# Deep Learning:
# Project presentation

## Deep Reinforcement Learning with JAVAgario

DEBES Baptiste

NELISSEN Louis

WEYDERS Pierre-Francois

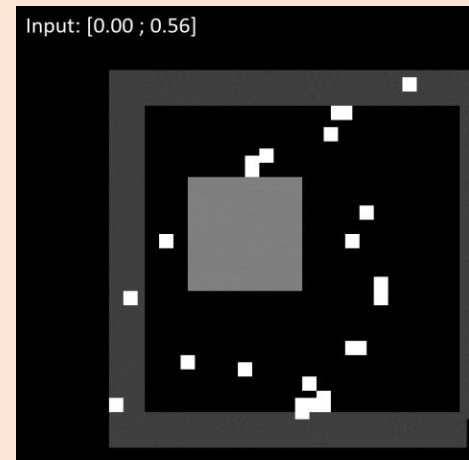LIÈGE université

August 2020

# The Game: Agar.io

- Popular Massively Multiplayer browser game

- Player is a cell that eats pellets

- Constantly shrinks, but grows when eating

  pellets
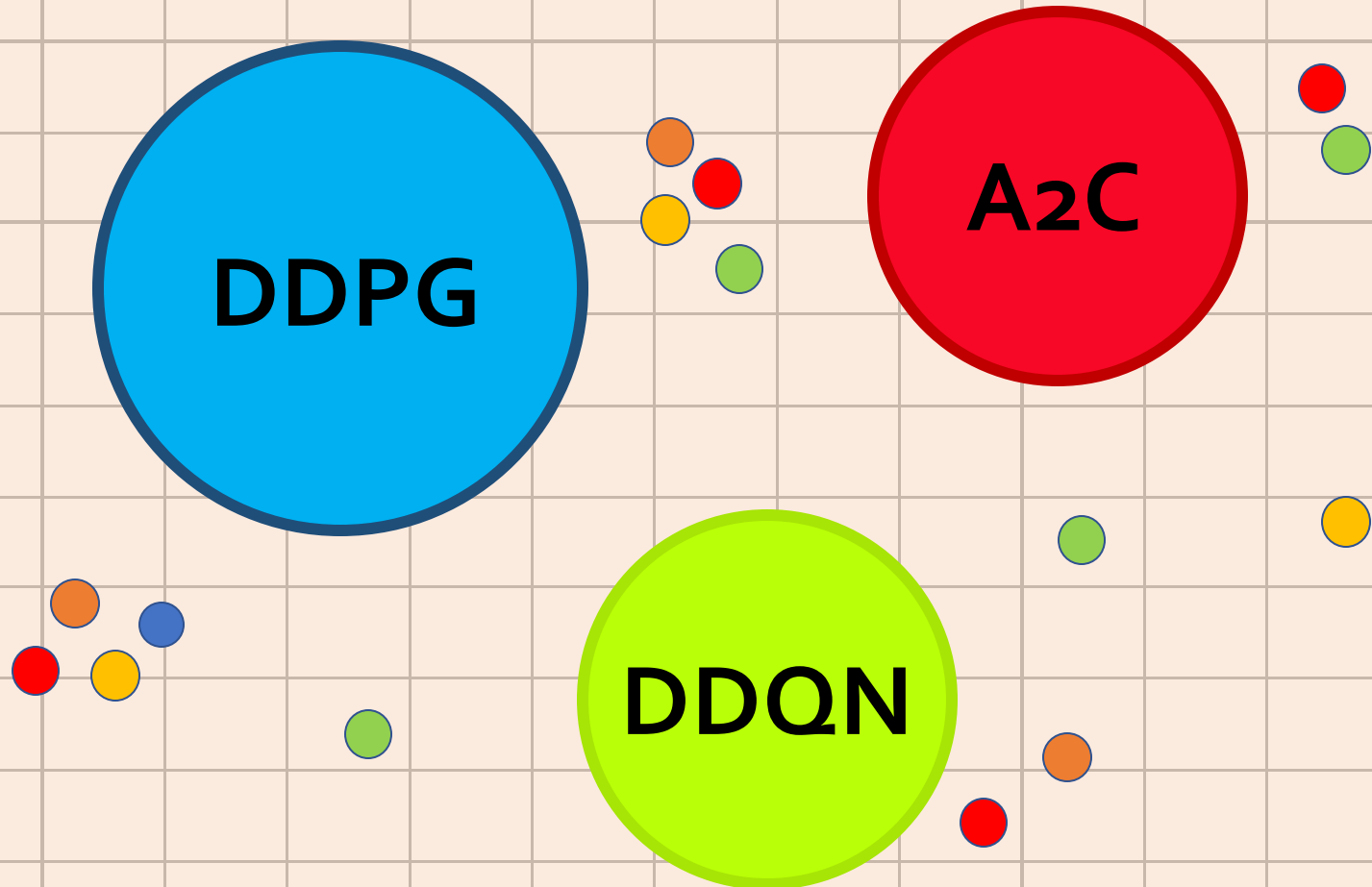
- Controlled with mouse

# The Game: JAVAgar.io

- Coded in JAVA

- Greyscale image

- Ability to launch our own servers

- Simplified rules
  - One agent
  - One box (small)
  - 20 pellets

- Objective: get all the pellets as fast as possible without dying



