

IA pour la Robotique

M2 ANDROIDE 2021 – 2022

Encadrant : Benoît Girard



**SORBONNE
UNIVERSITÉ**

CRÉATEURS DE FUTURS
DEPUIS 1257

Self-Improving Reactive Agents Based On Reinforcement Learning, Planning and Teaching

Long-Ji LIN (1992)

LIM Vincent
WANG Tianyu

SOMMAIRE



**SORBONNE
UNIVERSITÉ**
CRÉATEURS DE FUTURS
DEPUIS 1257

- 1 Introduction
- 2 Environnement
- 3 Réseau de neurones
- 4 Résultats

1

Introduction

- ❖ Etat de l'art en 1992 :
 - Adaptive Heuristic Critic (Sutton, 1984 et Barto, 1990)
 - Connectionist error backpropagation algorithm (Rumelhart, 1986)
 - Temporal difference learning (Sutton, 1988)
 - Q-Learning (Watkins, 1989)

- ❖ Testées sur des tâches simples

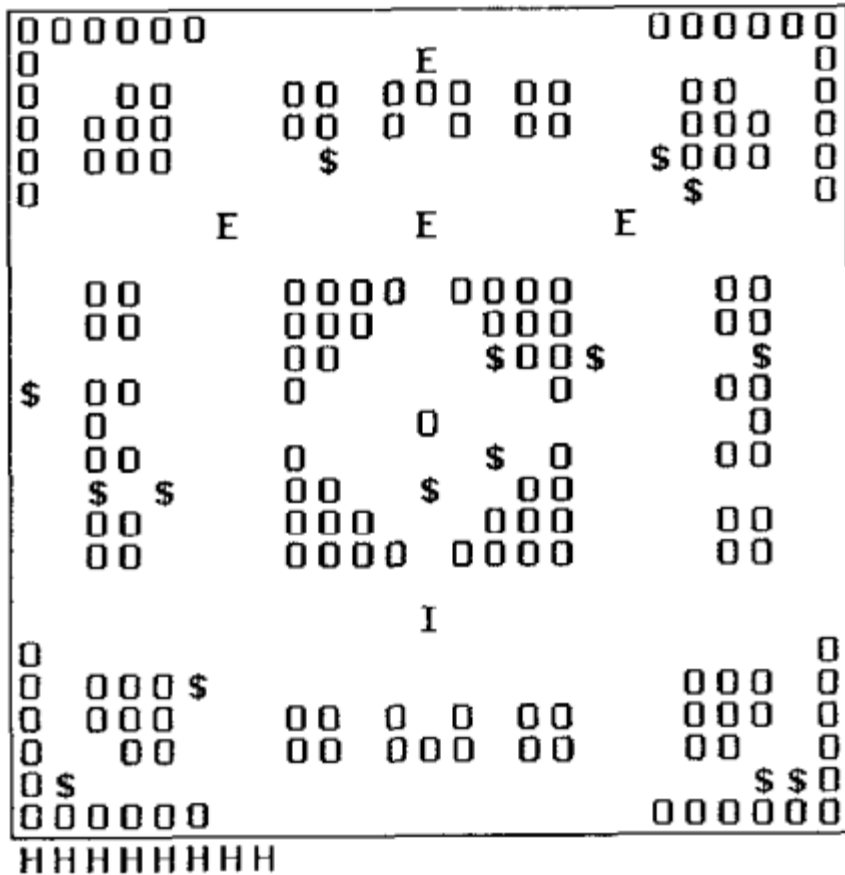
- ❖ Convergence lente

- ❖ Deux aspects :
 - Analyse des performances du RL sur des tâches plus complexes
 - Chercher des méthodes pour rendre la convergence plus rapide
- ❖ Comparer QCON et AHCON
- ❖ Trois méthodes pour améliorer la convergence :
 - Experience replay
 - Learning action model
 - Teaching

- ❖ Reproduire l'environnement de simulation
- ❖ Implémenter l'algorithme d'apprentissage QCON ainsi que la variante Experience Replay
 - Q-Learning
 - Réseau de neurones
- ❖ Reproduire le protocole expérimental et obtenir des résultats similaires à ceux de l'article

