

# DEEP REINFORCEMENT LEARNING 1

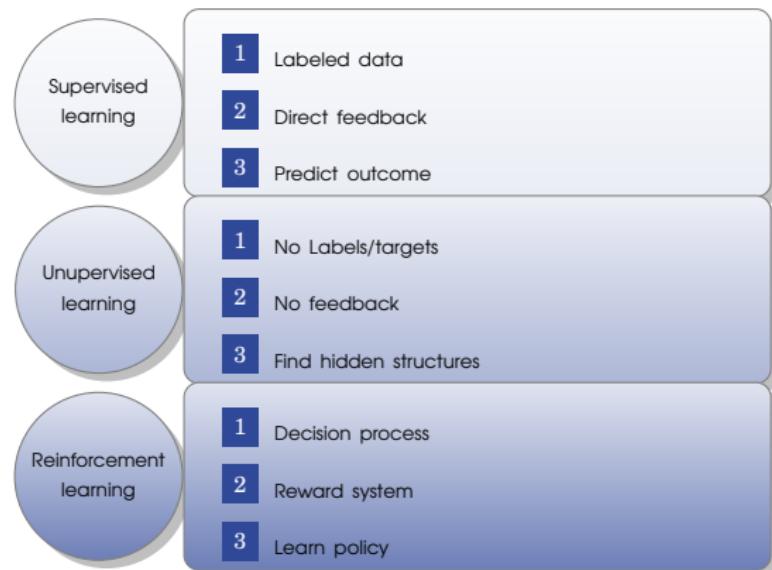
Vincent Barra  
LIMOS, UMR 6158 CNRS, Université Clermont Auvergne



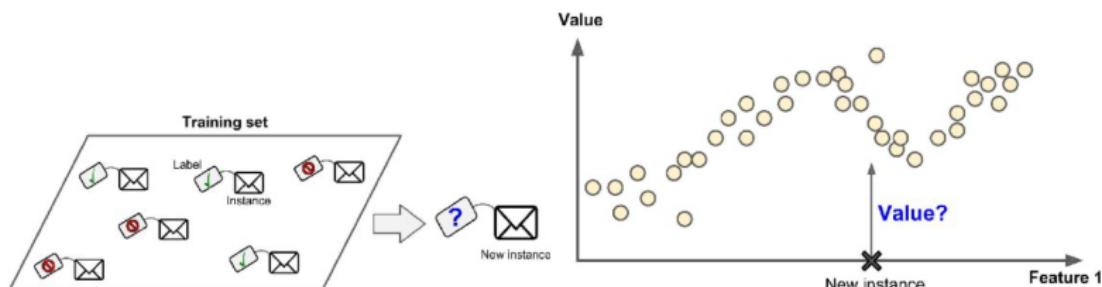
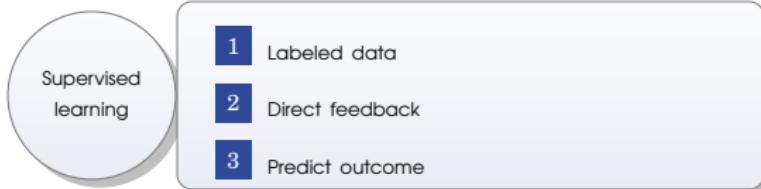
LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# A POSSIBLE TAXONOMY OF MACHINE LEARNING ALGORITHMS



