

## MSc on Intelligent Critical Infrastructure Systems

# Machine Learning

## Lecture 12: Reinforcement Learning (part II)

**Kleanthis Malialis**

Research Associate

KIOS Research and Innovation Center of Excellence

University of Cyprus

FUNDED BY:



# Course outline

- **Week 1**
  - Introduction and Preliminaries
- **Week 2**
  - Linear Regression
  - Regularisation, Logistic Regression, SVMs
- **Week 3**
  - Neural Networks and Deep Learning
- **Week 4**
  - Feature Engineering and Evaluation
  - Online Learning
- **Week 5**
  - Unsupervised Learning
- **Week 6**
  - Reinforcement Learning
- **Week 7**
  - Revision
- **Exam**
  - 25/05 – 01/06



# SARSA vs. Q-Learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$$

**SARSA (state-action-reward-state-action): on-policy control** compared to Q-Learning, which is an **off-policy control** algorithm.

# Convergence criteria

- The environment's states are Markov states.
- **An agent visits all state-action pairs infinitely often.**
- The learning rate reduces to zero.
- The exploration rate reduces to zero (for SARSA).
- Bounded rewards.



# Expected SARSA

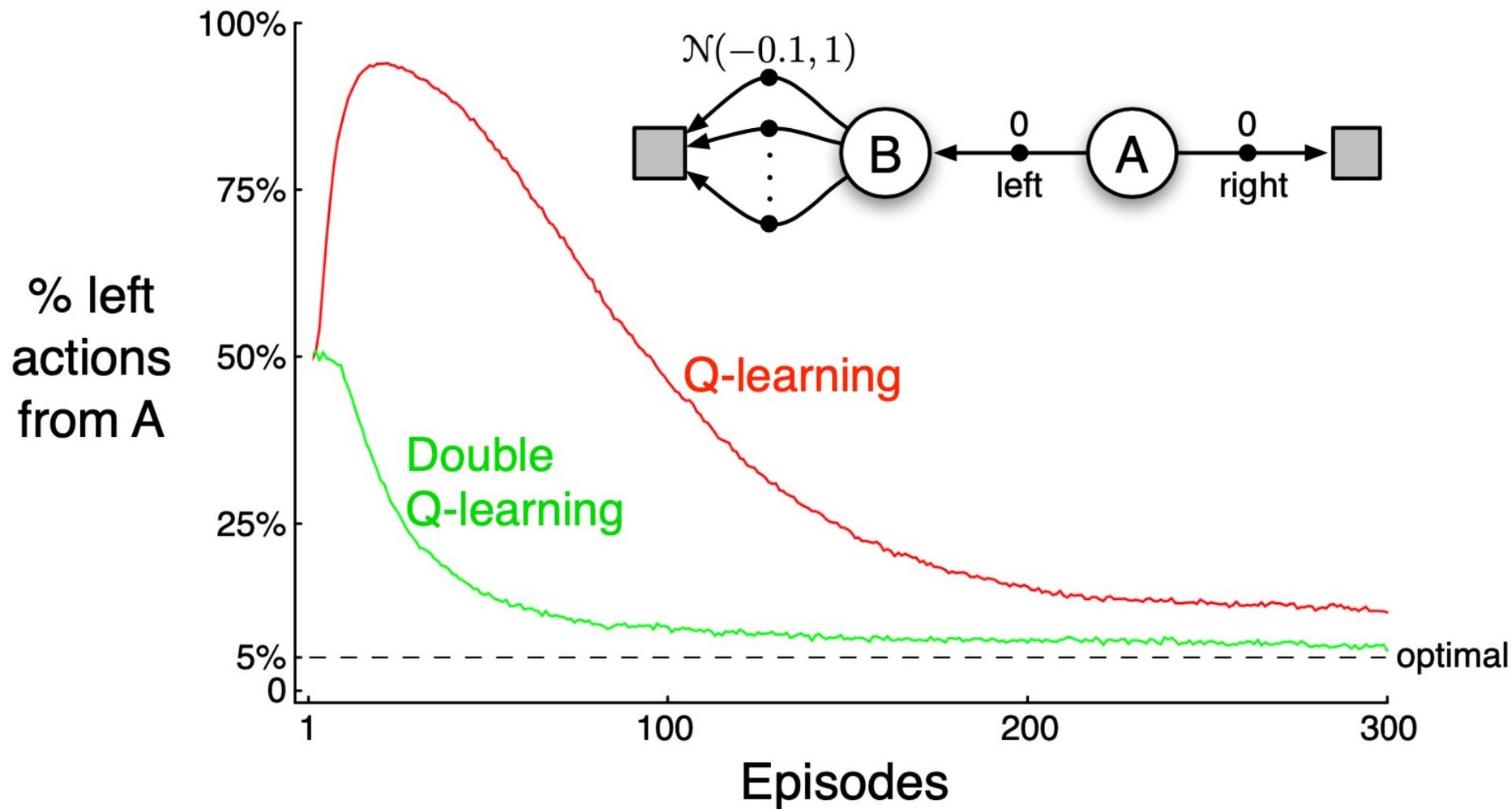
- Consider Q-learning except that instead of the maximum over next state–action pairs it uses the **expected value**, taking into account how likely each action is under the current policy.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- It generally outperforms SARSA.
- It can be an on-policy (like SARSA) or off-policy, for instance:
  - Use 0.1-greedy for exploration, and 0.01-greedy for the target.
  - Use  $\epsilon$ -greedy for exploration, and Greedy for the target. **Q-learning!**



# Maximisation bias



# Double Q-learning

- Learn **two independent** Q functions, on different sets of experience.
- The first Q-function is used for **action selection**.
- The second Q-function is used for **action evaluation**.

with 0.5 probability

$$Q_1(S_t, A_t) = Q_1(S_t, A_t) + \alpha [R_{t+1} + \gamma Q_2(S_{t+1}, \mathbf{argmax}_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t)]$$

else:

$$Q_2(S_t, A_t) = Q_2(S_t, A_t) + \alpha [R_{t+1} + \gamma Q_1(S_{t+1}, \mathbf{argmax}_a Q_2(S_{t+1}, a)) - Q_2(S_t, A_t)]$$

**Question:**

How do the (i) memory requirements, and (ii) computational time compare to the original Q-learning?



# Tabular solution methods

- So far, we have represented the value function  $Q$  as a **lookup table** with an entry for every state-action pairs  $(s, a)$
- Recall that one of the convergence criteria of Q-learning is for an agent to visit all state-action pairs infinitely often.
- **Problems:**
  - In some cases, it is impossible to satisfy the above criterion, e.g., the game of Go (19x91) has  $10^{170}$  states.
  - Some domains require continuous state and/or action space (e.g., mobile robot).
  - Too many states and/or actions to store in memory.
  - Too slow to learn the value for each state-action.