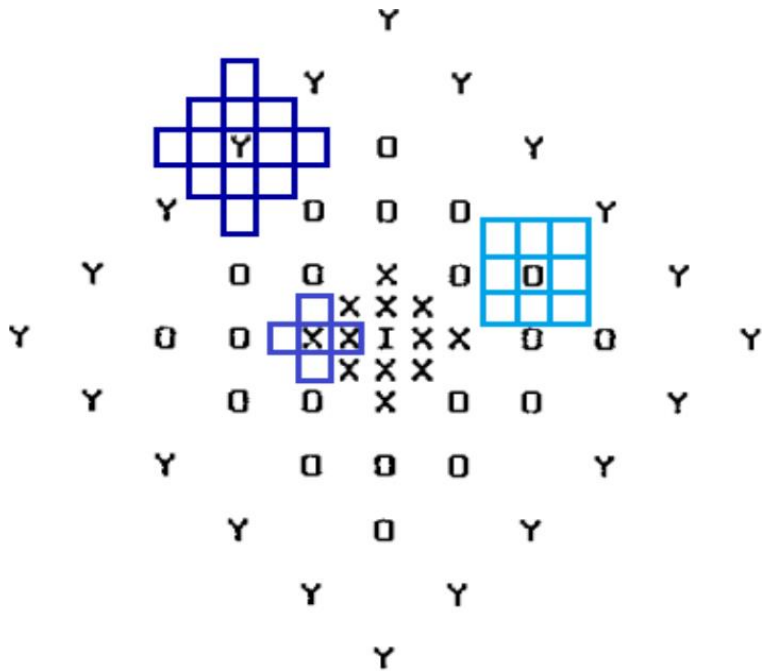
2

L'environnement

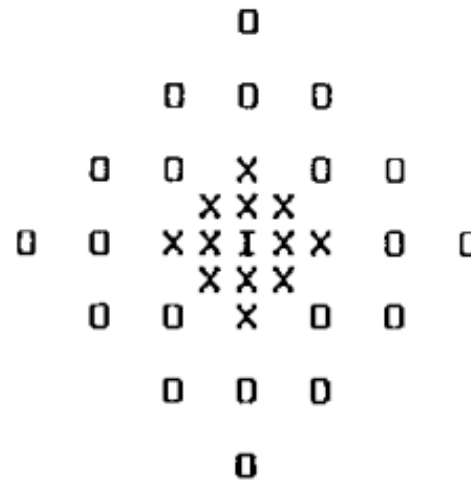


Représentation de l'environnement

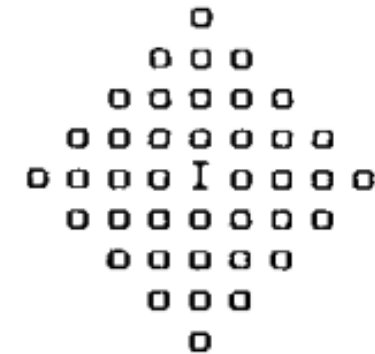
- ❖ Grille 25x25
- ❖ 1 agent (I)
- ❖ 4 ennemies (E)
- ❖ 15 nourritures (\$)
- ❖ Obstacles (O)
- ❖ Energie de l'agent (H)
- ❖ Déplacement possible dans 4 directions
- ❖ Trois rewards :
 - + 0.4 lorsque l'agent collecte de la nourriture
 - - 1.0 lorsque l'agent meurt
 - 0.0 sinon