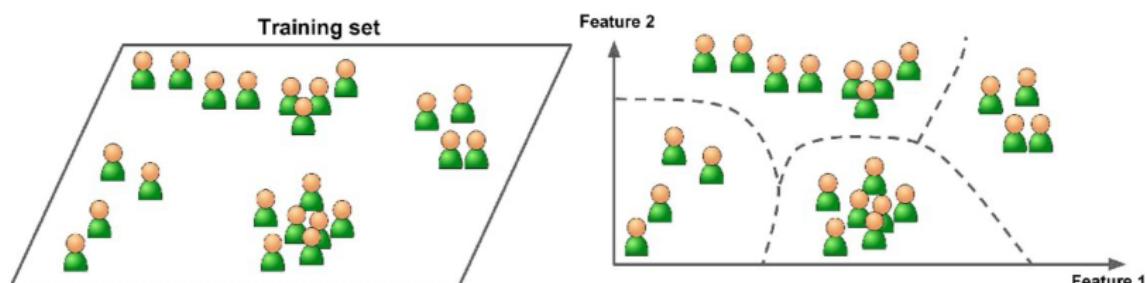
# SUPERVISED LEARNING



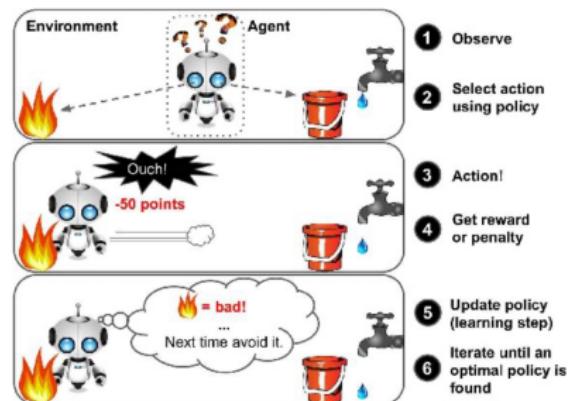
# UNSUPERVISED LEARNING



- 1 No Labels/targets
- 2 No feedback
- 3 Find hidden structures



# REINFORCEMENT LEARNING



# REINFORCEMENT LEARNING



Topic of the two next lectures

## AGENT: TAKES ACTIONS

Agent



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# ENVIRONMENT: THE "WORLD"



