

INTRODUCTION TO DEEP REINFORCEMENT LEARNING

IA FRAMEWORKS

TOOLS

ML Python Libraries



Python Environment



Viz' Python Libraries



Framework & Tool



TABLE OF CONTENTS

Introduction

Definitions

Iteration algorithm

Q-learning

Deep Q-Learning (DQN)

D3QN

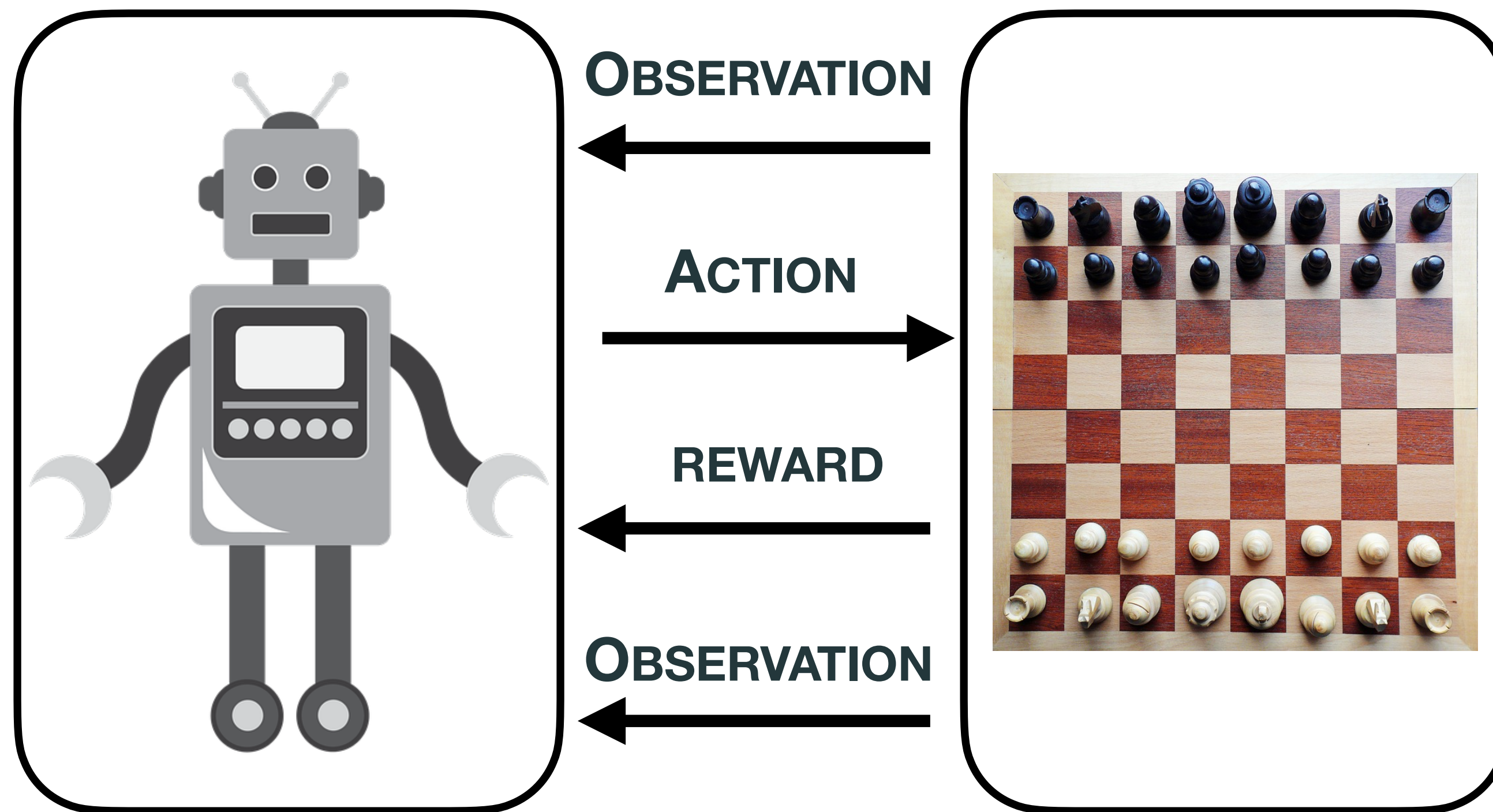
INTRODUCTION

DEFINITIONS

In **reinforcement learning**, an agent makes **OBSERVATIONS** of the **STATES** of an **ENVIRONMENT**.

It takes **ACTION** within the **ENVIRONMENT** and receives **REWARDS**.

Example: Chess Game



OBSERVATION: Locations of each pieces on the chessboard ($32 \times (id, x, y)$).

ACTION: Move piece p from (x_a, y_a) to (x_b, y_b) .

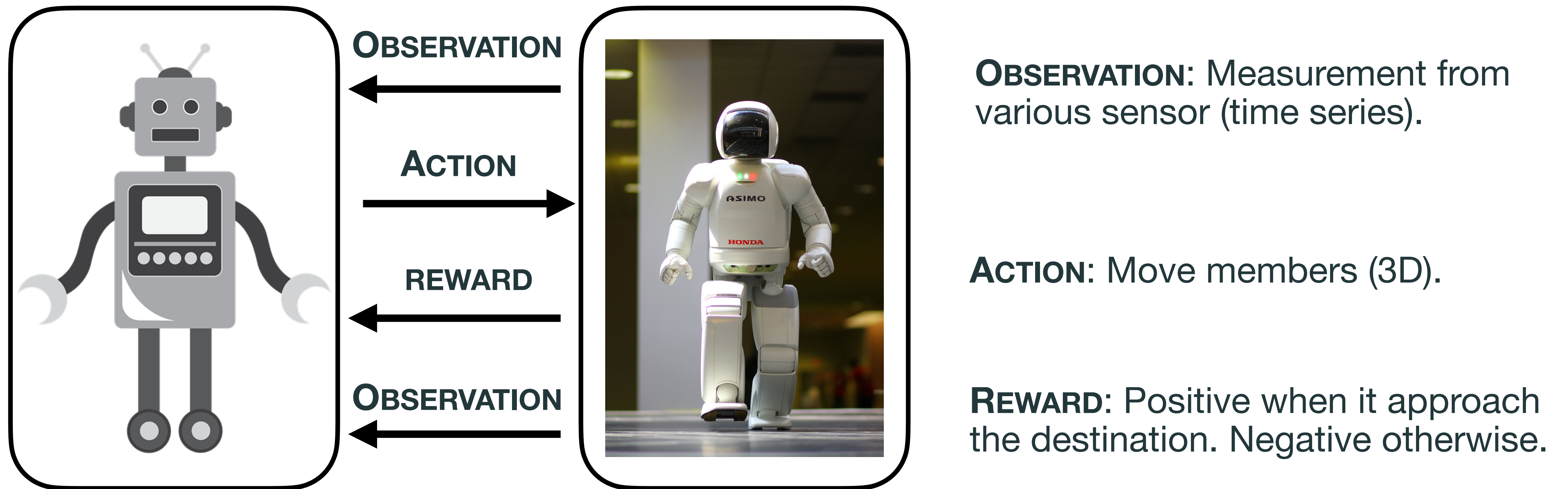
REWARD: Positive if a piece has been token.

DEFINITIONS

In **reinforcement learning**, an agent makes **OBSERVATIONS** of the **STATES** of an **ENVIRONMENT**.

It takes **ACTION** within the **ENVIRONMENT** and receives **REWARDS**.

Example: Walking Robot

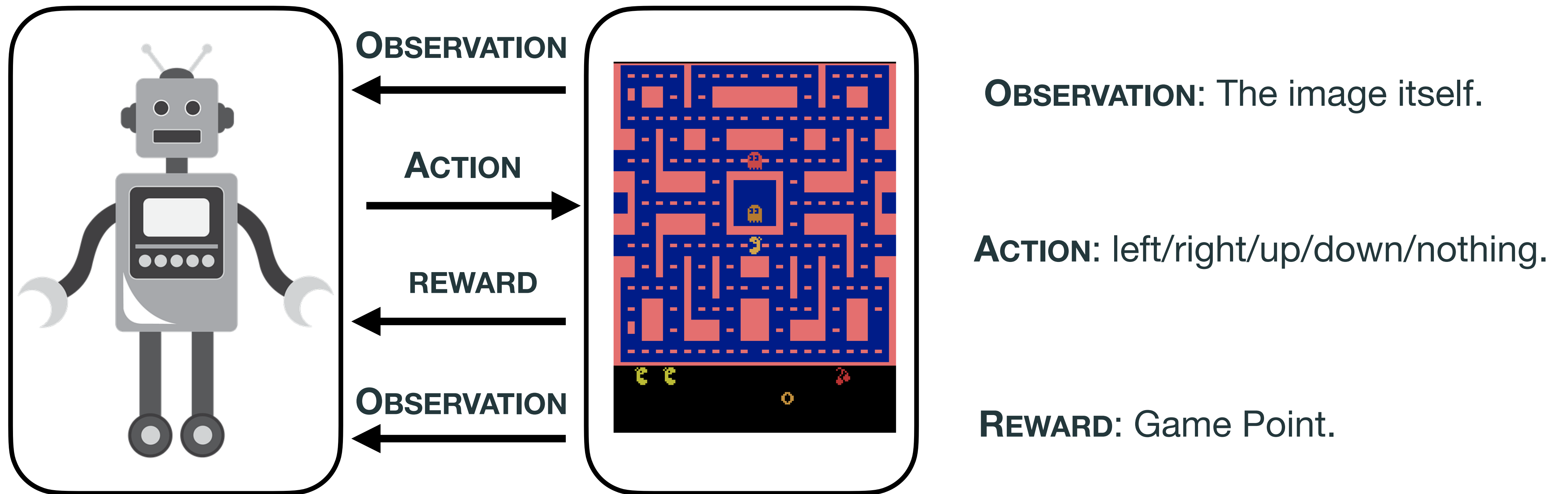


DEFINITIONS

In **reinforcement learning**, an agent makes **OBSERVATIONS** of the **STATES** of an **ENVIRONMENT**.

It takes **ACTION** within the **ENVIRONMENT** and receives **REWARDS**.

Example: Pac Man



DEFINITIONS

When using reinforcement learning you need to define:

An **AGENT**. The one who will take **action** within the **environment**.

An **ENVIRONMENT** where the **action** will be taken.

The **STATE** of the **environment** defines it after each **action**.

The **OBSERVATION** is what we will use from the **state** to decide the next **action**.

A list of possible **ACTIONS** actions that will affect the **environment** and produce a new **state**.

A **REWARD**. A numerical value which reveals how positive or negative the **action** is.

You can't apply reinforcement learning if you're not able to define these objects properly!

OBJECTIVE

Maximise the **long-term** rewards.

Do not pull all the effort on capturing the queen if it means losing all you pieces.

How ? There are two mains approaches:

VALUE-BASED. Look for the optimal **reward**.

- *Learn to estimate the expected rewards for each action in each state.*
- *Use this knowledge to choose the best action.*

POLICY-BASED. Look for the optimal **policy**.

- *Learn directly the best action to take for each observation.*

NB: Methods like Actor-Critic try to optimise both policy and rewards.

OBJECTIVE

Maximise the **long-term** rewards.

Do not pull all the effort on capturing the queen if it means losing all you pieces.

How ? There are two mains approaches:

VALUE-BASED. Look for the optimal reward.

- *Learn to estimate the expected rewards for each action in each state.*
- *Use this knowledge to choose the best action.*

POLICY-BASED. Look for the optimal **policy**.

- *Learn directly the best action to take for each observation.*

NB: Methods like Actor-Critic try to optimise both policy and rewards.