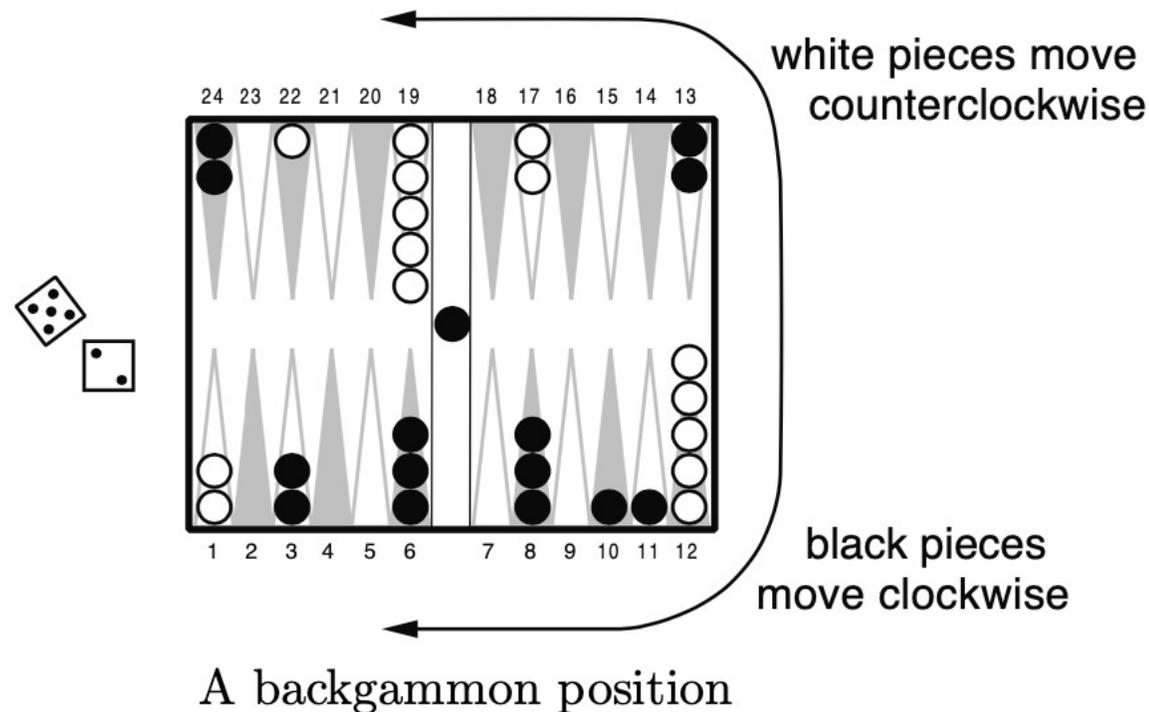


# APPROXIMATE SOLUTION METHODS

# IBM: Mastering the game of Backgammon (1995)



- One of the earliest and most impressive applications of RL to date is **TD-Gammon** by Gerald Tesauro.
- The algorithm combined TD-learning with non-linear function approximation using a multi-layered neural network.



# DeepMind: Human-level control of Atari games (2015)

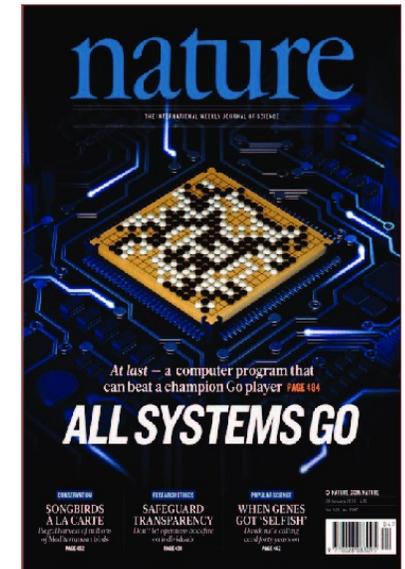
- DeepMind (U.K.) was bought by Google for \$500M in 2014!
- The **Deep Q-Network (DQN)** agent, receiving only the pixels and the game score as inputs, was able to achieve a level comparable to that of a professional human games tester across a set of 49 games, using the same algorithm, network architecture and hyperparameters.



[Paper](#), [video](#)

# DeepMind: Mastering the game of Go (2016)

- **AlphaGo** is the first computer program to defeat a professional human Go player, the first to defeat a Go world champion (Lee Sedol), and is arguably the strongest Go player in history.
- **Next generations:** AlphaGo Zero, AlphaZero, MuZero.

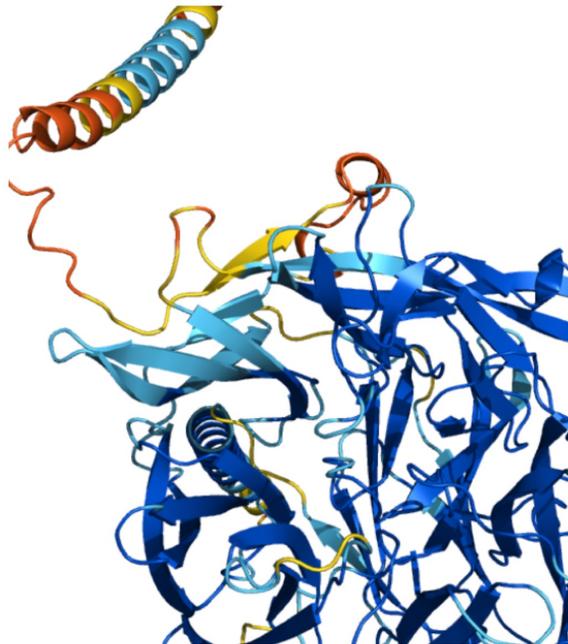


[Paper, documentary](#)

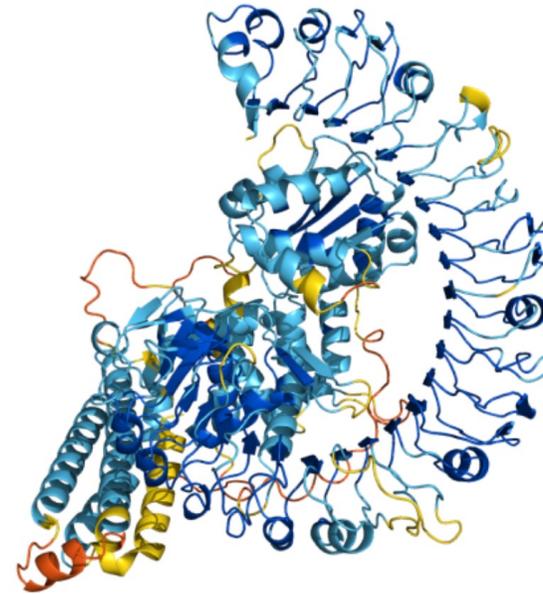
# DeepMind: Protein structure prediction (2021)



- **AlphaFold** can accurately predict 3D models of protein structures and has the potential to accelerate research in every field of biology.
- A new Alphabet company was launched (Isomorphic labs) to re-imagine the drug discovery process.



Q8I3H7: May protect the malaria parasite against attack by the immune system. Mean pLDDT 85.57.



Q8W3K0: A potential plant disease resistance protein. Mean pLDDT 82.24.

# Sony AI: Mastering Gran Turismo (2022)

- **Gran Turismo Sophy** is a racing RL agent that has learned to master the PlayStation game and driving simulator, Gran Turismo (GT) Sport.



[Paper](#), [video](#)

# RL with function approximation



- Use function approximation to estimate the value function, and generalise from seen states to unseen states.
- The approximate value function is not represented as a table, but as a parameterized functional form.
- Approximate solution methods:
  - Linear methods
  - Deep RL

# Linear methods

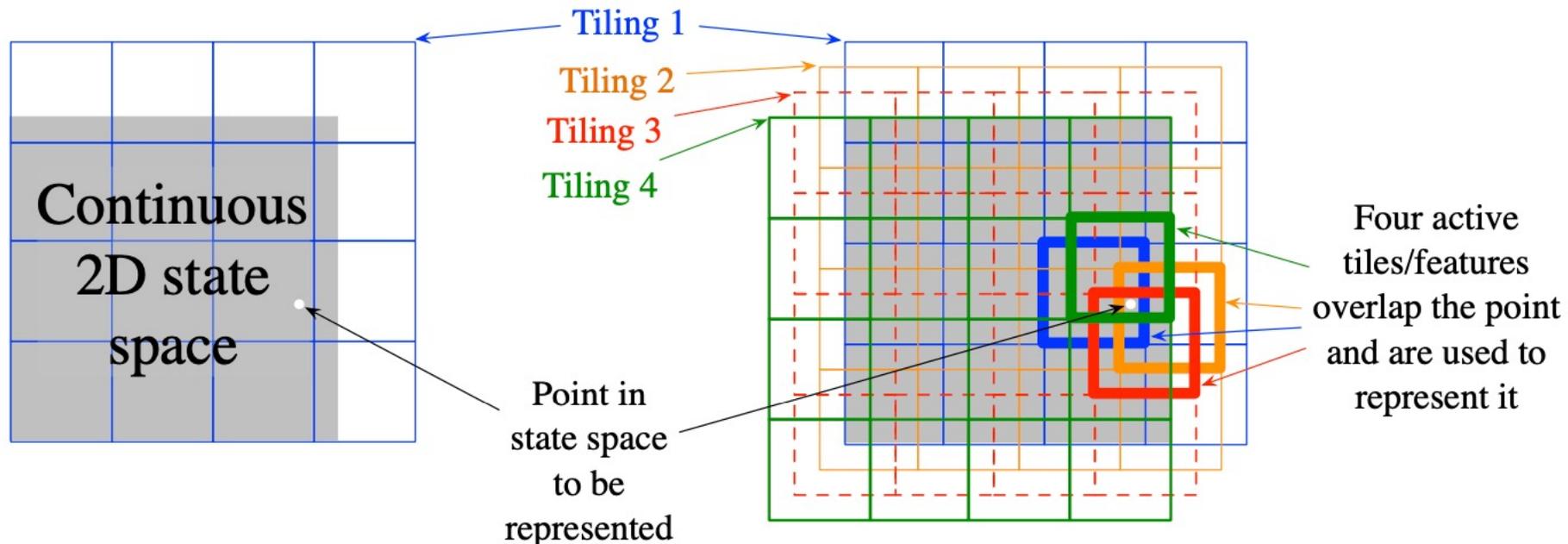


- The approximate function  $v(s; w)$  is a linear function of the weight vector  $w = (w_1, w_2, \dots, w_d)$ .
- Corresponding to each state  $s$ , there is a real-valued vector  $x(s) = (x_1(s), x_2(s), \dots, x_d(s))$ .
- Linear methods approximate the state value function as follows:

$$v(s; w) = w^T x(s) = \sum_{i=1}^d w_i x_i(s)$$

# Tile Coding

- Tile Coding partitions the state space into **tilings**, each tiling is further partitioned into **tiles** where state feature values are grouped into.
- The number of tilings indicate the degree of **resolution**.
- The size / shape of tiles indicate the nature of **generalisation**.



# Tile Coding (2)



- If the state is inside a tile, then the corresponding feature has the value 1, otherwise the feature is 0. This kind of 0/1-valued feature is called a **binary feature**.

- Tile Coding activates  $m < d$  tiles:  $TC(s) \rightarrow (x_1, x_2, \dots, x_m)$

- The Q-value is now calculated as:  $Q(s, a) = Q(TC(s), a) = \sum_{i=1}^m Q(x_i, a)$

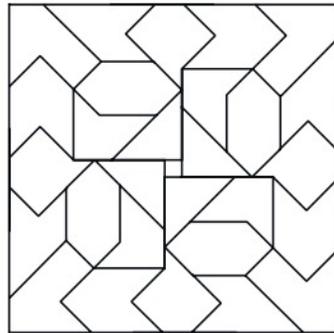
- The update rule is then:

$$\Delta Q = r + \gamma \max_a Q(s', a) - Q(s, a)$$

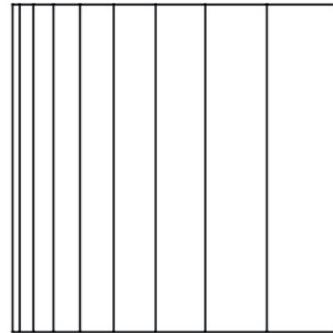
$$Q(x_i, a) = Q(x_i, a) + \frac{\alpha}{m} \Delta Q$$

# Tile Coding (3)

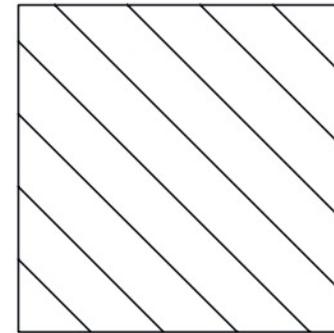
- Tilings need not be grids, they can be arbitrarily shaped and non-uniform. However, these are rarely used in practise as Tile Coding is flexible and computationally efficient.



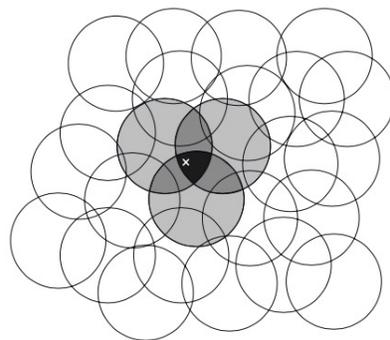
Irregular



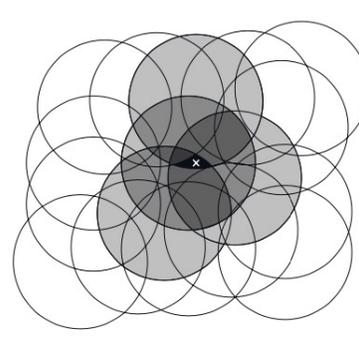
Log stripes



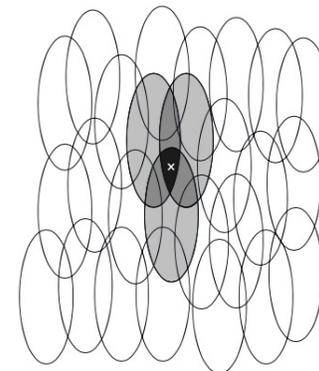
Diagonal stripes



Narrow generalization



Broad generalization



Asymmetric generalization



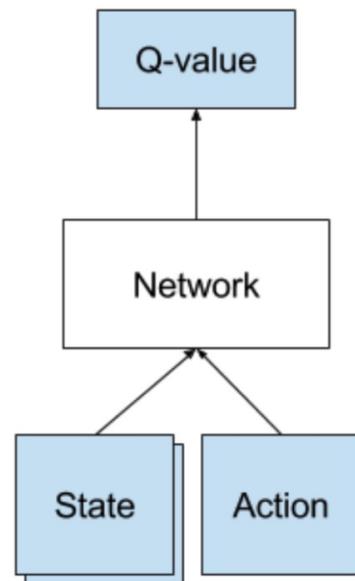
# Deep RL

- Neural networks as functions approximators have been around for a long time (e.g., TD-Gammon 1995).
- Key elements introduced by DQN:
  - Neural network architecture
  - Experience replay
  - Target network

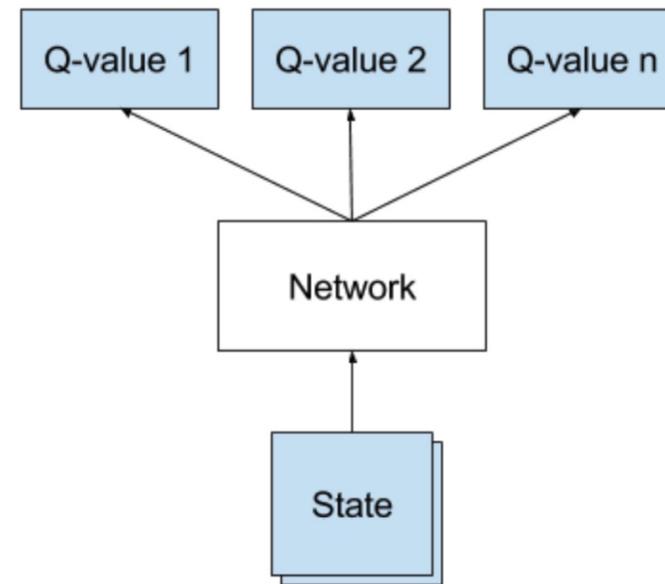


# DQN: Neural architecture

- **Input:** Current state vector
- **Output:** Contrary to the traditional RL setup, DQN produces a Q-value for each state-action pair in a **single forward pass**.