Input: [0.00 ; 0.56]

# The Theory: Q-Learning

- Branch of reinforcment learning

- Model-free

- Value based learning

  - Q → (**s**tate, **a**ction)

- $Q$-learning aims to approximate the optimal $Q^*$

- Deep Q-learning: use NN as approximators

# The Theory: Policy Gradient

- Branch of reinforcment learning

- Model-free

- Seek an approximator $\pi_\theta \ (a|s)$ of the optimal policy
  $\pi^*(a|s)$

# The Theory: The RL family

**Value based** DDQN
- No policy (implicit)
- Value function

**Policy based**
- Policy
- No value function
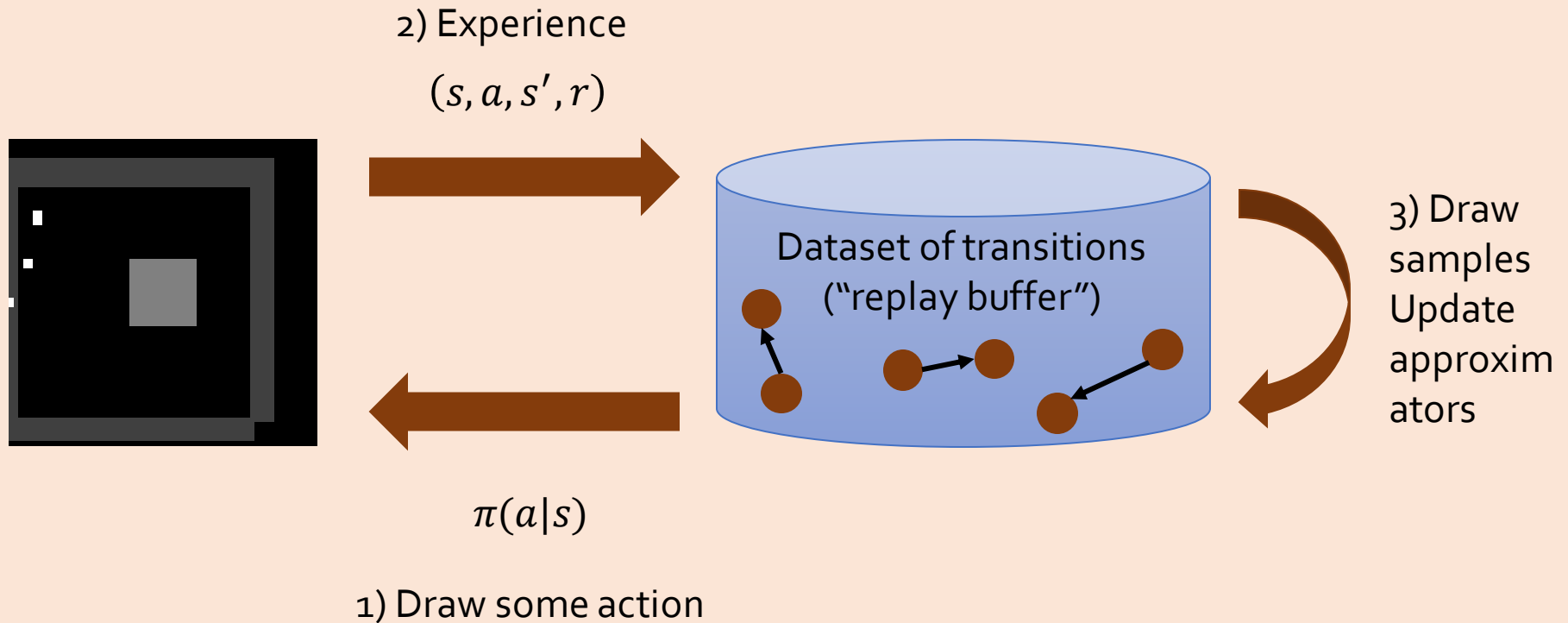
**Actor critic** DDPG A2C
- Policy
- Value function