DEFINITIONS

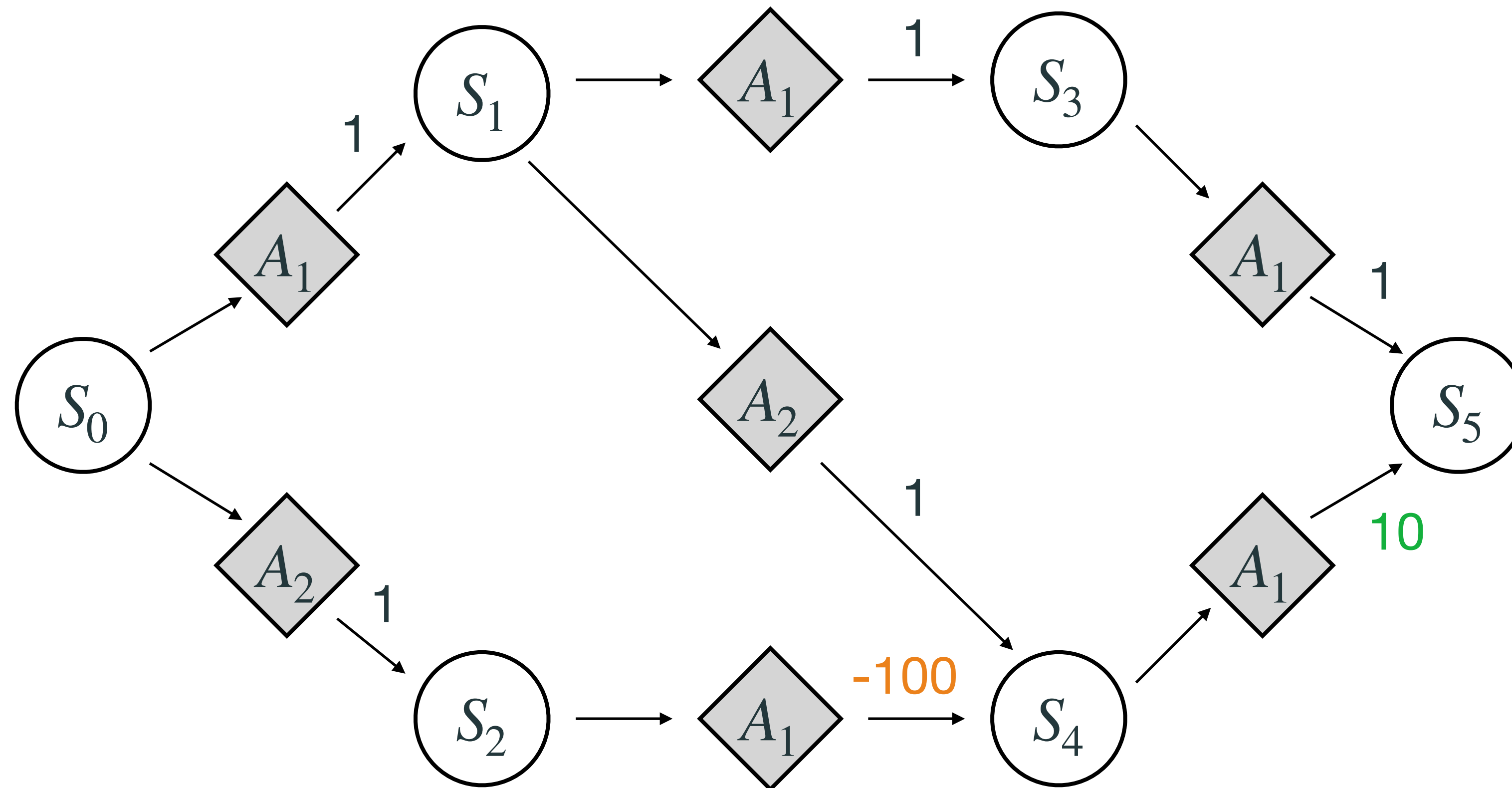
MARKOV DECISION PROCESS - DEFINITION

A **Markov Decision Process** is composed of :

- A set of state $S = \{s_1, s_2, \dots, s_n\}$.
- A set of action $A = \{a_1, a_2, \dots, a_m\}$
- A reward function : $R = S \times A \times S \rightarrow \mathcal{R}$
- A transition function: $P_{ij}^a = P(s_j | s_i, a_k)$

MARKOV DECISION PROCESS - EXAMPLE

- A set of state $S = \{1,2,3,4,5\}$.
- A set of action $A = \{1,2\}$
- A reward function : $R = [\{0,A,1\} = 1, \{0,B,2\} = 2, \dots]$
- A transition function: $P_{ij}^a = 1, \forall (i,j,a), 0$ otherwise



POLICIES

There are 3 policies for this **MDP**.

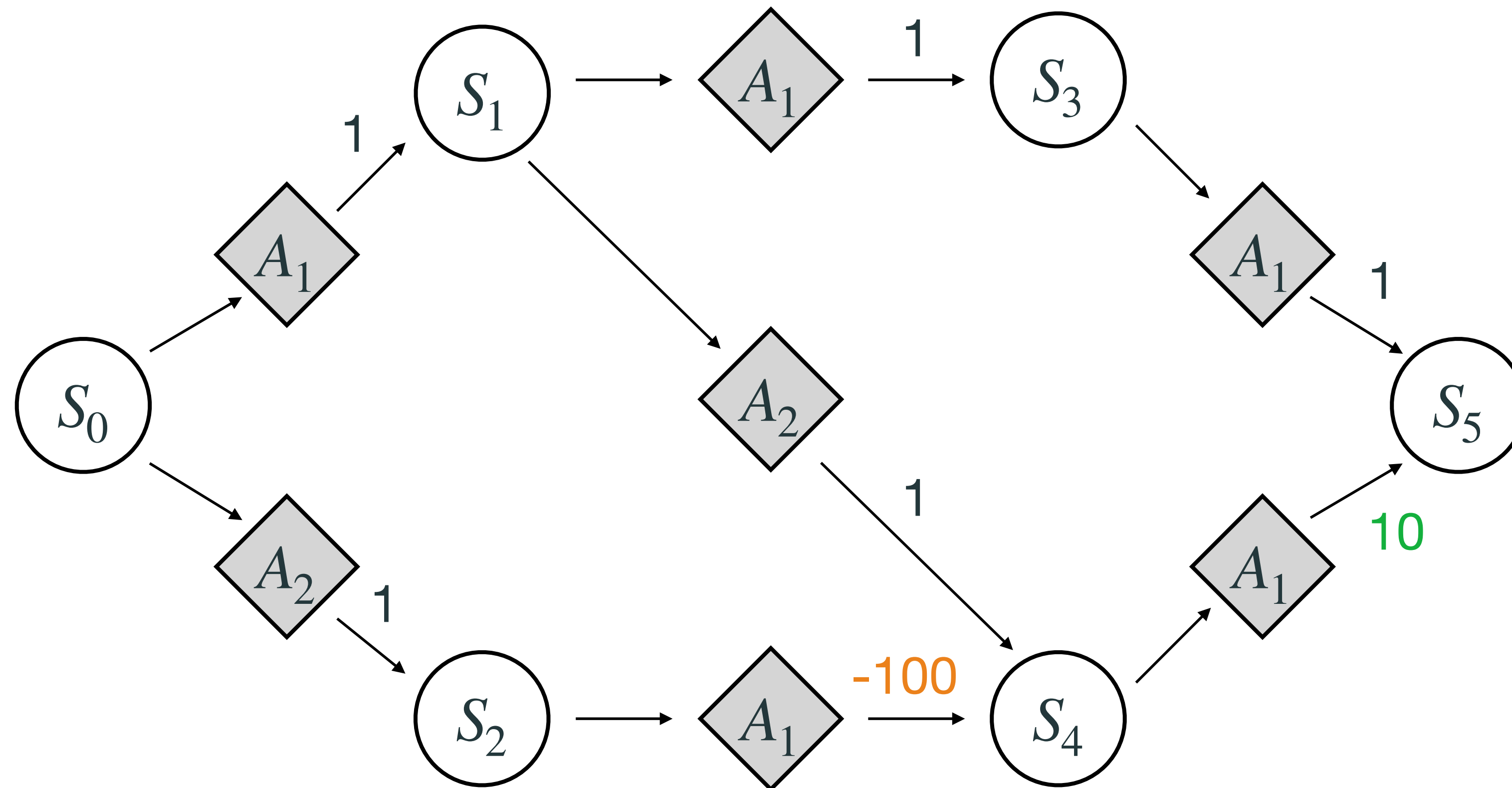
$$\Pi_1 = \{0 \rightarrow 1 \rightarrow 3 \rightarrow 5\} \quad R = 3$$

$$\Pi_2 = \{0 \rightarrow 1 \rightarrow 4 \rightarrow 5\} \quad R = 12$$

$$\Pi_3 = \{0 \rightarrow 2 \rightarrow 4 \rightarrow 5\} \quad R = -89$$

Which one is the best?

Compute their total reward !



POLICIES

There are 3 policies for this **MDP**.

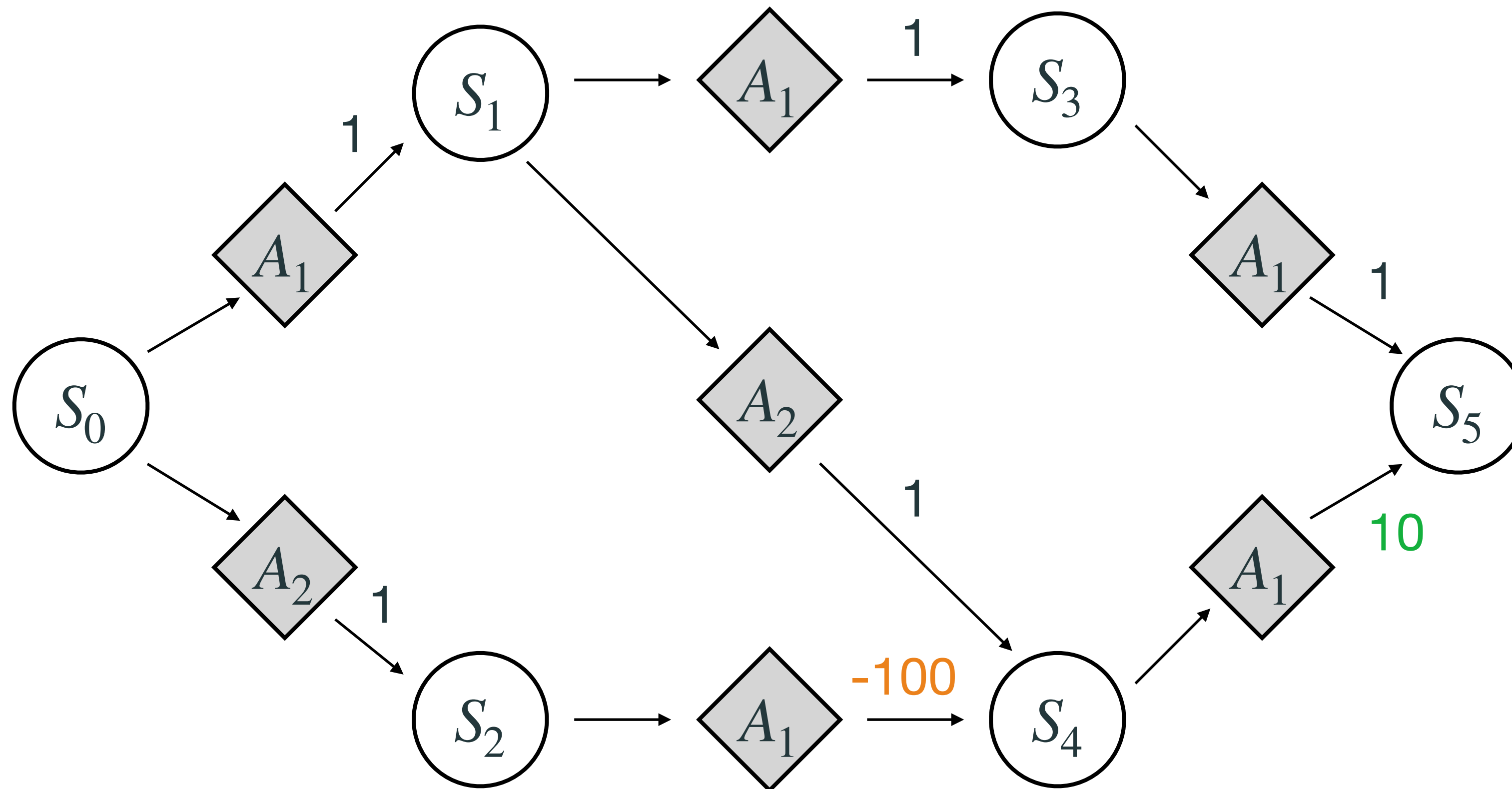
$$\Pi_1 = \{0 \rightarrow 1 \rightarrow 3 \rightarrow 5\} \quad R = 3$$