Capteurs de nourritures



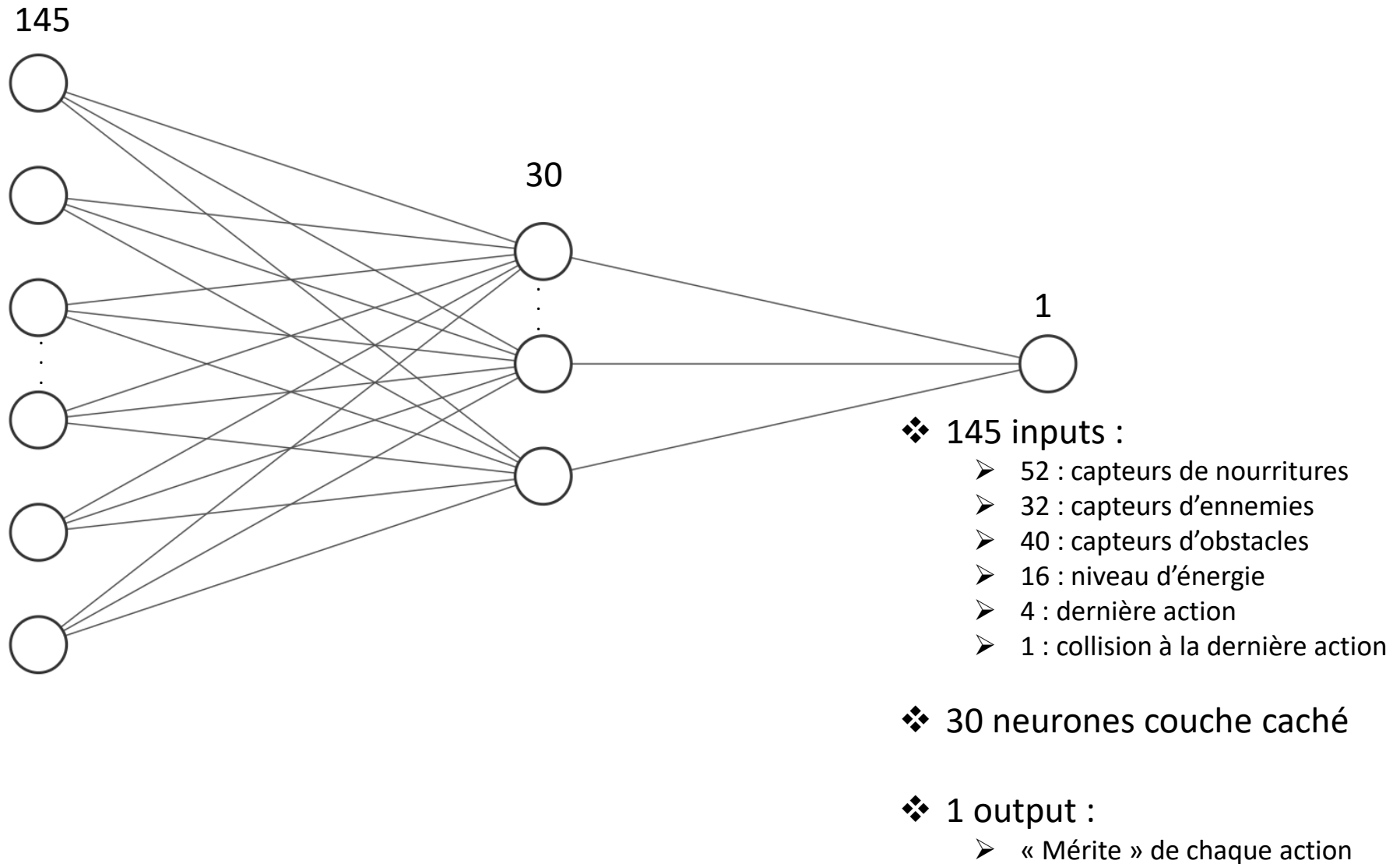
Capteurs d'ennemies