Environment

Agent



# ACTION SPACE



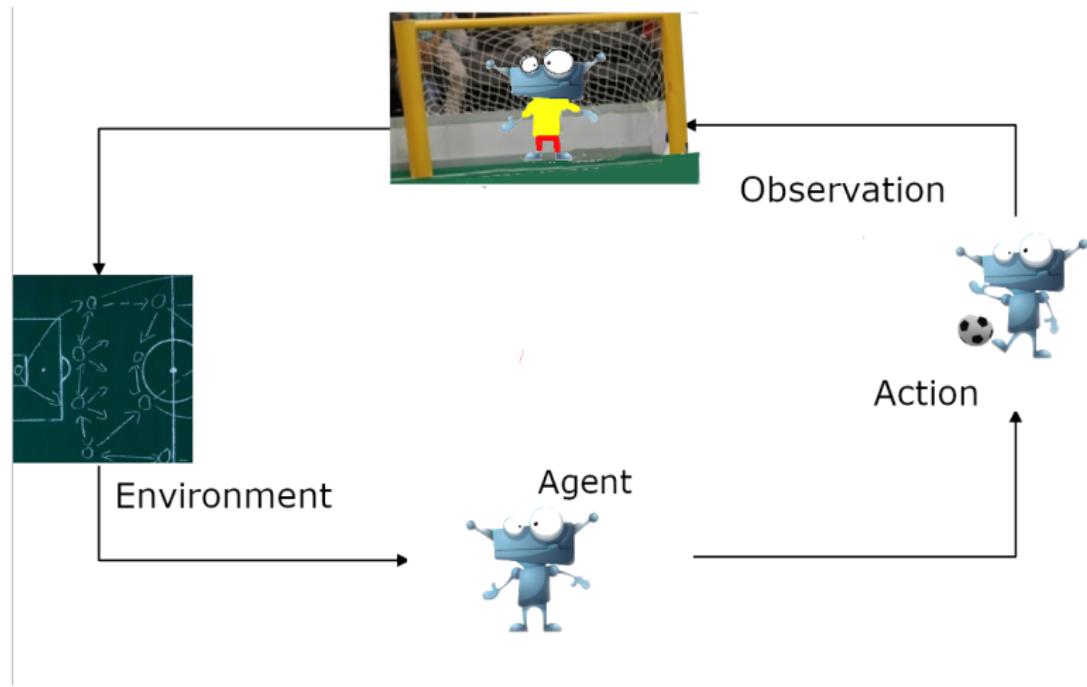
Environment



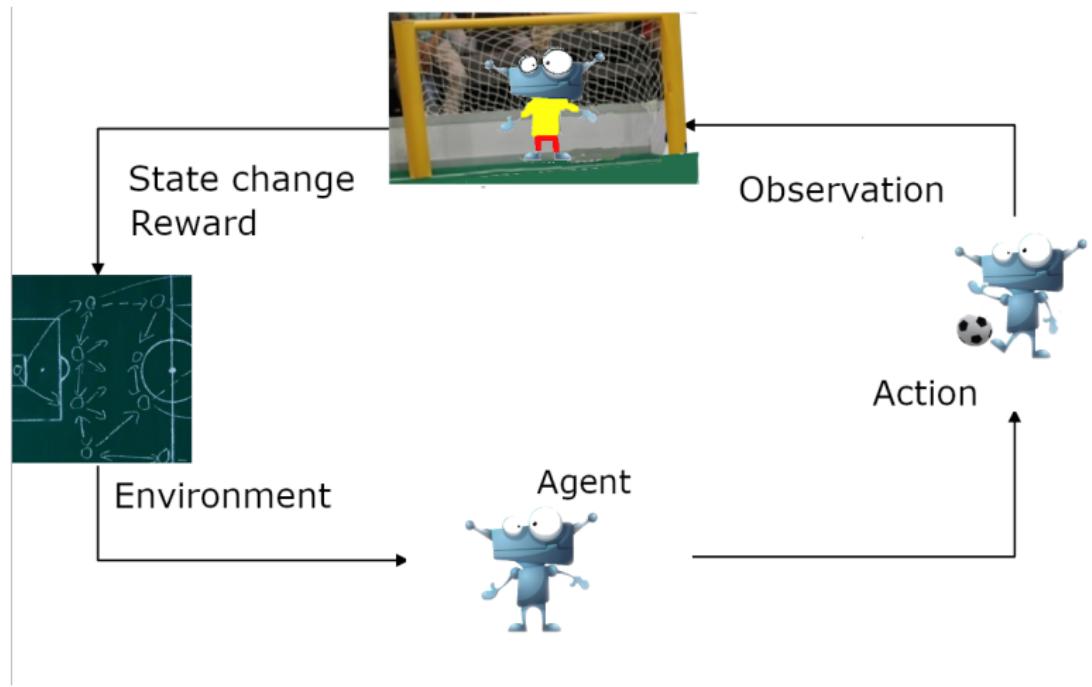
Action



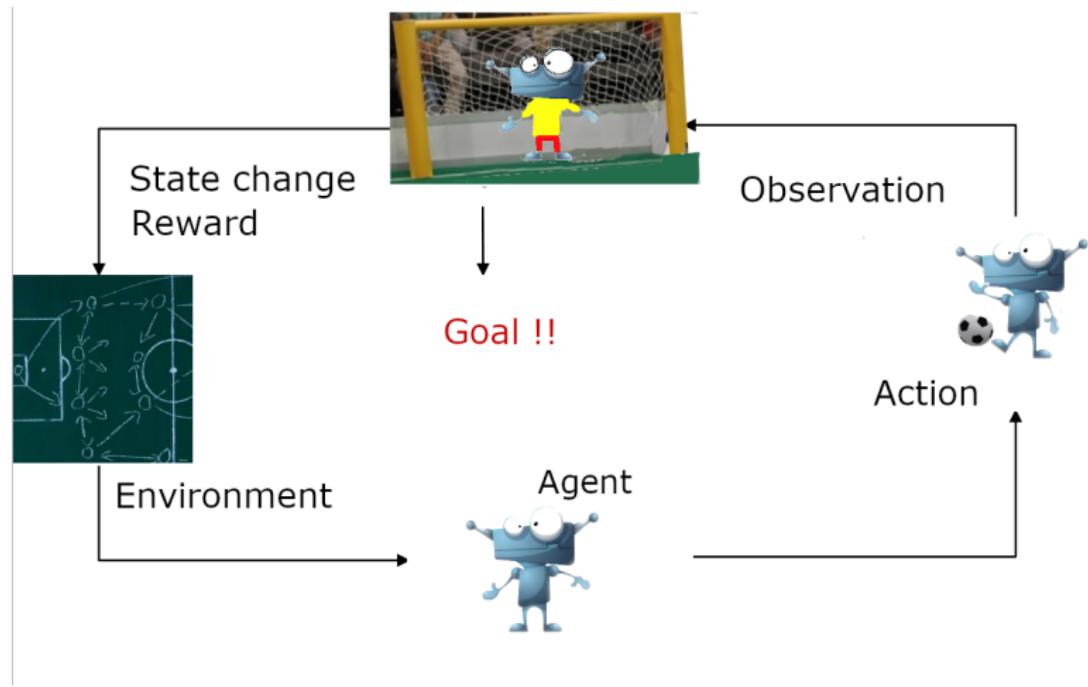
# CHANGES IN THE ENVIRONMENT



# STATE AND REWARD



# OBJECTIVE FUNCTION



# NOTATIONS

## Space, Action, Reward

For each timestep  $t$ :

- ▶  $a_t \in A$ : action
- ▶  $s_t \in S$ : state
- ▶  $r_t : S \times A \times S \rightarrow \mathbb{R}^+$ : reward obtained when arriving at  $s_{t+1}$  from  $s_t$ , taking action  $a_t$

$\forall t$ , the agent receives a state representation  $s_t$ , selects an action  $a_t$ , and receives a reward  $r_t$  after which it will find itself in a new state  $s_{t+1}$ .

Trajectory:  $\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$

# NOTATIONS

## Space, Action, Reward

For each timestep  $t$ :

- ▶  $a_t \in A$ : action
- ▶  $s_t \in S$ : state
- ▶  $r_t : S \times A \times S \rightarrow \mathbb{R}^+$ : reward obtained when arriving at  $s_{t+1}$  from  $s_t$ , taking action  $a_t$

$\forall t$ , the agent receives a state representation  $s_t$ , selects an action  $a_t$ , and receives a reward  $r_t$  after which it will find itself in a new state  $s_{t+1}$ .

$$\text{Trajectory: } \tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$$

## Hypothesis

Markovian environment:

- ▶  $P(s_{t+1}|s_t, a_t, \dots, s_0, a_0) = P(s_{t+1}|s_t, a_t)$
- ▶  $P(r_t|s_t, a_t, \dots, s_0, a_0) = P(r_t|s_t, a_t)$

# NOTATIONS

## Objective

Total reward:  $R_t = \sum_{i=t}^{+\infty} r_i$

Total discounted reward:  $R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$

$0 < \gamma < 1$ : discount factor, controls the trade-off between immediate and long-term rewards.

