

Introduction to Q-learning and Deep Q-networks

Homeproblem B: Playing Tetris using Reinforcement: Learning

Advanced machine learning with neural networks

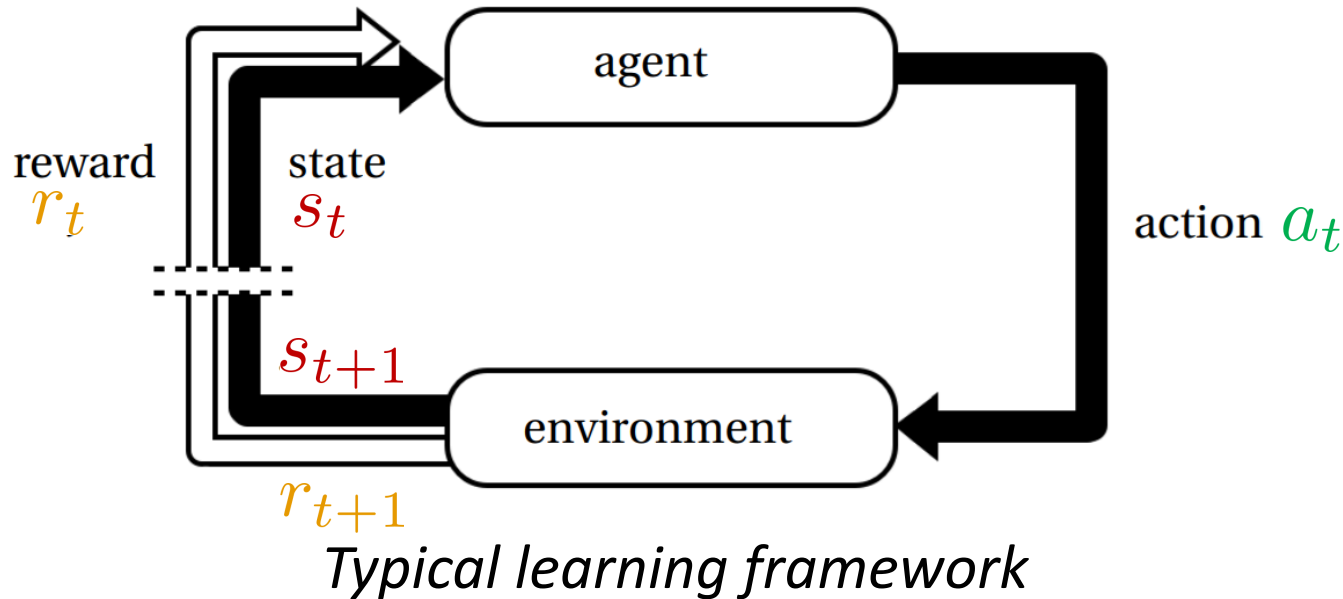
Kristian Gustavsson, Department of Physics, University of Gothenburg, Sweden

kristian.gustafsson@physics.gu.se

What is Reinforcement Learning?

Reinforcement Learning: An introduction, Sutton & Barto (2018)
Machine learning with neural networks, Mehlig (2021)

Scheme that uses trial and error to find optimal strategies to achieve a predetermined task.



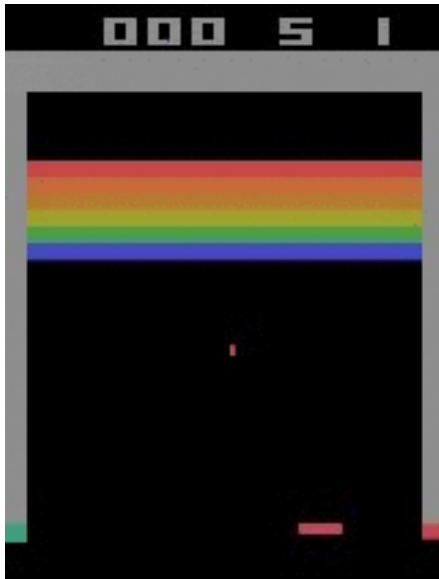
Learning uses a reward signal only (no supervisor).

Example applications

- Gaming strategies
- Control theory
- Industrial optimizations
- Understanding nature

Application: Gaming strategies

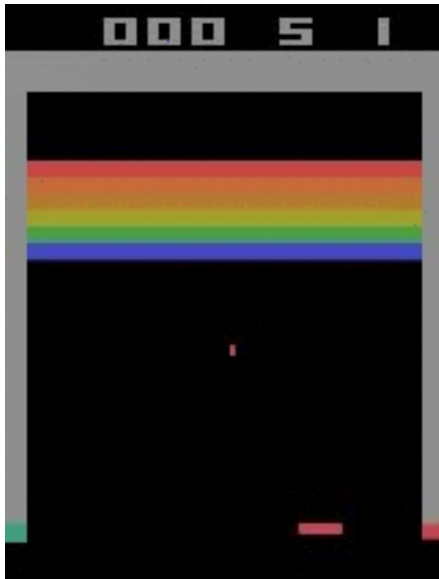
before training



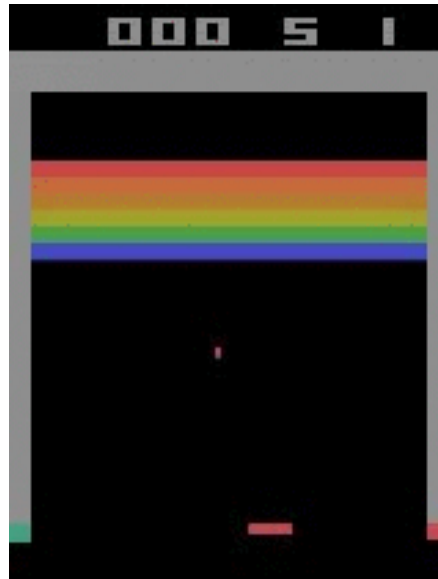
Learning to play 'Breakout'

Application: Gaming strategies

before training



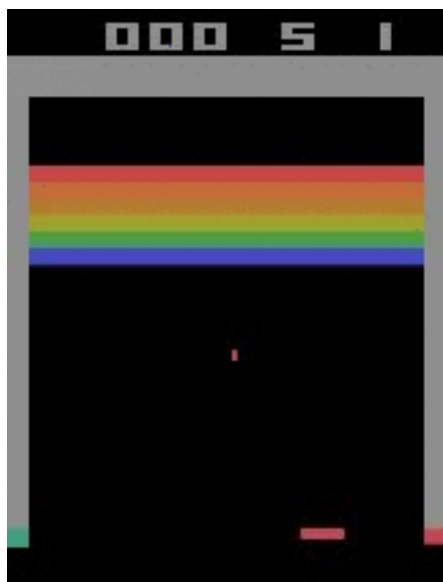
during training



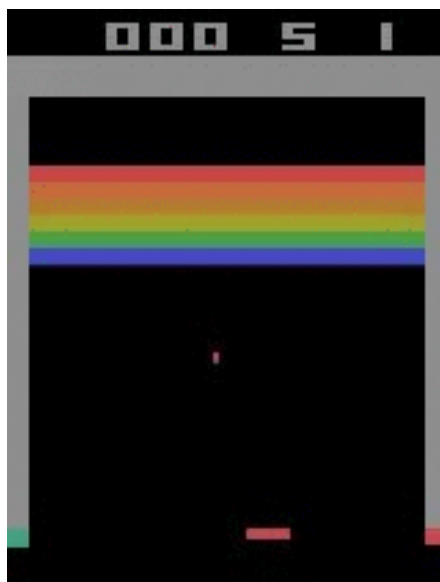
Learning to play 'Breakout'

Application: Gaming strategies

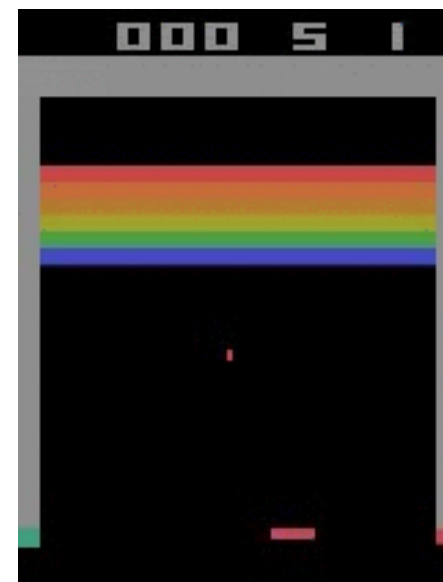
before training



during training



after training



Learning to play 'Breakout'

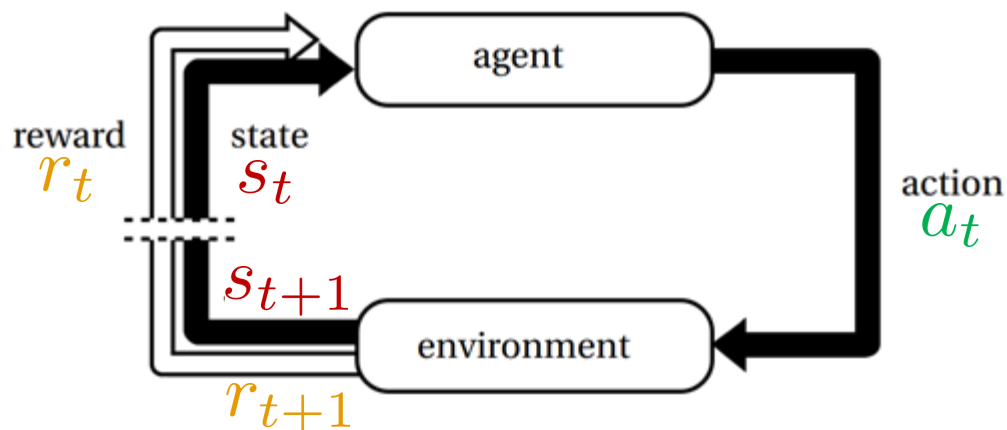
Agent The player

Environment The game (and player)