$$\Pi_2 = \{0 \rightarrow 1 \rightarrow 4 \rightarrow 5\} \quad R = 12$$

$$\Pi_3 = \{0 \rightarrow 2 \rightarrow 4 \rightarrow 5\} \quad R = -89$$

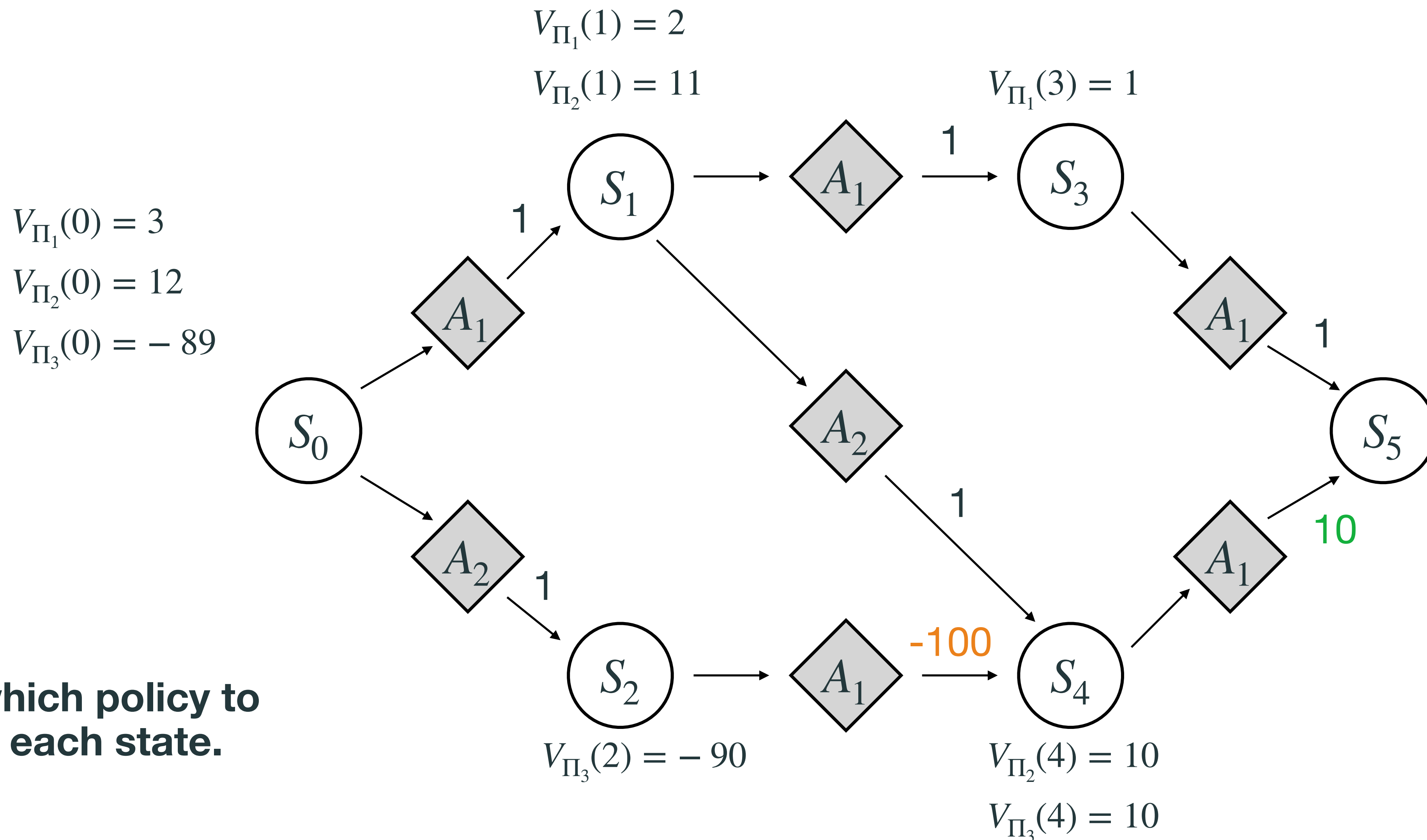
Which one is the best?

Compute their total reward !



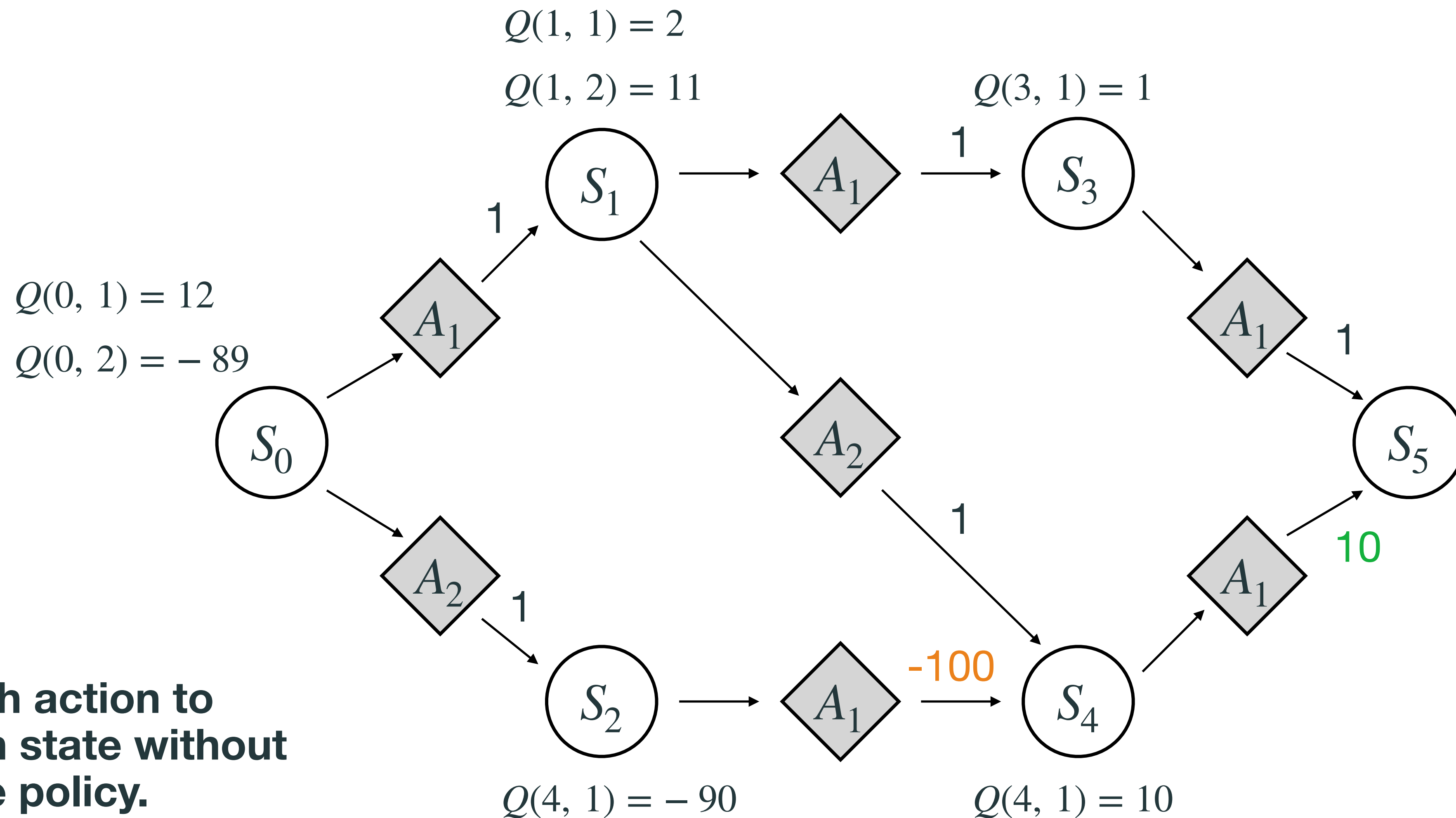
STATE VALUE

A **state value** $V_{\Pi}(s)$ describes how good is it to run policy Π from state s .



ACTION-STATE VALUE

An **action-state value** $Q(s, a)$ describes the value of taking and action a from state s .



GENERAL DEFINITION

State value of a policy π :

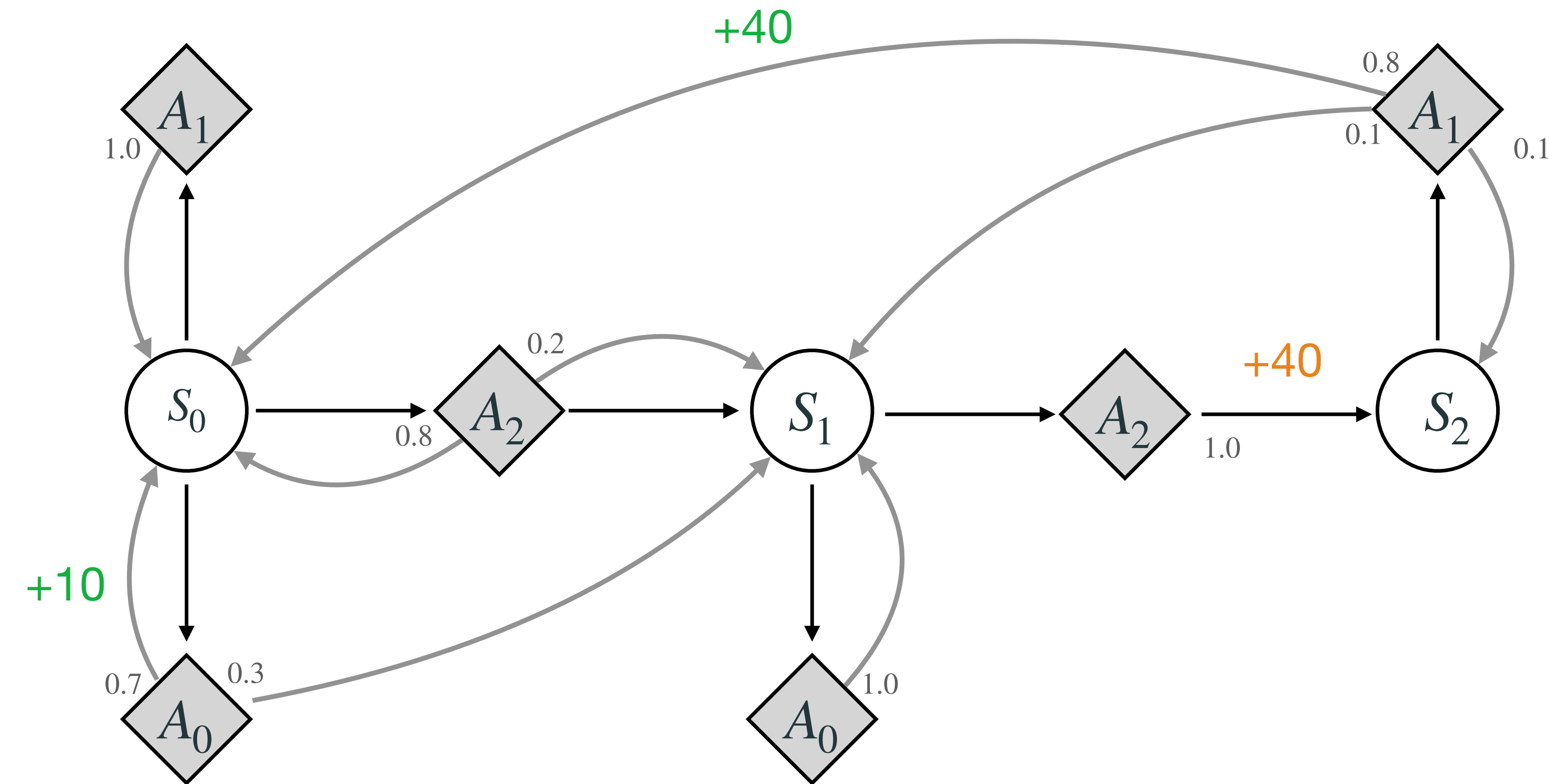
$$V^\pi(s) = \mathbb{E} \left(\lim_{H \rightarrow \infty} \sum_{t=0}^H \gamma^t r_t \middle| s_0 = s, \pi \right)$$

Action-state value function of a policy

$$Q^\pi(s, a) = \mathbb{E} \left(\lim_{H \rightarrow \infty} \sum_{t=0}^H \gamma^t r_t \middle| s_0 = s, a_0 = a, \pi \right)$$

ACTION-STATE VALUE - PROBABILISTIC ACTION

What if your action can lead you to different state ?



From state = S_0

$$P_{S_0, S_0}^{A_0} = 0.7 \quad P_{S_0, S_1}^{A_0} = 0.3$$

$$P_{S_0, S_0}^{A_1} = 1$$

$$P_{S_0, S_0}^{A_2} = 0.8 \quad P_{S_0, S_1}^{A_2} = 0.2$$

From state = S_1

$$P_{S_1, S_1}^{A_0} = 1$$

$$P_{S_0, S_2}^{A_2} = 1$$

From state = S_2

$$P_{S_2, S_0}^{A_1} = 0.8 \quad P_{S_2, S_1}^{A_1} = 0.1 \quad P_{S_2, S_2}^{A_1} = 0.1$$

SOLUTION: The Bellman Equation

BELLMAN OPTIMALITY EQUATION

For all s :

$$V^*(s) = \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V^*(s')]$$

- V^* is the **optimal state value** at state s .
- $P_{s,s'}^a$ is the **transition probability** from state s to state s' given that the agent choice action a .
- $R(s, a, s')$ is the **reward** that the agent gets when it goes from state s to state s' given that the agent choose action a .
- γ is the **discount rate** (*the importance we give to future rewards*).