Naive

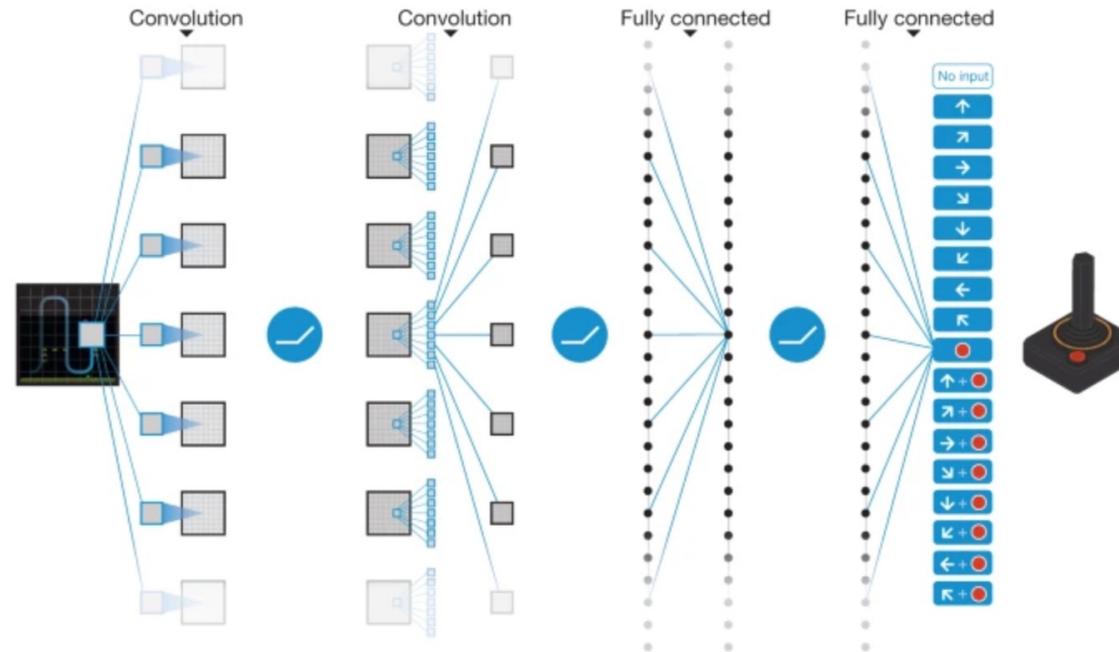


DQN



# DQN: Neural architecture (2)

- **Input:** Current state vector
- **Output:** Contrary to the traditional RL setup, DQN produces a Q-value for each state-action pair in a **single forward pass**.



# DQN: Experience replay



- Inspired from biological systems
  - Neuroscientists have discovered that spontaneous recollections, measured directly in the brain, often occur as very fast sequences of multiple memories. Neural “replay” sequences were originally discovered by studying the hippocampus in rats.
- Experience replay
  - Contrary to the traditional RL setup, all experiences  $\mathbf{e}_t = \langle \mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_{t+1}, \mathbf{s}_{t+1} \rangle$  are stored in a memory / buffer.

$\langle \mathbf{s}_1, \mathbf{a}_1, \mathbf{r}_2, \mathbf{s}_2 \rangle$
$\langle \mathbf{s}_2, \mathbf{a}_2, \mathbf{r}_3, \mathbf{s}_3 \rangle$
$\langle \mathbf{s}_3, \mathbf{a}_3, \mathbf{r}_4, \mathbf{s}_4 \rangle$
...
$\langle \mathbf{s}_{t-1}, \mathbf{a}_{t-1}, \mathbf{r}_t, \mathbf{s}_t \rangle$
$\langle \mathbf{s}_t, \mathbf{a}_t, \mathbf{r}_{t+1}, \mathbf{s}_{t+1} \rangle$

# DQN: Experience replay (2)



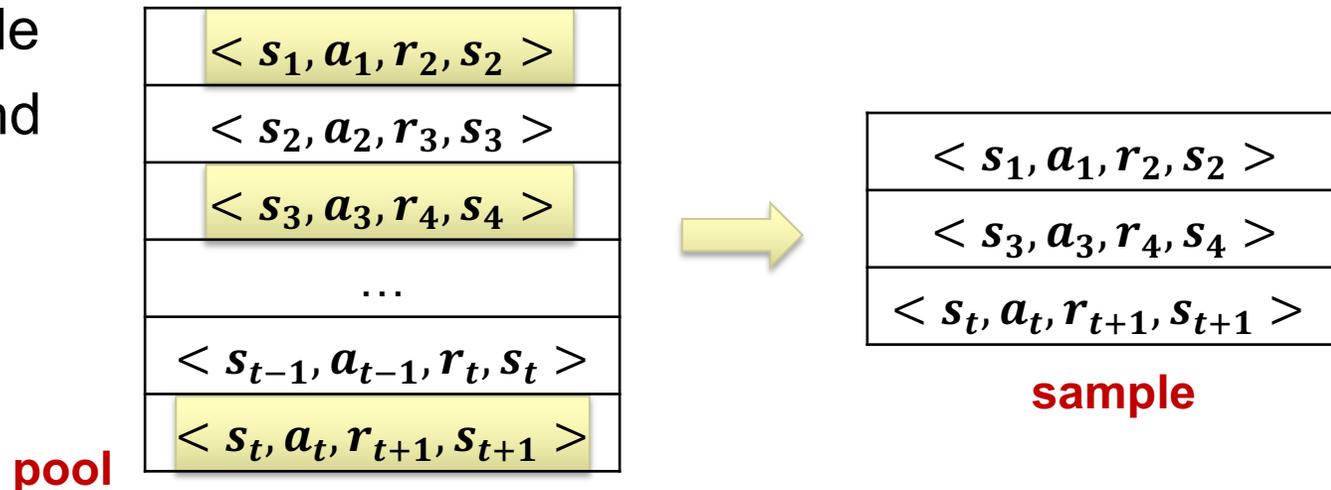
- **Advantages**

- Each experience can potentially be used for many updates, thus being more data efficient.
- Some experiences may be rare; it'd be useful if we could recall them.

- **Problem:** Sample correlations

- It causes instability in the NN's performance.

- **Solution:** Draw a random sample from the pool of experiences, and use that to train the network.



# DQN: Loss function



- DQN brings Q-learning closer to supervised learning (backpropagation, stochastic / mini-batch GD) while still allowing it to bootstrap.

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \underbrace{[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]}_{\text{TD error}}$$

$$L = \frac{1}{2} \left[ \underbrace{R_{t+1} + \gamma \max_a Q(S_{t+1}, a; w)}_{\text{target}} - \underbrace{Q(S_t, A_t; w)}_{\text{prediction}} \right]^2$$

$$\Delta w = \alpha \left( R_{t+1} + \gamma \max_a Q(S_{t+1}, a; w) - Q(S_t, A_t; w) \right) \nabla_w Q(S_t, A_t; w)$$

# DQN: Target network



- **Problem: Non-stationary target**
  - The loss function's dependence on  $w$  complicates the process compared to the simpler supervised-learning situation in which the targets do not depend on the parameters being updated. This is another source of instability.
- **Solution: Fix Q-targets**
  - To address the non-stationary target issue, DQN uses a separate target network ( $Q^{\sim}$ ) to estimate the target. The target network has the same architecture as the function approximator but with frozen parameters. The parameters from the prediction network are copied to the target network at every  $C$  iterations.

$$L = \frac{1}{2} \left[ \underbrace{R_{t+1} + \gamma \max_a Q^{\sim}(S_{t+1}, a; w)}_{\text{target}} - \underbrace{Q(S_t, A_t; w)}_{\text{prediction}} \right]^2$$

# DQN: other elements



- **State**

- A human playing any Atari game sees 210x160 pixel image frames with 128 colors at 60Hz. To reduce memory and processing requirements, the authors preprocessed each frame to produce an 84x84 array of luminance values.
- The authors “stacked” the four most recent frames so that the input dimension is 84x84x4. This did not eliminate partial observability for all of the games, but it was helpful in making many of them more Markovian.

