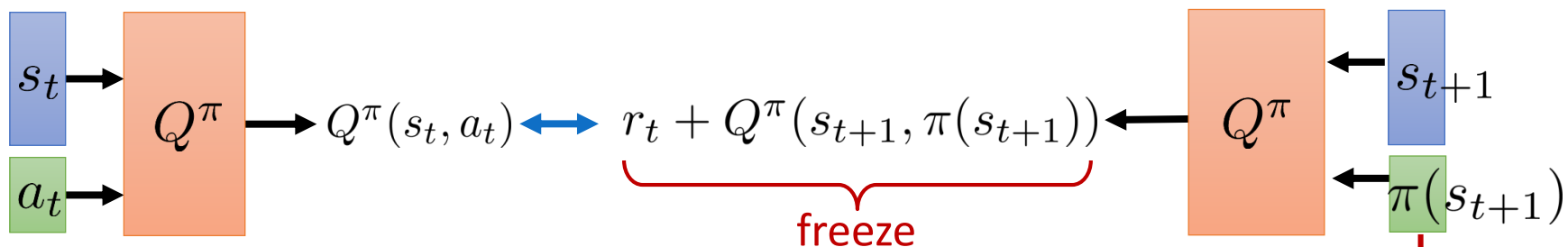$$P(a \mid s) = \frac{\exp(Q(s, a))}{\sum_a \exp(Q(s, a))}$$

# Replay Buffer

put the experience into buffer

$\pi$ interacts with the environment

the experience in the buffer comes from different $\pi$

drop the old one if full

exp
exp
exp
exp

$\pi = \pi'$

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

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

# Replay Buffer

put the experience into buffer

$\pi$ interacts with the environment

exp
exp
exp
exp

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

In each iteration:
1. Sample a batch
2. Update Q-function

$\pi = \pi'$

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

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

Off-policy

# Stable Solution 2: Fixed Target Q-Network

To avoid oscillations, fix parameters used in Q-learning target



- Compute Q-learning targets w.r.t. old, fixed parameters $w^-$

$$r + \gamma \max_{a'} \hat{Q}(s', a', w^-)$$

- Optimize MSE between Q-network and Q-learning targets

$$\mathcal{L}(w) = \mathbb{E}_{s,a,r,s' \sim D}\left[ \left( r + \gamma \max_{a'} \hat{Q}(s', a', w^-) - Q(s, a, w) \right)^2 \right]$$