**Model based**
- Transition and reward model

**Model free**
- No transition and no reward model (implicit)

2) Experience

$(s, a, s', r)$



Dataset of transitions ("replay buffer")

3) Draw samples Update approxim ators

$\pi(a|s)$

1) Draw some action

# The Theory: The losses

A2C: Advantage Actor Critic

- $L^{PG}(\theta) = E\left[-\log(\pi_\theta(a_n|s_n)A(s_n,a_n))\right]$

DDPG: Double DPG

- $L^{DDPG}(\theta) = -E\left[Q_\tau(s,a)\,|_{s=s_n,a_n=\pi_\theta(s_n)}\right]$
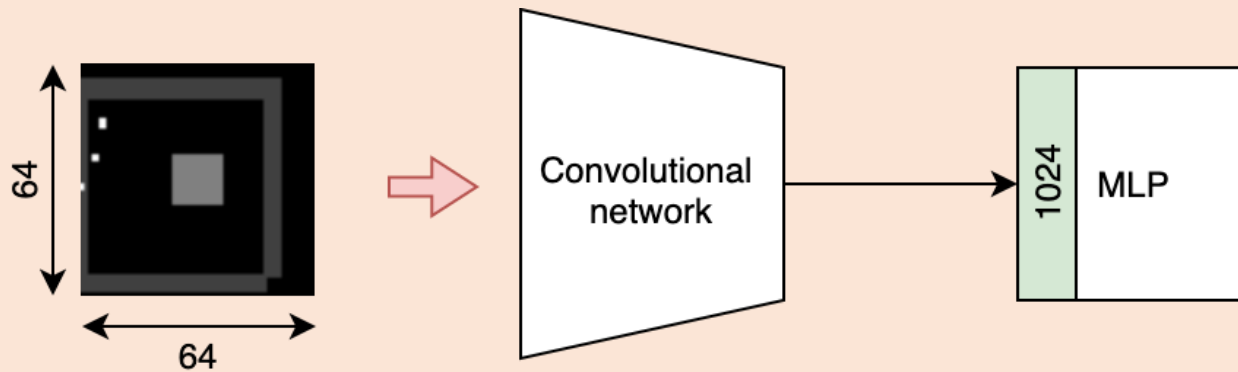
Caveat in the report

DDQN: Double Deep Q Network

- $L^{TD}(\theta) = \frac{1}{2}\left[Q_\theta(s,a) - r - \gamma \max_{\{a'\}} Q_{\overline{\theta}}(s',a')\right]^2$

- Batch normalization:
  - Observed in the DDPG paper
- Target networks:
  - Stabilizes training
  - Soft update equation