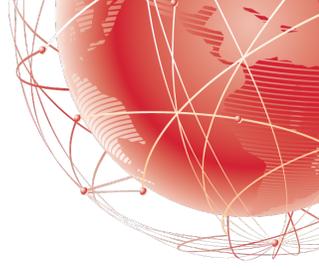
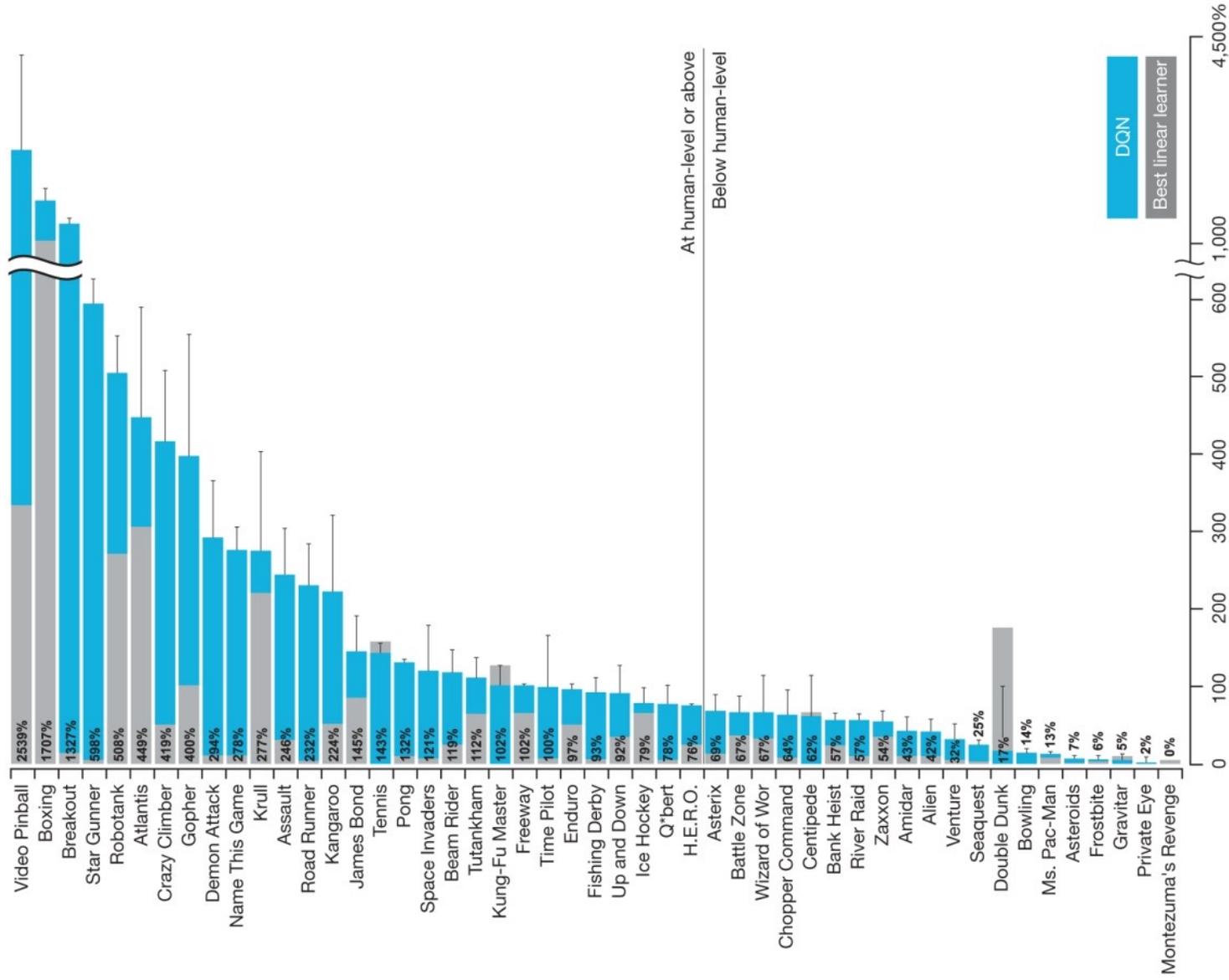
- **Reward function**

- The reward indicated how a game’s score changed from one time step to the next: +1 whenever it increased, -1 whenever it decreased, and 0 otherwise. This standardized the reward signal across the games.

- **Error clipping**

- The authors clipped the TD error in  $[-1, +1]$ .

# DQN: Results





# Prioritised Experience Replay

- Experience Replay does not consider the importance of an experience  $e_j = \langle s_j, a_j, r_{j+1}, s_{j+1} \rangle$ .

- **Idea:** sample experiences based on a priority function.

- Priority is proportional to TD error.

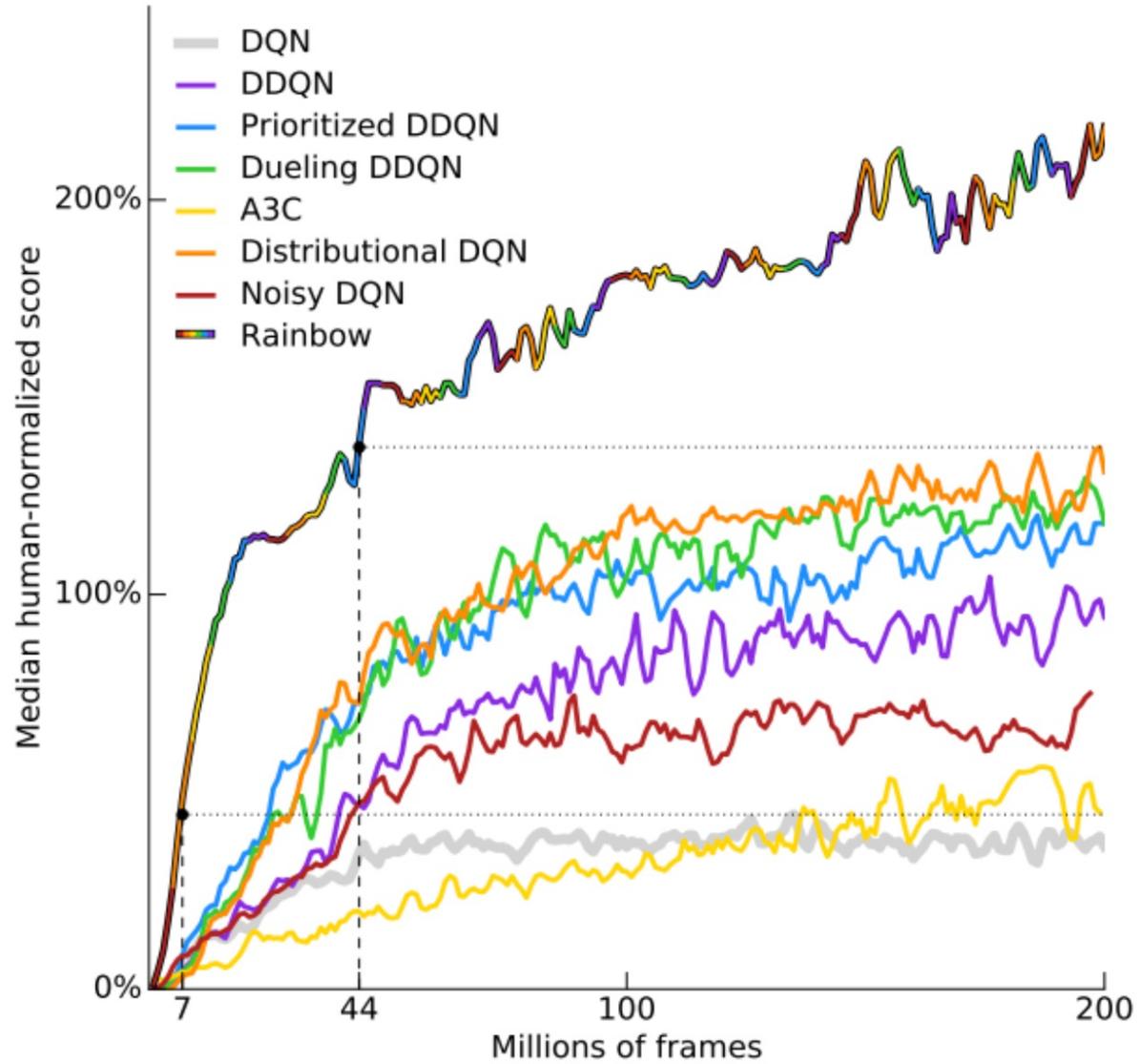
$$p_j = |R_{j+1} + \gamma \max_a Q^\sim(S_{j+1}, a; w) - Q(S_j, A_j; w)|$$

- Priority function:

- Parameter  $\alpha$  determines “how much” prioritisation is used.