Capteurs d'obstacles

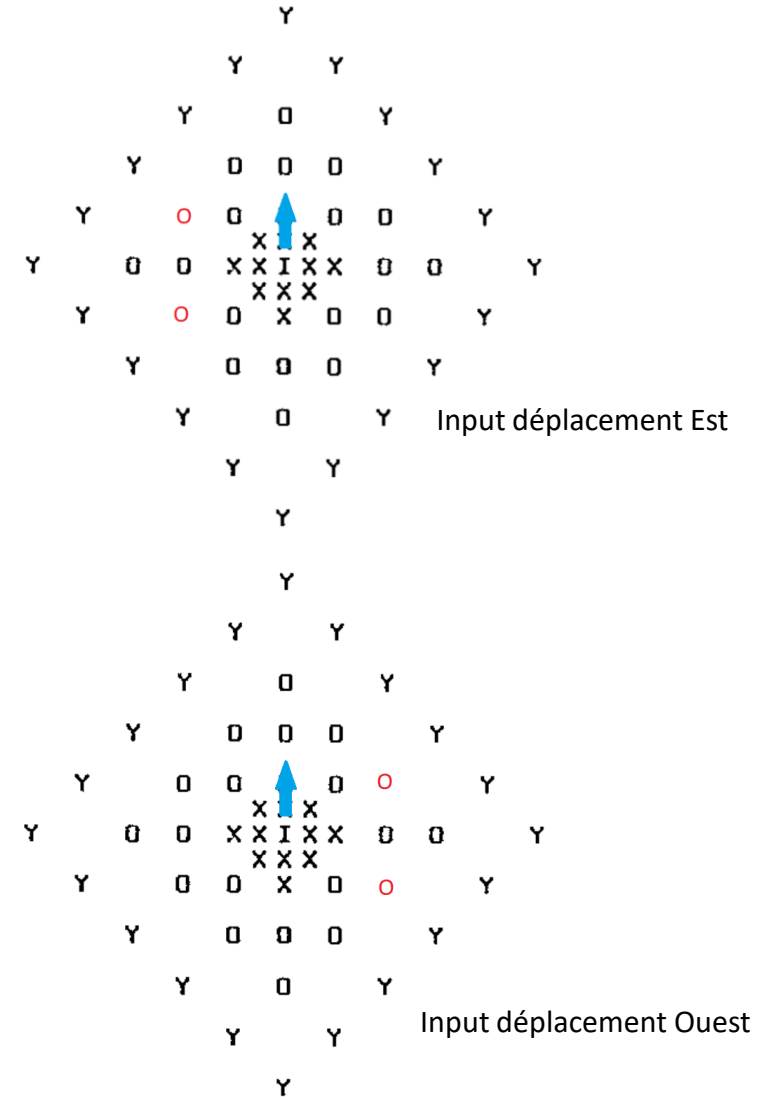
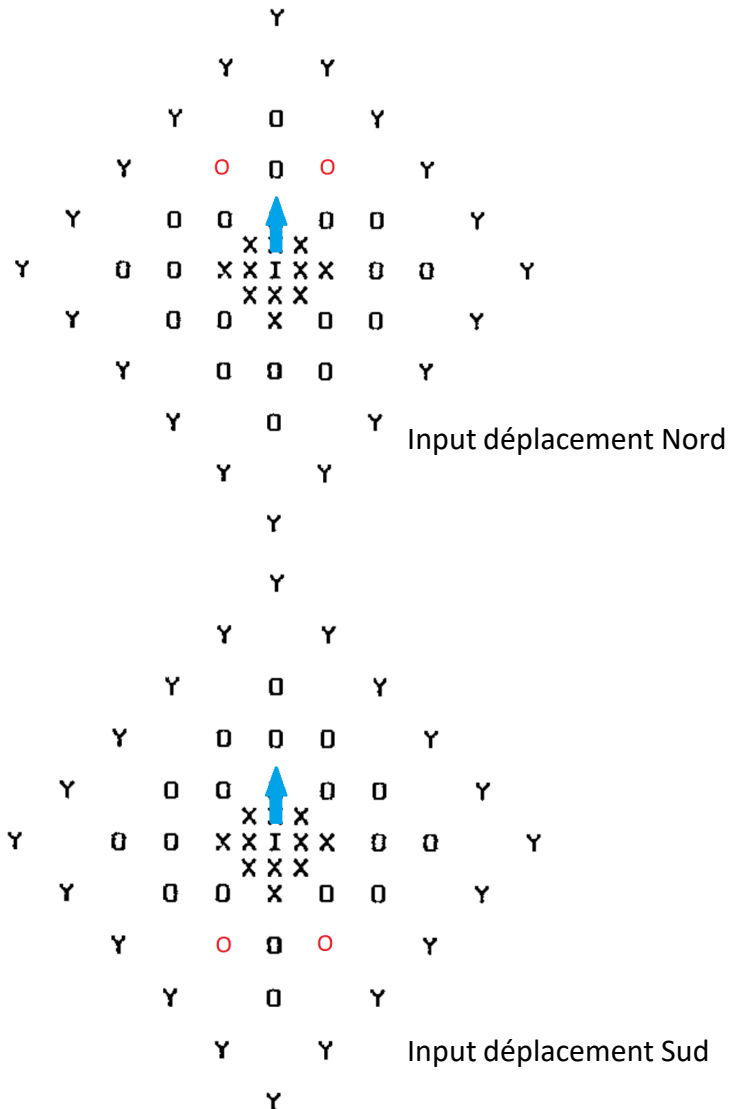
3