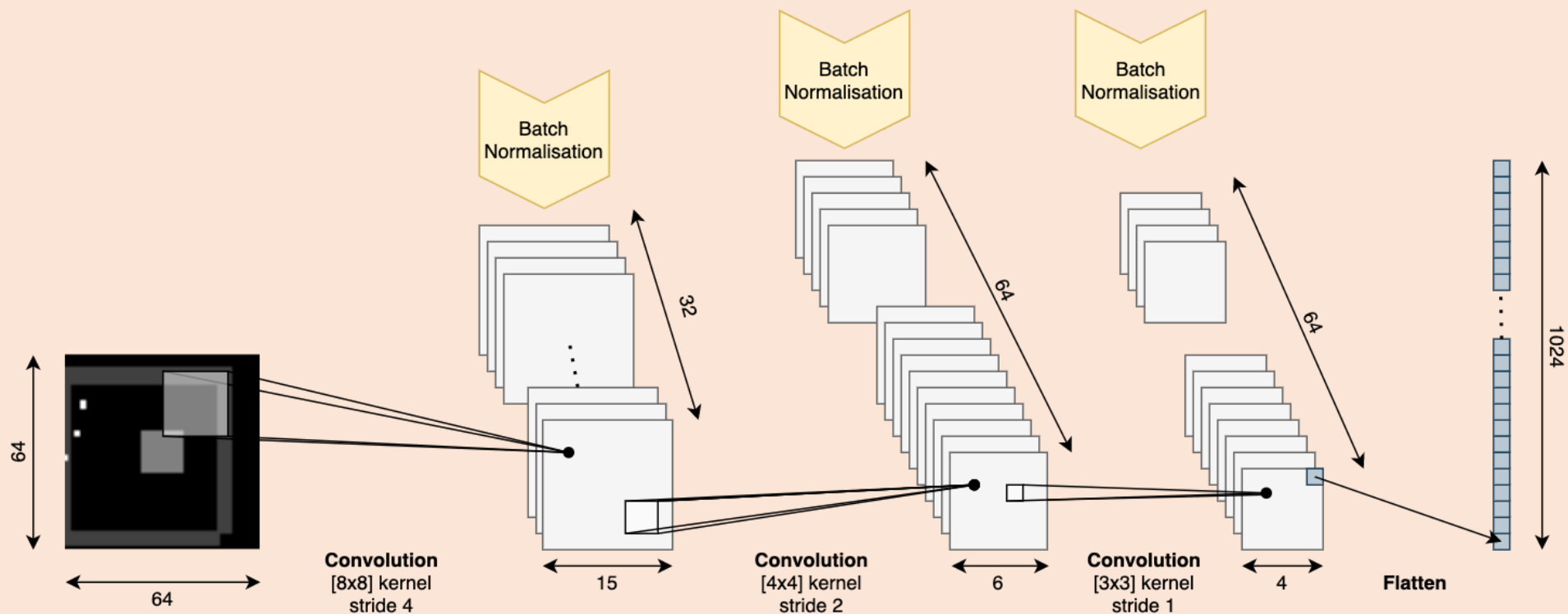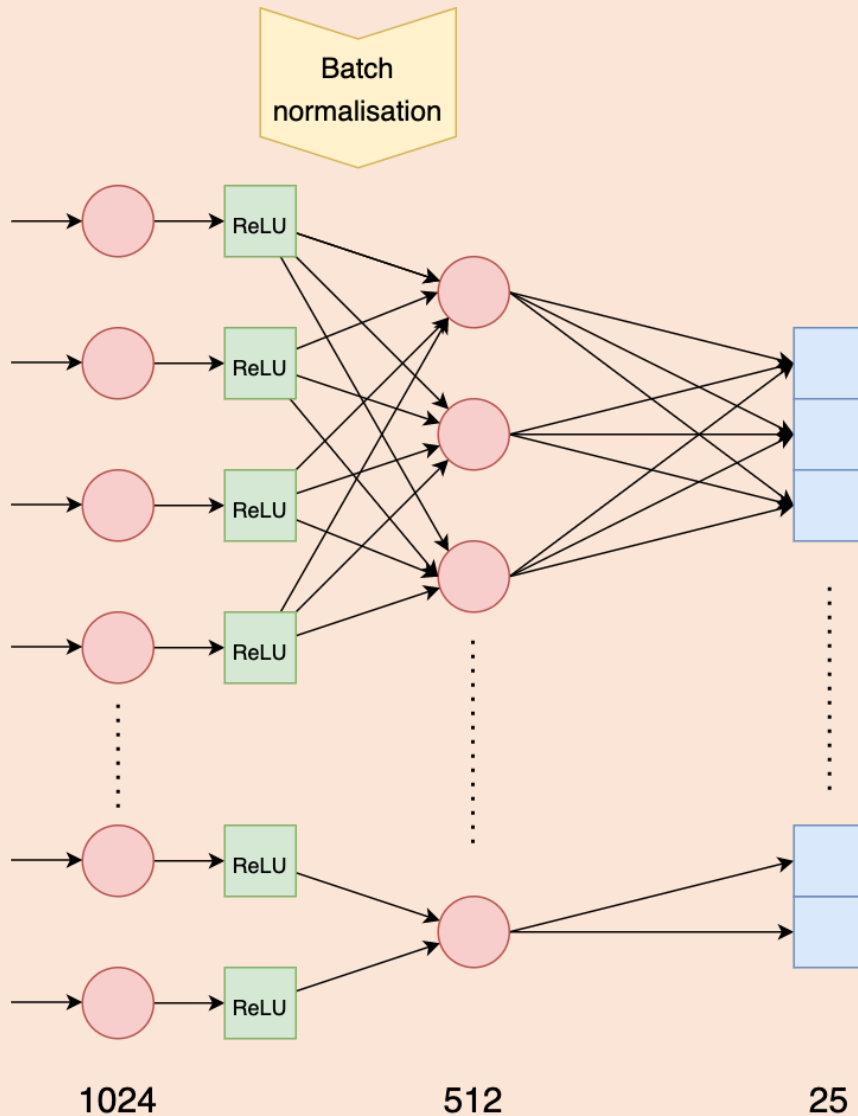$$\theta' \leftarrow \tau * \theta + (1 - \tau) * \theta' \text{ with } \tau \in [0,1]$$

