

DQN

Deep Q-Networks

Angela Carnevale

MathLearn - Galway, June 2026

## Q. Brief Recap

Q-table :  $Q(s, a)$   
state-action function

The diagram shows the notation  $Q(s, a)$  for a Q-table. An arrow points from the word 'state' above to the variable 's'. Another arrow points from the word 'action' above to the variable 'a'. A horizontal bracket is drawn under the entire expression  $Q(s, a)$ , with the text 'state-action function' written below it. A long arrow points from the text 'state-action function' back to the 'Q' in the notation.

governed by Bellman equation (together with  $V(s)$ )

use TD-target & TD-control to update Q-values

That is: Q-Learning learns a table

# 0. Brief Recap

Q-table :  $Q(s, a)$  update according to

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

learning rate

discount factor

TD-target

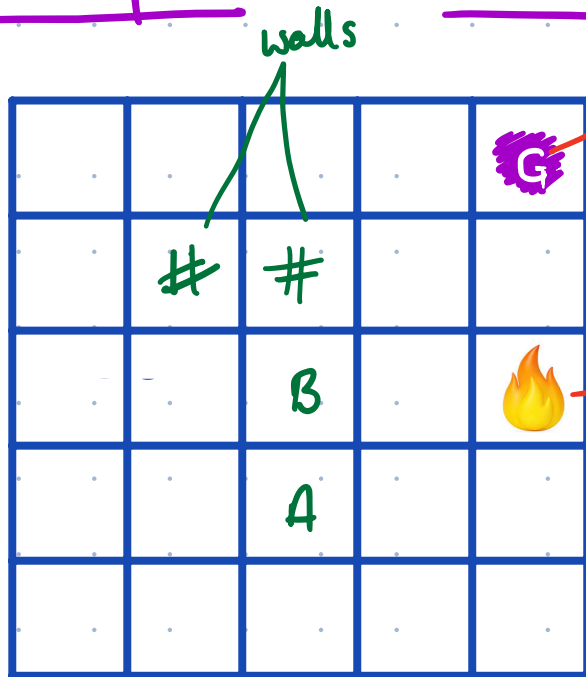
TD-error

So each update improves one number in our

$|\hat{S}| \times |\hat{A}|$  table...

# Example

# Gridworld



Actions:  $\{\uparrow, \downarrow, \leftarrow, \rightarrow\}$

Reward -1 for each move

So our Q table is a

$(5.5) \times 4 = 25 \times 4$  table ✓

Example

Matrix World

$$\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix} \in \text{Mat}_{3 \times 3}(\mathbb{F}_2)$$

Goal: turn into upper triangular form using

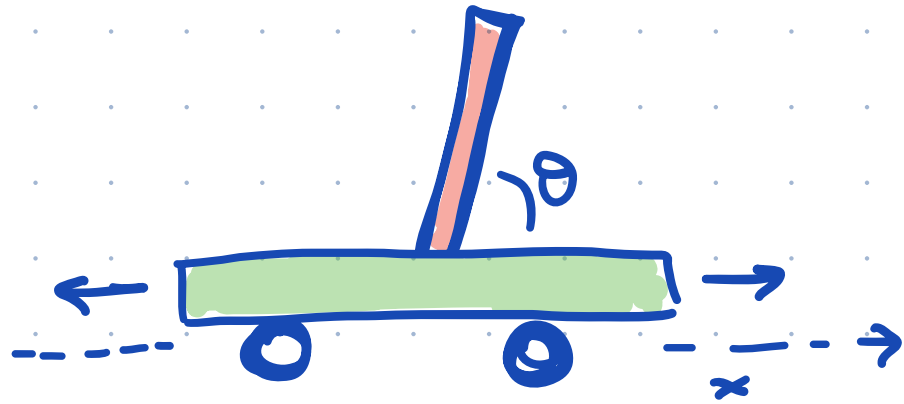
Actions := { row operations }

$$|\mathcal{Q}\text{-tabel}| = \underbrace{2^3}_{|S|} \cdot \underbrace{\binom{3}{2} \cdot 2 \cdot \binom{3}{2}}_{|\hat{A}|}$$

becomes huge very quickly!