BELLMAN OPTIMALITY EQUATION

For all s :

$$V^*(s) = \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V^*(s')]$$

- V^* is the **optimal state value** at state s .
- $P_{s,s'}^a$ is the **transition probability** from state s to state s' given that the agent choice action a .
- $R(s, a, s')$ is the **reward** that the agent gets when it goes from state s to state s' given that the agent choose action a .
- γ is the **discount rate** (*the importance we give to future rewards*).

BELLMAN OPTIMALITY EQUATION

For all s :

$$V^*(s) = \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V^*(s')]$$

- V^* is the **optimal state value** at state s .
- $P_{s,s'}^a$ is the **transition probability** from state s to state s' given that the agent choice action a .
- $R(s, a, s')$ is the **reward** that the agent gets when it goes from state s to state s' given that the agent choose action a .
- γ is the **discount rate** (*the importance we give to future rewards*).

BELLMAN OPTIMALITY EQUATION

For all s :

$$V^*(s) = \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V^*(s')]$$

- V^* is the **optimal state value** at state s .
- $P_{s,s'}^a$ is the **transition probability** from state s to state s' given that the agent choice action a .
- $R(s, a, s')$ is the **reward** that the agent gets when it goes from state s to state s' given that the agent choose action a .
- γ is the **discount rate** (*the importance we give to future rewards*).

How to estimate this optimal sates? The **VALUE ITERATION ALGORITHM**.

ITERATION ALGORITHM

ITERATION ALGORITHM

VALUE ITERATION ALGORITHM.

- Initialise: $V_0(s) = 0, \forall s$.
- Iterate, $\forall s$:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V_k(s')]$$

- $V_k(s)$ is the estimated value of state s at the k^{th} iteration of the algorithm.
- Algorithm is guaranteed to converge to the optimal state value (given enough time)

We can now evaluate the optimal policy...

... but we don't know which action to take at each state !

ITERATION ALGORITHM

Q-VALUE ITERATION ALGORITHM.

- Initialise: $Q_0(s, a) = 0, \forall (s, a)$.
- Iterate, $\forall (s, a)$:

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot \max_{a'} Q_k(s', a')]$$

Once you have the optimal **Q-value**, Q^* , the **optimal policy**, $\Pi^*(s)$, is defined as:

$$\Pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

REAL APPLICATION

PROBLEM: The **Q-value iteration algorithm** required to know **the complete MDP**:

- complete list of possible states $s \in S$ and actions $a \in A$,
- all transition probabilities $T(s, a, s')$
- all possible reward $R(s, a, s')$

In real case, the agent has only **partial knowledge of the MDP**:

- complete list of possible states $s \in S$ and actions $a \in A$,
- $T(s, a, s')$ and $R(s, a, s')$ are unknown

PROPOSITION: Estimate those values

- $R(s, a, s')$ requires to see each transition at least once.
- $T(s, a, s')$ requires to experience it multiple times.

Q-LEARNING

TEMPORAL DIFFERENCE LEARNING

TD-LEARNING ITERATION ALGORITHM

$$V_{k+1}(s) \leftarrow (1 - \alpha)V_k(s) + \alpha V_{k+1}^*(s)$$

- Initialise: $V_0(s) = 0, \forall s$.
- At iteration k , we have an estimation of the optimal state value $V_k(s)$, then $\forall s$:
 - Choose a reachable state s' (*we'll come back to that later*).
 - Compute an estimation of the next state value $V_{k+1}^*(s) = R(s, a, s') + \gamma \cdot V_k(s')$
 - Update V_{k+1} as a combination of the previous estimation $V_k(s)$ and the estimation $V_{k+1}^*(s)$

Where the learning rate $\alpha \in [0,1]$

TEMPORAL DIFFERENCE LEARNING VS VALUE ITERATION ALGORITHM

TD-LEARNING ITERATION ALGORITHM.

$$V_{k+1}(s) \leftarrow (1 - \alpha)V_k(s) + \alpha[R(s, a, s') + \gamma \cdot V_k(s')]$$

$$V_{k+1}(s) \leftarrow (1 - \alpha)V_k(s) + \alpha V_{k+1}^*(s)$$

VALUE ITERATION ALGORITHM.

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot V_k(s')]$$

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P_{s,s'}^a V_{k+1}^*(s)$$

- The **discount rate** γ represent how much you can trust the features reward.
- The **learning rate** α represent how much you can trust the next state value estimation.

The Q-value iteration algorithm can be adapted the same way.

Q-LEARNING

Q-LEARNING ITERATION ALGORITHM

$$Q_{k+1}(s, a) \leftarrow (1 - \alpha)Q_k(s, a) + \alpha \left[R(s, a, s') + \gamma \cdot \max_{a'} Q_k(s', a') \right]$$

- We build a memory table to store **Q-value** for all possible combinations of s and a .
- As we keep playing we update this table.
- And it will converge...
- ...with some tricks!

OFF POLICY

OFF-POLICY ALGORITHM

The **TD-Learning** and **Q-learning** iteration algorithm used with random policy exploration.

Π_l : The policy you are learning (*by updating either the Q-value or the state value*).

Π_e : The policy you are executing (*to choose the action you're taking*).

- Π_e is random when using off-policy algorithm
- You are learning how to act by observing someone doing random action.
- Q-learning will learn only if the exploration policy explores the MDP enough.
- It may take an extremely long time to do so..

Can't we do better?

EXPLORATION POLICY : ϵ — *greedy* POLICY

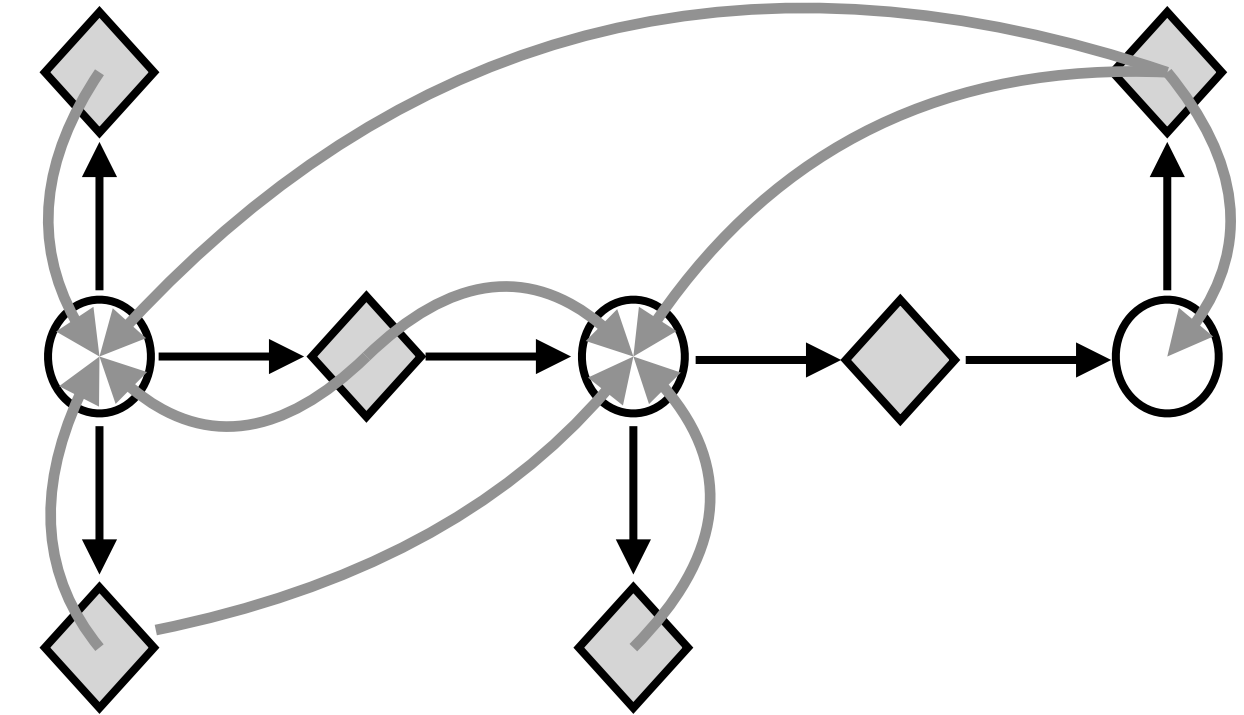
At each step you will choose either:

- The exploration policy (random one) Π_e with a probability ϵ
- The learned policy Π_l with a probability $(1 - \epsilon)$

You start with $\epsilon = 1$: purely random exploration of the environment.

You decrease the value of epsilon over iteration to explore less and exploit more the interesting parts of the environment

- *Q_learning.ipynb* : Implement **Q-VALUE ITERATION ALGORITHM** and **Q-LEARNING ITERATION ALGORITHM** on the toy MDP of this presentation



DEEP Q-LEARNING (DQN)

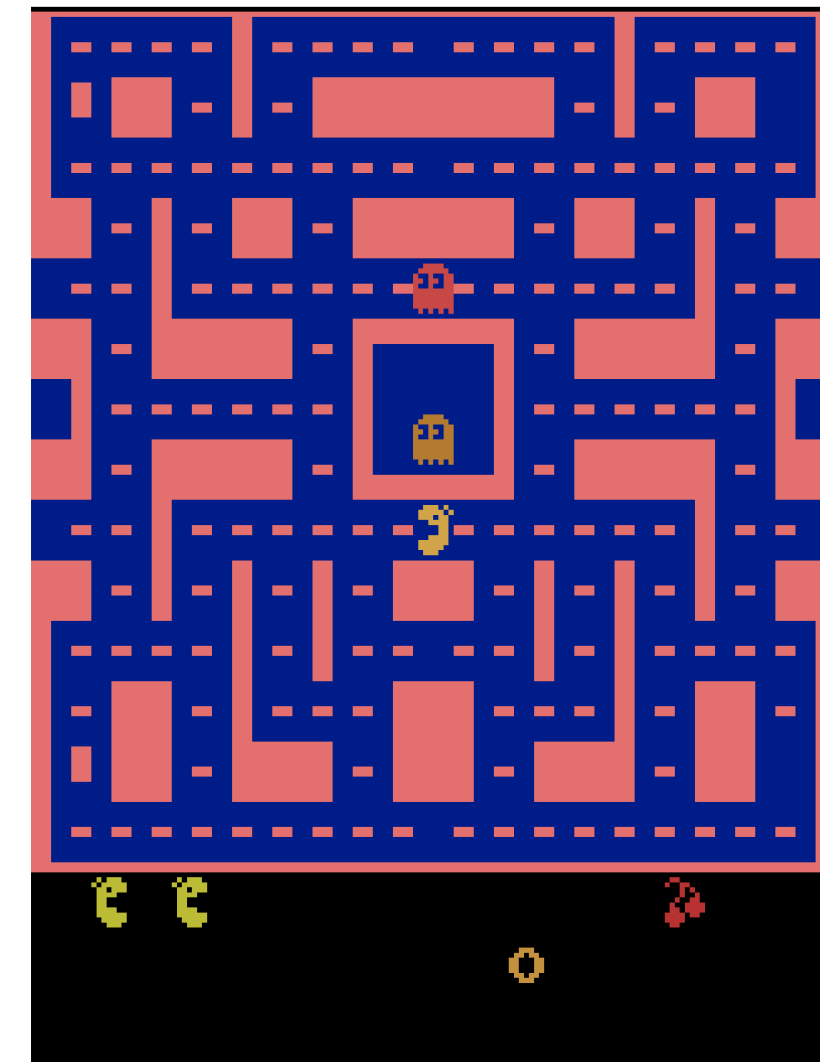
APPROXIMATE Q-LEARNING

PROBLEM: Q-learning still does not scale to large MDPs

- $T(s, a, s')$ and $R(s, a, s')$ are not computed for every (s, a) couple.
- But $Q(s, a)$ are!

EXAMPLE: Pacman

- 240 pellets that can be eaten or here or not.
- It represent $2^{240} \approx 10^{72}$ configuration. *Number of atoms in the universe.*
- And we need to consider ghost and Pacman position...



We can't compute nor store Q-value for every (s, a) couple.

APPROXIMATE AND DEEP Q-LEARNING

OBJECTIVE: Use a function to evaluate the Q-value

Q-LEARNING

Store and update all Q-values in a big table

Pick value from it.