- Periodically update fixed parameters $w^- \leftarrow w$

# Stable Solution 3: Reward / Value Range

To avoid oscillations, control the reward / value range
- DQN clips the rewards to [−1, +1]
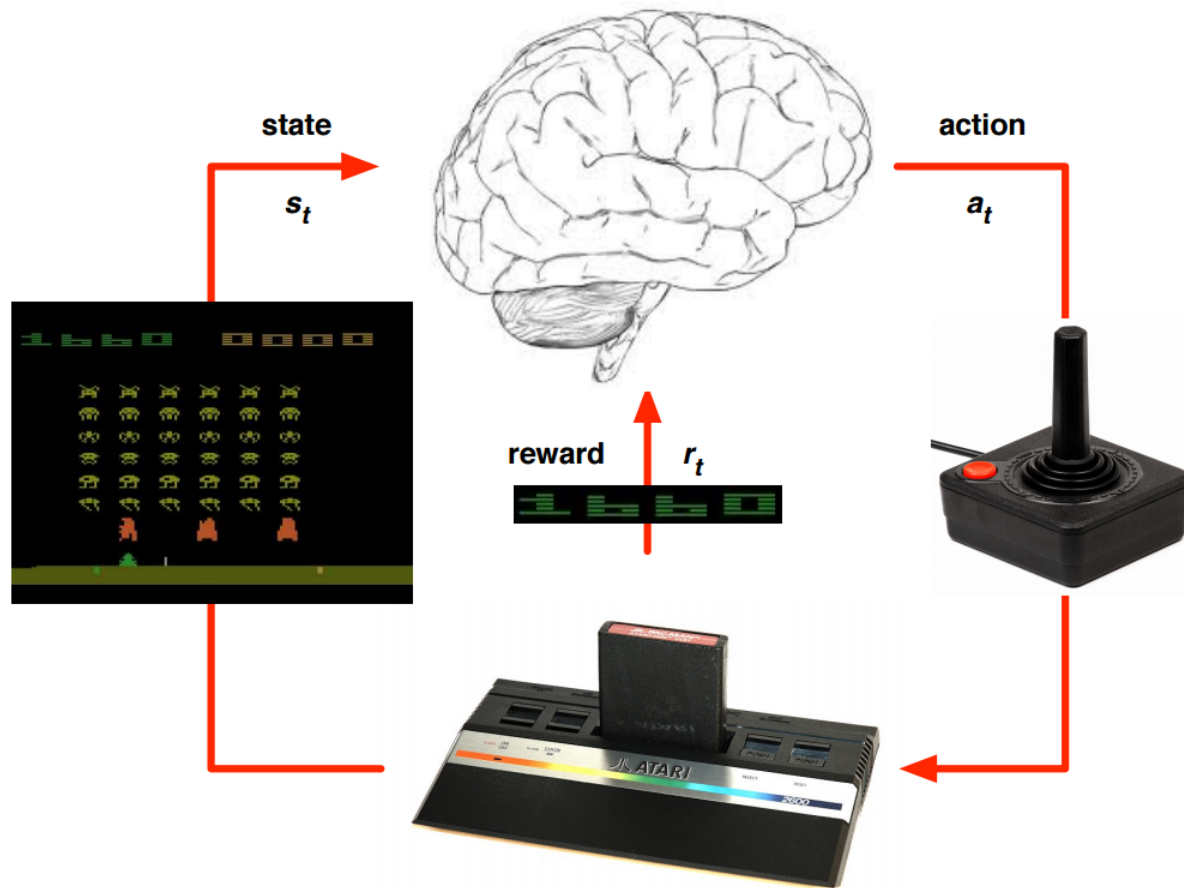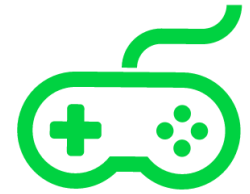  - Prevents too large Q-values
  - Ensures gradients are well-conditioned
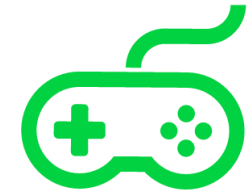
# Typical Q-Learning Algorithm

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 $Q$ (epsilon greedy)
  - 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)
  - Update the parameters of $Q$ to make $Q(s_i, a_i) \approx r_i + \max_a \hat{Q}(s_{i+1}, a)$
  - Every $C$ steps reset $\hat{Q} = Q$

