# Lab meeting

## Deep dive into deep RL

Andrea Pierré

February 12th, 2024

Brown University

# Outline
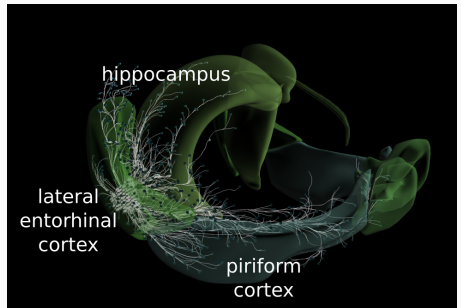
# Outline

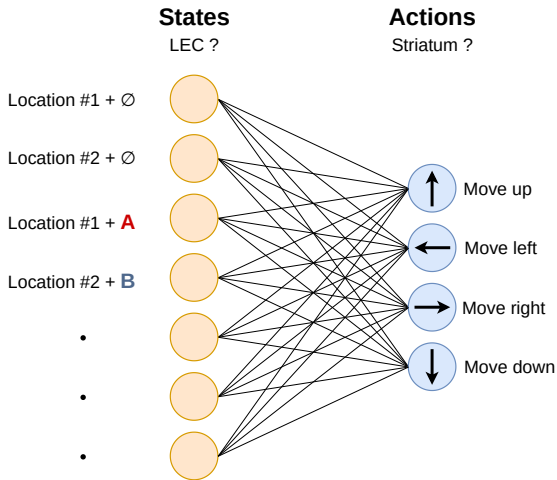# How do neurons in the LEC integrate sensory and spatial information?

- Piriform Cortex encodes olfactory information
- Hippocampus encodes spatial information
- Lateral Entorhinal Cortex (LEC) encodes both olfactory & spatial information
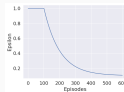
# The half-triangle task

# Mapping states to action



**States**
LEC ?

**Actions**
Striatum ?

Location #1 + ∅

Location #2 + ∅

Location #1 + **A**

Location #2 + **B**

↑ Move up

← Move left
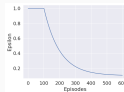
→ Move right

↓ Move down

# Outline

# Deep RL lessons learned

- Deep RL is different from supervised learning → the data to optimize on is not fixed (moving target)

- DQN tricks

# Deep RL lessons learned

- Deep RL is different from supervised learning → the data to optimize on is not fixed (moving target)

- DQN tricks:

  - Replay buffer → save the data experienced in a buffer and sample from it to break the temporal correlation of the data
  - Exploration warm-up with ε-greedy → experience more diverse data
  - Batching → update the weights of the network based on several data examples at the same time instead of only one
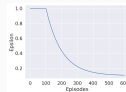  - Target network → use 2 networks to stabilize learning (one of them is updated with a lag)
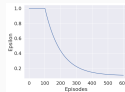
# Deep RL lessons learned

- Deep RL is different from supervised learning → the data to optimize on is not fixed (moving target)

- DQN tricks:

    - Replay buffer → save the data experienced in a buffer and sample from it to break the temporal correlation of the data

    - Exploration warm-up with ε-greedy → experience more diverse data

    - Batching → update the weights of the network based on several data examples at the same time instead of only one

    - Target network → use 2 networks to stabilize learning (one of them is updated with a lag)
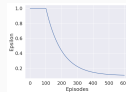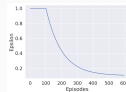
# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

  - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data

  - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data

  - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one

  - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)
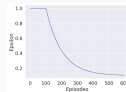
# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

  - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data

  

  - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data

  - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one

  - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)
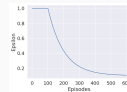
# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

    - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data
    - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data



    - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one
    - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)
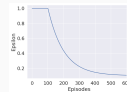
# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

  - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data
  - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data
  - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one
  - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)

# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

    - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data
    - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data
    - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one
    - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)

# Deep RL lessons learned

- Deep RL is different from supervised learning $\rightarrow$ the data to optimize on is not fixed (moving target)

- DQN tricks:

    - Replay buffer $\rightarrow$ save the data experienced in a buffer and sample from it to break the temporal correlation of the data
    - Exploration warm up with $\epsilon$-greedy $\rightarrow$ experience more diverse data
    - Batching $\rightarrow$ update the weights of the network based on several data examples at the same time instead of only one
    - Target network $\rightarrow$ use 2 networks to stabilize learning (one of them is updated with a lag)