APPROXIMATE Q-LEARNING

We train a function Q_θ that will generalise the approximation of the Q-values table

Use this function to estimate value

Compute hand-crafted features (localisation of paceman, distance to the ghost, localisation of remaining pellets...)

DEEP Q-LEARNING

Use a CNN network as a Q_θ function and use image as features.

DeepMind makes it work (very well) on various Atari Game in 2014.

APPROXIMATE AND DEEP Q-LEARNING

Q-VALUE ITERATION ALGORITHM.

$$Q_{k+1}(s, a) \leftarrow \sum_{s'} P_{s,s'}^a [R(s, a, s') + \gamma \cdot \max_{a'} Q_k(s', a')]$$

Q-LEARNING

$$\begin{aligned} target &= R(s, a, s') + \gamma \cdot \max_{a'} Q_k(s', a') \\ Q_{k+1}(s, a) &\leftarrow (1 - \alpha)Q_k(s, a) + \alpha [target] \end{aligned}$$

APPROXIMATE/DEEP Q-LEARNING

$$target = R(s, a, s') + \gamma \cdot \max_{a'} Q_k(s', a')$$

At step k , generate $batch_size$ *targets* and train model on this batch

$$\theta_{k+1} \leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s \sim P(s'|s,a)} [(Q_{\theta}(s, a) - target(s'))^2]_{\theta=\theta_k}$$

- θ_k are the weights of the (Deep) Q function computed at step k (batch)
- α is the learning rate and can be interpreted the same way.

PSEUDO CODE

- Create Q network.
- Iterate while *max_number_episode* or *goal_reward* is not achieved:
 - Play one episode - Explore randomly or Exploit with Q according to *random_probability*.
 - If we train more than *min_pre_trained_episode*
 - Create target using Q .
 - Train the Q network.
 - Start decreasing *random_probability*.

DEFINITIONS

EPISODE: One complete run of a game.

Pacman: all the move until the pacman died.

STEP: One move of the game.

Pacman: one action taken in the environment (up, down, left, right)

EXPERIENCE: All the information that resume one step: .

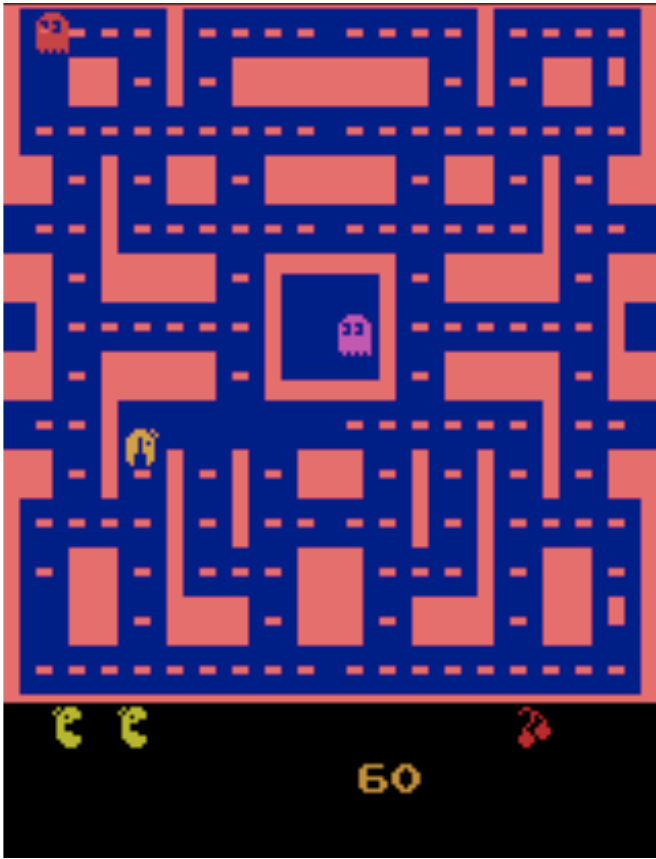
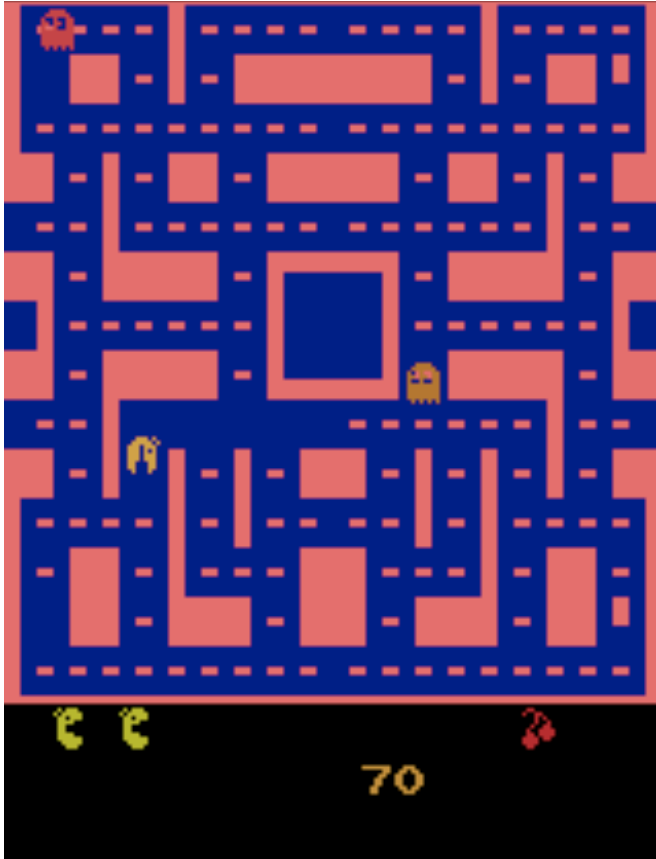
s : from state

a : action chosen

s' : to state

r : reward obtained

d : is game over

<i>Pacman</i>	<i>From state</i>	<i>Action</i>	<i>To state</i>	<i>Reward</i>
		DOWN		10

EXPERIENCES FOR TRAINING DATASET

<i>from state</i>	<i>action</i>	<i>to state</i>	<i>reward</i>
$s_{0,0}$	$a_{0,0}$	$s'_{0,0}$	$r_{0,0}$
$s_{0,1}$	$a_{0,1}$	$s'_{0,1}$	$r_{0,1}$
\dots	\dots	\dots	\dots
s_{0,ns_0}	a_{0,ns_0}	s'_{0,ns_0}	r_{0,ns_0}
$s_{1,0}$	$a_{1,0}$	$s'_{0,1,0}$	$r_{1,0}$
\dots	\dots	\dots	\dots
$s_{ne,0}$	$a_{ne,0}$	$s'_{ne,0}$	$r_{ne,0}$
\dots	\dots	\dots	\dots
$s_{ne,ns_{ne}}$	$a_{ne,ns_{ne}}$	$s'_{ne,ns_{ne}}$	$r_{ne,ns_{ne}}$

- $\{s, a, s', r\}_{i,j}$, from state, action, to state and reward of **episode** i at **step** j .
- ne number of **episodes** in the training dataset.
- ns_i number of **steps** in **episode** i .
- Number of experience in the training dataset : $\sum_{i=0}^{ne} ns_i$

PROBLEM: The input of the supervised learning problem are not $i . i . d$.

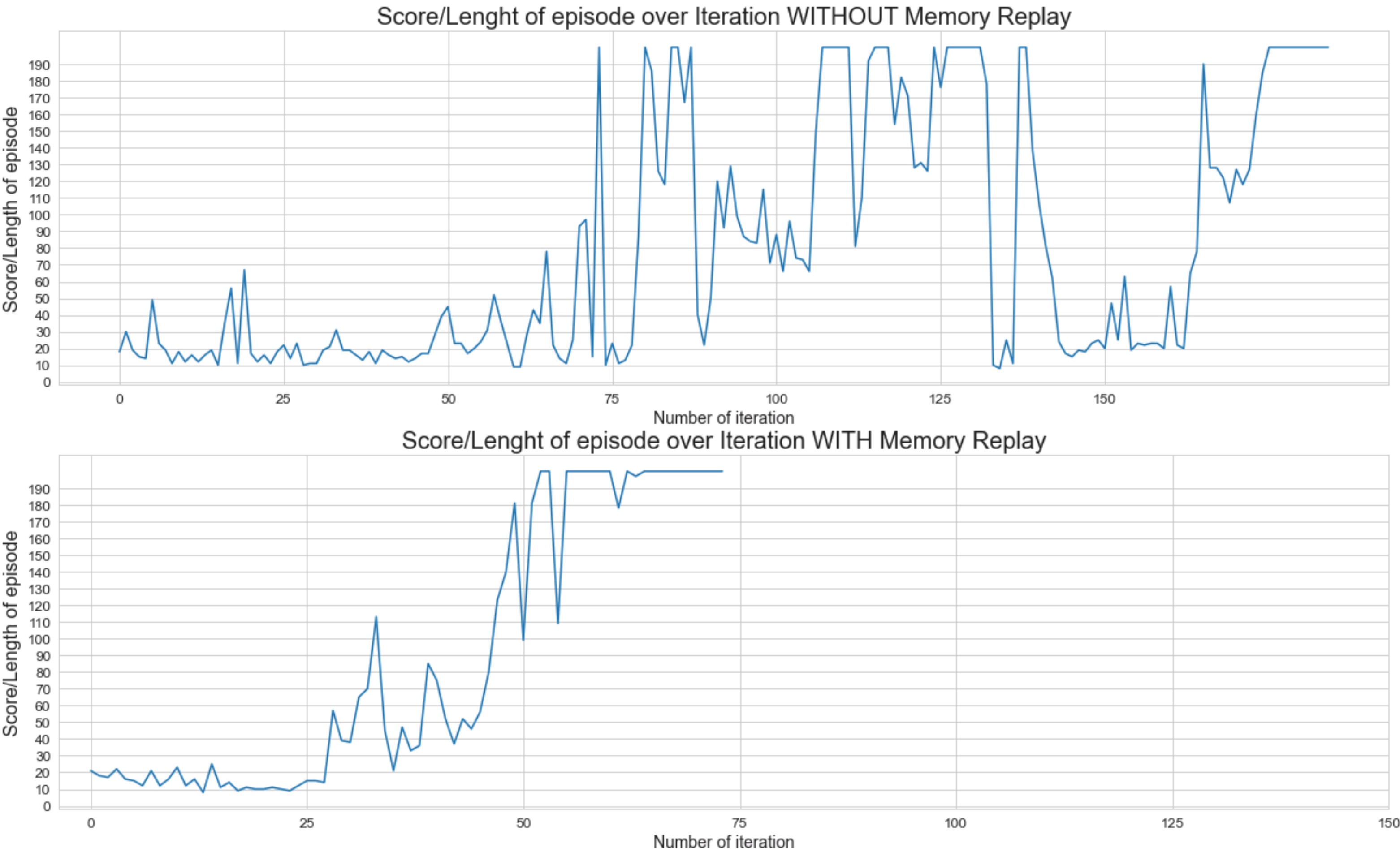
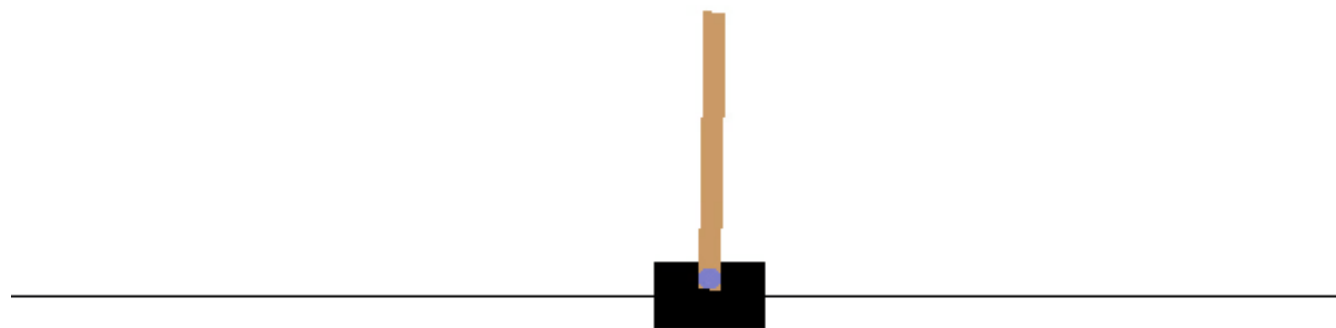
EXPERIENCE REPLAY