Example

CartPole



State:  $(x, \dot{x}, \theta, \dot{\theta})$

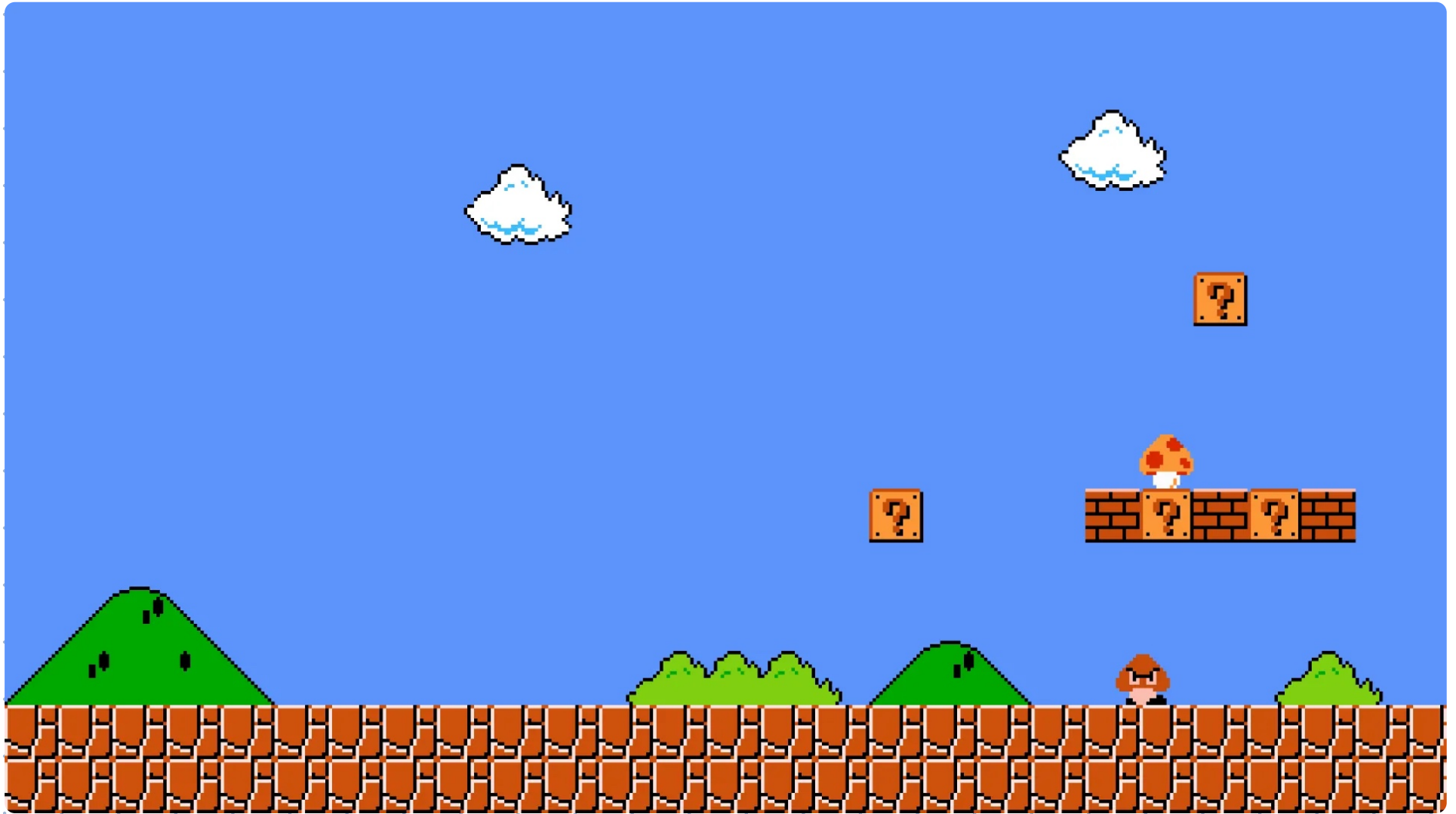
Actions:  $\{\leftarrow, \rightarrow\}$

Only 2 columns but continuous states

# Example Mario!



# Example Maui!





We need to:

- be able to deal with large size
- make more efficient updates
- be able to generalise across states

Idea: bring in a neural network ...  
("Learn the table's pattern")

Q-table  $\rightarrow$   $Q_{\theta}(s,a)$   $\rightarrow$  DQN

... and then make it work

Towards  $Q_{\theta}(s,a)$

# 1. Features and weights

We can view our table lookup for the entry  $Q(s,a)$  in

| $Q(s,a)$ | $a_1$ | $a_2$ | ... |
|----------|-------|-------|-----|
| $s_1$    |       |       |     |
| $s_2$    |       |       |     |
| $\vdots$ |       |       |     |

as:

$$Q(s,a) = \underbrace{w^T}_{\text{weights}} \underbrace{\phi(s,a)}_{\text{feature vector } \phi(s,a)}$$

$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ \vdots \end{bmatrix}$

but we can generalise the "state feature" to a set of features describing the state

# 1. Features and weights

Example GridWorld: from  $(s, a) = ((x, y), a)$  to

$$\phi(s, a) = \begin{bmatrix} \text{distance to goal} \\ \text{is near wall} \\ \text{is high reward nearby} \\ \text{is near lava} \text{ 😱} \\ \text{action is down} \\ \vdots \end{bmatrix}$$

and use this to approximate

$$Q_w(s, a) = w^T \phi(s, a)$$

# 1. Features and weights

From here:

→ can use same Q-Learning target,

TD error etc

→ update weights

↗  
already improved our update: now it's  
not a single entry but all weights.

However, hand-designed features are not yet enough  
to deal with more complex states.

## 2. Introducing $Q_{\theta}(s,a)$

NN parameters

Introducing a neural network  $Q_{\theta}(s,a)$

we move from hand-designed features to learned features

For discrete actions:

$$Q_{\theta}(s, \cdot) = \begin{bmatrix} Q_{\theta}(s, a_1) \\ \vdots \\ Q_{\theta}(s, a_k) \end{bmatrix}$$

with action choice  
 $a = \underset{a'}{\operatorname{argmax}} Q(s, a')$

## 2. Introducing $Q_{\theta}(s,a)$

First attempt: we want to predict  $Q_{\theta}(s,a)$

$$\text{Target: } y = r + \gamma \max_{a'} Q_{\theta}(s',a')$$

$$\text{Loss: } L(\theta) = (y - Q_{\theta}(s,a))^2$$

Looks like supervised learning! But it is not, and this  
this creates problems.

# Problems & fixes:

## DQN

- highly correlated data
- chasing moving targets
- ( • large function approximation)

"deadly triad"

### 3. Data correlation: the problem

In RL, we get highly correlated, serial data:

$$(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$$

and then

$$(S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, A_{t+2})$$

and so on.

Example Walking on a stretch of very similar images in a Mario game, or along a corridor in GridWorld.

⚠ NNs become unstable when training on highly correlated data (standard assumption is iid)