- 3 Convolution layers
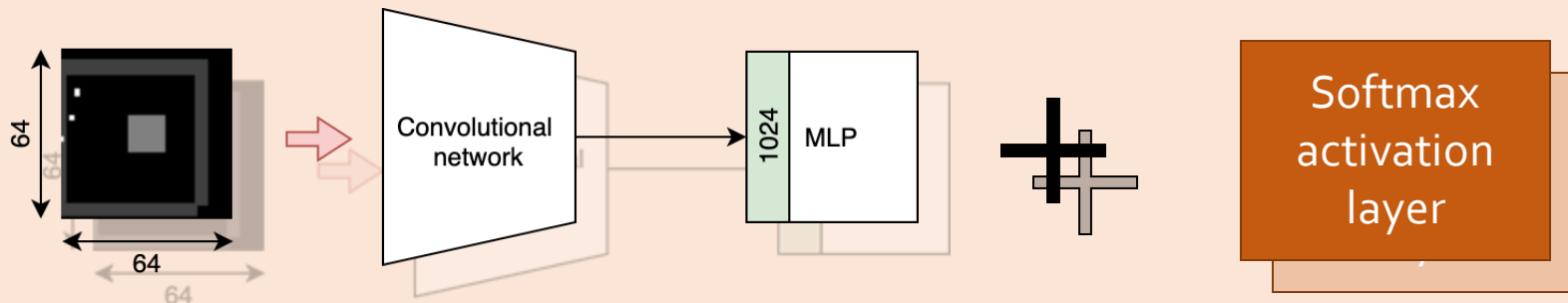- Batch normalisation

