ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA
DEPARTMENT OF ELECTRICAL, ELECTRONIC AND INFORMATION ENGINEERING
MASTER'S DEGREE IN AUTOMATION ENGINEERING

# A Deep Reinforcement Learning approach based on policy gradient for Mobile Robot Navigation

*Candidate:*
Enrico Pianazzi
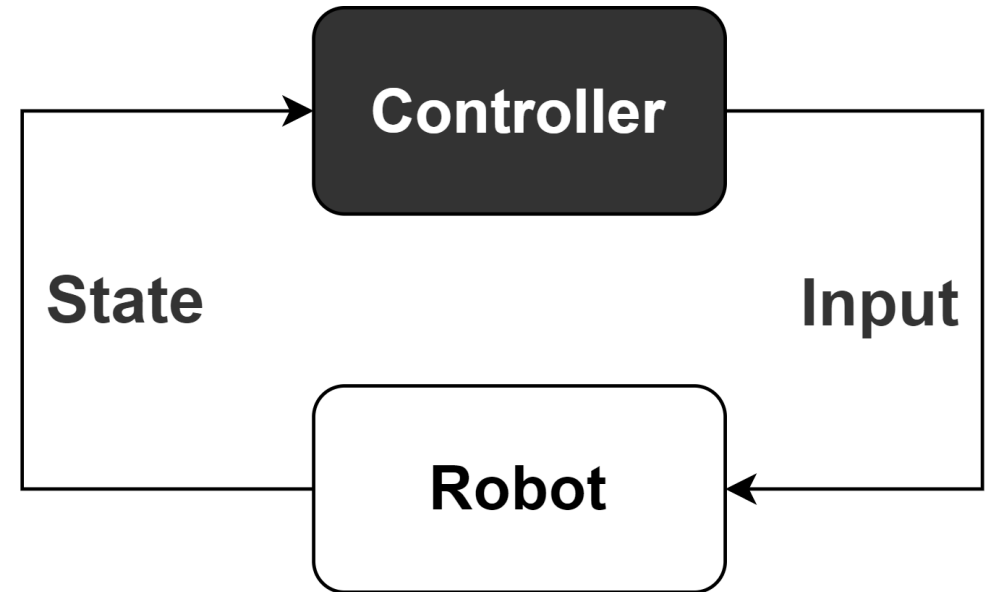
*Advisor:*
Prof. Giuseppe Notarstefano

*Co-Advisors:*
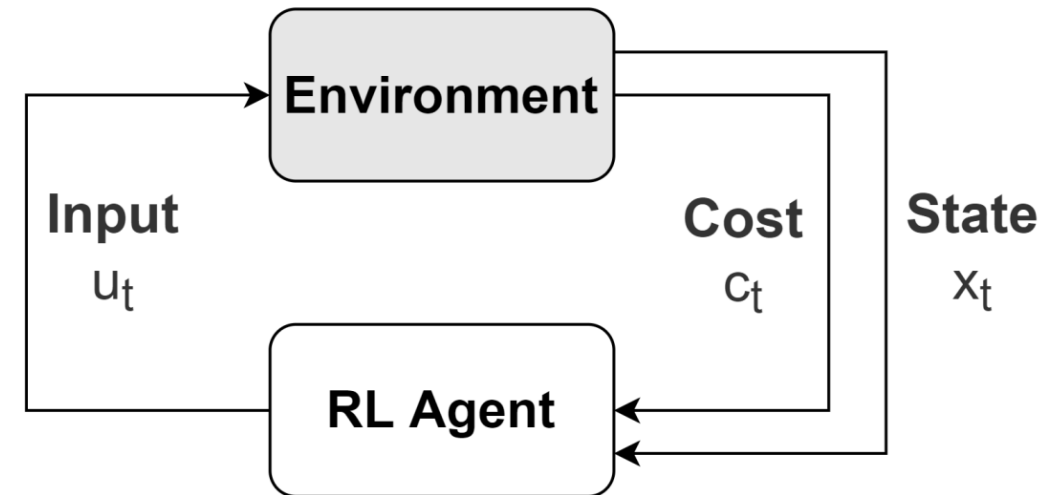Dott. Andrea Camisa
Ing. Lorenzo Sforni

# Motivations

- Autonomous navigation features in mobile robots are highly valuable

- Reinforcement Learning can achieve complex data-driven control

- State feedback control without the need of complex sensory data

# Reinforcement Learning

- Optimize a policy (controller)