States Four most recent game images

Actions Joystick direction

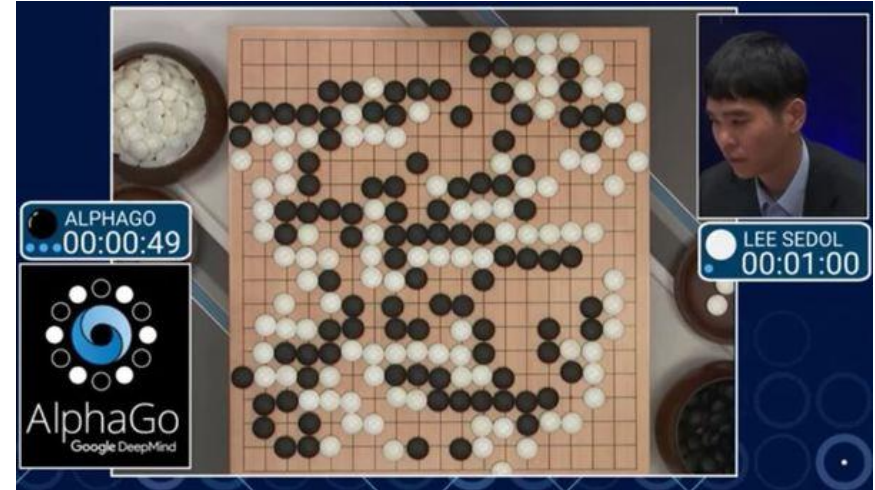
Reward Increment in game score



Application: Gaming strategies



Learning to play 'Go'



Top view

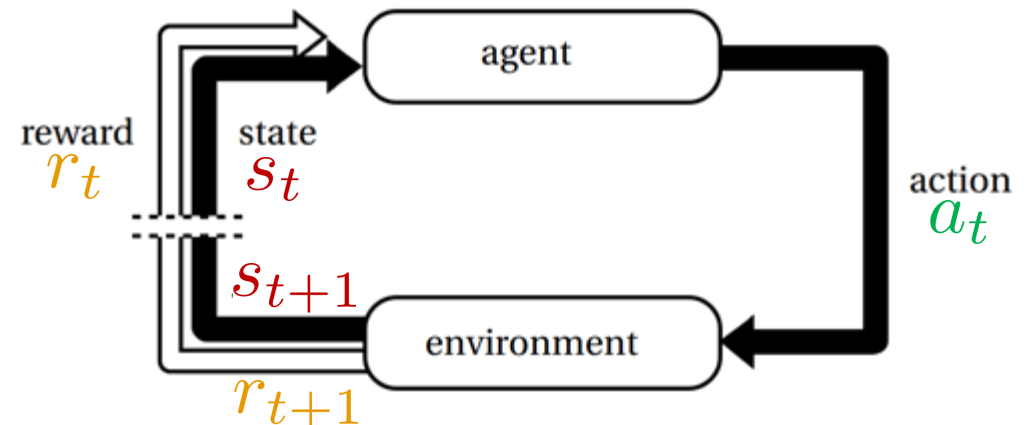
Agents Both players

Environment The game (and players)

States Positions of game markers

Actions Place new game marker

Reward + on win; - on lose



Application: Control theory



Helicopter aerobatics

Agent The helicopter operator

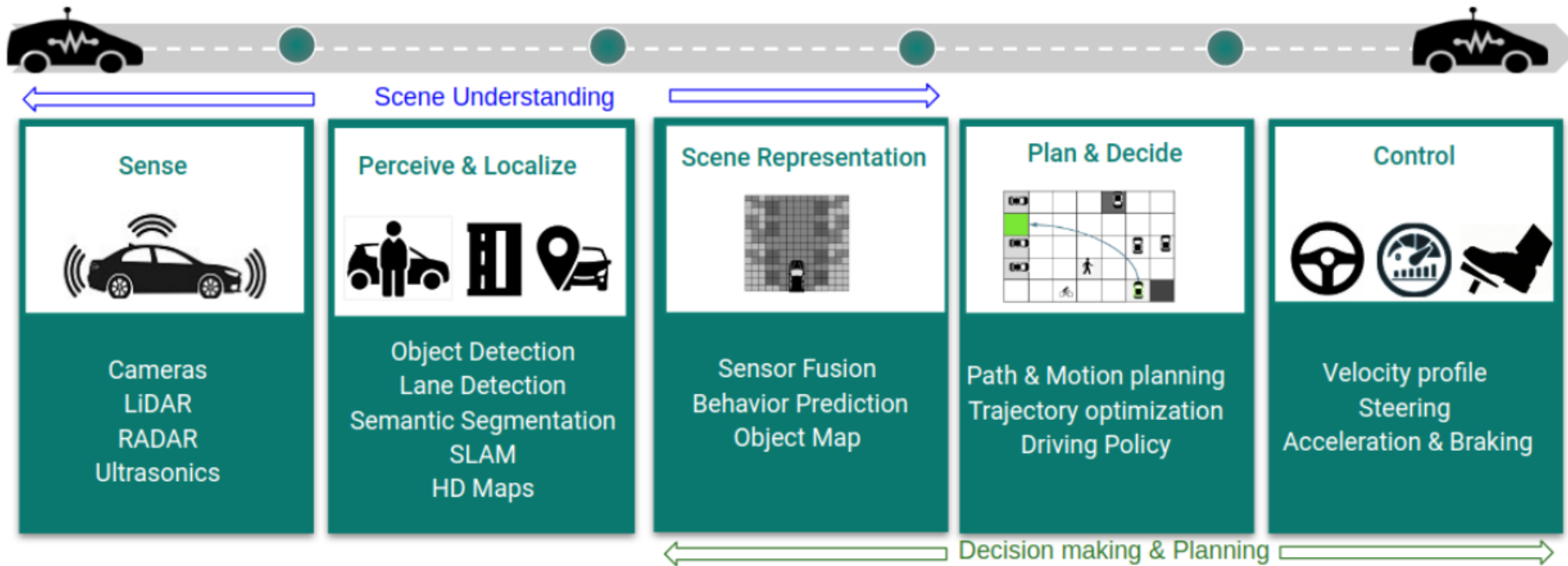
Environment The helicopter and the surroundings

States Position, orientation, velocity, and angular velocity

Actions Change main rotor tilt or angle of attack, change tail rotor thrust

Reward Penalty for deviations from target trajectory

Application: Control theory



Decision making and planning in self-driving cars

Agent The driver

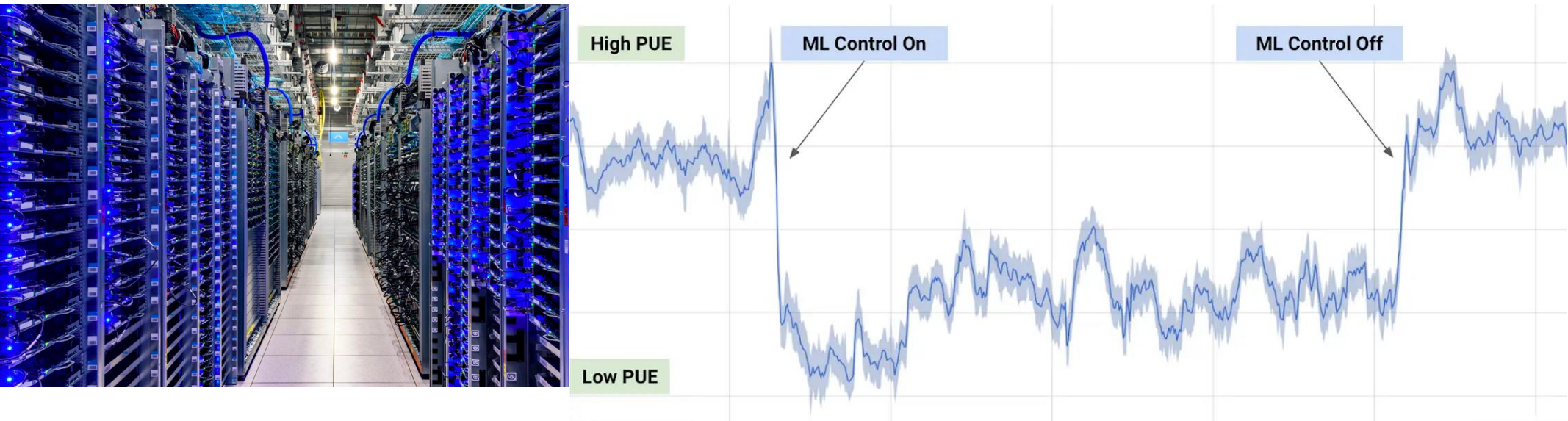
Environment The car and the surroundings

States Position and velocity. Scene representation.

Actions Steering, accelerating, braking. Control blinkers, lights, horn etc.