# NOTATIONS

## Objective

Total reward:  $R_t = \sum_{i=t}^{+\infty} r_i$

Total discounted reward:  $R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$

$0 < \gamma < 1$ : discount factor, controls the trade-off between immediate and long-term rewards.

Objective: maximise the long-term rewards.

$$R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$$

# VALUE FUNCTION

$$R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$$

- ▶ Estimate how good/bad it is for an agent to be in a particular state
- ▶ Goodness of state  $s_t$  defined in terms of future rewards that can be expected by from  $s_t$ .



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# VALUE FUNCTION

$$R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$$

- ▶ Estimate how good/bad it is for an agent to be in a particular state
- ▶ Goodness of state  $s_t$  defined in terms of future rewards that can be expected by from  $s_t$ .

Policy function  $\pi : S \rightarrow A$



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# VALUE FUNCTION

$$R_t = \sum_{i=t}^{+\infty} \gamma^i r_i$$

- ▶ Estimate how good/bad it is for an agent to be in a particular state
- ▶ Goodness of state  $s_t$  defined in terms of future rewards that can be expected by from  $s_t$ .

Policy function  $\pi : S \rightarrow A$

- ▶ state-value function  $V^\pi(s) = \mathbb{E}(R_t | s_t = s, \pi)$
- ▶ state-action function  $Q^\pi(s_t, a_t) = \mathbb{E}(R_t | s_t, a_t, \pi)$



LIMOS  
LABORATOIRE DINFORMATIQUE,  
DE MODELISATION ET OPTIMISATION DES SYSTÈMES



## VALUE FUNCTION

- ▶ The environment is unknown.
- ▶  $V^\pi, Q^\pi \sim$  knowledge representation of an agent, accumulating knowledge through its interactions with the environment, stored by updating its value function.
- ▶ If a value function is accurate an agent will know everything he needs to know for interacting with an environment.

# $Q$ FUNCTION

## $Q$ function

$$Q(s_t, a_t) = \mathbb{E}(R_t | s_t, a_t)$$

Captures the expected reward of a agent being in state  $s_t$  at time  $t$ , and taking the action  $a_t$ .

Question: how to decide actions to take, given a  $Q$  function ?

# POLICY

## Policy

- ▶ Policy function  $\pi : S \rightarrow A$
- ▶ Infer the best action to take when being in state  $s \in S$

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



LIMOS  
LABORATOIRE DINFORMATIQUE,  
DE MODELISATION ET OPTIMISATION DES SYSTÈMES



# POLICY

## Policy

- ▶ Policy function  $\pi : S \rightarrow A$
- ▶ Infer the best action to take when being in state  $s \in S$

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

## Value Learning

Find  $Q(s, a)$   
 $a = \arg \max_{a' \in A} Q(s, a')$



LIMOS  
LABORATOIRE DINFORMATIQUE,  
DE MODELISATION ET OPTIMISATION DES SYSTÈMES



# POLICY

## Policy

- ▶ Policy function  $\pi : S \rightarrow A$
- ▶ Infer the best action to take when being in state  $s \in S$

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

## Value Learning

Find  $Q(s, a)$   
 $a = \arg \max_{a' \in A} Q(s, a')$

## Policy Learning

Find  $\pi(s)$   
Sample  $a \sim \pi(s)$



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# *Q*-LEARNING

In classical reinforcement learning, problem related to Markov Decision Processes (out of scope).

## *Q*-learning algorithm

- 1 build a memory table to store  $Q$ -values for all possible combinations of  $s \in S$  and  $a \in A$
- 2 as the algorithm evolves, the table is updated...
- 3 ... until convergence...
- 4 ... using some tricks ( $\epsilon$ -greedy policy)

## *Q*-LEARNING

In classical reinforcement learning, problem related to Markov Decision Processes (out of scope).

### *Q*-learning algorithm

- 1 build a memory table to store *Q*-values for all possible combinations of  $s \in S$  and  $a \in A$
- 2 as the algorithm evolves, the table is updated...
- 3 ... until convergence...
- 4 ... using some tricks ( $\epsilon$ -greedy policy)

### *Q*-learning does not scale to large problems

If  $A$  and/or  $S$  are large.