# The Methods: Multi Layer Perceptron



- 2 dense layers
- Batch normalisation

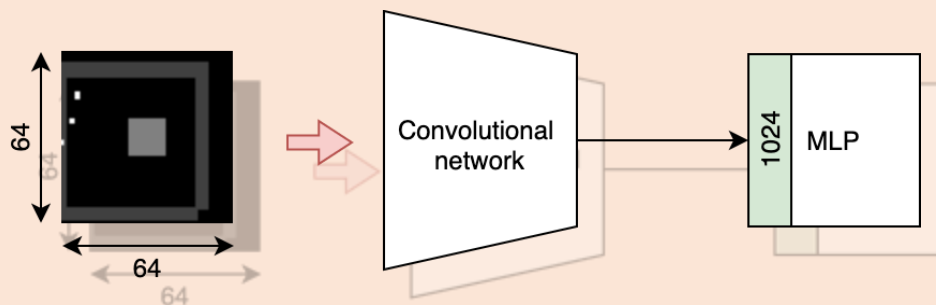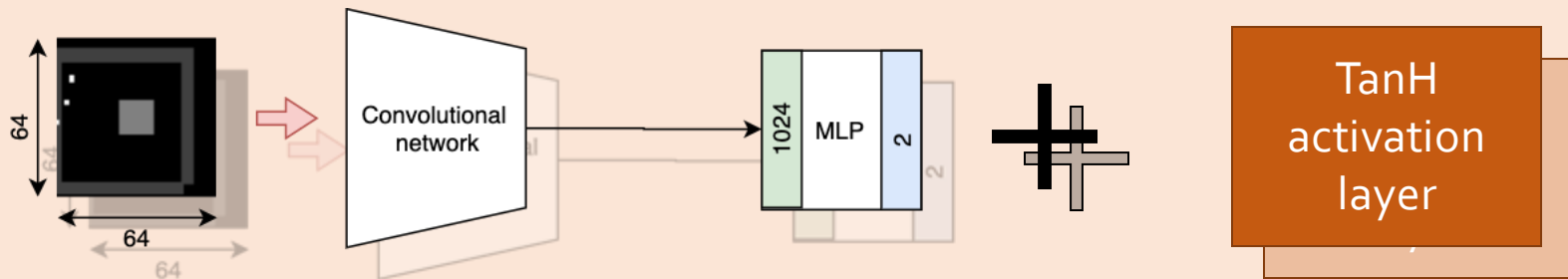**Different versions for the different algorithms**

## Actor:



## Critic:



| Parameters | Value |
|---|---|
| Learning rate | 0.0001 |
| Tau | 0.001 |
| Discount factor | 0.99 |

14

# The Methods: DDPG

## Actor:



TanH activation layer

## Critic:



| Parameters | Value |
|---|---|
| Learning rate critic | 0.0005 |
| Learning rate actor | 0.00001 |
| Tau | 0.001 |
| Discount factor | 0.99 |

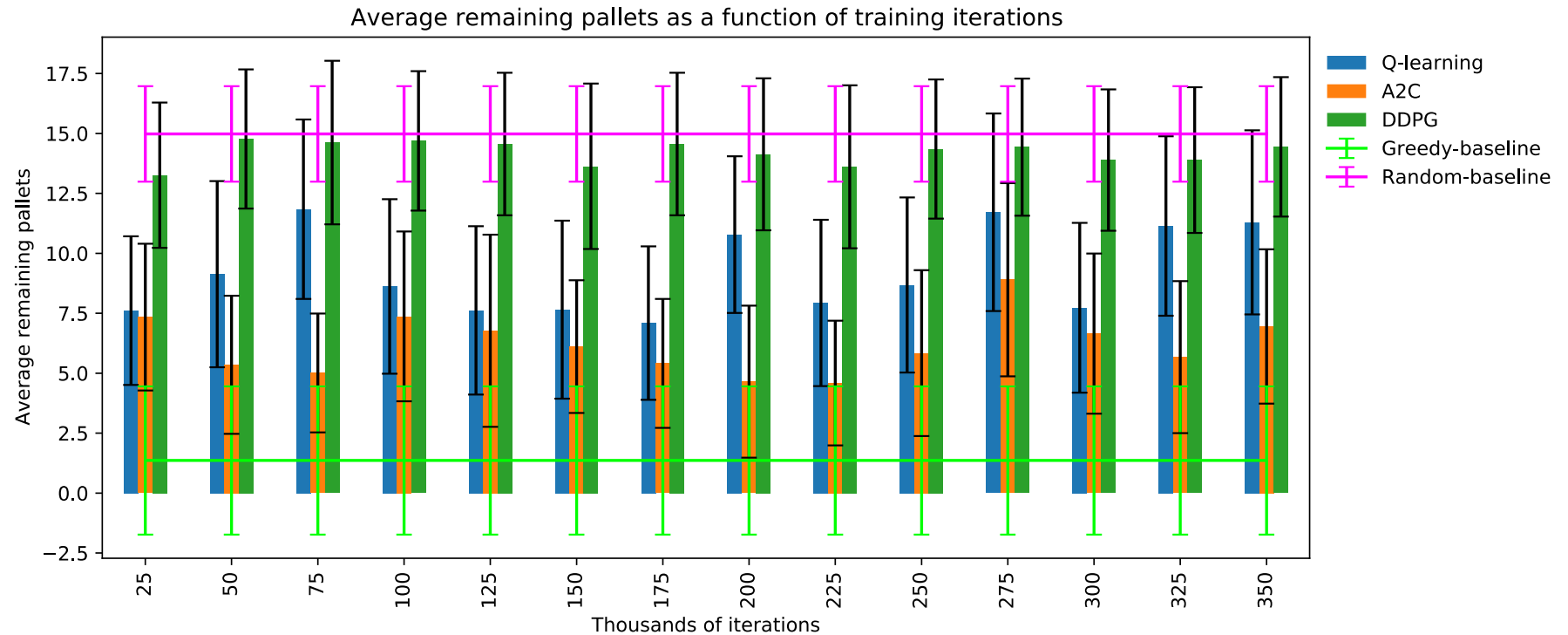| Parameters | Value |
| --- | --- |
| Learning rate | 0.0001 |
| Tau | 0.001 |
| Discount factor | 0.99 |