  $\mu(x_t) = u_t$ to minimize future cost

- Only experience cost and state signals
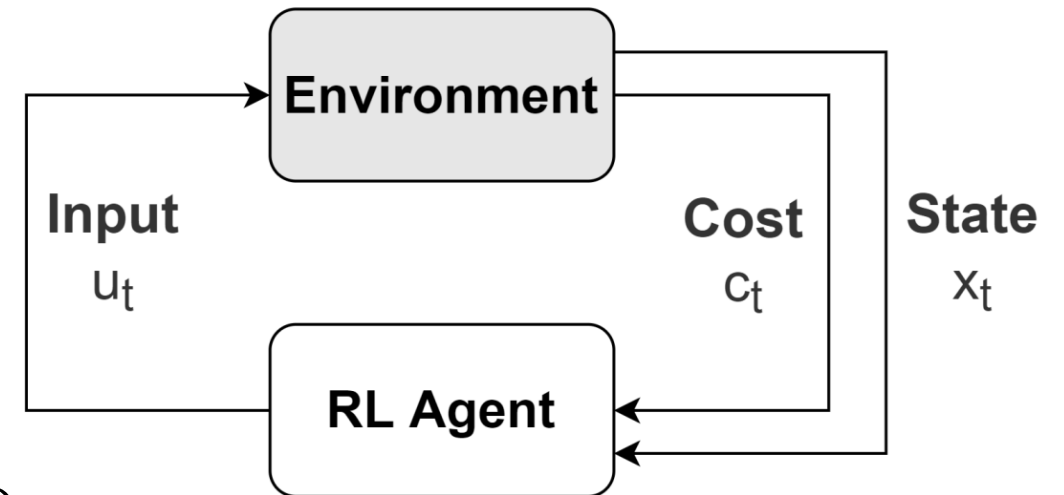
- Division in discrete time steps and

  episodes

# Q-function and Bellman's equation

- Q-function: estimates future cost

- Bellman's equation:

$$Q(x_t, u_t) = c(x_t, u_t) + \gamma Q(x_{t+1}, \mu(x_{t+1}))$$

- Temporal Difference error:

$$TD_{er} = c(x_t, u_t) + \gamma Q(x_{t+1}, \mu(x_{t+1})) - Q(x_t, u_t)$$
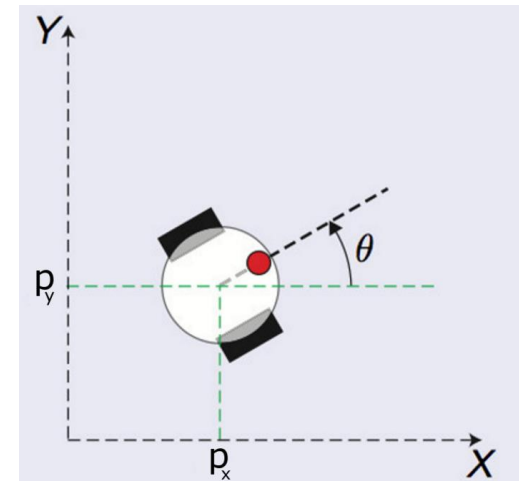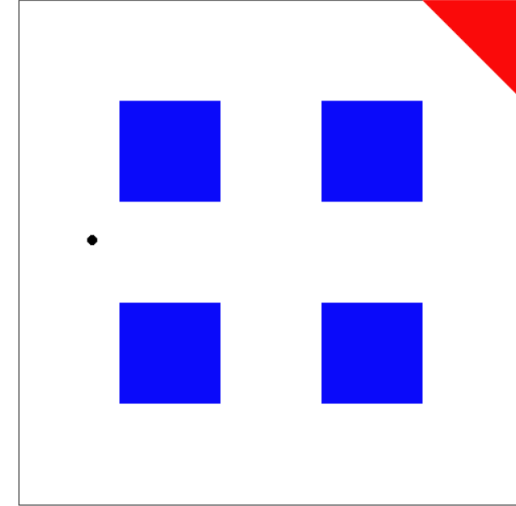
# Problem Statement