Reward + for following road and reaching target; Neg. for collision or crash

Application: Industrial optimizations



Reducing cooling energy consumption by 40% in Google servers

Agent Controller of cooling system

Environment Data center

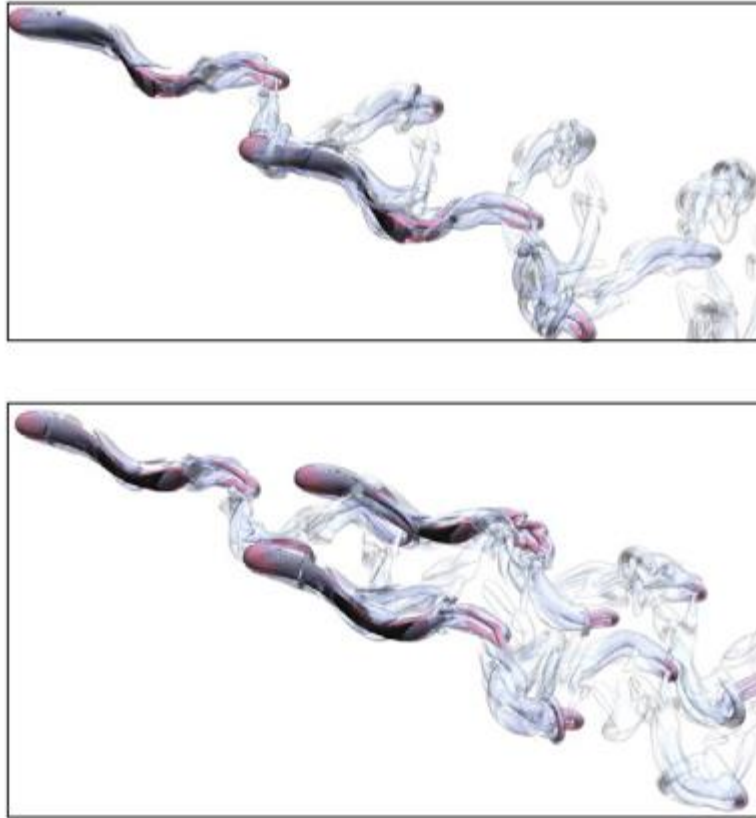
States Workload and ambient temperature

Actions Control cooling system

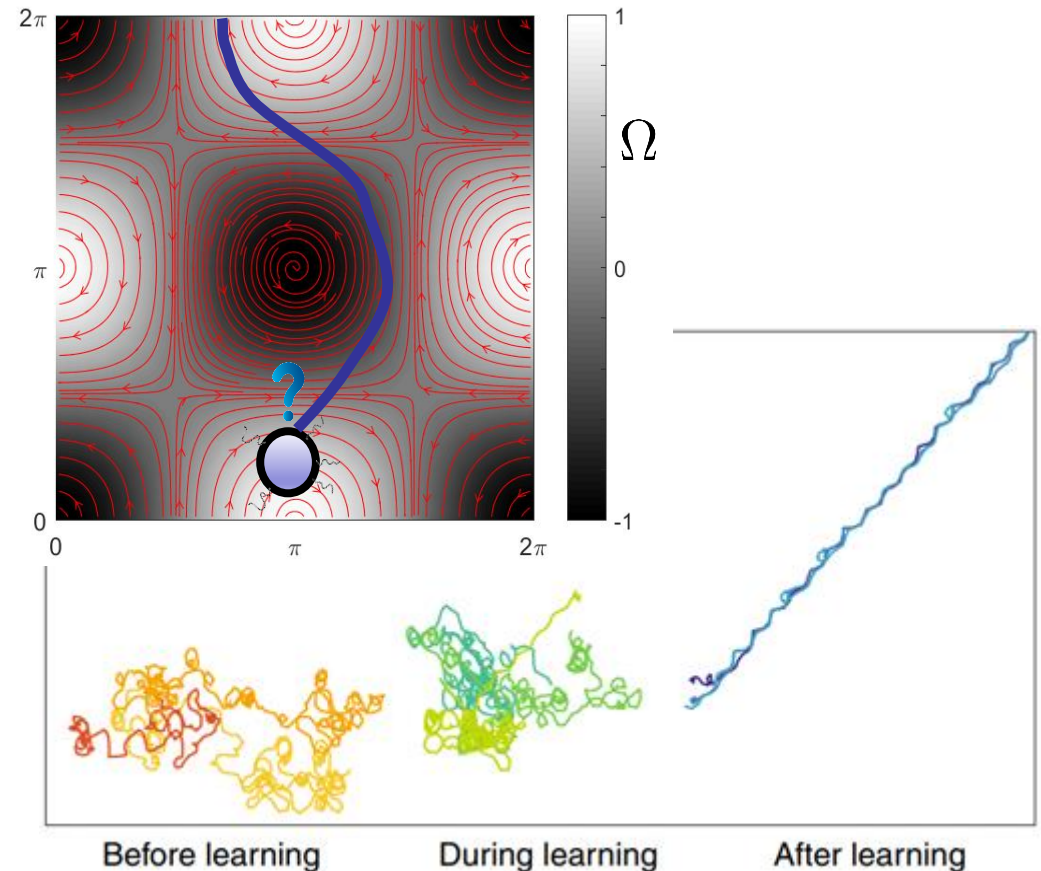
Reward Neg. for power usage or outlet temperature outside specified range

Application: Understanding nature

Verma et al, PNAS **115**, 5849–5854 (2018)



Colabrese et al, PRL **118**, 158004 (2017)



Swimming strategies in turbulent flows

Agent Swimmer

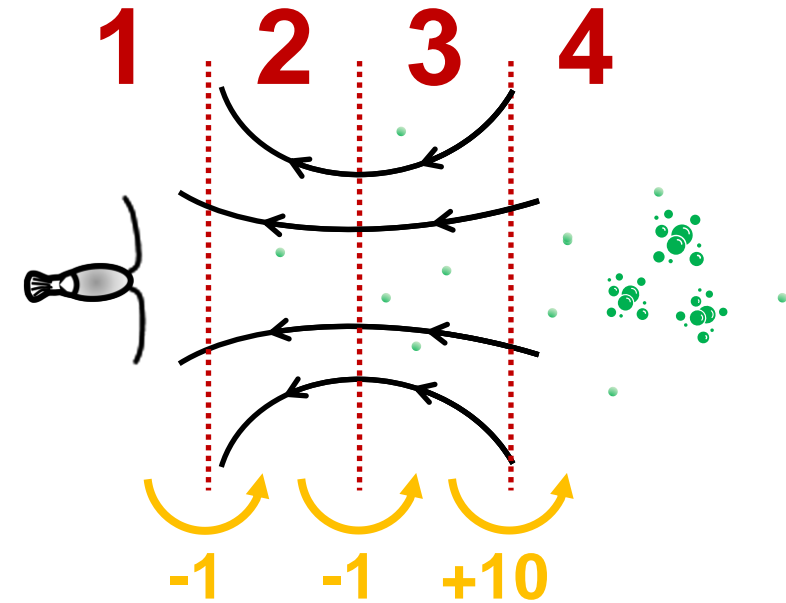
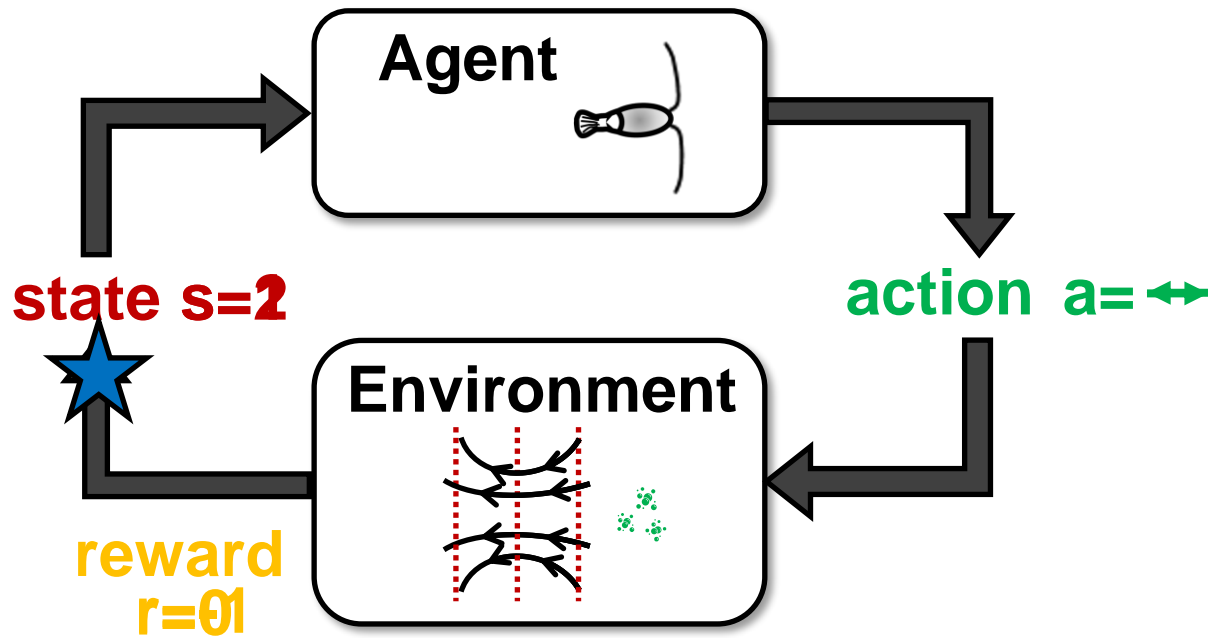
Environment Turbulent flow

States Flow properties

Actions Swimming behavior

Reward Spend little energy, swimming in target direction

Example: Swimming against the flow



Choose 'greedy' **action**

$$a = \operatorname{argmax}_{a'} Q(s, a')$$

action \rightarrow
 \leftarrow

	state			
	1	2	3	4
\rightarrow	-1	0	0	0
\leftarrow	0	0	0	0

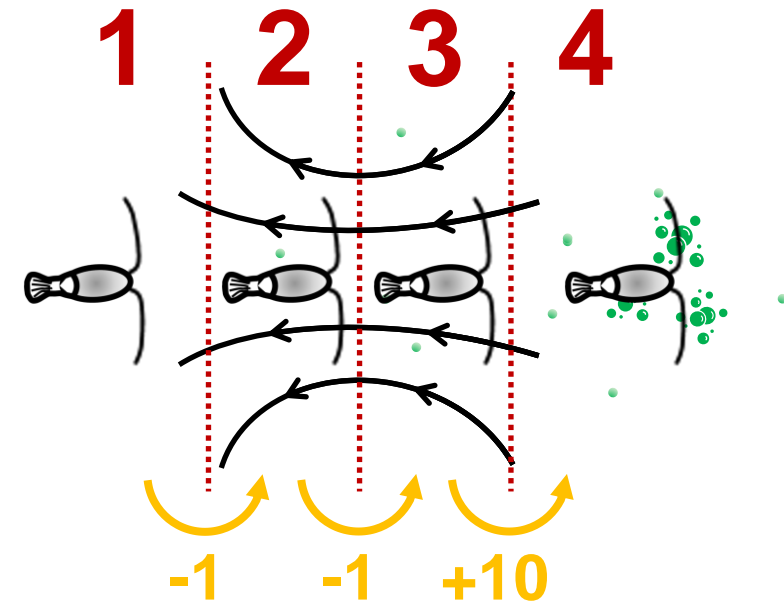
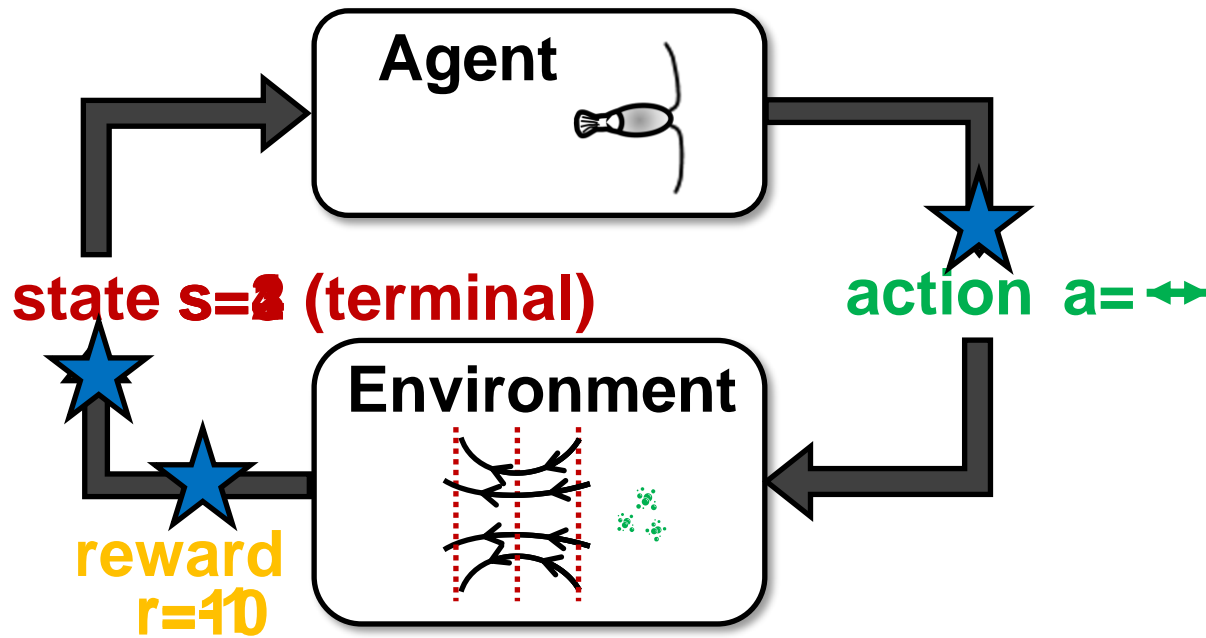
Q-table