SOLUTION: The Experience Replay

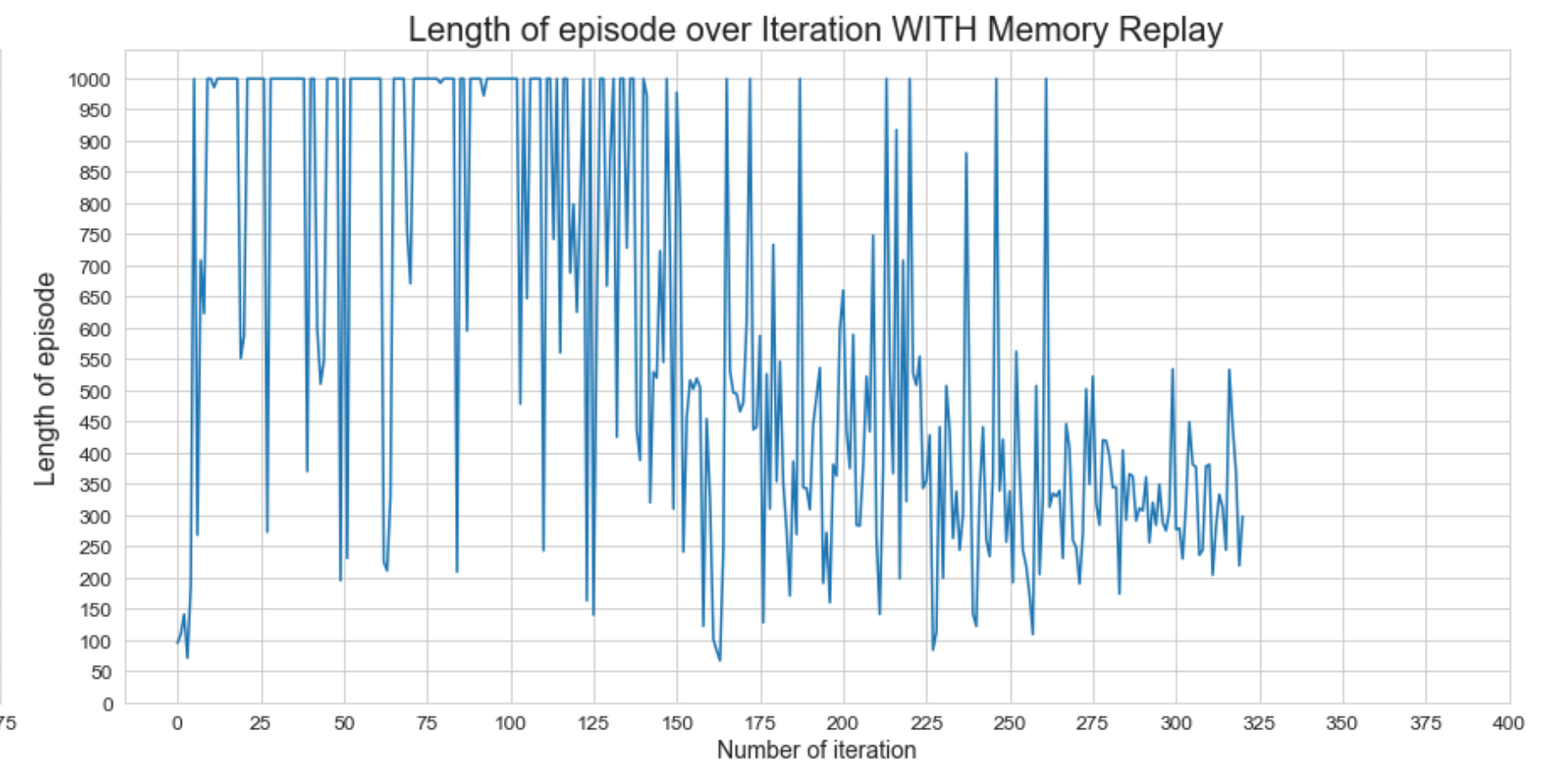
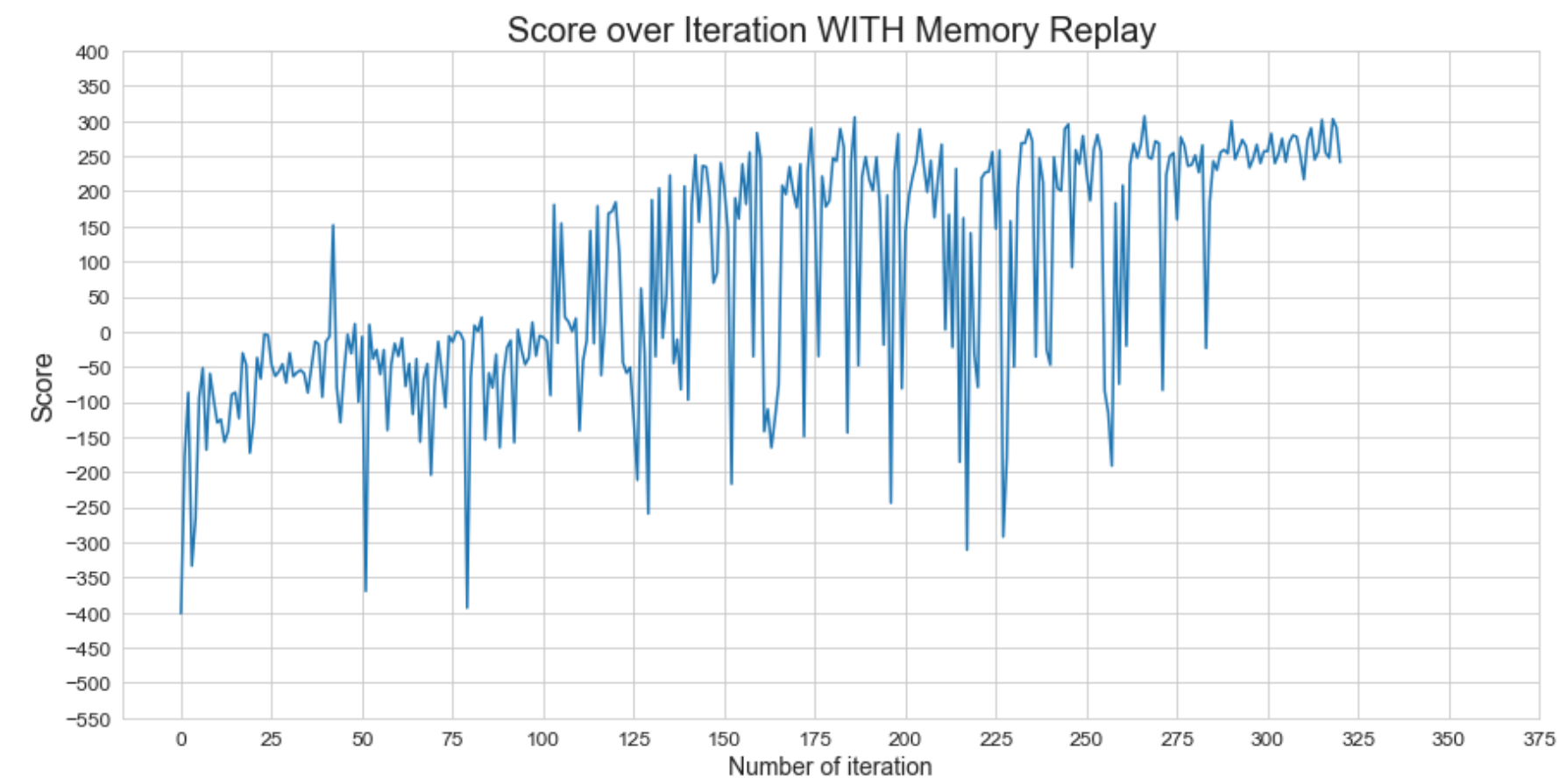
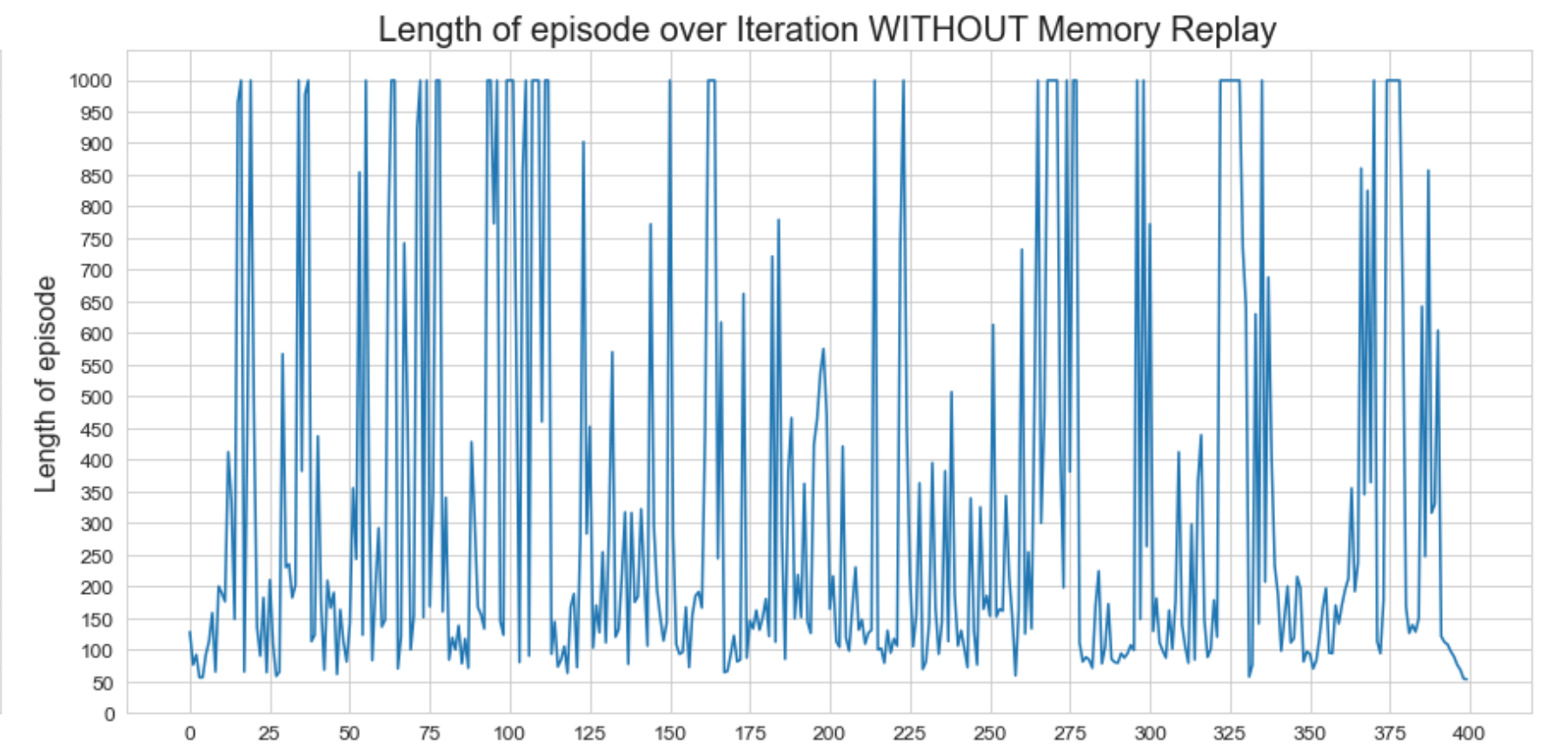
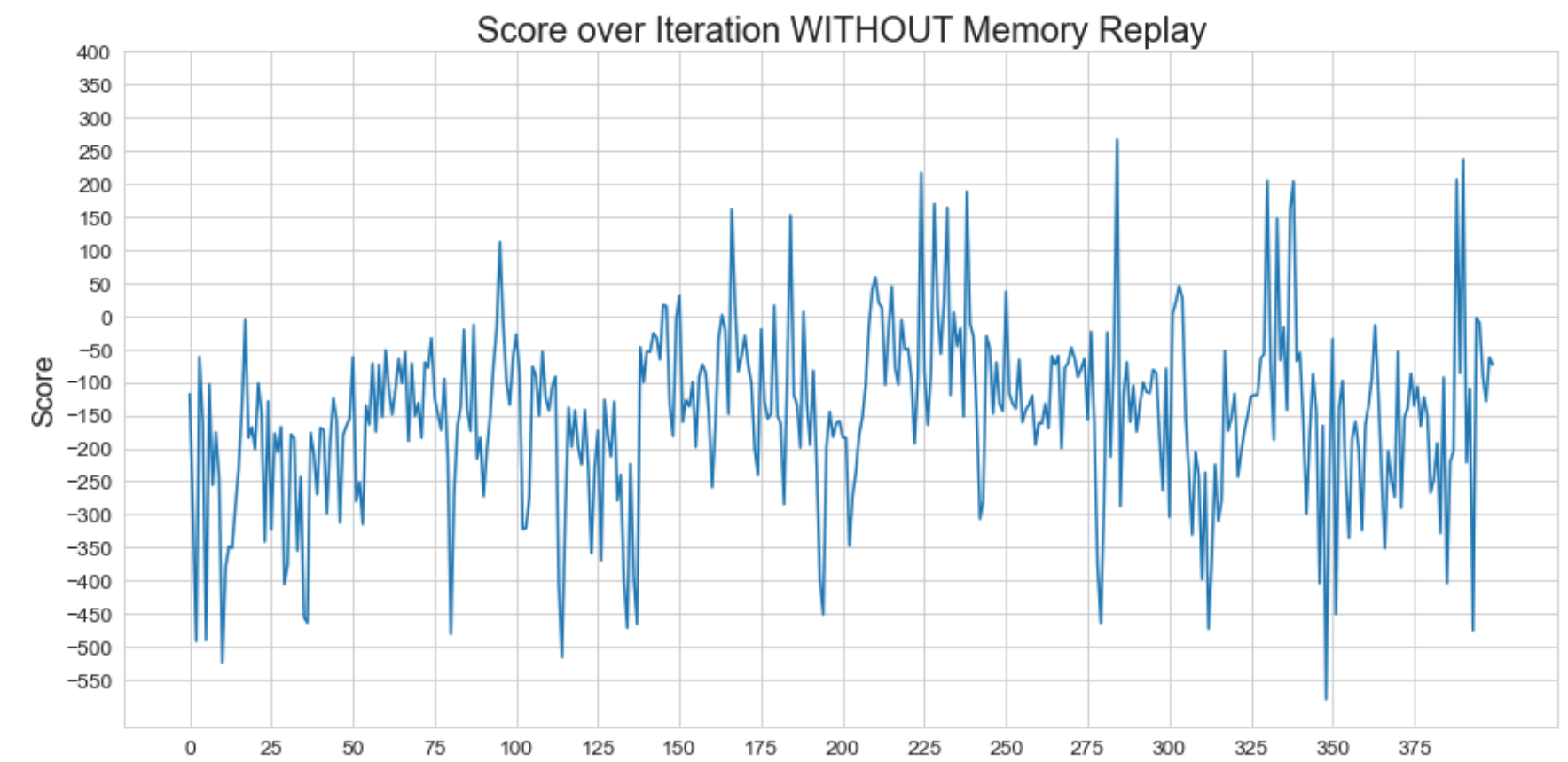
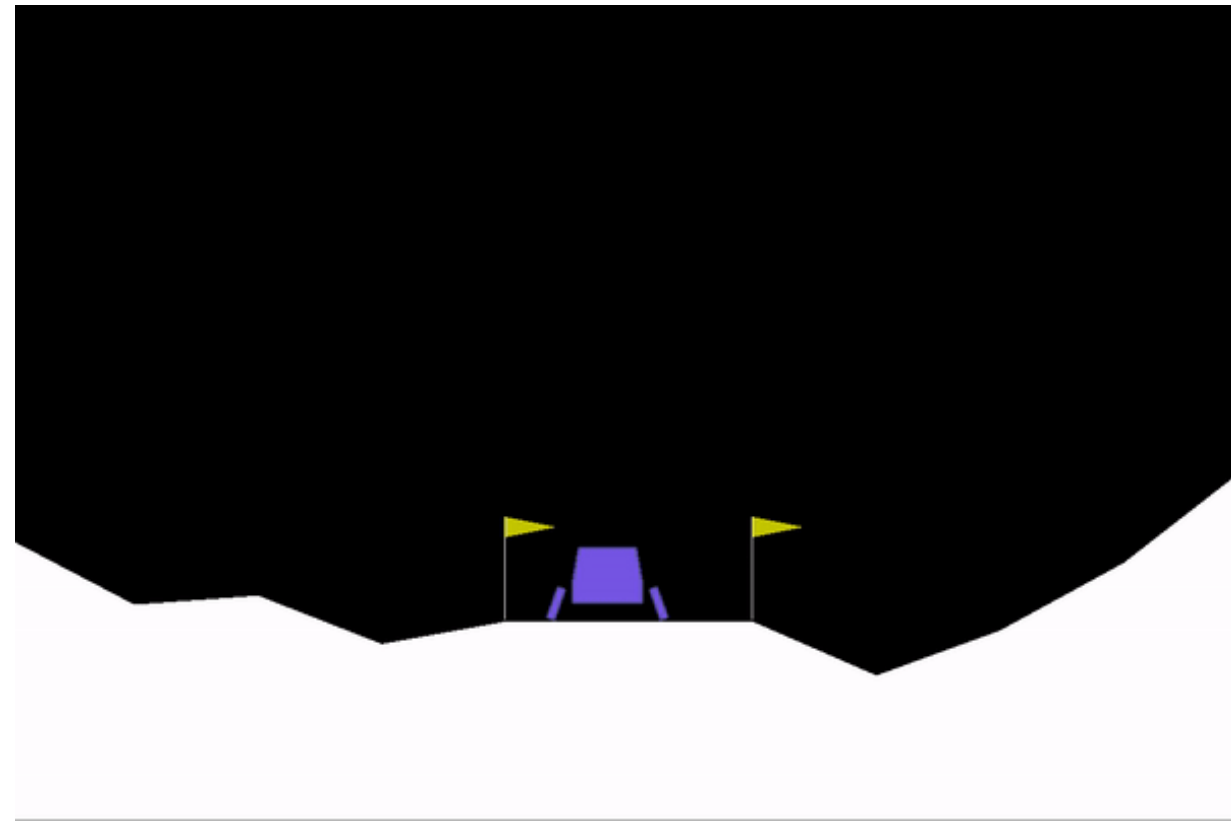
<i>from state</i>	<i>action</i>	<i>to state</i>	<i>reward</i>
$s_{10,3}$	$a_{10,3}$	$s'_{10,3}$	$r_{10,3}$
$s_{8,37}$	$a_{8,37}$	$s'_{8,37}$	$r_{8,37}$
...
$s_{3,17}$	$a_{3,17}$	$s'_{3,17}$	$r_{3,17}$
$s_{5,38}$	$a_{5,38}$	$s'_{5,38}$	$r_{5,38}$
...
$s_{8,6}$	$a_{8,6}$	$s'_{8,6}$	$r_{8,6}$
...
$s_{3,10}$	$a_{3,10}$	$s'_{3,10}$	$r_{3,10}$