$$P(j) = \frac{p_j^\alpha}{\sum_k p_k^\alpha} \quad \boxed{\alpha = 0?}$$

# Rainbow



# MULTI-AGENT RL





# Multi-agent systems and RL

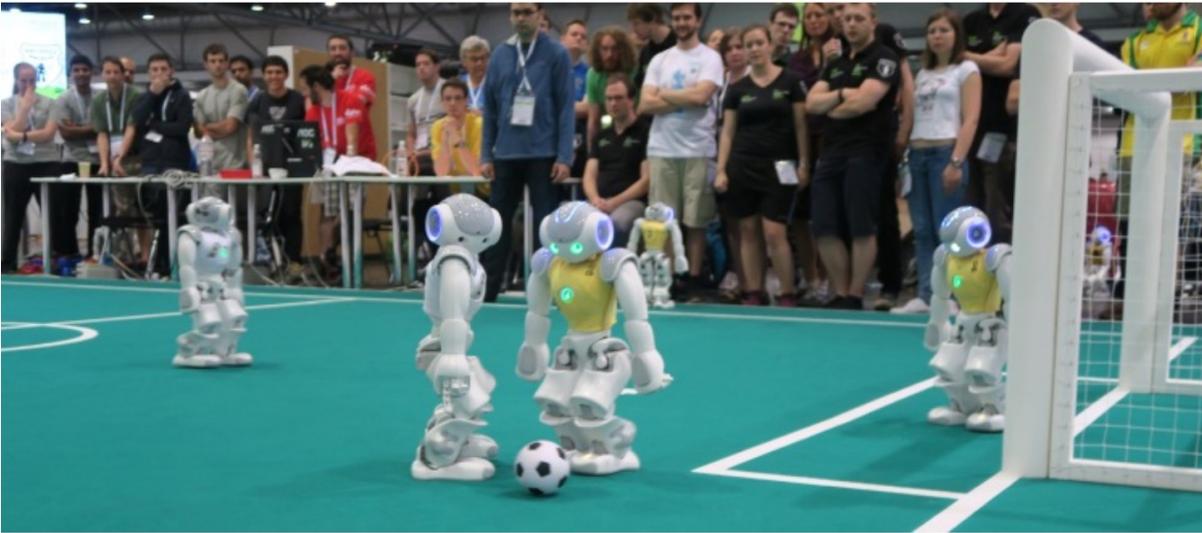
- “A **multiagent system** deals with the construction of a system involving multiple autonomous and interacting agents, which are situated in a common environment that can perceive through sensors, and act upon it through actuators.” (Busoniu et al., 2008)
- “Multiagent systems often need to be very complex, and **multiagent reinforcement learning** is a promising candidate for dealing with this emerging complexity.” (Stone & Veloso 2000).



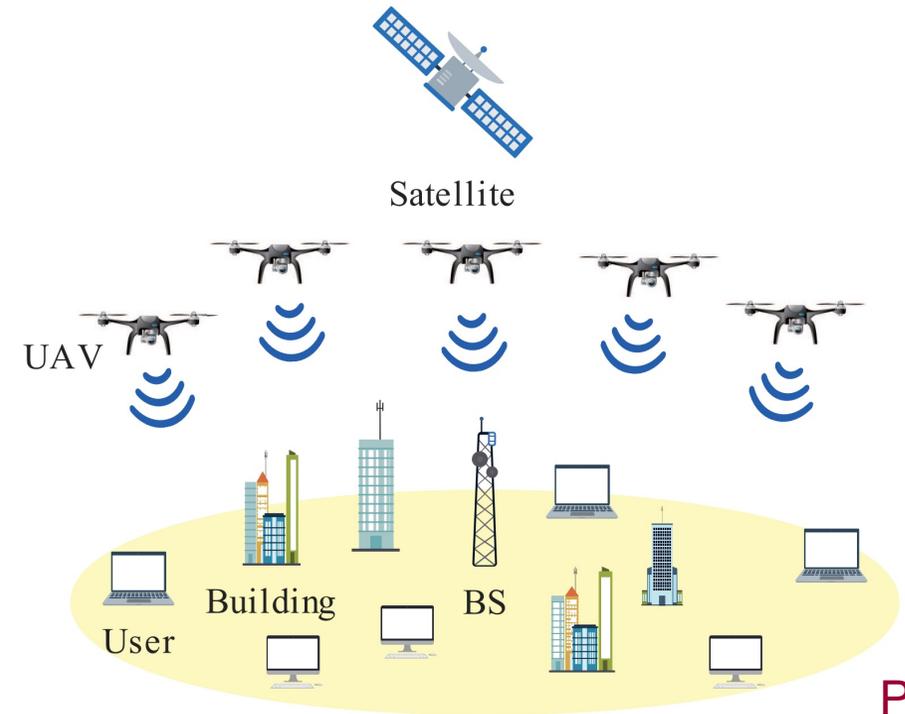
# Multi-agent RL: Applications

- **Robocup**

- *“By the middle of the 21st century, a team of fully autonomous humanoid robot soccer players shall win a soccer game, complying with the official rules of FIFA, against the winner of the most recent World Cup.” [Website](#)*



- **UAV-aided communication network**

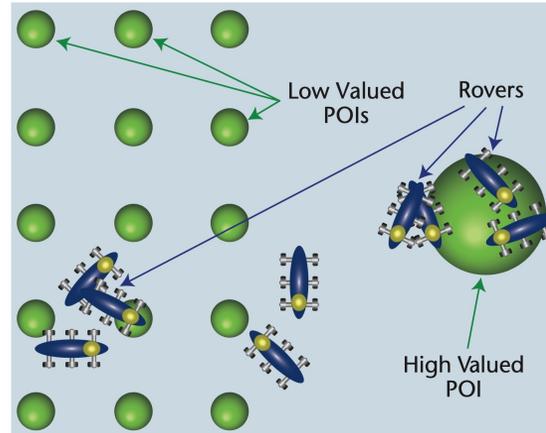
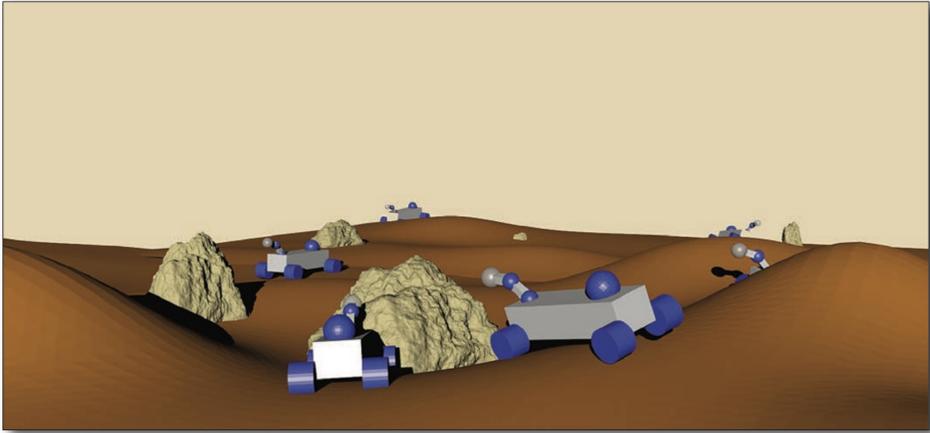


[Paper](#)

# Multi-agent RL: Applications (2)



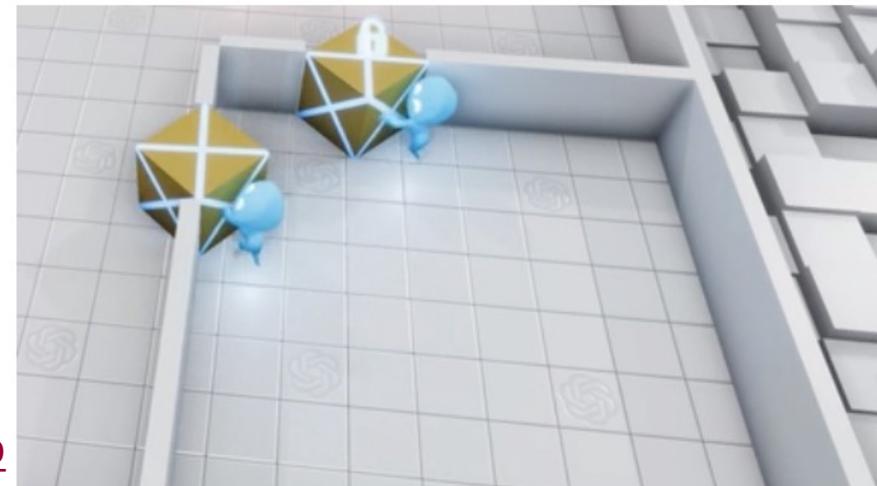
- **Space exploration**



[Paper](#)

- **Hide-and-Seek**

- Agents discover progressively more complex tool use while playing a simple game of hide-and-seek.