Le réseau de neurones



Rotation des inputs

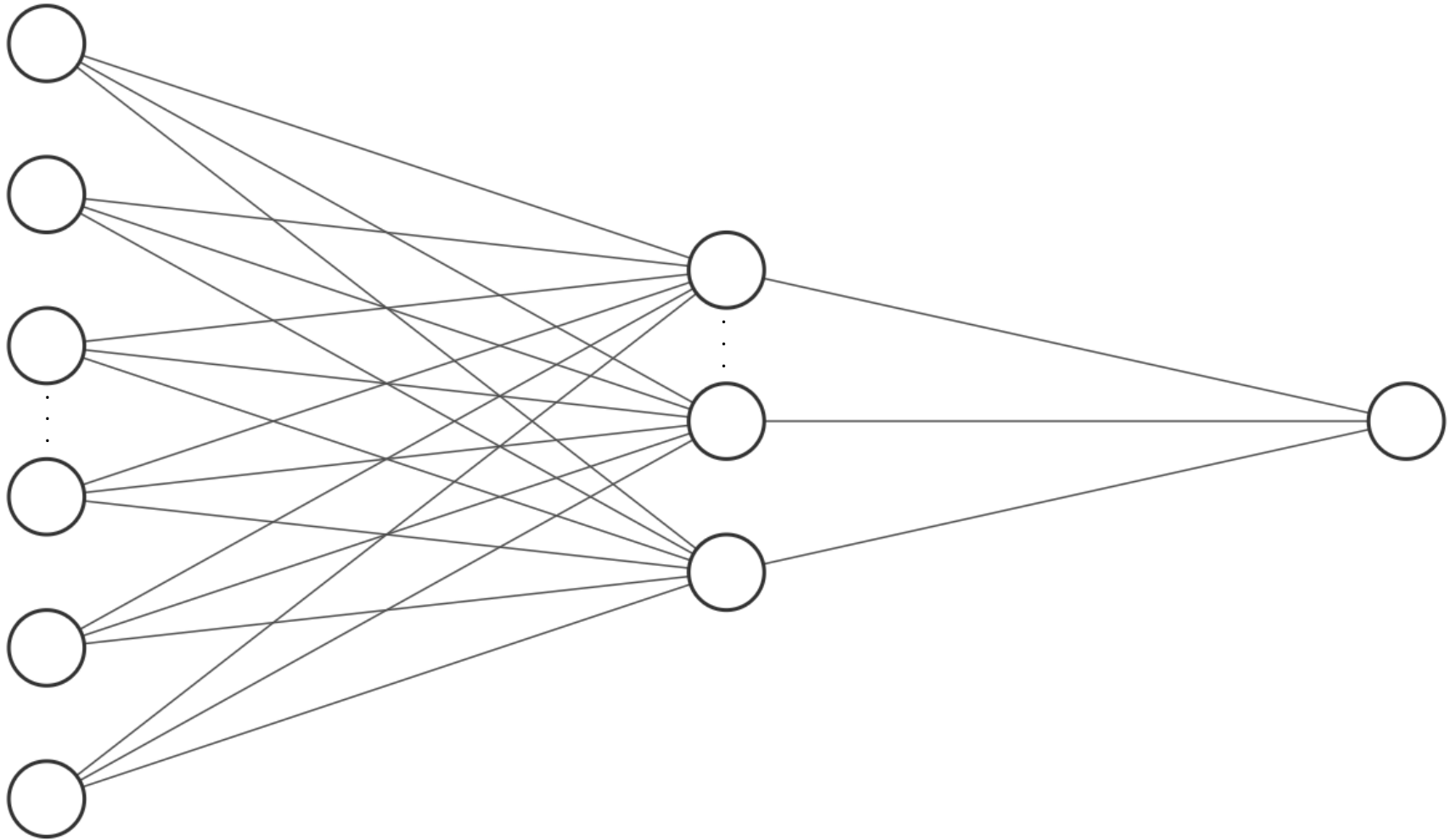
2



La propagation

3

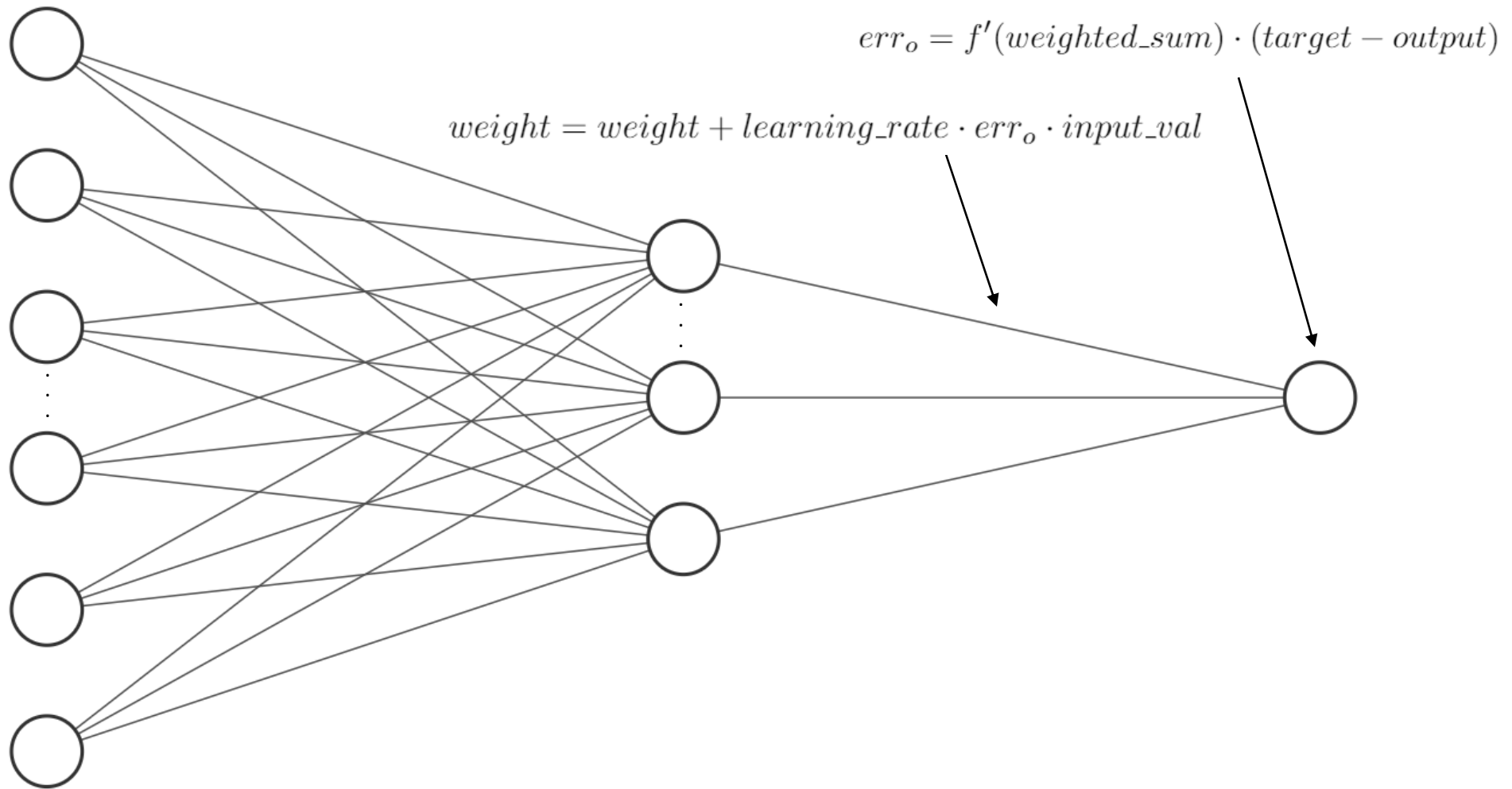
Fonction d'activation : $f(x) = \left(\frac{1}{1 + e^{-x}} - 0.5\right) \cdot 2$