- Task: reach goal area from any initial configuration in the environment

- Environment limits, obstacles are unknown

- Continuous states and inputs
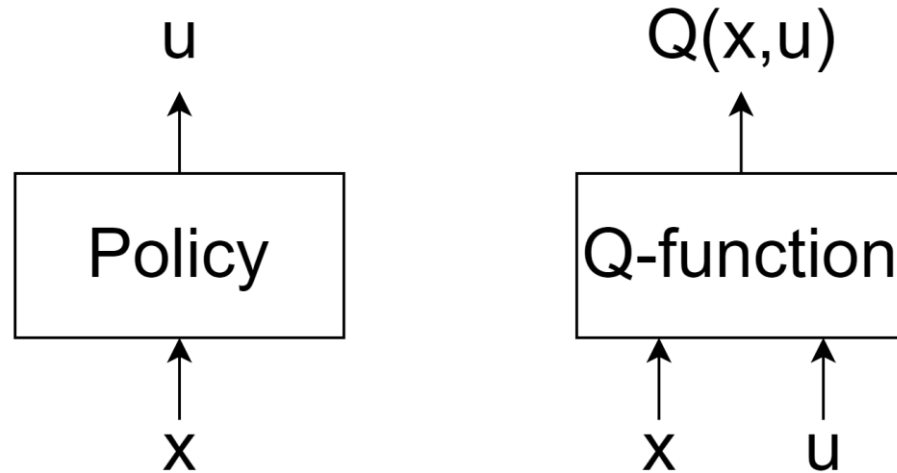
- Robot simulated as a unicycle, but model is unknown

State: $x = [p_x \ p_y \ \theta]$

Model: $\dot{x} = \begin{bmatrix} \cos(\theta) & 0 \\ \sin(\theta) & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} v \\ \omega \end{bmatrix}$

# Deep Deterministic Policy Gradient (DDPG)

- Policy: controller
- Q-function: measures the cost of a state-input pair
- Implemented as deep neural networks

- Experience Replay: keeps a buffer of past transitions {state, input, cost, next state}
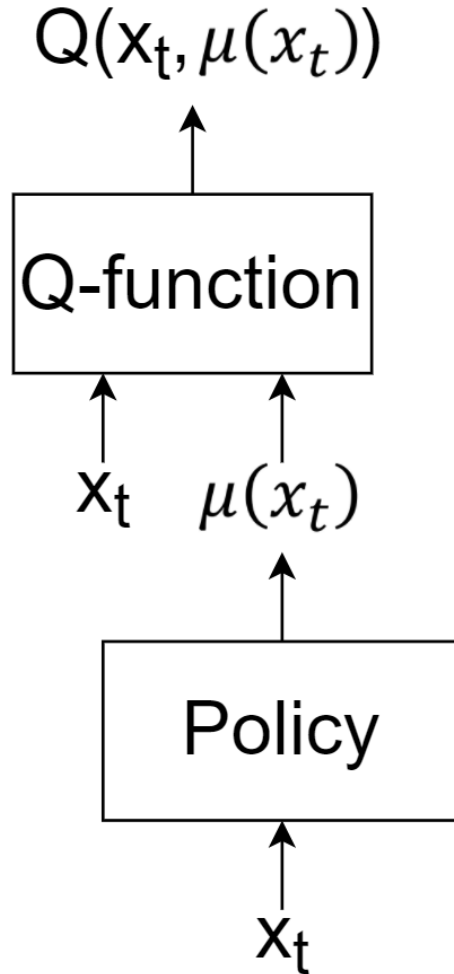


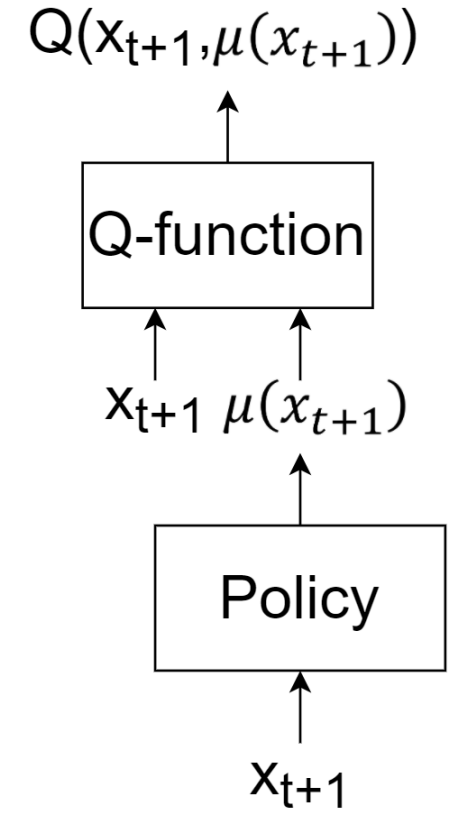| Replay Buffer |
|---|
| $\{x_1, u_1, c(x_1, u_1), x_2\}$ |
| $\{x_2, u_2, c(x_2, u_2), x_3\}$ |
| $\cdots$ |
| $\{x_\tau, u_\tau, c(x_\tau, u_\tau), x_{\tau+1}\}$ |