### 3. Data correlation: The Fix

Replay Memory

Store transitions in a Replay Buffer

$$D = \{(s, a, r, s', done)\}$$

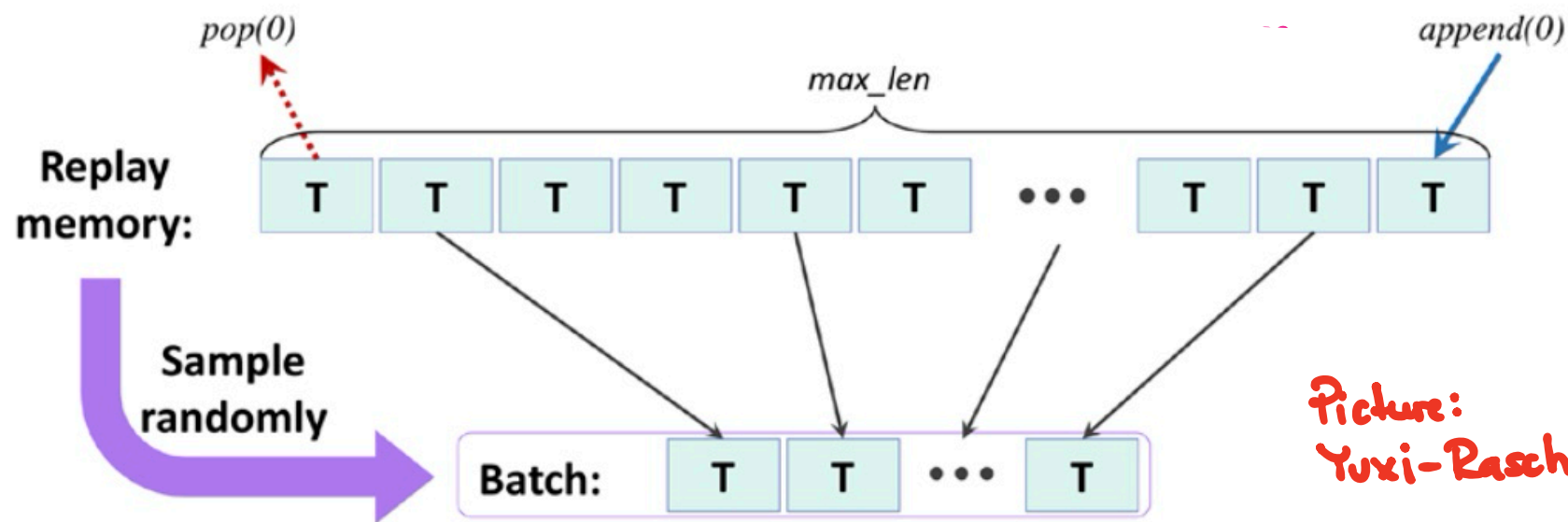
Then sample random mini-batches and train on those

Good!   
 → breaks correlation   
 → reuses experience

RL env → "mini-batch supervised learning".

### 3. Data correlation: The Fix

## Replay Memory (Flashcards!)



Picture:  
Yuxi-Raschka

Good! → breaks correlation  
→ reuses experience

RL env → "mini-batch supervised learning".

## 4. Moving targets: the problem

Assuming we introduced a NN  $Q_{\theta}(s,a)$ , we have

prediction:  $Q_{\theta}(s,a)$

target:  $y = r + \gamma \max_{a'} Q_{\theta}(s',a')$

That is: the agent is chasing  
a moving target



## 4. Moving targets: the fix

DQN fixes this issue by introducing an auxiliary network. So it uses

online network  $Q_{\theta}$   $\leftarrow$  trained every update

target network  $Q_{\theta^-}$   $\leftarrow$  used to build targets

$$y = r + \gamma \max_{a'} Q_{\theta^-}(s, a')$$

Every  $C$  steps, we copy the online weights into target:

$$\theta^- \leftarrow \theta$$

## 4. Moving targets: the fix

Recall this means that DQN is off-policy,

Behaviour policy  $\neq$  target policy

we also act according to an  $\epsilon$ -greedy policy  
(note this works nicely with experience replay!)

# 5. Putting the pieces together

## Full DQN algorithm

1. Initialise online network  $Q_{\theta}$
2. Initialise target network  $Q_{\theta^-} \leftarrow Q_{\theta}$
3. Initialise replay buffer  $\mathcal{D}$  decay  $\epsilon$  over time
4. Step:
  - choose action ( $\epsilon$ -greedy)
  - observe  $r, s', done$
  - store  $(s, a, r, s', done)$  in  $\mathcal{D}$
  - sample mini-batch from  $\mathcal{D}$
  - compute target

$$y = \begin{cases} r & \text{if done} \\ r + \gamma \max_{a'} Q_{\theta^-}(s', a') & \text{o/w} \end{cases}$$