Update Q-table

$$Q(s, a) = r$$

The scheme gets stuck at state **1**
=> Need to **add exploration**

Example: Swimming against the flow



Choose ' ϵ -greedy' **action**

probability ϵ probability $1 - \epsilon$
 random a $\text{argmax}_{a'} Q(s, a')$

action \rightarrow
 \leftarrow

state			
1	2	3	4
-1	-1	10	0
0	0	0	0

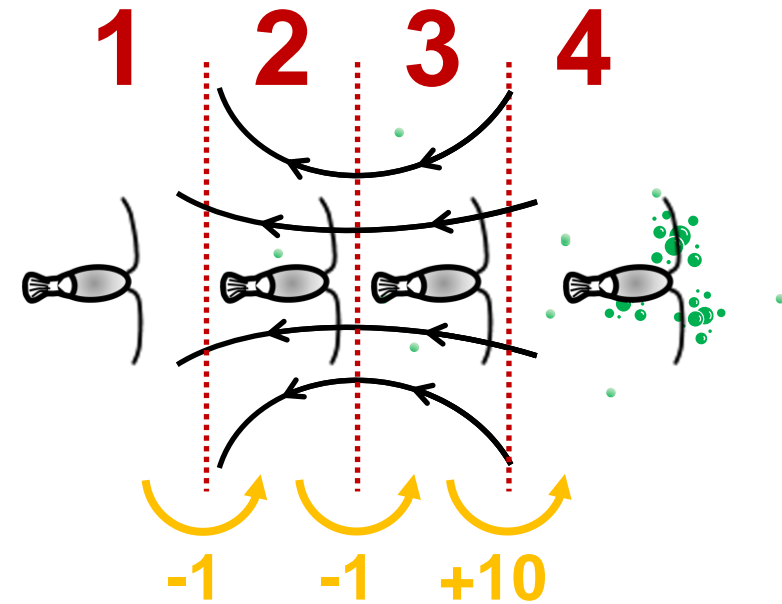
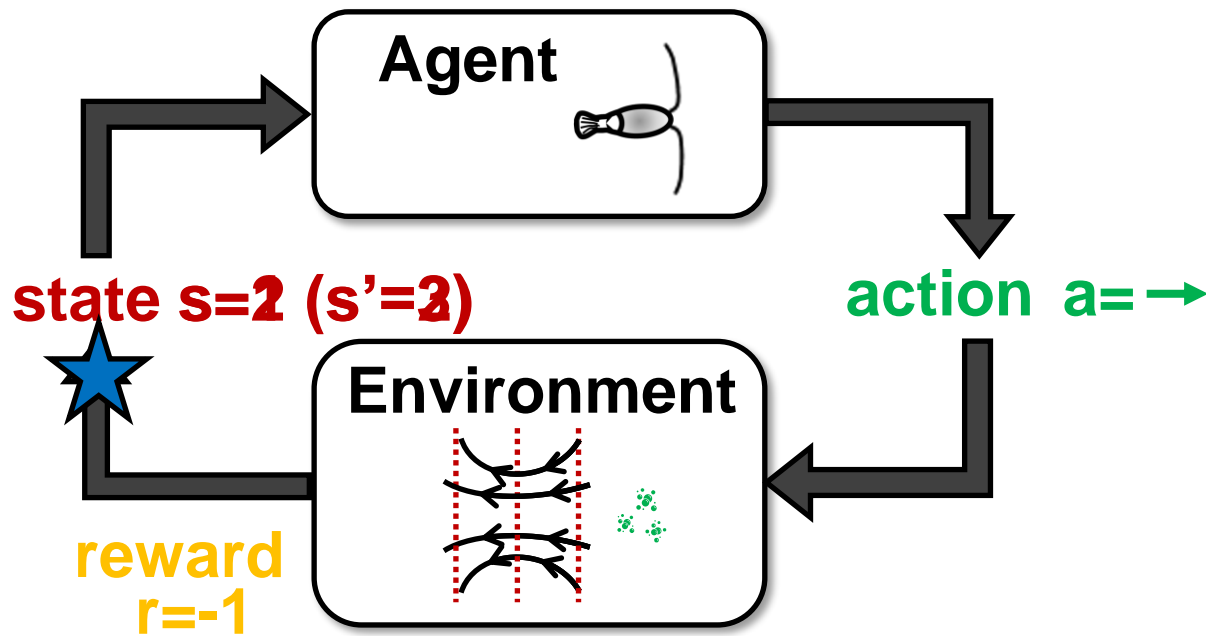
Q-table

Update Q-table

$$Q(s, a) = r$$

Exploration \Rightarrow all states are reached.
 But Q-table does not give a good strategy.
 \Rightarrow Need to improve update of Q

Example: Swimming against the flow



Choose ' ϵ -greedy' action

probability ϵ probability $1 - \epsilon$
 random a $\operatorname{argmax}_{a'} Q(s, a')$

Update Q-table (Q-learning)

$Q(s, a) = r + \max_{a'} Q(s', a')$
 here s' is the next state.

state

1 2 3 4

8	9	10	0
0	0	0	0

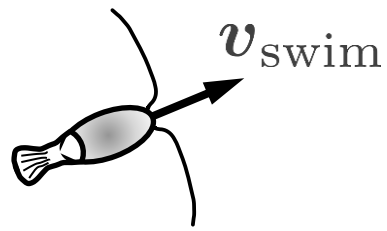
action \rightarrow \leftarrow

Q-table

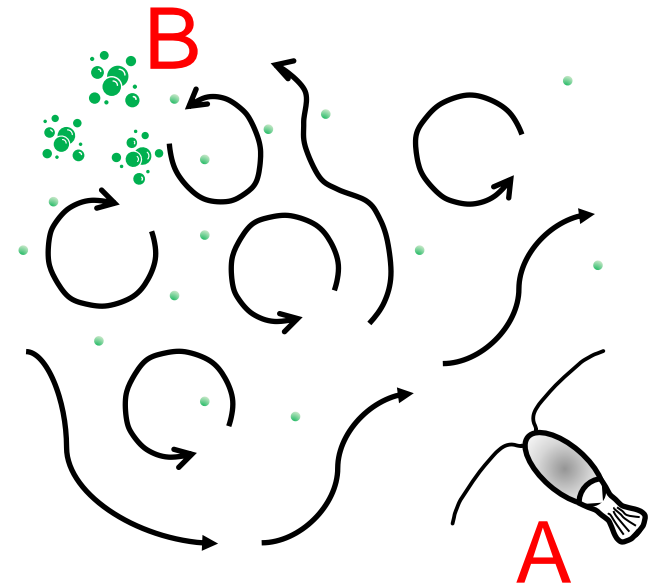
Q-table gives successful strategy

Example: Zermelo's problem

- Navigation from point **A** to point **B** in a background flow
- Constant swimming speed v_{swim}

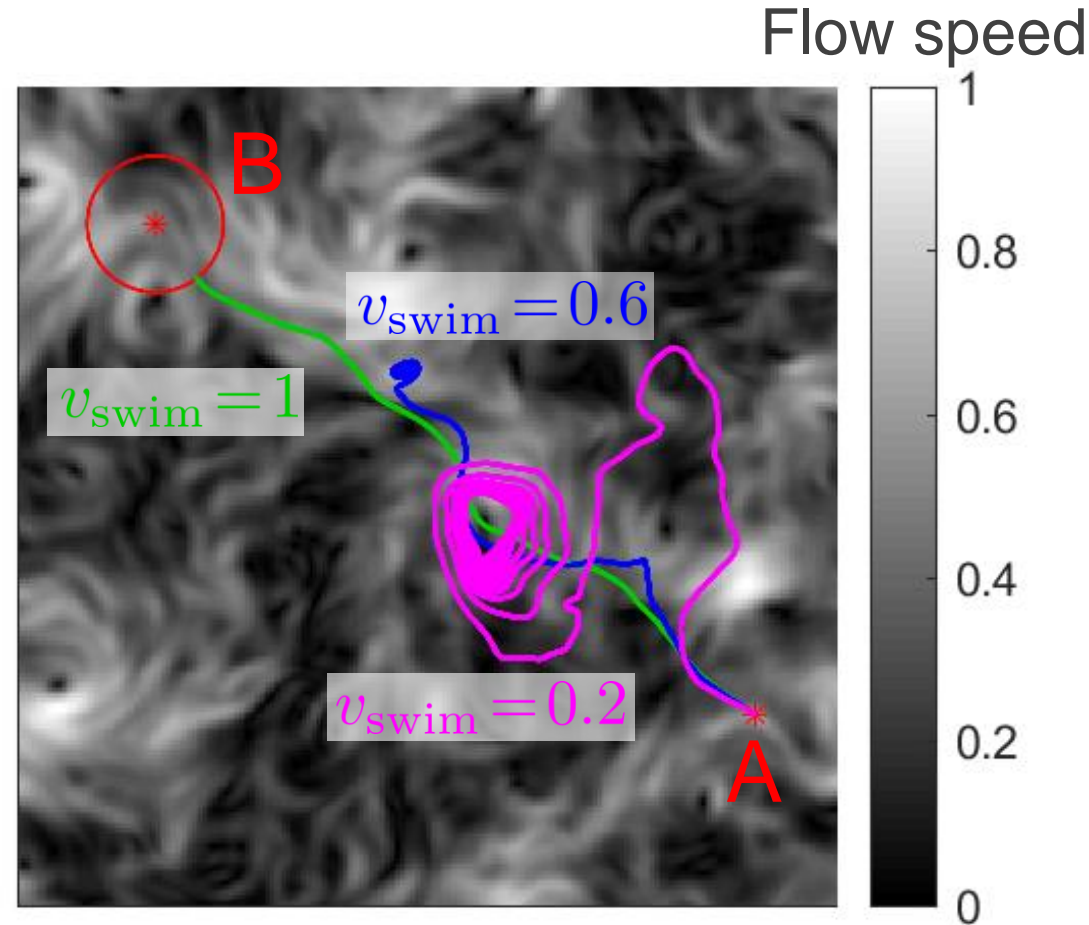


- How to choose swimming direction?