# Deep RL implementation

## Algorithm 1: Deep Q-Network (DQN)

Initialize replay memory D to capacity N
Initialize action-value network Q with random weights $\theta$
Initialize target action-value network $\hat{Q}$ with random weights $\theta^- = \theta$
for $episode \leftarrow 1 \ldots M$ do
    $state \leftarrow reset(env)$
    $done \leftarrow False$
    while $done \neq True$ do
        $Q \leftarrow forward\_pass(state)$           /* 4 action values vector */
        $action \leftarrow \epsilon_{greedy}(action\_space, state, Q)$
        $state_{new}, reward, done \leftarrow env.step(action, state)$
        Store transition (state, action, reward, next_state, done) in D
        Sample random minibatch of transitions from D
        $Q \leftarrow forward\_pass(state_{new})$     /* 4 action values vector */
        $Q_{new} \leftarrow reward + \gamma max(\hat{Q})$          /* scalar */
        $y \leftarrow max(Q)$          /* scalar */
        if $done = True$ then
          |  $\hat{y}_{pred} \leftarrow reward$          /* scalar */
        else
          |  $\hat{y}_{pred} \leftarrow Q_{new}$          /* scalar */
        end
        $Loss \leftarrow (y - \hat{y}_{pred})^2$
        Perform a gradient descent step on Loss with respect to the network parameters $\theta$
        Every C steps reset $\hat{Q} = Q$
    end
end

Inspired from (Mnih et al., 2015)

# Outline

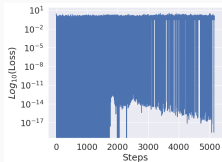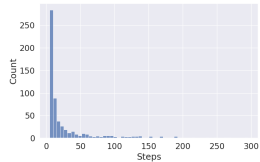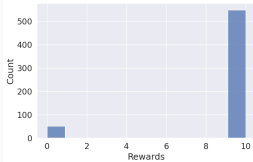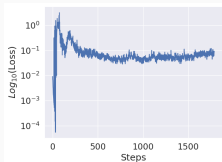# How deep RL feels like

# Networks architectures
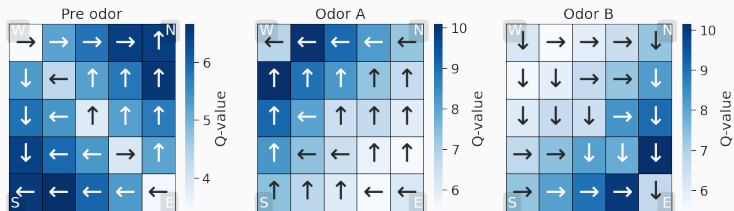
$$28 \rightarrow 54 \rightarrow 54 \rightarrow 54 \rightarrow 4$$

# Loss, rewards & steps distributions

# Policy learned (when it works)

## Outline

# Adding memory into the task

## Current environment

| step | location | cue | reward |
|------|----------|---------|--------|
| 1 | 2 | No odor | 0 |
| 2 | 3 | No odor | 0 |
| 3 | 4 | Odor A | 0 |
| 4 | 3 | Odor A | 0 |
| 5 | 2 | Odor A | 0 |
| 6 | 1 | Odor A | 0 |
| 7 | 0 | Odor A | 10 |

## With memorization needed

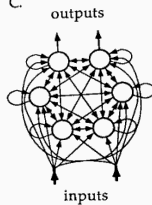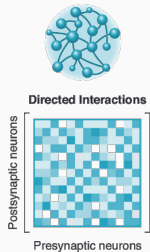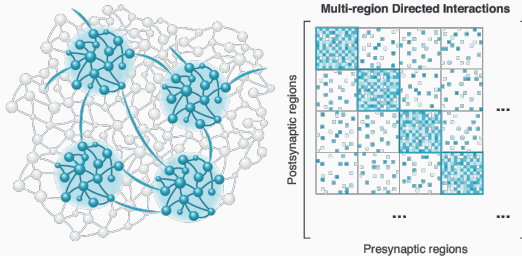| step | location | cue | reward |
|------|----------|--------|--------|
| 1 | 2 | $\emptyset$ | 0 |
| 2 | 3 | $\emptyset$ | 0 |
| 3 | 4 | Odor A | 0 |
| 4 | 3 | $\emptyset$ | 0 |
| 5 | 2 | $\emptyset$ | 0 |
| 6 | 1 | $\emptyset$ | 0 |
| 7 | 0 | $\emptyset$ | 10 |

# Network of networks



(a) Single-region RNN

(b) Multi-region RNN (mRNN)

Directed Interactions

Multi-region Directed Interactions

Questions ?