La rétropropagation

3

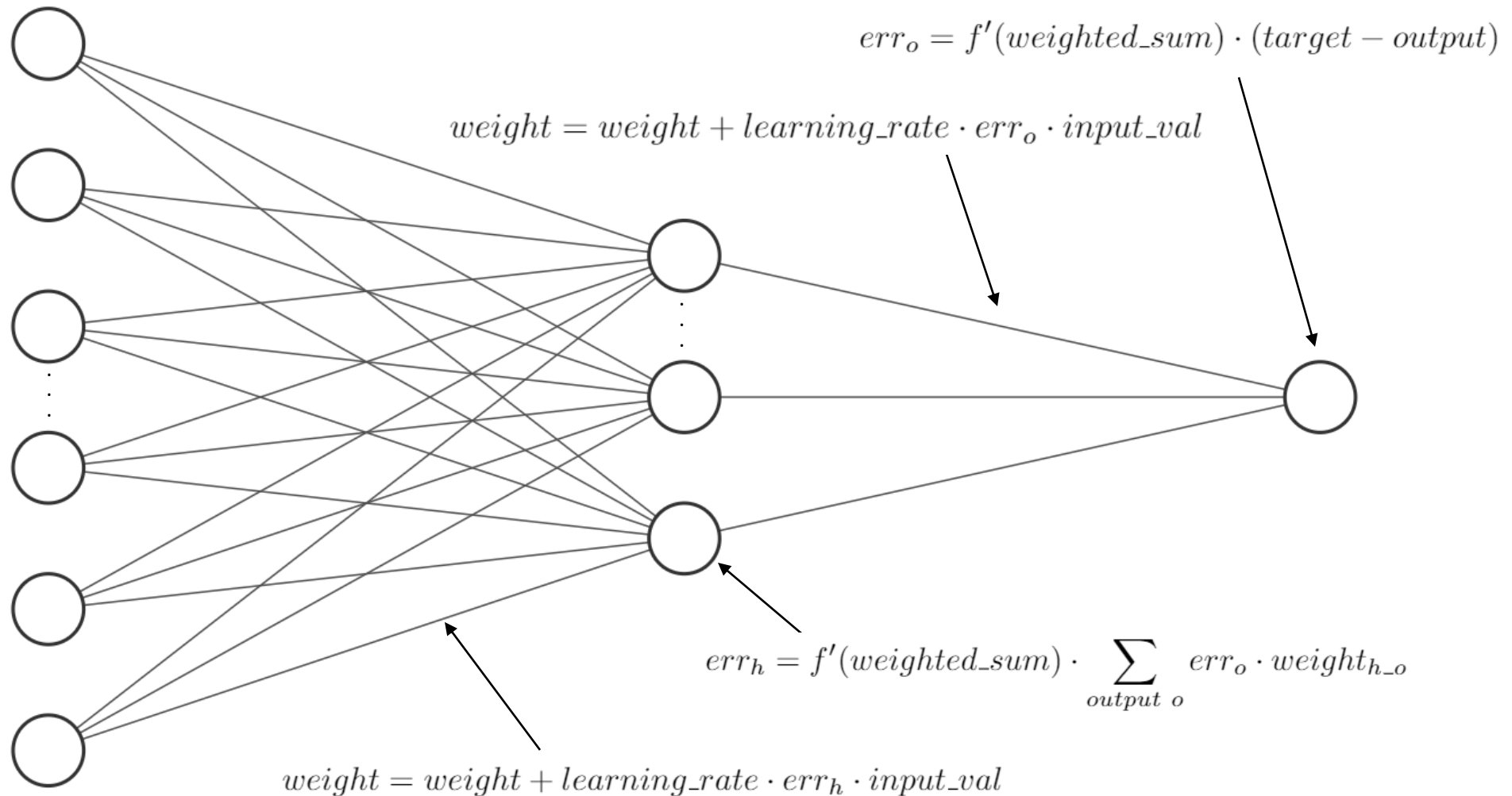
Fonction d'activation : $f(x) = \left(\frac{1}{1 + e^{-x}} - 0.5\right) \cdot 2$