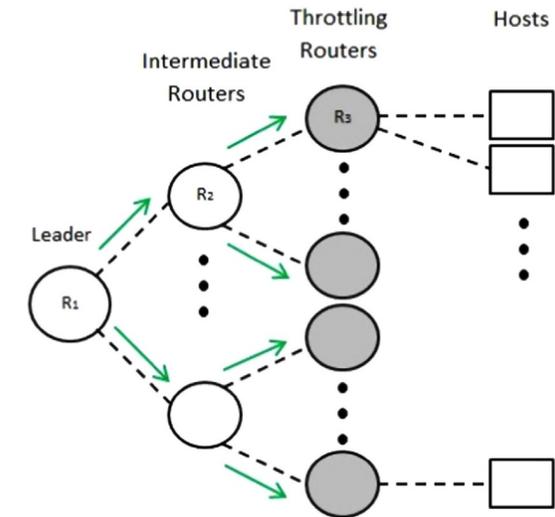
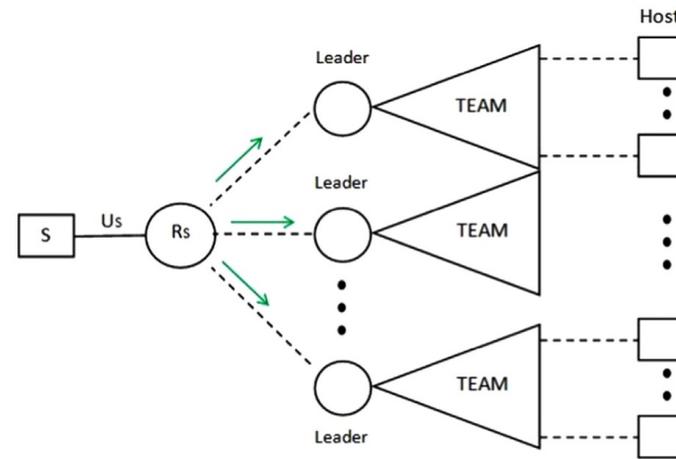
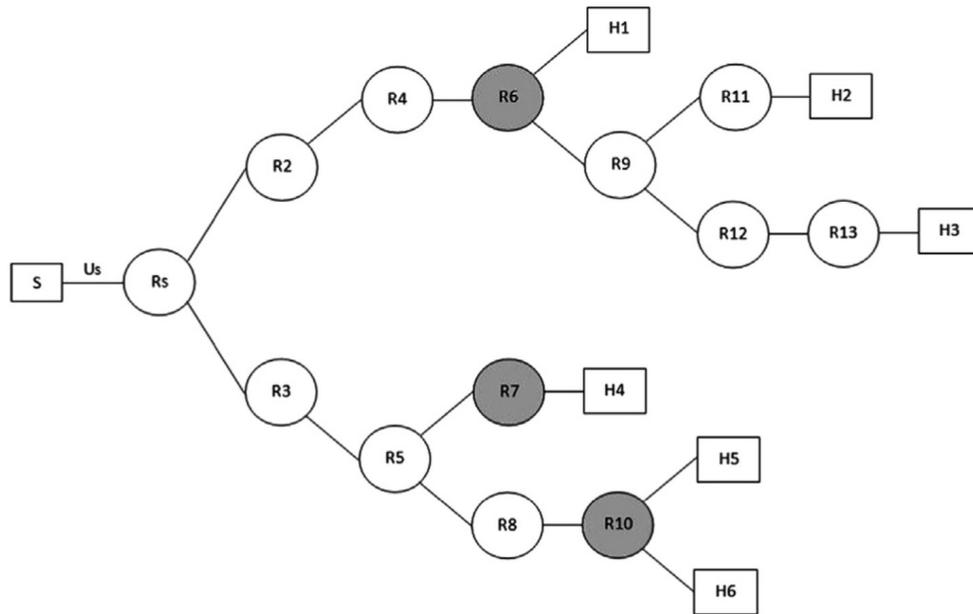


[Video](#)



# Multi-agent RL: Applications (3)

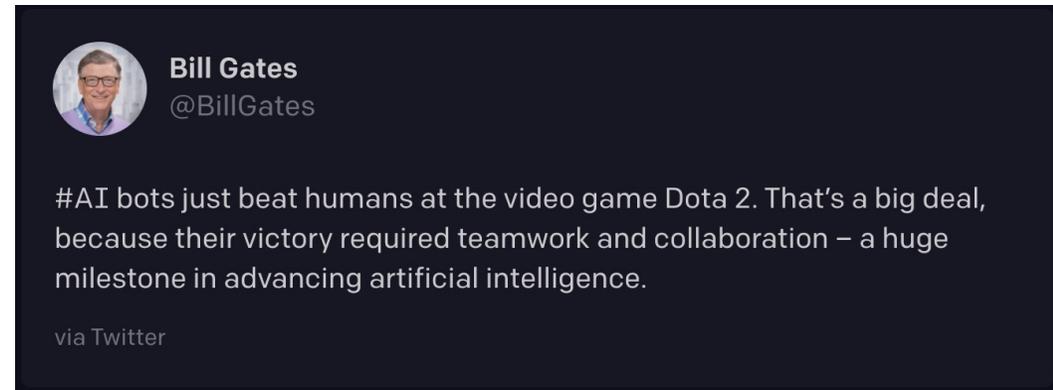
- **Network intrusion and response**
  - Multiple RL agents are installed on a set of routers and learn to throttle or rate-limit traffic towards a victim server.



# Multi-agent RL: Applications (4)



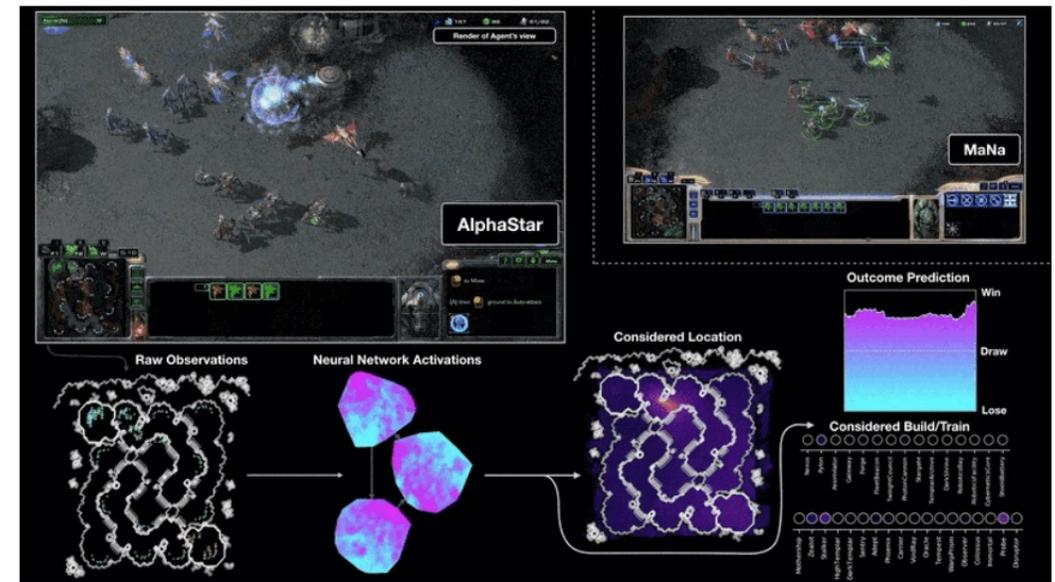
- **OpenAI: Expert-level performance in Dota 2**
  - In 2019, OpenAI and Microsoft forms exclusive partnership, and a \$1B investment from Microsoft!
  - **OpenAI Five** learned by playing over 10,000 years of games against itself. It demonstrated the ability to achieve expert-level performance, and learn human–AI cooperation.



# Multi-agent RL: Applications (5)

## DeepMind: Mastering StarCraft II

- “Despite the recent successes (Atari, Dota 2), AI techniques have struggled to cope with the complexity of StarCraft II. **AlphaStar** is the first AI to defeat a top professional player.”