Example: Zermelo's problem

- Naive approach: always swim towards the target



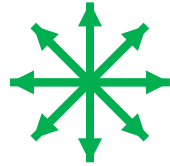
Example: Zermelo's problem (Q-learning)

States

Flow grid

Actions

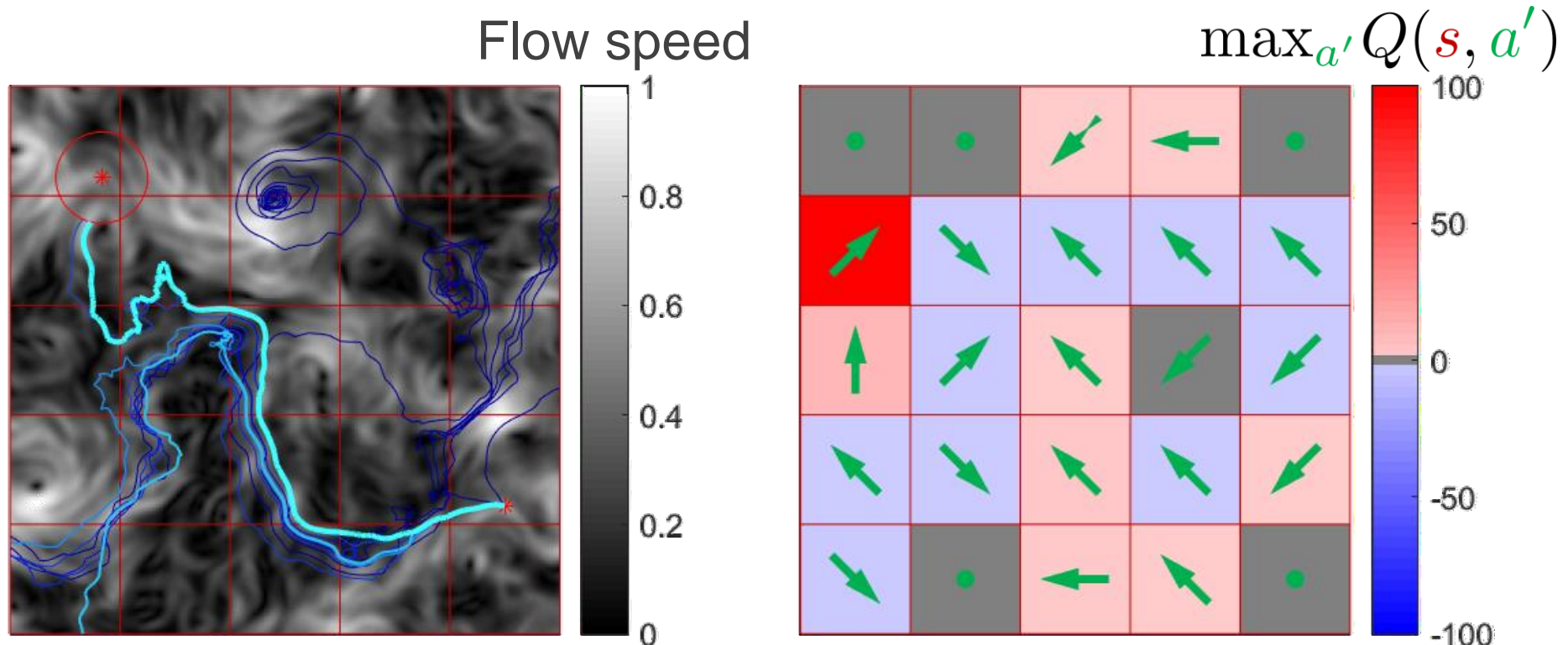
Swim direction



Rewards

- For state update/crossing bounds
- + For moving towards/reaching target

Learning rate α : $Q(s, a) = \alpha [r + \max_{a'} Q(s', a')] + (1 - \alpha) Q(s, a)$



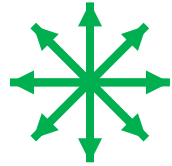
Example: Zermelo's problem (Q-learning)

States

Flow grid

Actions

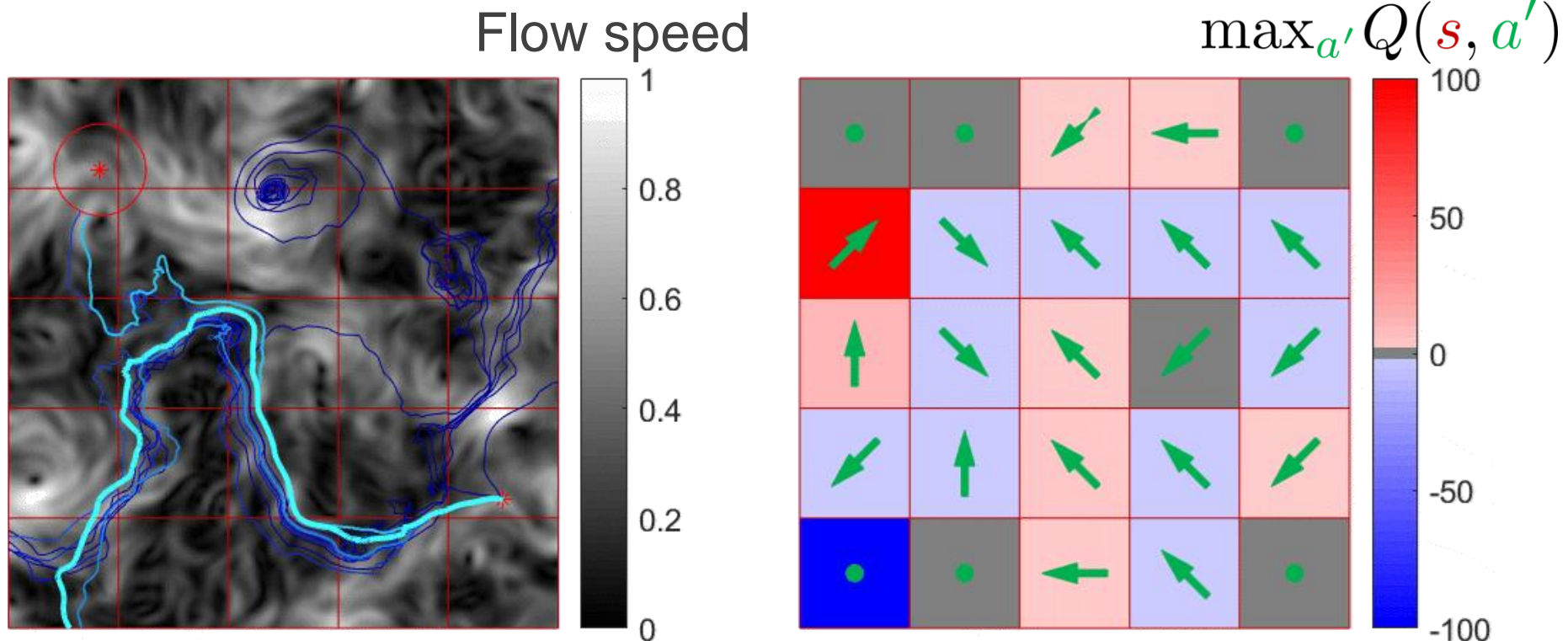
Swim direction



Rewards

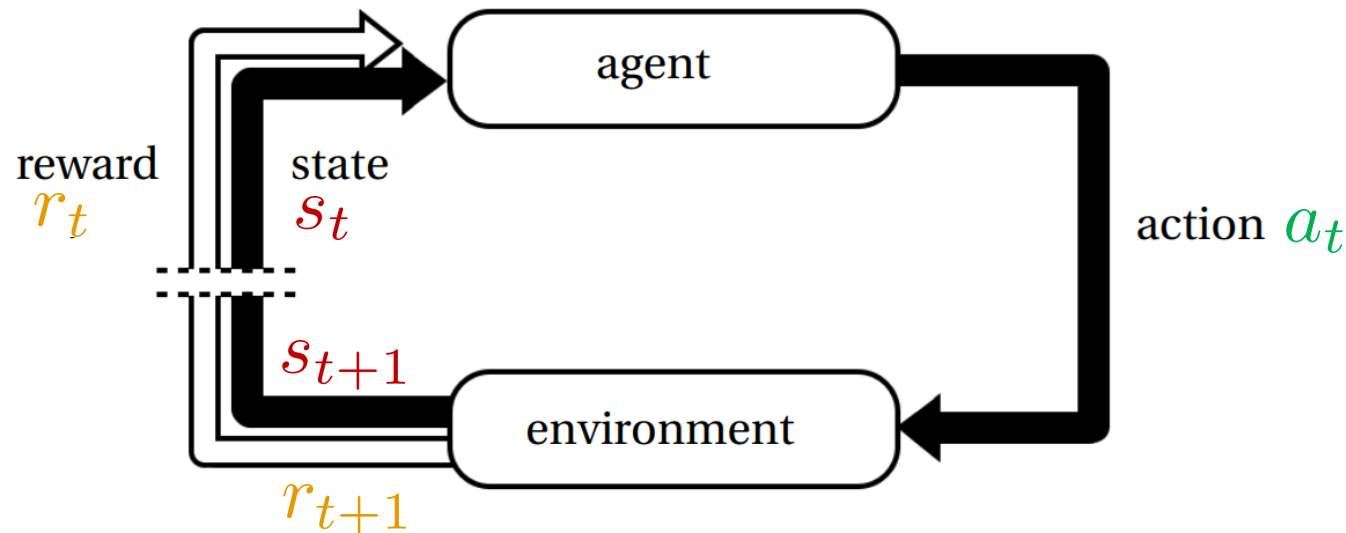
- For state update/crossing bounds
- + For moving towards/reaching target

Learning rate α : $Q(s, a) = \alpha [r + \max_{a'} Q(s', a')] + (1 - \alpha) Q(s, a)$