Pacman:

- ▶ 240 pellets can be present/absent/eaten
- ▶  $2^{240} \approx 10^{72}$  configurations
- ▶ + ghost and Pacman positions....



# APPROXIMATING $Q$ FUNCTIONS

Aim: build an approximation of the  $Q$ -function

- ▶ Approximate  $Q$ -learning: train  $Q_W$  that will generalize the memory table of the  $Q$ -learning algorithm

# APPROXIMATING $Q$ FUNCTIONS

Aim: build an approximation of the  $Q$ -function

- ▶ Approximate  $Q$ -learning: train  $Q_W$  that will generalize the memory table of the  $Q$ -learning algorithm
- ▶ **Deep  $Q$ -learning:** uses a Deep (convolutional) Network as the  $Q_W$  function

# APPROXIMATING $Q$ FUNCTIONS

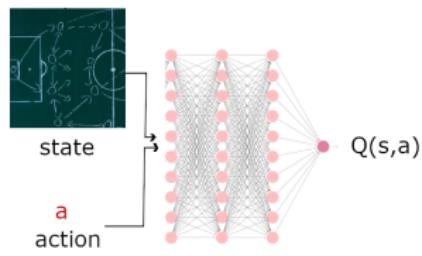
Aim: build an approximation of the  $Q$ -function

- ▶ Approximate  $Q$ -learning: train  $Q_W$  that will generalize the memory table of the  $Q$ -learning algorithm
- ▶ **Deep  $Q$ -learning:** uses a Deep (convolutional) Network as the  $Q_W$  function

$V^\pi$  can be approximated in the same way

# USING DEEP NETWORKS TO MODEL $Q$ FUNCTIONS

## Deep Q-Networks (DQN)

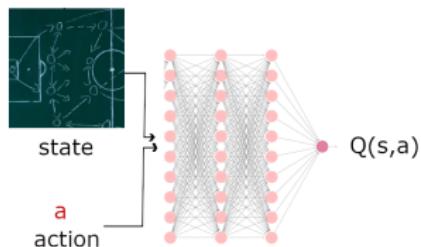


Input              Agent              Output

$$S \times A \rightarrow Q(s, a)$$

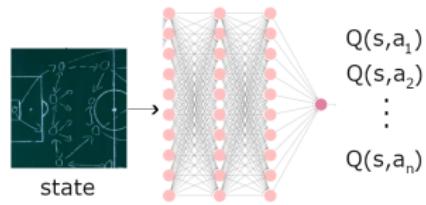
# USING DEEP NETWORKS TO MODEL $Q$ FUNCTIONS

## Deep Q-Networks (DQN)



Input      Agent      Output

$$S \times A \rightarrow Q(s, a)$$



Input      Agent      Output

$$S \rightarrow Q(s, a)$$

# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET OPTIMISATION DES SYSTÈMES



# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

Let put it in equation

# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

$$\underbrace{(r + \gamma \max_{a'} Q(s, a'))}_{\text{target}}$$



LIMOS  
INSTITUT  
DE  
MODÉLISATION  
ET  
OPTIMISATION  
DES  
SYSTÈMES



# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

$$\underbrace{(r + \gamma \max_{a'} Q(s, a'))}_{\text{target}} \underbrace{Q(s, a)}_{\text{predicted}}$$

# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

$$\underbrace{(r + \gamma \max_{a'} Q(s, a'))}_{\text{target}} - \underbrace{Q(s, a)}_{\text{predicted}}$$



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

$$\mathcal{L} = \mathbb{E}[\underbrace{\|(r + \gamma \max_{a'} Q(s, a'))}_{\text{target}} - \underbrace{Q(s, a)}_{\text{predicted}}\|^2] \quad Q\text{-loss}$$

# TRAINING

## Strategy

- ▶ take all the best actions
- ⇒ Maximize expected reward
- ⇒ Train the agent

$$\mathcal{L} = \mathbb{E}[\underbrace{\|(r + \gamma \max_{a'} Q(s, a'))}_{\text{target}} - \underbrace{Q(s, a)}_{\text{predicted}}\|^2] \quad Q\text{-loss}$$

And so

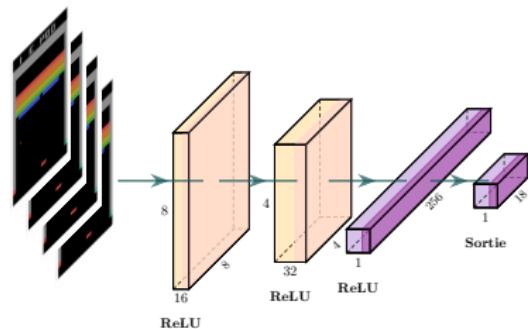
- 1 Use a (deep) neural network to learn  $Q$  function and use it to infer the optimal policy  $a_t = \pi^*(s_t)$ .
- 2 Send action back to the environment:  $s_t \xrightarrow{a_t} s_{t+1}$



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



## EXAMPLE - DQN



Deep Q-Network (DQN) on 49 ATARI games (Mnih et al., 2015).