[Paper](#), [demo](#)



# Multi-agent RL: Challenges

- **Curse of dimensionality**
  - Table-based methods are also affected by this. In MARL cases, the problem is worse as the complexity is now exponential in the number of agents in the system. **(function approx., organisation)**
- **Partial observability and Non-stationarity**
  - Each agent cannot observe the entire environment, but has only knowledge about its local environment.
  - The effects of an agent's action on the environment, also depend on the other agents' actions. Therefore, the Markov property does not hold if an agent cannot observe the joint state / action.  
**(communication)**



# Multi-agent RL: Challenges (2)

- **Multiagent credit assignment**
  - Typically, in a multiagent system an agent is provided with a reward at the global / system level. Given that other agents also interact with the environment, an agent may be rewarded for taking a bad action, or punished for taking a good action. (**reward shaping**)

$r$	-10	0	10
$t_1$	10 A	10* B	10 C
$t_2$	6* A	9 B	15 C

# Markov Models

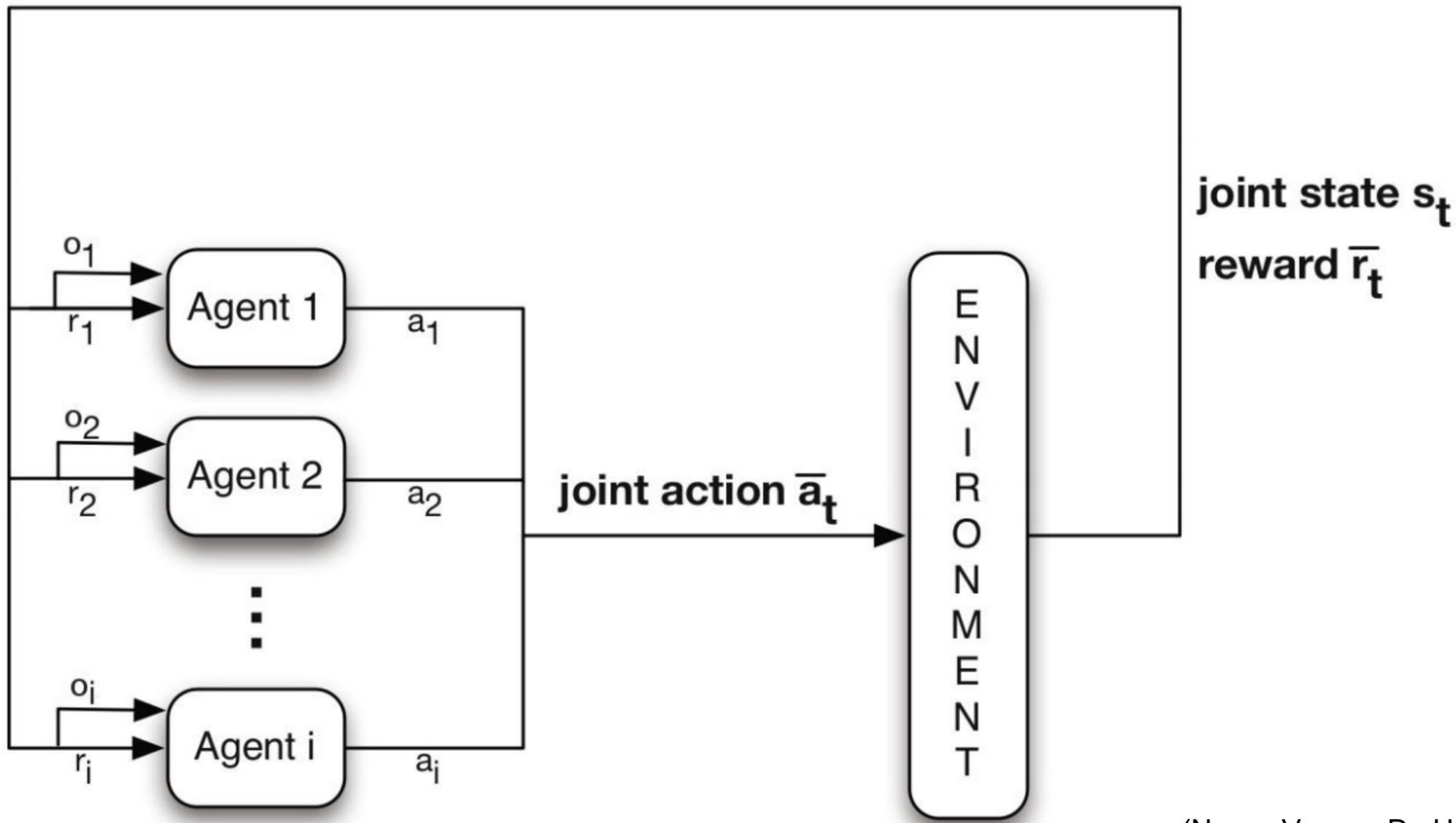


	Single-agent	Multi-agent
State known	Markov Decision Process (MDP)	Stochastic Game (SG)
State observed indirectly	Partially Observable MDP (POMDP)	Partially Observable SG (POSG)

## ■ Stochastic Game (SG)

- $n$ : number of agents
- $S$ : set of states
- $A = A_1 \times A_2 \times \dots \times A_n$ : joint action set, where  $A_i$  is the set of actions available to each agent  $i$
- $p: S \times R \times S \times A \rightarrow [0,1]$ : transition probability function
- $r_i: S \times R \times S \rightarrow [0,1]$ : reward function of agent  $i$

# Multi-agent RL



(Nowe, Vrancx, De Hauwere, 2012)

# Independent Learners (ILs)



- This involves the deployment of multiple agents each using a single-agent RL algorithm.
- Multiple ILs assume any other agents to be a part of the environment and so, as the others simultaneously learn, the environment appears to be dynamic as the probability of transition when taking action  $a$  in state  $s$  changes over time.

$$Q_i(S_t, A_t) = Q_i(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q_i(S_{t+1}, a) - Q_i(S_t, A_t)]$$

# Joint Action Learners (JALs)

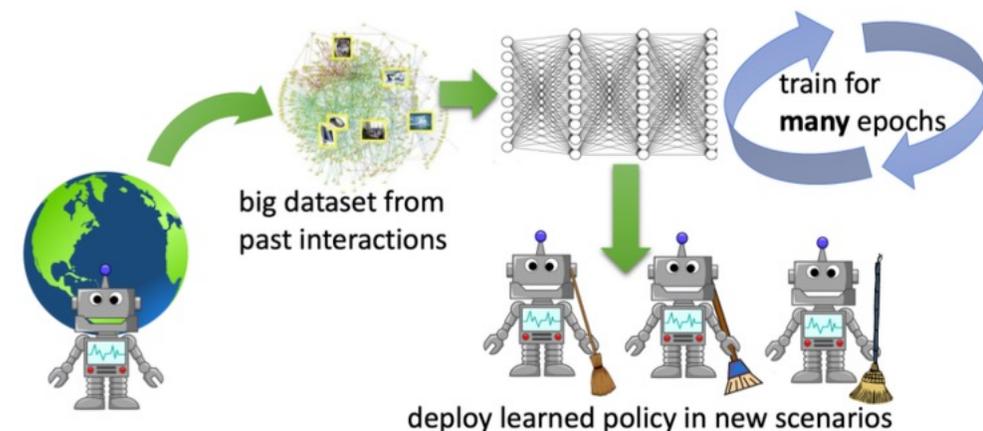


- This approach considers the existence of other agents. Specifically, an agent observes the actions of the other agents, or each agent communicates its action to the others.
- To overcome the appearance of a dynamic environment, JALs extend their value function to consider for each state the value of each possible combination of actions by all agents.
- If a JAL agent  $i$  uses SARSA, then its update rule is:

$$Q_i(s, a^{\rightarrow}) \leftarrow Q_i(s, a^{\rightarrow}) + \alpha[r + Q_i(s', a'^{\rightarrow}) - Q_i(s, a^{\rightarrow})]$$

# Trends in RL

- **Imitation learning**
  - An agent tries to learn the optimal policy by following or imitating the expert's decisions.
- **Offline RL**
  - The agent learns skills solely from previously collected datasets, without any active environment interaction. Datasets might be of the form of human demonstrations, prior experiments, domain-specific solutions and even data from different but related problems, to build complex decision-making engines.



# Trends in RL (2)



- **Safe RL**
  - It is defined as the process of learning policies that maximise the expectation of the return in problems in which it is important to ensure reasonable system performance and/or respect safety constraints during the learning and/or deployment processes

# Some Bibliography on Reinforcement Learning



- R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*, 2<sup>nd</sup> Edition, The MIT Press, 2018.
- D. Bertsekas, *Reinforcement Learning and Optimal Control*, Athena Scientific, 2019.
- C. Szepesvari, *Algorithms for Reinforcement Learning*, Morgan and Claypool, 2010.
- David Silver Lectures on RL: <https://www.davidsilver.uk/teaching>