Summary of Q-learning

$0 \leq \varepsilon \leq 1$ Exploration rate
 $0 < \alpha \leq 1$ Learning rate



ε -greedy policy

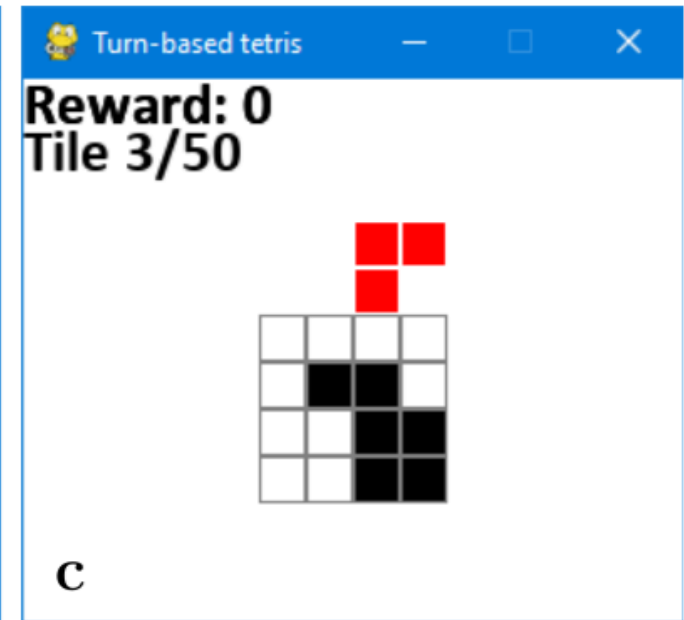
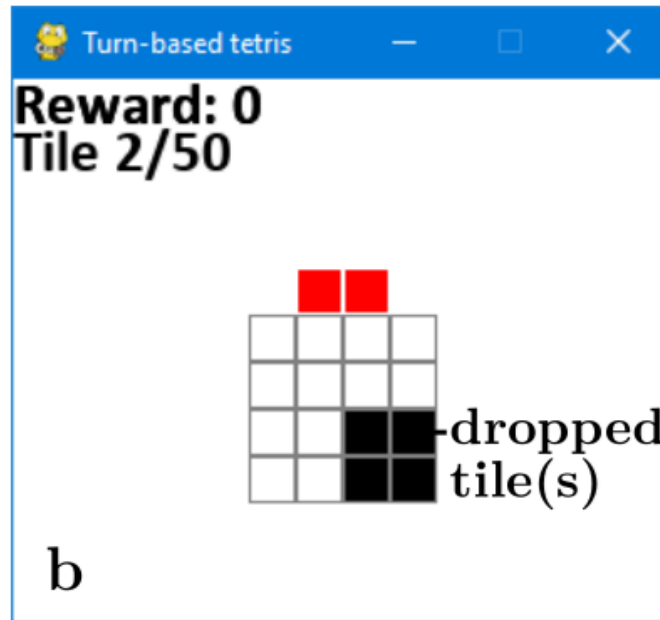
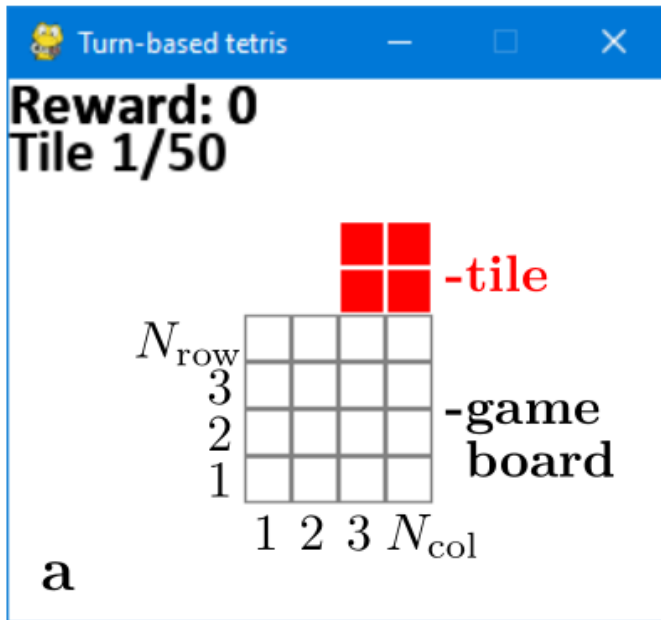
probability ε probability $1 - \varepsilon$
 random a $\operatorname{argmax}_{a'} Q(s, a')$

Use Q-table to store experience (expected future reward)

$$Q(s, a) = \alpha [r + \max_{a'} Q(s', a')] + (1 - \alpha) Q(s, a)$$

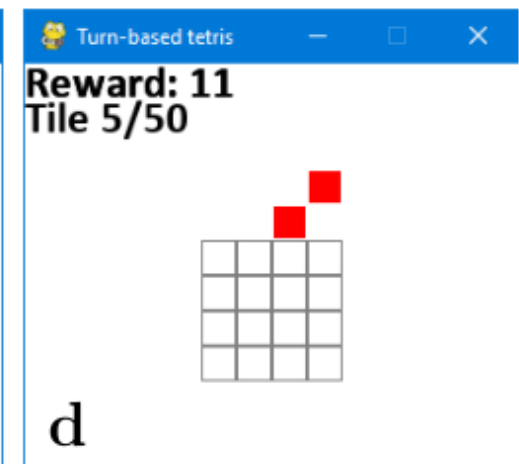
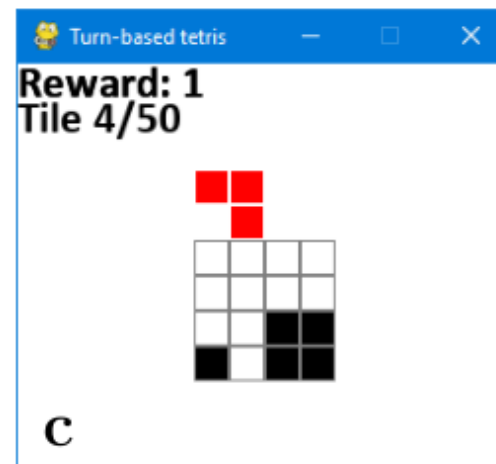
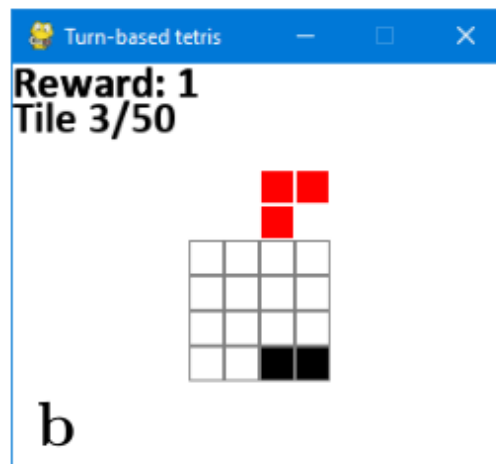
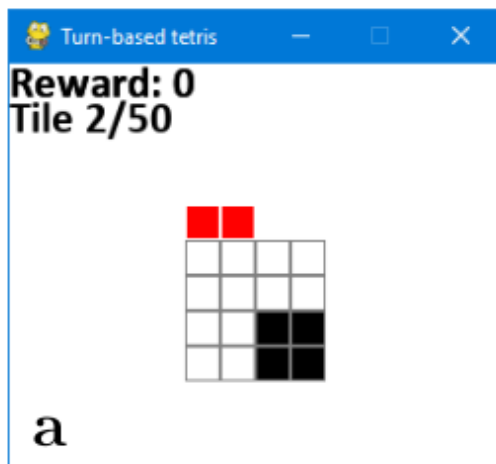
Home problem B: Introduction

Task To learn to play simplified Tetris



Tile set ( ,  ,  , )

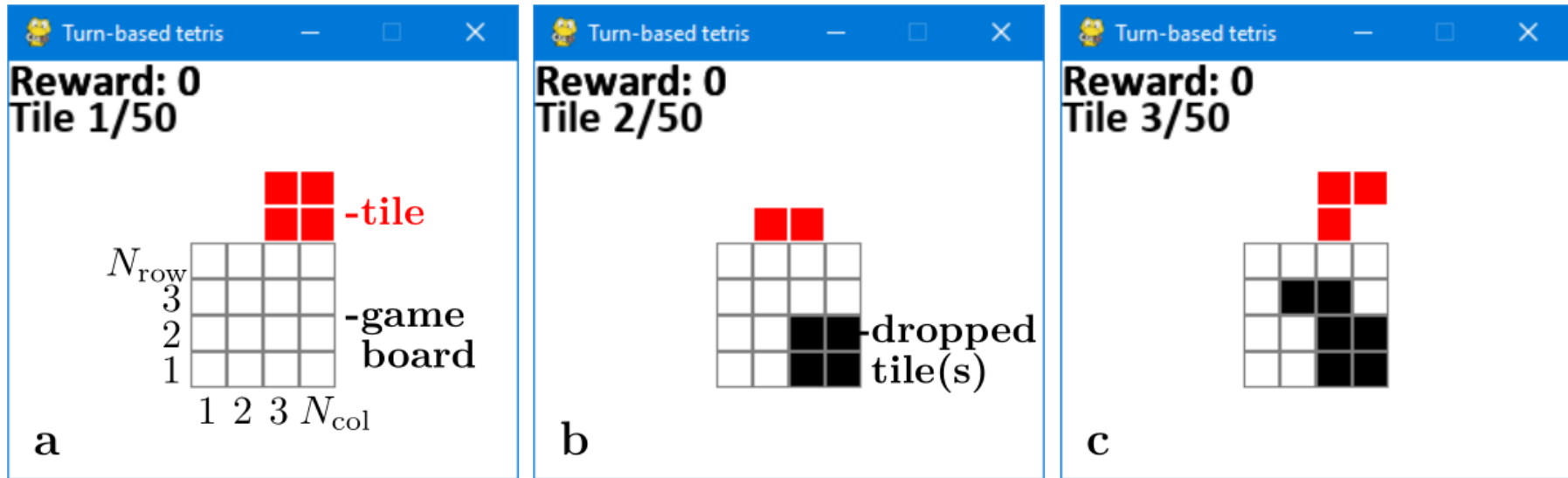
Reward is given for completing rows



Home problem B: Getting started

- **Download the following files**
 - [Canvas/files/Homeworks/HW B Kristian/tetris.py](#)
 - [Canvas/files/Homeworks/HW B Kristian/gameboardClass.py](#)
 - [Canvas/files/Homeworks/HW B Kristian/agentClass.py](#)
- **Make sure you have all required packages installed**
 - numpy
 - pygame
 - h5py
 - torch/tensorflow
- **Run tetris.py**
 - Use original parameters to test game
 - Set `human_player=0` and `param_set` to one of `PARAM_TASK1a`, `PARAM_TASK1b`, `PARAM_TASK1c`, `PARAM_TASK1d`, `PARAM_TASK2a` (and optional `PARAM_TASK2b`) to address the corresponding task in [HomeworkB_RL.pdf](#)
- **Complete the code in [agentClass.py](#) to address the tasks in [HomeworkB_RL.pdf](#)**

Home problem B: Discussion



Tile set ( ,  ,  , )

Think about/Discuss with your neighbors

What are possible choices of states and actions?

What is good/bad with different choices of states and actions?

Agent The player

Environment The game (and player)

States ???


Actions ???