- ▶ Input  $4 \times 84 \times 84$  images
- ▶ Conv(16)-ReLU
- ▶ Conv(32)-ReLU
- ▶ FCC(256)
- ▶ Output: FCC(18): joystick positions

## EXAMPLE - DQN

### Add-ons

- ▶ Experience replay: previous agent experiences  $(s_t, a_t, s_{t+1}, r_t)$  are stored in a replay memory and randomly drawn during training
  - ⇒ suppress the correlation between training data
- ▶ Use of a target network  $Q_{\mathbf{W}^-}$ , keeping its own parameters and frequently updating them
  - ⇒ reduce the correlations between  $Q$ -values and the target
$$r + \gamma \max_{a'} Q(s, a')$$

$$\mathcal{L} = \mathbb{E} \left[ \|(r + \gamma \max_{a'} Q_{\mathbf{W}_t^-}(s, a')) - Q_{\mathbf{W}_t}(s, a)\|^2 \right]$$

# EXAMPLE - DQN

Episode = a complete run of game;  $S$ : target network update step;  $T$ : length of a game :  $x_t$  input at time  $t$  ;

$\phi$ : preprocessing function

**begin**

**Input** : Image of the game; score,  $T$ ,  $S$

**Output** :  $Q$  value

    Initialize Experience replay memory  $D$

    Initialize  $Q$  with random weights  $W$

    Initialize target  $\hat{Q}$  with weights  $W^- = W$

**foreach** episode **do**

        Initialize a sequence  $s_1 = \{x_1\}$ .  $\phi_1 = \phi(s_1)$   
        **for**  $t=1$  **to**  $T$  **do**

            Select an action  $a_t = \begin{cases} \text{randomly} & \text{with probability } \epsilon \\ \operatorname{Arg} \max_a Q(\phi(s_t), a, \theta) & \text{otherwise} \end{cases}$

            Execute  $a_t$  and observe  $r_t$  and  $x_{t+1}$

$s_{t+1} = s_t, a_t, x_{t+1}$

$\phi_{t+1} = \phi(s_{t+1})$

$D \leftarrow D \cup (\phi_t, a_t, r_t, \phi_{t+1})$ ;

            /\* Experience replay

            Sample a batch  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$ ;

$y_j = \begin{cases} r_j & \text{if episode finishes in } j + 1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a', \theta^-) & \text{otherwise} \end{cases}$

            Gradient descent step w.r.t.  $\theta$  on  $(y_j - Q(\phi_j, a_j, \theta))^2$

            /\* Update target network

            Each  $S$  step do  $\theta^- = \theta$

        \*/

**end**

**end**

**end**



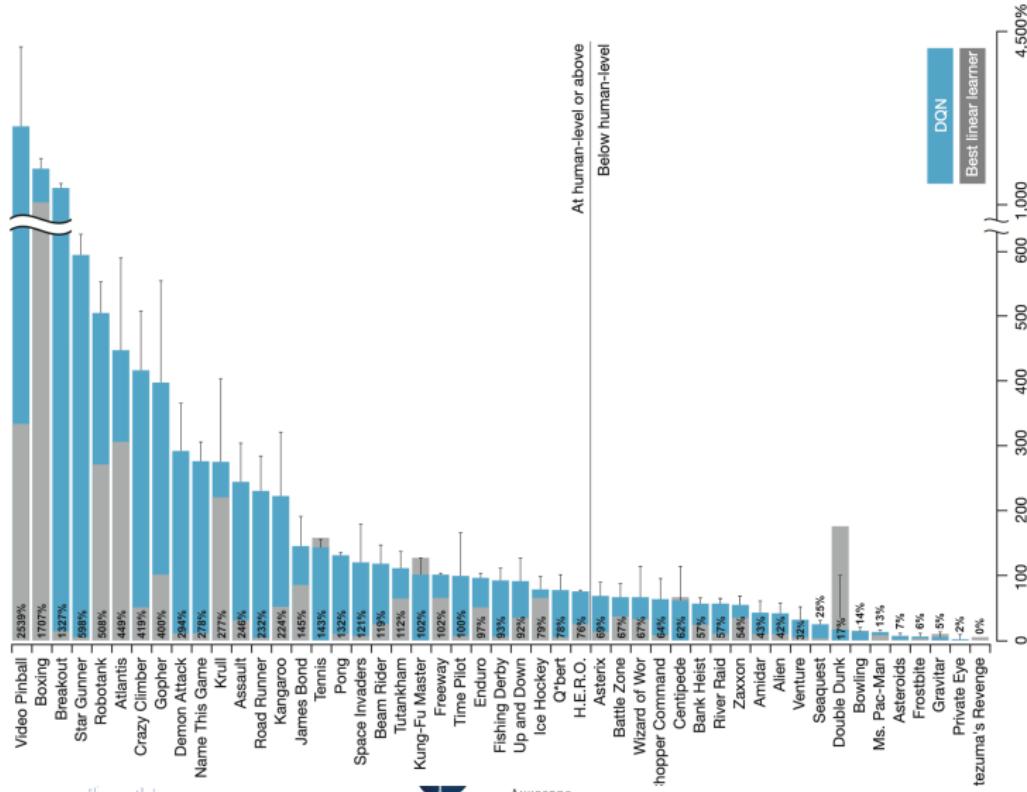
LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