# The Methods: Training

- Random environments
  - Random pellet position
  - Random init position
- Constants:
  - Size of map
  - Number of pellets (20)
- Images of 64 x 64
- Episodes of 2500 steps
- Total of 350 000 frames
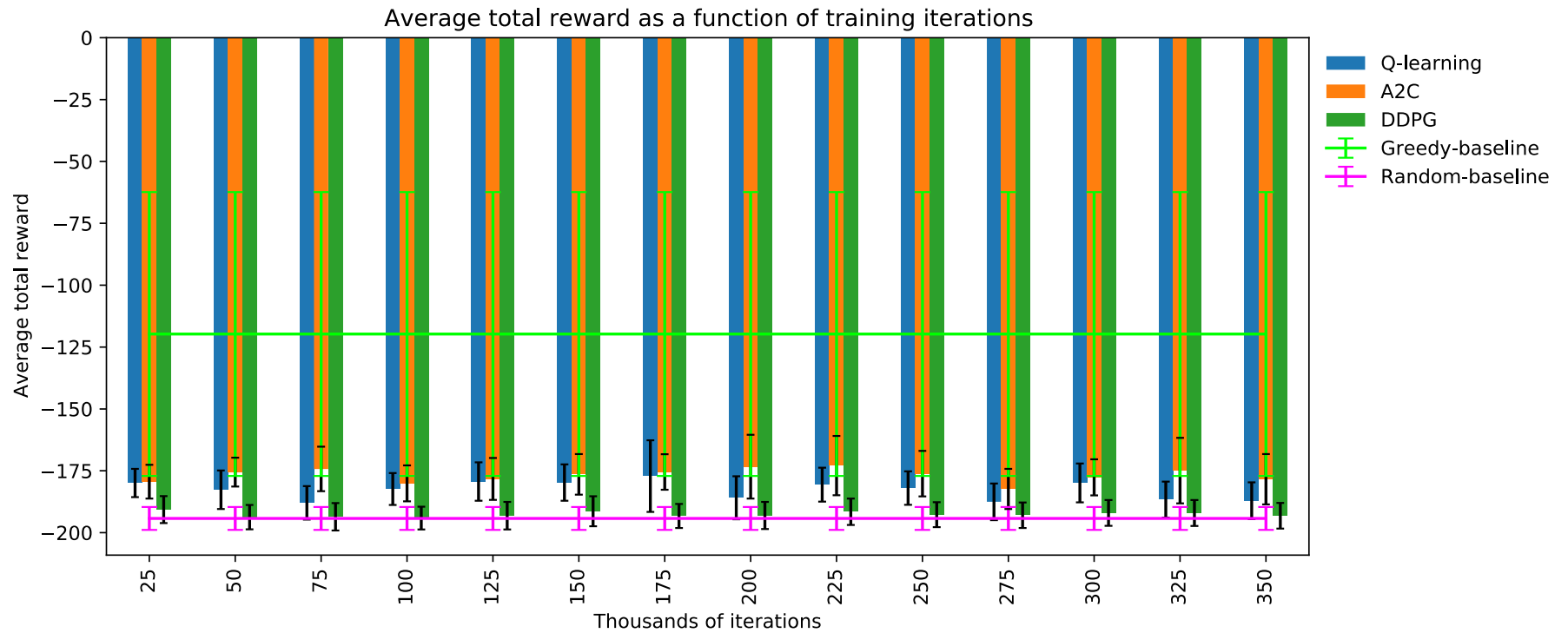- Replay buffer of size 65 000
- Adam optimizer