- $\{s, a, s', r\}_{i,j}$, from state, action, to state and reward of **episode** i at **step** j .
- Fix the number of **experiences** in the training dataset ($batch_size$).
- Play ne **episodes** such that
$$\sum_{i=0}^{ne} ns_i > batch_size.$$
- Sample $batch_size$ **experiences** from the ne episodes generated.

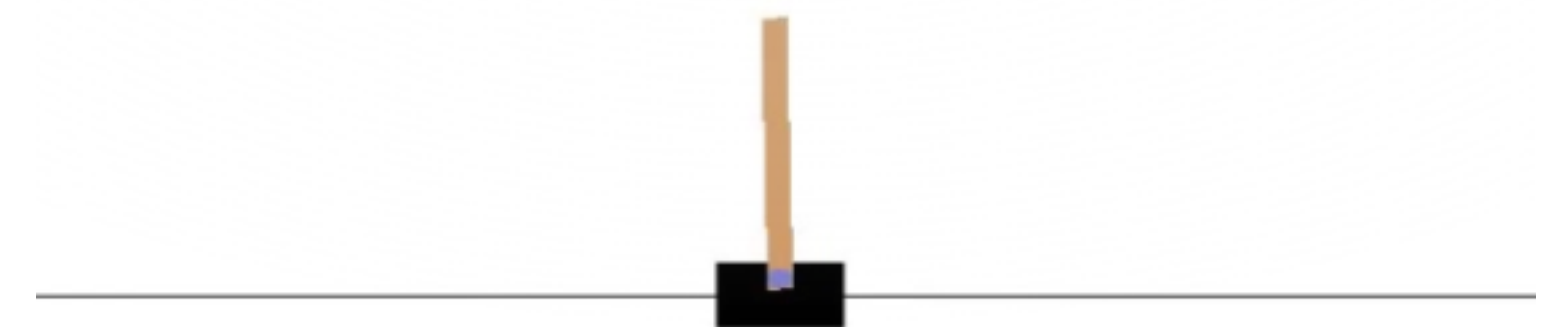
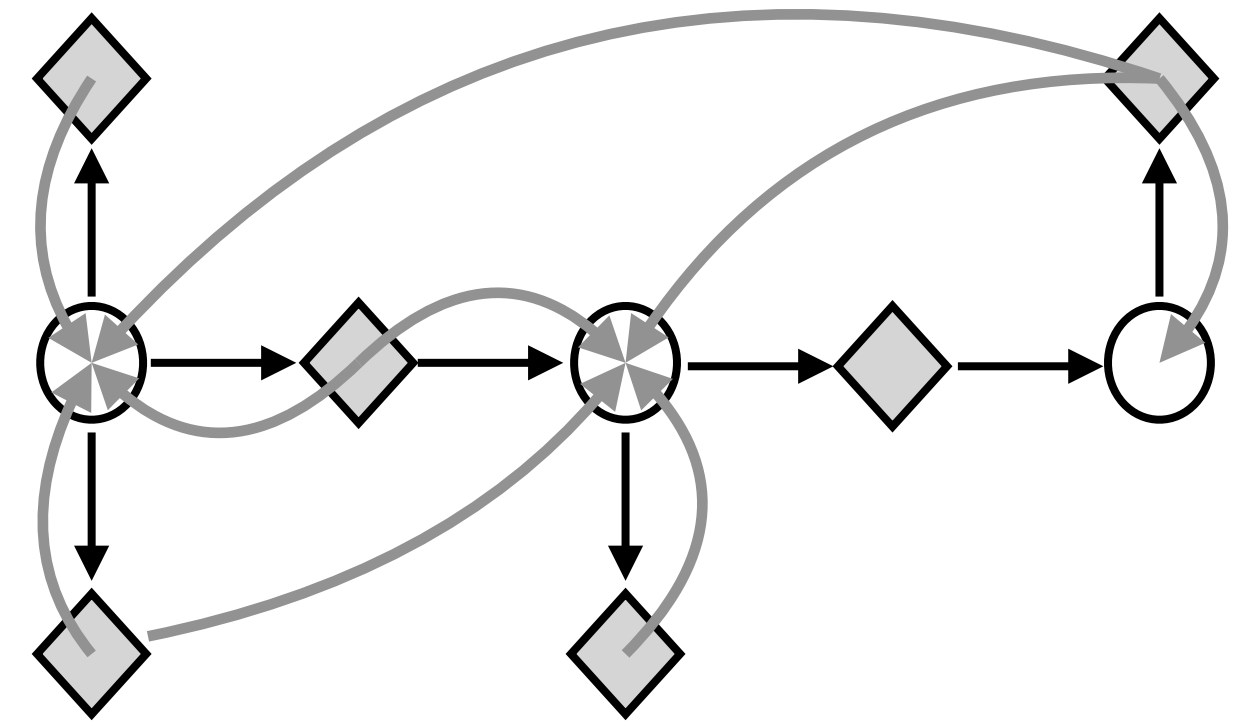
EXPERIENCE REPLAY - CART POLE



EXPERIENCE REPLAY – LUNAR LANDING



- *Q_learning.ipynb* : Implement **Q-VALUE ITERATION ALGORITHM** and **Q-LEARNING ITERATION ALGORITHM** on the toy **MDP** of this presentation.
- *Deep_Q_learning_CartPole.ipynb* : Implement **Deep Q Learning** with and without Replay memory buffer on **CartPole**.



TARGET NETWORK

PROBLEM: The targets are unstable and depend of Q itself.

$$\begin{aligned} \text{targets} &= [R(s, a, s') + \gamma \max_{a'} Q(s', a')] \\ \theta_{k+1} &\leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s' \sim P(s'|s,a)} [(Q_{\theta}(s, a) - \text{target}(s'))^2]_{\theta=\theta_k} \end{aligned}$$

At each train iteration:

The Q values shift to get closer to the target which shift the same way.

It is chasing a **non-stationary target**

TARGET NETWORK

PROBLEM: The targets are unstable and depend of Q itself.

$$\begin{aligned} \text{targets} &= [R(s, a, s') + \gamma \max_{a'} Q_{\text{target}}(s', a')] \\ \theta_{k+1} &\leftarrow \theta_k - \alpha \nabla_{\theta} \mathbb{E}_{s \sim P(s'|s,a)} [(Q_{\text{main}\theta}(s, a) - \text{target}(s'))^2]_{\theta=\theta_k} \end{aligned}$$

At each train iteration:

The Q values shift to get closer to the target which shift the same way.