# Deep RL in Atari Games
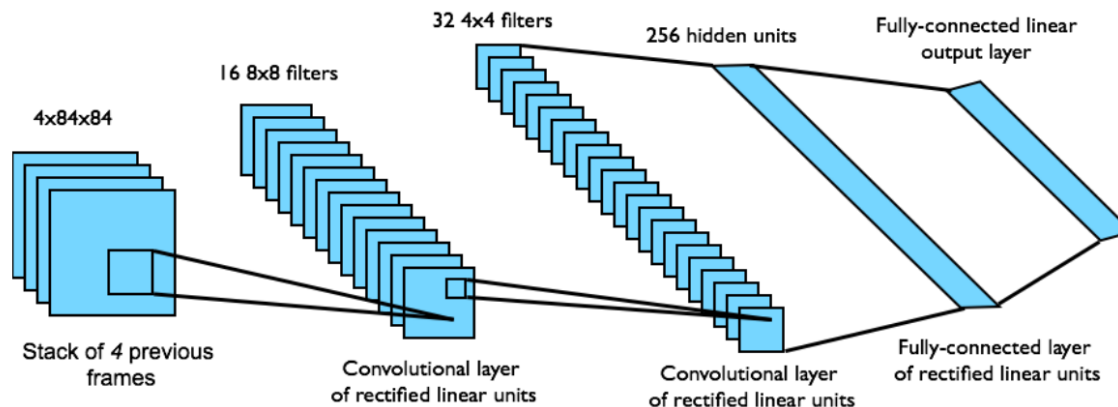
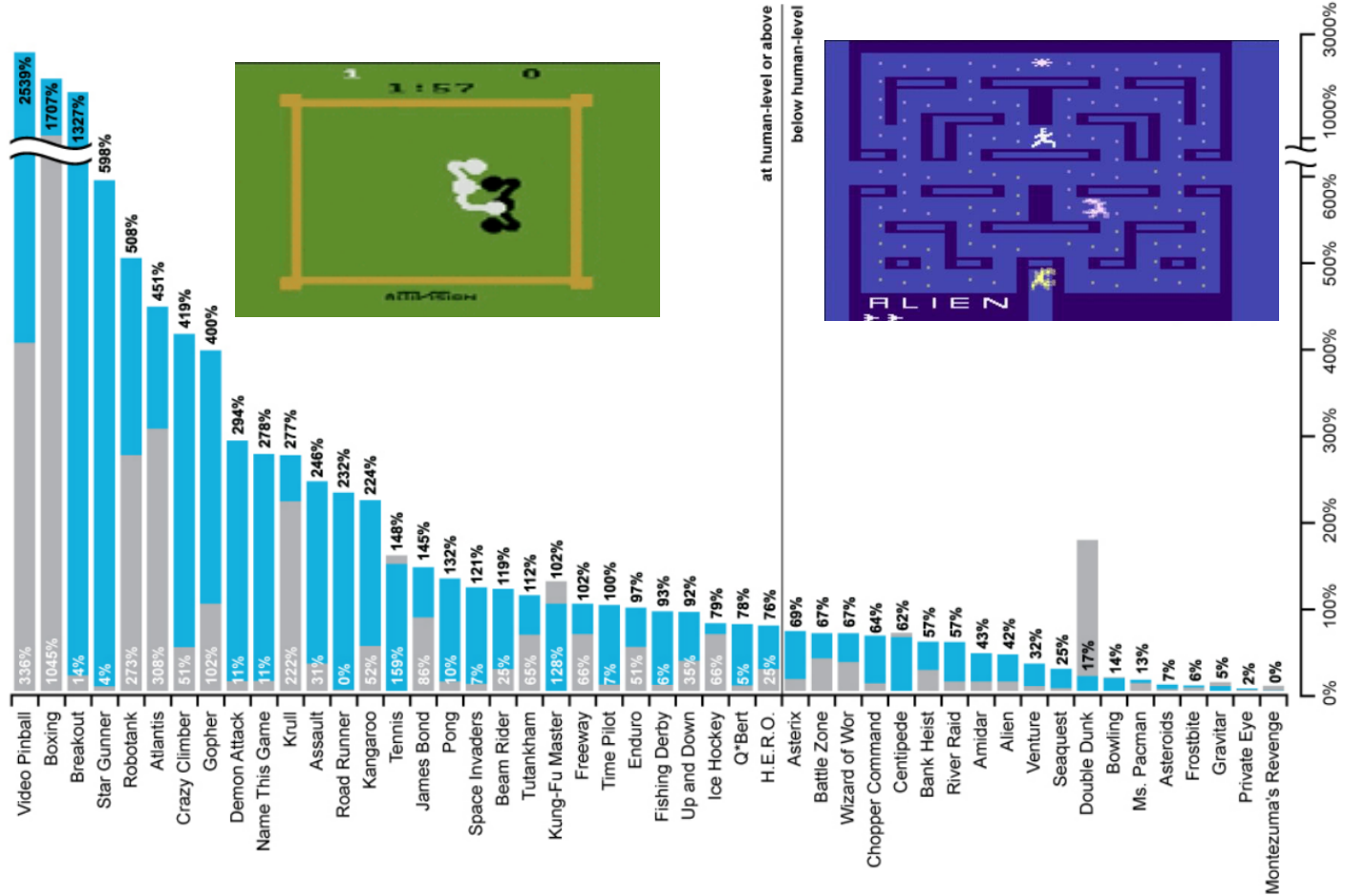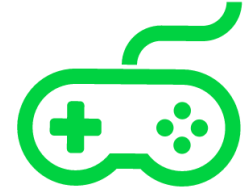

state $s_t$

action $a_t$

reward $r_t$

# DQN in Atari

Goal: end-to-end learning of values $Q(s, a)$ from pixels

$$\mathcal{L}(w) = \mathbb{E}_{s,a,r,s' \sim D}\left[\left(r + \gamma \max_{a'} Q(s', a', w^-) - Q(s, a, w)\right)^2\right]$$

◦ Input: state is stack of raw pixels from last 4 frames
◦ Output: $Q(s, a)$ for all joystick/button positions $a$
◦ Reward is the score change for that step

# DQN in Atari

# Concluding Remarks

RL is a general purpose framework for **decision making** under interactions between agent and environment

A value-based RL measures how good each state and/or action is via a value function
◦ Monte-Carlo (MC) v.s. Temporal-Difference (TD)