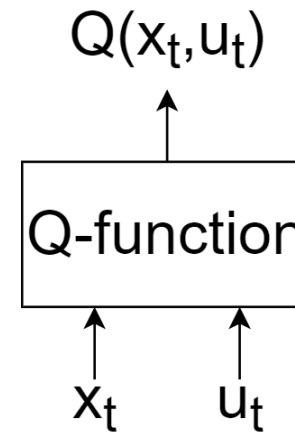# Policy Training

$$Q(x_t, \mu(x_t))$$



- Trained to choose inputs that minimize cost

- Policy network parameters: $\theta_k^\mu$

- Policy gradient update:

$$\theta_{k+1}^\mu = \theta_k^\mu - \alpha_\mu \nabla_{\theta^\mu} Q(x_t, \mu(x_t))$$

# Q-function Training

- Sampled transition: $\{x_t, u_t, c(x_t, u_t), x_{t+1}\}$

- Q-function network parameters: $\theta_k^Q$

- Gradient descent update minimizing the Temporal Difference error:

$$\theta_{k+1}^Q = \theta_k^Q - \alpha_Q \nabla_{\theta^Q} [\underbrace{c(x_t, u_t) + \gamma Q(x_{t+1}, \mu(x_{t+1}))}_{\textbf{Target of the update}} - \underbrace{Q(x_t, u_t)}_{\textbf{Prediction}}]^2$$

$Q(x_t, u_t)$

Q-function

$x_t \quad u_t$

$Q(x_{t+1}, \mu(x_{t+1}))$

Q-function

$x_{t+1} \; \mu(x_{t+1})$

Policy

$x_{t+1}$

# Target Networks

- Copies of the main networks: $Q'(x, u)$ and $\mu'(x)$, with parameters: $\theta^{Q'}$ and $\theta^{\mu'}$
- Slowly track the main networks parameters:

$$\begin{cases} \theta^{Q'} = \tau\theta^Q + (1-\tau)\theta^{Q'} \\ \theta^{\mu'} = \tau\theta^\mu + (1-\tau)\theta^{\mu'} \end{cases}$$