## EXAMPLE - DQN



# EXAMPLE - DQN

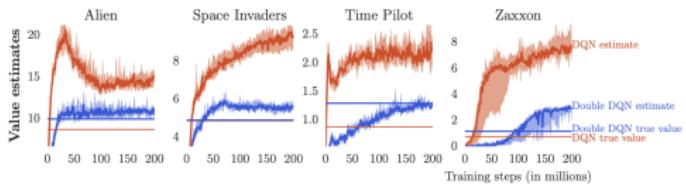


source: Mnih et al., 2015

# IMPROVEMENTS

## Improvements

- ▶ Double DQN (2016)

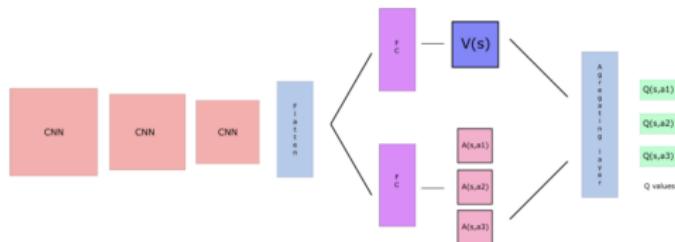


- ▶ Q Learning has been shown to overestimate its targets, because it uses a single estimator for estimating the Q values and choosing the maximum over actions.
- ▶ Solution: to overcome this, Double Q Learning uses a double estimator technique:
  - The target network is used for estimation
  - The online network for choosing the greedy action

# IMPROVEMENTS

## Improvements

- ▶ Double DQN (2016)
- ▶ Dueling DQN (2016)



Dueling architecture:

- ▶ Decompose the  $Q$  value into a state-dependent part and a state-action dependent part
- ▶ This allows to share the state-dependent estimations (good for actions that are less chosen), focuses in estimating the state-action dependent part (good for relative ranking of the actions).

# PROS AND CONS

## Flexibility

- ▶  $Q$  determines  $\pi^*$  in a deterministic way
- ⇒ cannot manage stochastic policies

## Complexity

- ▶ OK if  $|A|$  small, and  $A$  discrete
- ▶ KO if  $A$  is continuous



LIMOS  
LABORATOIRE INFORMATIQUE,  
MODELISATION ET OPTIMISATION DES SYSTÈMES



# PROS AND CONS

## Flexibility

- ▶  $Q$  determines  $\pi^*$  in a deterministic way
- ⇒ cannot manage stochastic policies

## Complexity

- ▶ OK if  $|A|$  small, and  $A$  discrete
- ▶ KO if  $A$  is continuous

See next lecture

Policy gradient methods

# GYM



Gym: toolkit for developing and comparing reinforcement learning algorithms



LABORATOIRE D'INFORMATIQUE,  
DE MODÉLISATION ET D'OPTIMISATION DES SYSTÈMES



# ENVIRONMENTS



BipedalWalker-v2  
Train a bipedal robot to walk.



BipedalWalkerHardcore-v2  
Train a bipedal robot to walk over rough terrain.



CarRacing-v0  
Race a car around a track.



LunarLander-v2  
Navigate a lander to its landing pad.



LunarLanderContinuous-v2  
Navigate a lander to its landing pad.



Acrobot-v1  
Swing up a two-link robot.



CartPole-v1  
Balance a pole on a cart.



MountainCar-v0  
Drive up a big hill.