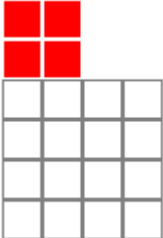
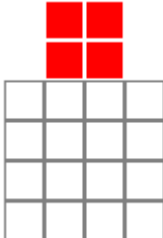
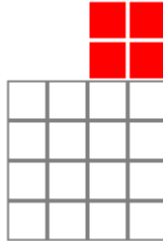
Reward + 10^{N-1} for completing N rows, -100 for game over


Home problem B: Actions

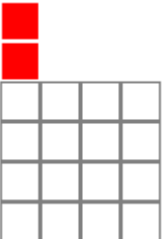
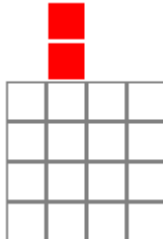
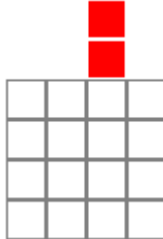
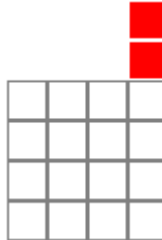
Use all allowed ways to place a tile as actions

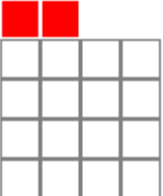
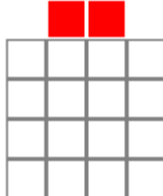
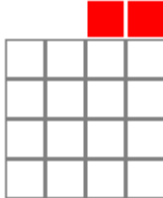
- Tile orientation
- Horizontal position of left tile edge

Example  has 3 actions:

Example  has 7 actions:

 has 12 actions and  has 6 actions

Home problem B: States

States that should be used are a combination of the following

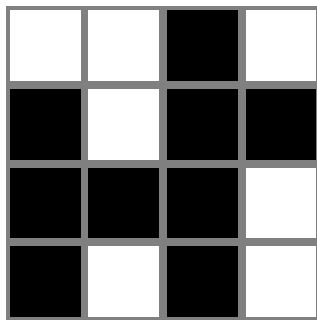
- Tile identifier of tile to be placed



Tile identifier 0 1 2 3

- Binary matrix representation of game board occupation

Example

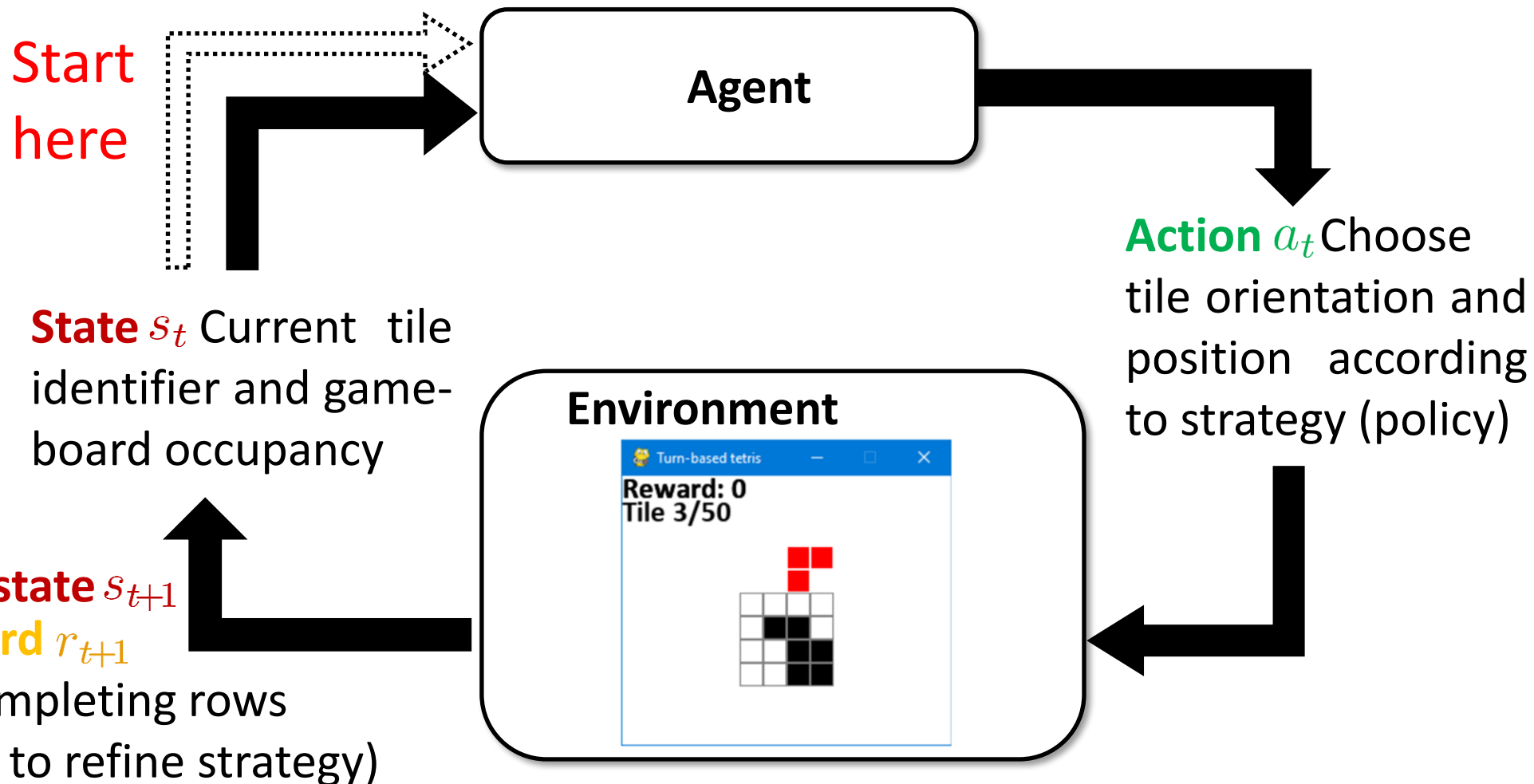


Game board

$$\begin{pmatrix} -1 & -1 & 1 & -1 \\ 1 & -1 & 1 & 1 \\ 1 & 1 & 1 & -1 \\ 1 & -1 & 1 & -1 \end{pmatrix}$$

Matrix representation

Home problem B: Reinforcement learning



Repeat until game over or a predefined number of moves (episodic task)

$s_0, a_0, r_1, s_1, a_1, \dots, s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}, \dots, r_T, s_T$

Reinforcement learning algorithms aim to find strategy/policy (choice of action given state) that maximizes long term reward.

Optimal policy in Markovian systems

Policy Method of choosing a given s

Defined as probability $\pi(a|s)$

Examples

- Always choose the same action $a^{(1)}$: $\pi(a|s) = \delta_{a,a^{(1)}}$
- Choose action randomly: $\pi(a|s) = \frac{1}{N_a}$ with $a \in \{a^{(1)}, a^{(2)}, \dots, a^{(N_a)}\}$
- Use the optimal policy: $\pi^*(a|s) = \delta_{a,a_*}$ where $a_* = \operatorname{argmax}_{a'} Q^*(s, a')$ and $Q^*(s_t, a_t) = \langle r_{t+1} + r_{t+2} + \dots + r_T \rangle$ is the expected future reward when following the optimal policy.

Goal Find $Q^*(s, a) \Rightarrow$ optimal policy $\pi^*(a|s)$

Brute force Find $Q^*(s, a)$ by evaluating $Q_\pi(s, a)$ of all policies π

$$\max_a Q^*(s, a) \geq \max_a Q_\pi(s, a)$$

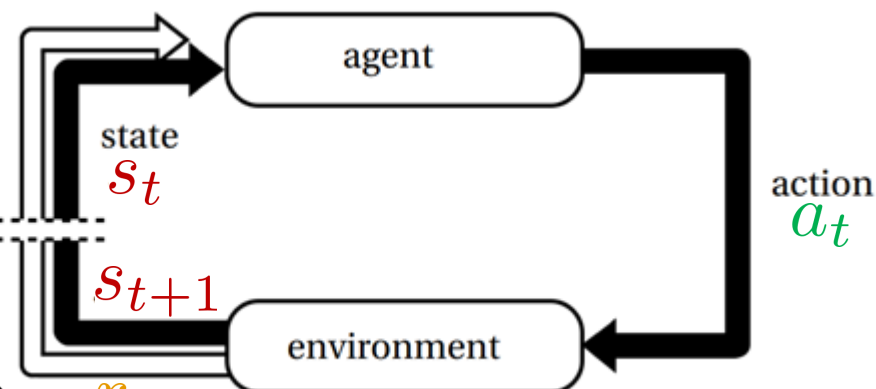
Problem: $N_a^{N_s}$ different policies if N_s states and N_a actions

Optimal solution using Q-learning

Optimal policy $a_t = \operatorname{argmax}_{a'} Q^*(s_t, a')$

Expected future reward

$$\begin{aligned} Q^*(s_t, a_t) &= \langle r_{t+1} + r_{t+2} + \dots + r_T \rangle^{r_t} \\ &= \langle r_{t+1} \rangle + \langle r_{t+2} + \dots + r_T \rangle \\ &= \langle r_{t+1} \rangle + \max_{a'} Q^*(s_{t+1}, a') \end{aligned}$$



Q-learning algorithm (aims to find approximate $Q(s, a) \approx Q^*(s, a)$)

1. Start with arbitrary Q-table $Q_0(s, a)$ (size $N_s \times N_a$)
2. Evaluate initial state s_t (with $t = 0$)
3. Choose action a_t according to policy, e.g. ϵ -greedy
4. Apply a_t and evaluate s_{t+1} and r_{t+1}
5. Update $Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha \Delta Q$

$$\Delta Q = r_{t+1} + \max_{a'} Q_t(s_{t+1}, a') - Q_t(s_t, a_t)$$

6. Repeat from 3. until episode is finished
7. Repeat from 2. (keeping Q) until convergence

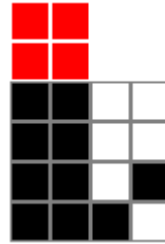
Q-learning Example I

$$\Delta Q = r_{t+1} + \max_{a'} Q_t(s_{t+1}, a') - Q_t(s_t, a_t)$$

- First time state \hat{s} is encountered

$$Q(\hat{s}, a) = (0, 0, 0)$$

Choose random action on tie



Chosen action a_1



Game over

Reward $r = -100$

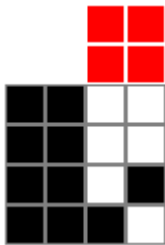
$$Q(\hat{s}, a_1) = -100\alpha = -20$$

(assume $\alpha = 0.2$)