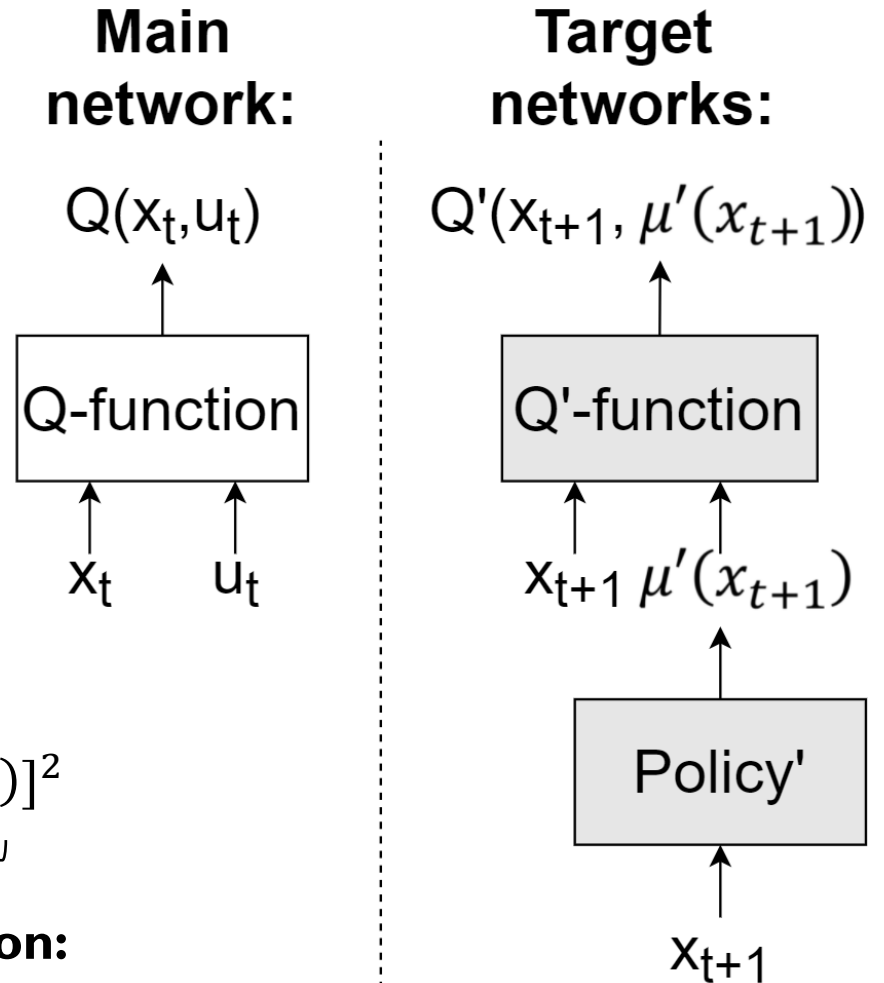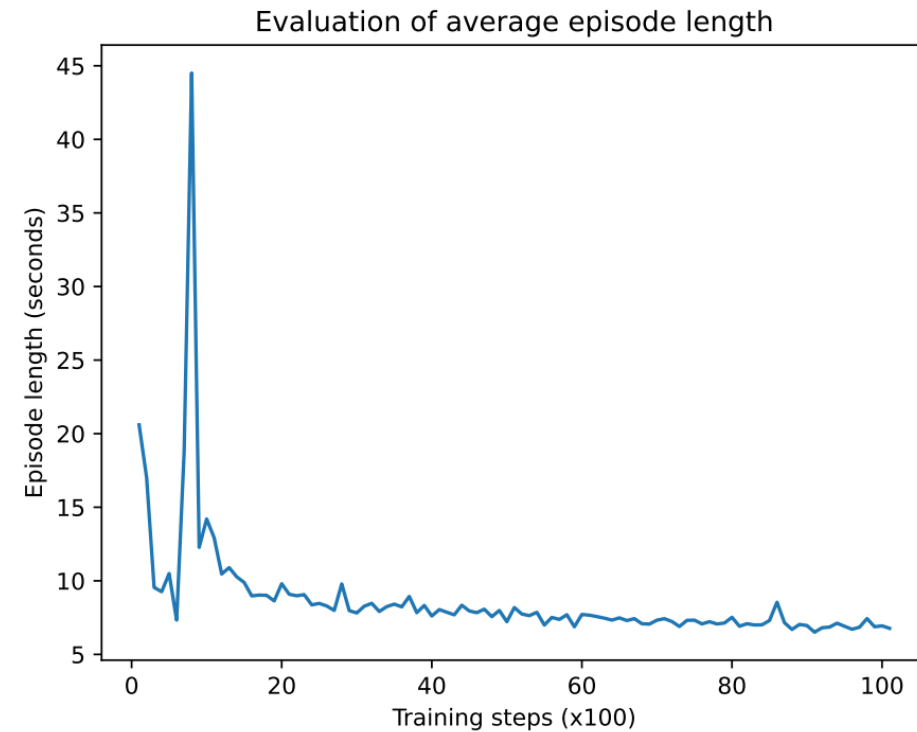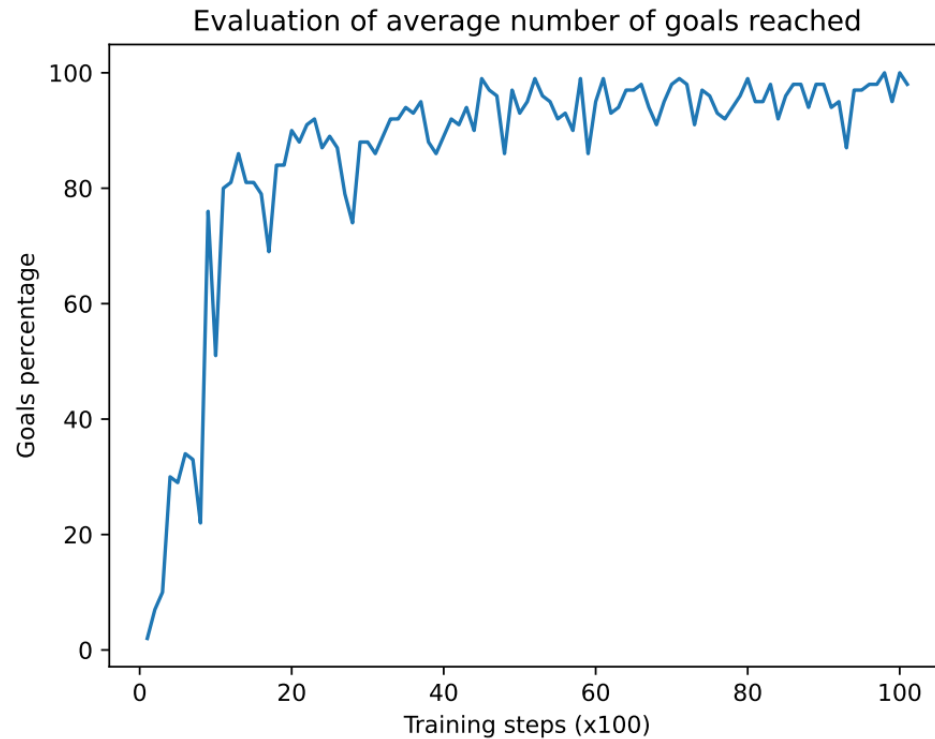With $\tau \ll 1$