Hardware:
- *CPU*: I7-6700k 4GHZ (4 cores)
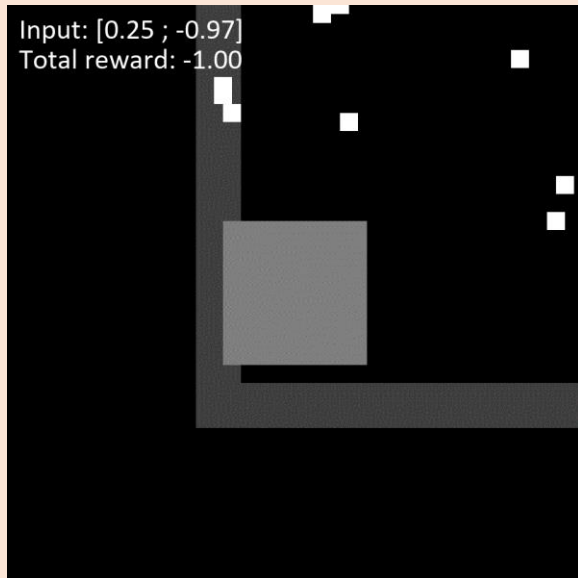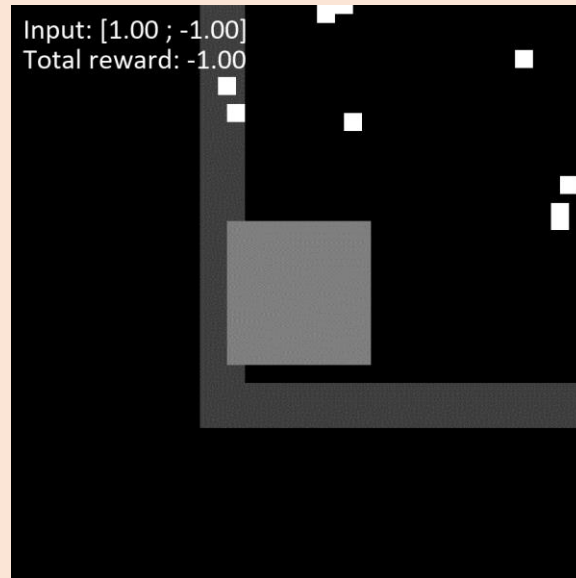- *GPU*: NVIDIA GTX 1070 (training was ~ 30 % faster)

17

# The Results



Average remaining pallets as a function of training iterations

# The Results



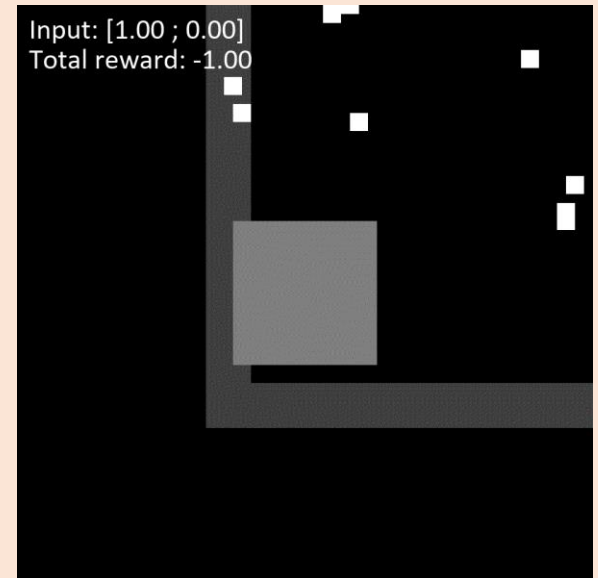Average total reward as a function of training iterations

# The Results: example trajectories



Greedy agent



A2C after 3000 steps



A2C after 4000 steps

# The Results: Discussion

- Lack of training

- Simplicity of RL Algorithms

- Complexity of the input

- Size of the replay buffer

# Thank you for your attention