It is chasing a **non-stationary target**

SOLUTION: Use a Q_{target} network to generate the target (different from the Q_{main} network) .

Q_{target} network is used to generate the target: $\text{target} = [R(s, a, s') + \gamma \max_{a'} Q_{\text{target}}(s', a')]$

Q_{main} network is the network train and used to generate experiences.

Q_{target} 's weights is updated with Q_{main} 's weights from time to time

REMARKS

MODEL DIMENSION

The train model Q actually take only s as an input and generate all $Q(s, a)$ value:

$$\text{Pacman example: } Q : s \rightarrow [Q(s, a_1), Q(s, a_2), Q(s, a_3), Q(s, a_4)]$$

Hence during training, the target dataset is generate with fixed valued for not tested a and $target$ for tested a .

Pacman example:

experience: $[s, a_2, s', R(s, a_2, s')]$

input: s

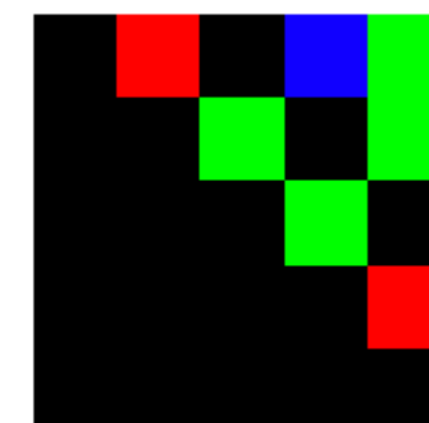
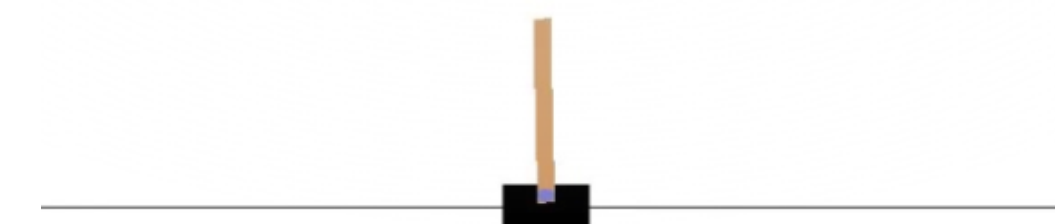
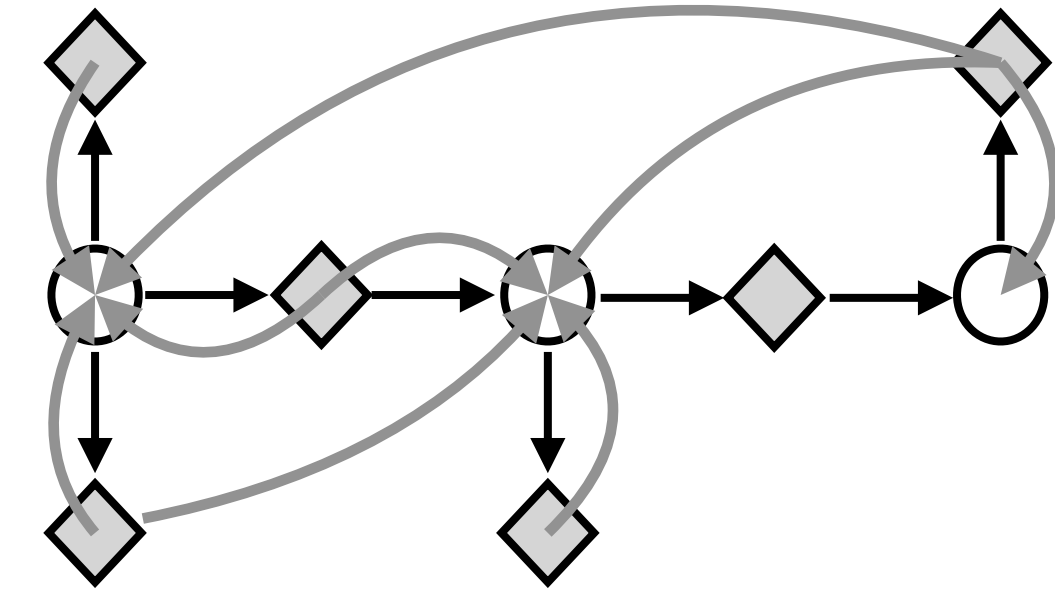
target:

$$\begin{bmatrix} Q_{main}(s, a_1) \\ R(s, a_2, s') + \gamma \cdot \max_{a'} Q_{target}(s', a') \\ Q_{main}(s, a_3) \\ Q_{main}(s, a_4) \end{bmatrix}$$

DEEP-Q NETWORK MNIH ET AL.[2015] - PSEUDO CODE

- Create Q_{main} and Q_{target} networks and initiate the replay memory M .
- Iterate while $max_number_episode$ or $goal_reward$ is not achieved:
 - Play one episode - Explore randomly or Exploit with Q_{main} according to $random_probability$.
 - Add the episode (*train data*) to M .
 - If we train more than $min_pre_trained_episode$
 - Start decreasing $random_probability$.
 - Every $train_frequency$:
 - Sample $batch_size$ experience from M .
 - Create target using Q_{target}
 - Train the Q_{main} network
 - Update Q_{target} weights with Q_{main} weights

- *Q_learning.ipynb* : Implement **Q-VALUE ITERATION ALGORITHM** and **Q-LEARNING ITERATION ALGORITHM** on the toy **MDP** of this presentation.
- *Deep_Q_learning_CartPole.ipynb* : Implement **Deep Q Learning** with and without **Replay memory buffer** on **CartPole**.
- *Deep_Q_learning_Gridworld.ipynb* : Implement **Deep Q Learning** with **separated target network** on **Gridworld**.



D3QN

DEEP-Q NETWORK IMPROVEMENT

Various improvement of the **DQN** algorithm lead to greater performance and stability.

- **Dueling DQN**. Wang et al. [2016]
- **Double DQN**. Van Hasselt et al [2016].
- **Prioritised Experience Replay**. Shaul et al. [2015]

DEEP-Q NETWORK IMPROVEMENT

Various improvement of the **DQN** algorithm lead to greater performance and stability.

- **Dueling DQN**. Wang et al. [2015]
- **Doube DQN**. Van Hasselt et al [2016].
- **Prioritised Experience Replay**. Shaul et al. [2015]

DUELING - WANG ET AL. [2016]

Main intuition

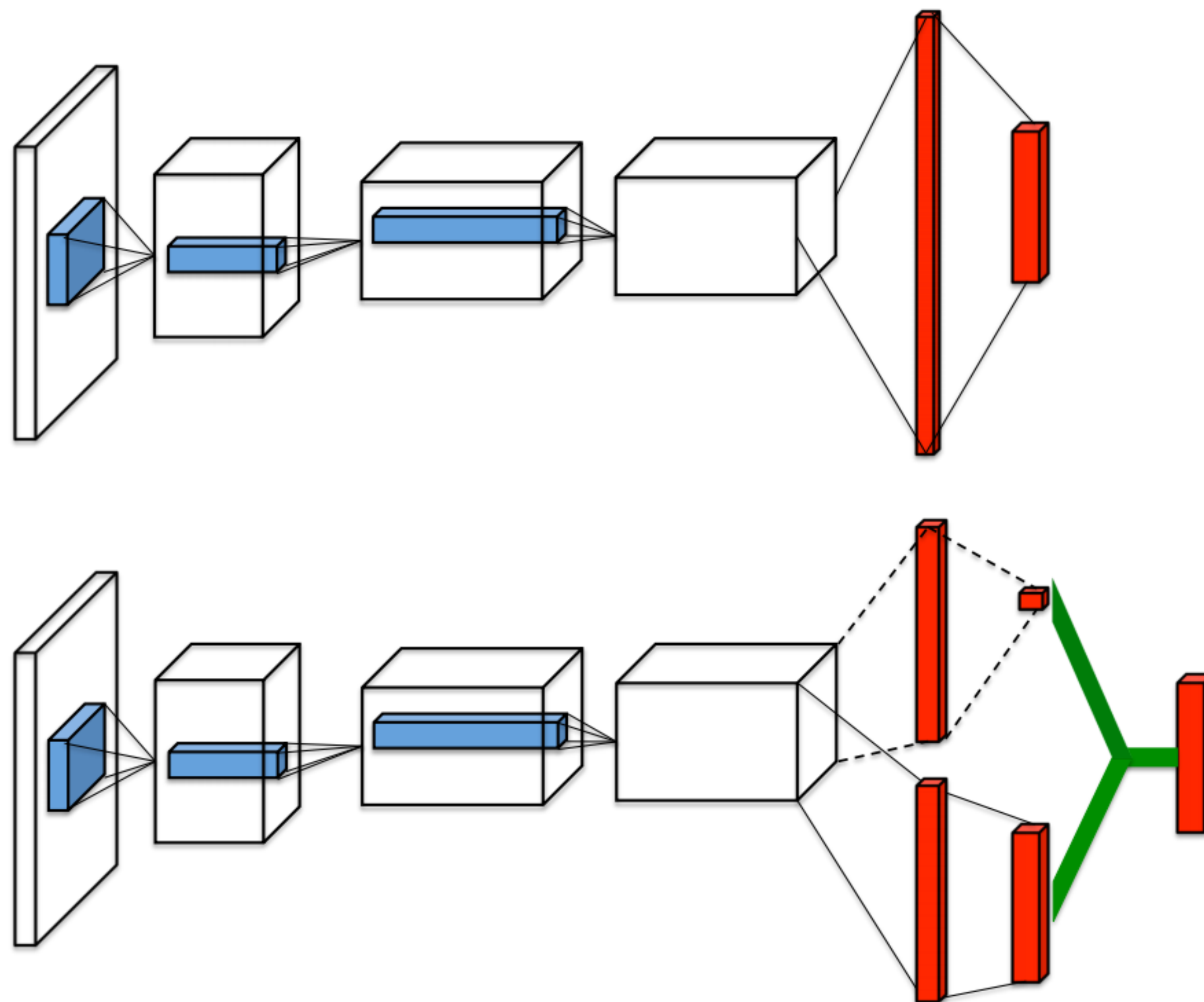
- **Q-value** represent how good it is to take an action a at a state s .
- It can be decomposed into two notions : $Q(s, a) = V(s) + A(s, a)$
 - The **value function** $V(s)$: *how good it is to be in this state?*
 - The **advantage function** $A(s, a)$: *how good it is to take this action in this state?*
- > It leads to poor performance.
- Let's force the Q value for the maximising action to equal V.

$$Q(s, a) = V(s) + (A(s, a) - \max_{a' \in |\mathcal{A}|} A(s, a'))$$

- Which is approximate in practice by:

$$Q(s, a) = V(s) + \left(A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right)$$

DUELING



DOUBLE DQN

Main intuition:

DQN overestimates the Q-values of the potential action to take

SEPARATED TARGET NETWORK

$$target = R(s, a, s') + \gamma \max_{a'} Q_{target}(s', a')$$

DOUBLE DQN Use Q_{main} to chose the action

$$target = R(s, a, s') + \gamma Q_{target}(s', \operatorname{argmax}_{a'}(Q_{main}(s', a')))$$

BLOG & CODE

<https://medium.com/@awjuliani/simple-reinforcement-learning-with-tensorflow-part-4-deep-q-networks-and-beyond-8438a3e2b8df>

<https://medium.com/emergent-future/simple-reinforcement-learning-with-tensorflow-part-0-q-learning-with-tables-and-neural-networks-d195264329d0>

<https://towardsdatascience.com/reinforcement-learning-tutorial-part-3-basic-deep-q-learning-186164c3bf4>

<https://sergioskar.github.io/Reinforcement\ learning/>

<http://www2.econ.iastate.edu/tesfatsi/RLUsersGuide.ICAC2005.pdf>

REFERENCES

Géron, A. (2017). Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts. *Tools, and Techniques to build intelligent systems*.

Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., ... & Petersen, S. (2015). Human-level control through deep reinforcement learning. *nature*, 518(7540), 529-533.

Schaul, T., Quan, J., Antonoglou, I., & Silver, D. (2015). Prioritized experience replay. arXiv preprint arXiv:1511.05952.

Van Hasselt, H., Guez, A., & Silver, D. (2015). Deep reinforcement learning with double q-learning. *arXiv preprint arXiv:1509.06461*.

Wang, Z., Schaul, T., Hessel, M., Hasselt, H., Lanctot, M., & Freitas, N. (2016, June). Dueling network architectures for deep reinforcement learning. In International conference on machine learning (pp. 1995-2003).