- Only used for target computation of the TD error:

$$\theta^Q_{k+1} = \theta^Q_k - \alpha_Q \nabla_{\theta^Q}[c(x_t, u_t) + \gamma \boldsymbol{Q}'(x_{t+1}, \boldsymbol{\mu}'(x_{t+1})) - Q(x_t, u_t)]^2$$

**Target of the update: Target networks**

**Prediction: main network**

**Main network:**

$Q(x_t, u_t)$

Q-function

$x_t$ $u_t$

**Target networks:**

Q'$(x_{t+1}, \mu'(x_{t+1}))$

Q'-function

$x_{t+1}$ $\mu'(x_{t+1})$

Policy'

$x_{t+1}$

# Controller Performance Evolution

- Experiments from a set of fixed initial conditions
- 10000 episodes training



Evaluation of average number of goals reached
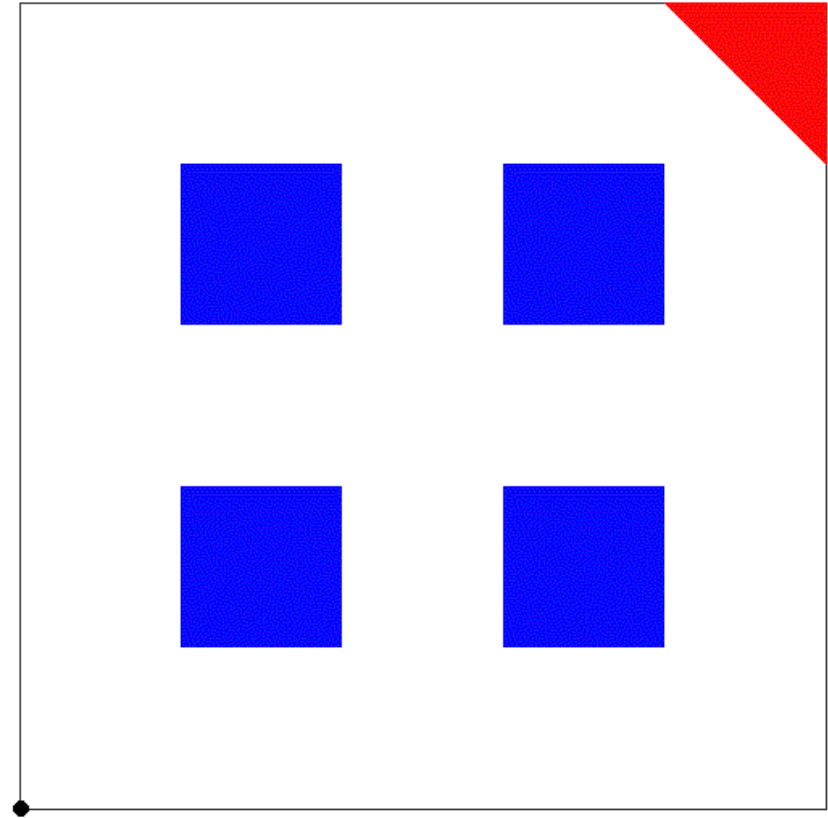


Evaluation of average episode length

# Resulting Robot Trajectories

10Hz simulation:

100Hz simulation:

# Conclusions

- Achieved above 95% success rate in the autonomous goal-reaching task

- Provided a custom Python implementation of DDPG and the unicycle simulation

- Drawback: length of training

- Future work: algorithm tuning; distributed setting; integration with model-based techniques