- Second time \hat{s} is encountered

$$Q(\hat{s}, a) = (-20, 0, 0)$$

Choose a_2 or a_3 randomly



Chosen action a_3



Two lines completed

Reward $r = +10$

$$Q(\hat{s}, a_1) = 10\alpha = 2$$

- Third time \hat{s} is encountered $Q(\hat{s}, a) = (-20, 0, 2) \longrightarrow$ action a_3

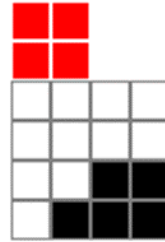
Q-learning Example 2

$$\Delta Q = r_{t+1} + \max_{a'} Q_t(s_{t+1}, a') - Q_t(s_t, a_t)$$

- First time state \hat{s} is encountered

$$Q(\hat{s}, a) = (0, 0, 0)$$

Choose random action on tie



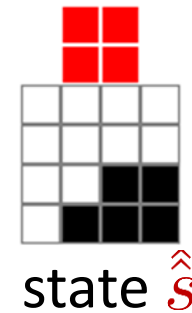
Chosen action a_1



One line is completed

Reward $r = 1$

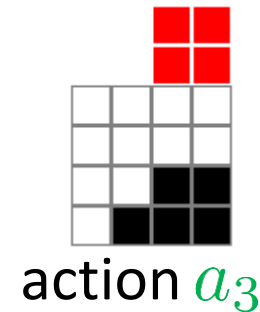
$$Q(\hat{s}, a_1) = 1\alpha = 0.2$$



- Second time \hat{s} is encountered

$$Q(\hat{s}, a) = (0.2, 0, 0) \longrightarrow \text{action } a_1$$

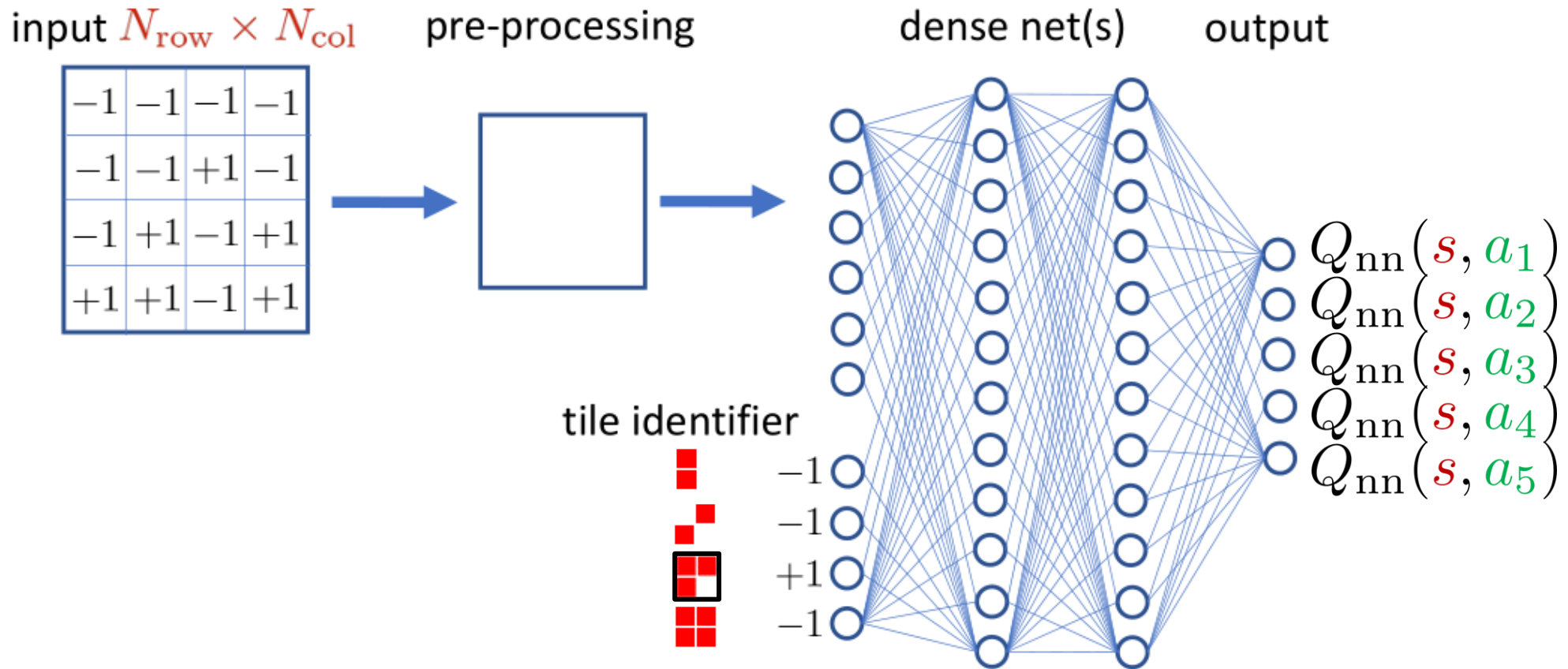
Problem action a_3 is better in the long run



Solution Add exploration, e.g. ϵ -greedy policy

Optimal solution using Deep Q-networks

Idea Replace Q-table by deep neural network



Policy $a_t = \begin{cases} \text{random valid action} & \text{with prob. } \epsilon_E \\ \operatorname{argmax}_{a'} Q_{\text{nn}}(s_t, a') & \text{otherwise} \end{cases}$

$$\epsilon_E = \max(\epsilon, 1 - E/E_0), \text{ episode } E$$

Optimal solution using Deep Q-networks

Minimize $|\Delta Q(T_t)| = |r_{t+1} + \max_{a'} Q_{nn}(s_{t+1}, a') - Q_{nn}(s_t, a_t)|$
 for transitions $T_t = \{s_t, a_t, r_{t+1}, s_{t+1}\}$

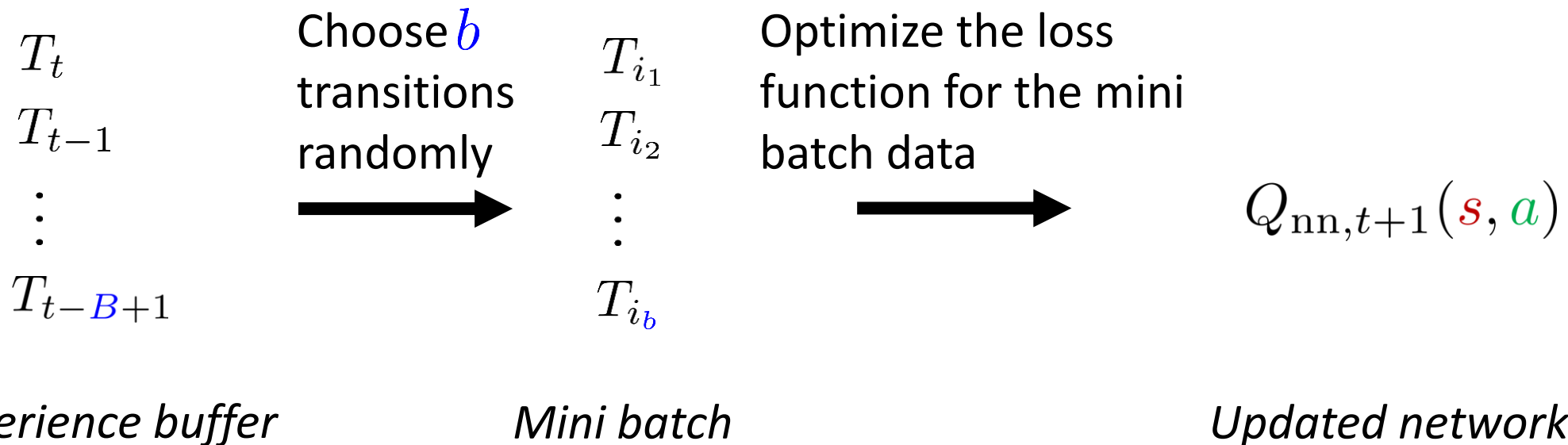
Loss function $\sum_{\tau} \Delta Q(T_{\tau})^2$

Problem 1

Data is not independent and identically distributed (needed for optimizing the loss function)

Solution

Use 'experience replay buffer' containing B last state transitions



Optimal solution using Deep Q-networks

Minimize $|\Delta Q(T_t)| = |r_{t+1} + \max_{a'} Q_{nn}(s_{t+1}, a') - Q_{nn}(s_t, a_t)|$
for transitions $T_t = \{s_t, a_t, r_{t+1}, s_{t+1}\}$

Loss function $\sum_{\tau} \Delta Q(T_{\tau})^2$

Problem 2

Potential instability due to ‘bootstrapping’: we use Q_{nn} to estimate the expected future reward when updating Q_{nn} .

Solution

Introduce ‘target network’ $\hat{Q}_{nn}(s, a)$ when estimating future reward

$$\Delta Q(T_t) = r_{t+1} + \max_{a'} \hat{Q}_{nn}(s_{t+1}, a') - Q_{nn}(s_t, a_t)$$

$\hat{Q}_{nn}(s, a)$ is a copy of Q_{nn} , but only occasionally updated.

Optimal solution using Deep Q-networks

Deep Q-network algorithm

1. Start with arbitrary and identical Q-networks $Q_{nn}(s, a)$ and $\hat{Q}_{nn}(s, a)$
2. Evaluate initial state
3. Choose action based on $Q_{nn}(s, a)$ and ϵ_E
4. Apply a_t and evaluate s_{t+1} and r_{t+1}
5. Store transition $T_t = \{s_t, a_t, r_{t+1}, s_{t+1}\}$ in replay buffer
6. Sample mini-batch of transitions from the replay buffer
7. For each transition in the mini-batch, use $\hat{Q}_{nn}(s, a)$ to calculate target value: $y = r$ if s_{t+1} is terminal,
 $y = r + \max_{a'} \hat{Q}_{nn}(s_{t+1}, a')$ otherwise.
8. Calculate loss $\mathcal{L} = (y - Q_{nn}(s, a))^2$
9. Update $Q_{nn}(s, a)$ using for example 'Adam' (built-in) to minimize loss
10. Every 100 episode copy weights from $Q_{nn}(s, a)$ to $\hat{Q}_{nn}(s, a)$
11. Repeat from 3. until convergence

Assignments

Complete the classes for Q-learning and Deep Q-networks in [agentClass.py](#)

Address the following tasks

- 1.a) [3p] Use Q-learning to train an artificial player on deterministic tile sequence. Use greedy policy.
- 1.b) [1p] Same as 1.a) with ϵ - greedy policy.
- 1.c) [1p] Same as 1.b) with random tile sequence.
- 1.d) [1p] Discuss possibility to scale up method to a larger game board

- 2.a) [4p] Solve task 1.c) using deep Q-networks.
- 2.b) [Optional task, 2p] Solve the problem on larger game board using deep Q-networks