La rétropropagation

3

Fonction d'activation : $f(x) = \left(\frac{1}{1 + e^{-x}} - 0.5\right) \cdot 2$

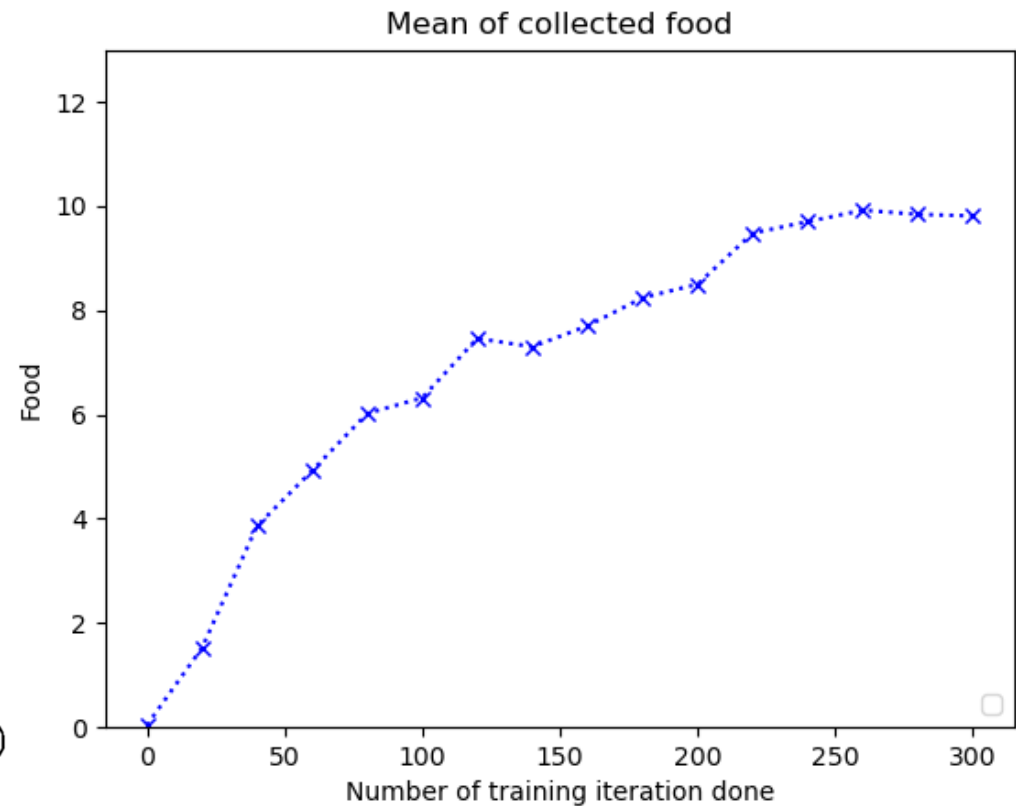
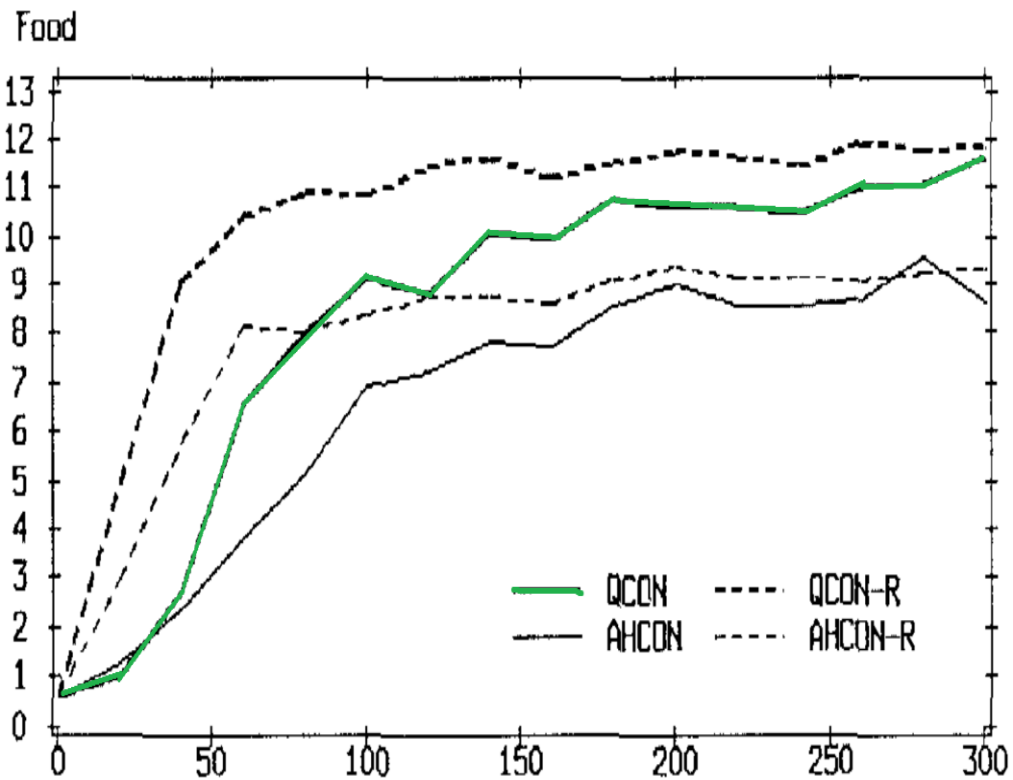