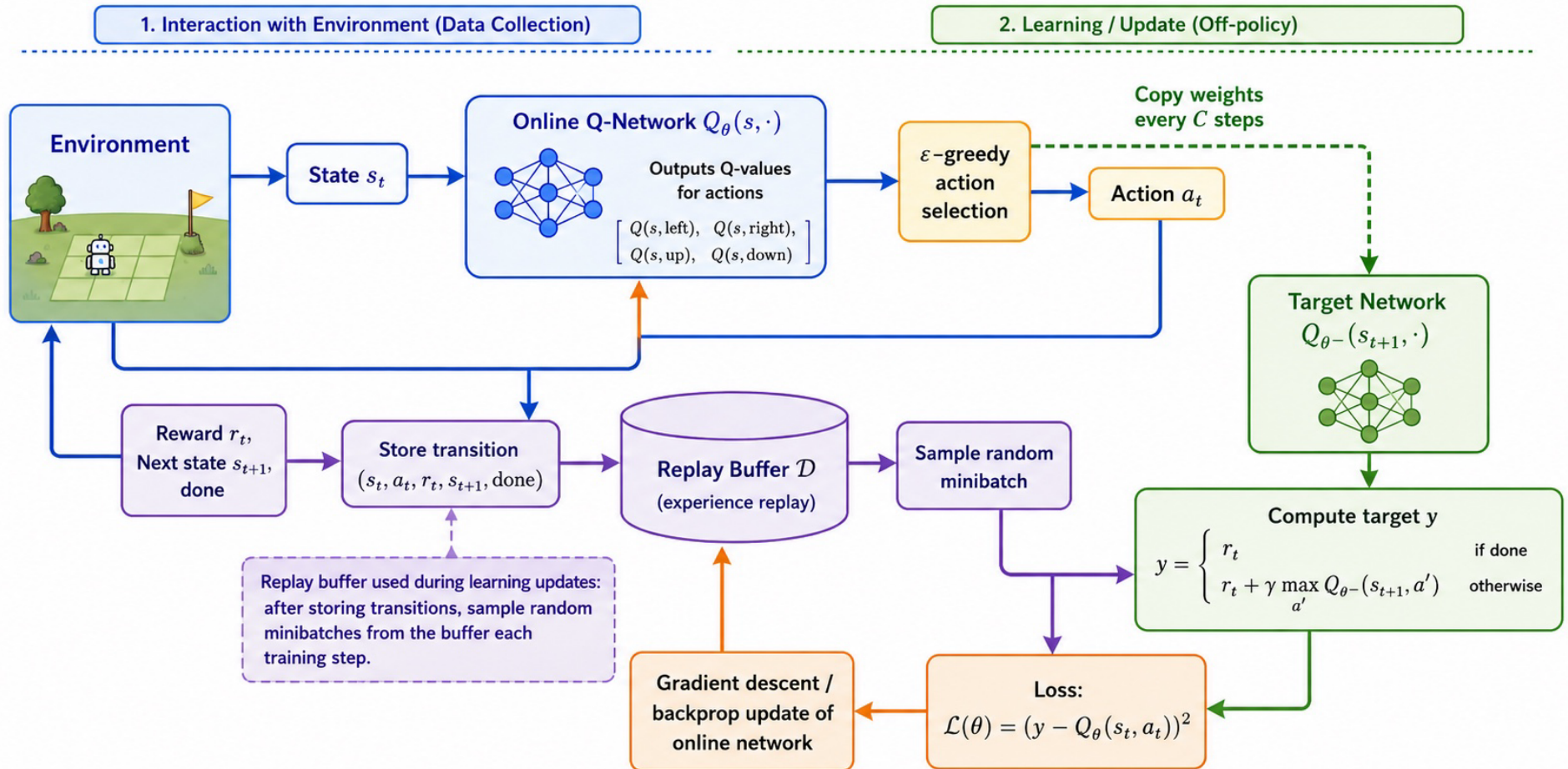
## 5. Putting the pieces together

Full DQN algorithm

5. Minimise  $(y - Q_{\theta}(s,a))^2$

6. [every C steps] Update  $Q_{\theta^-} \leftarrow Q_{\theta}$

# Deep Q-Network (DQN) Workflow



Picture: GPT5.5

## 6. What DQN is (not) good for

DQN can handle discrete action spaces, learns  $Q$  and the policy is derived from that by taking argmax

In contrast, it cannot handle continuous / stochastic action spaces. For that (and other limitations) it might work to learn the policy (by optimizing expected returns)

~> POLICY GRADIENT METHODS & PPO