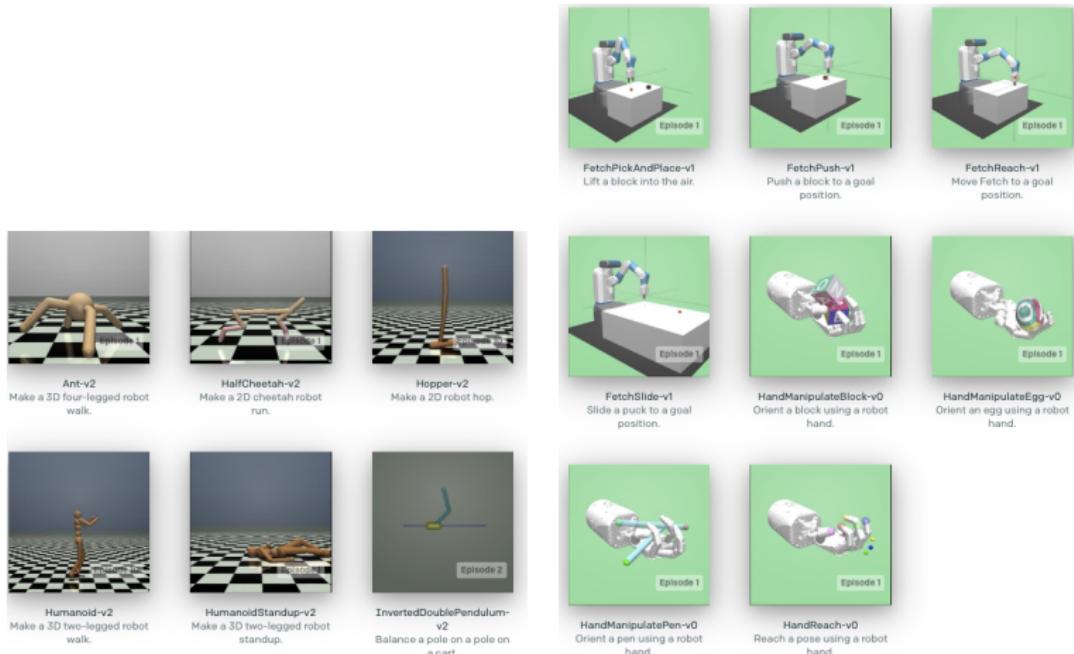


MountainCarContinuous-v0  
Drive up a big hill with continuous control.

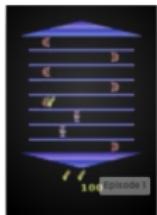


Pendulum-v0  
Swing up a pendulum.

# ENVIRONMENTS



# ENVIRONMENTS



**Asterix-v0**  
Maximize score in the game  
Asterix, with screen Images  
as input



**Asteroids-ram-v0**  
Maximize score in the game  
Asteroids, with RAM as  
input



**Asteroids-v0**  
Maximize score in the game  
Asteroids, with screen  
images as input



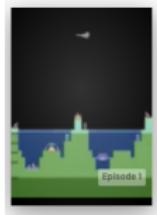
**BankHeist-v0**  
Maximize score in the game  
BankHeist, with screen  
images as input



**BattleZone-ram-v0**  
Maximize score in the game  
BattleZone, with RAM as  
input



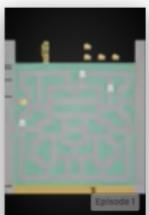
**BattleZone-v0**  
Maximize score in the game  
BattleZone, with screen  
images as input



**Atlantis-ram-v0**  
Maximize score in the game  
Atlantis, with RAM as input



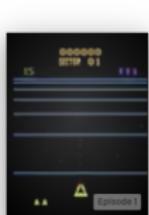
**Atlantis-v0**  
Maximize score in the game  
Atlantis, with screen Images  
as input



**BankHeist-ram-v0**  
Maximize score in the game  
BankHeist, with RAM as  
input



**BeamRider-ram-v0**  
Maximize score in the game  
BeamRider, with RAM as  
input

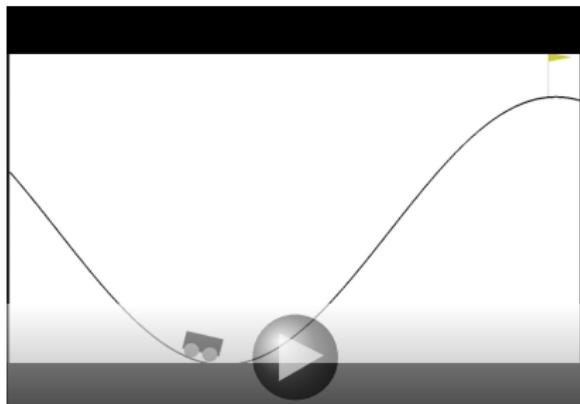


**BeamRider-v0**  
Maximize score in the game  
BeamRider, with screen  
images as input



**Berzerk-ram-v0**  
Maximize score in the game  
Berzerk, with RAM as input

# CODE



```
import gym
# Game to play
env = gym.make("MountainCar-v0")

# initial observation of the environment
observation = env.reset()

# A
print(env.action_space)

# S
print(env.observation_space)

for t in range(1000):

    env.render()

    # The agent chooses an action a(t)
    action = env.action_space.sample()

    # return s(t+1), r(t+1), time to reset, debugging info
    observation, reward, done, info = env.step(action)

    if done:
        # episode terminated
        observation = env.reset()
        env.close()
```