4

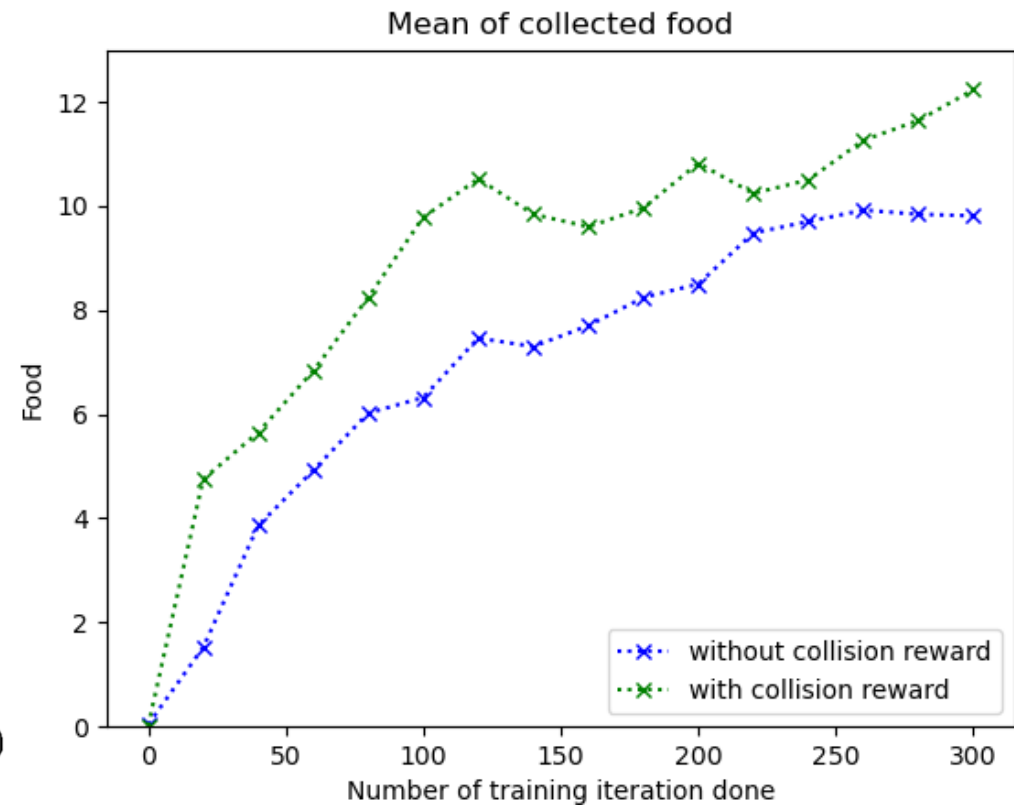
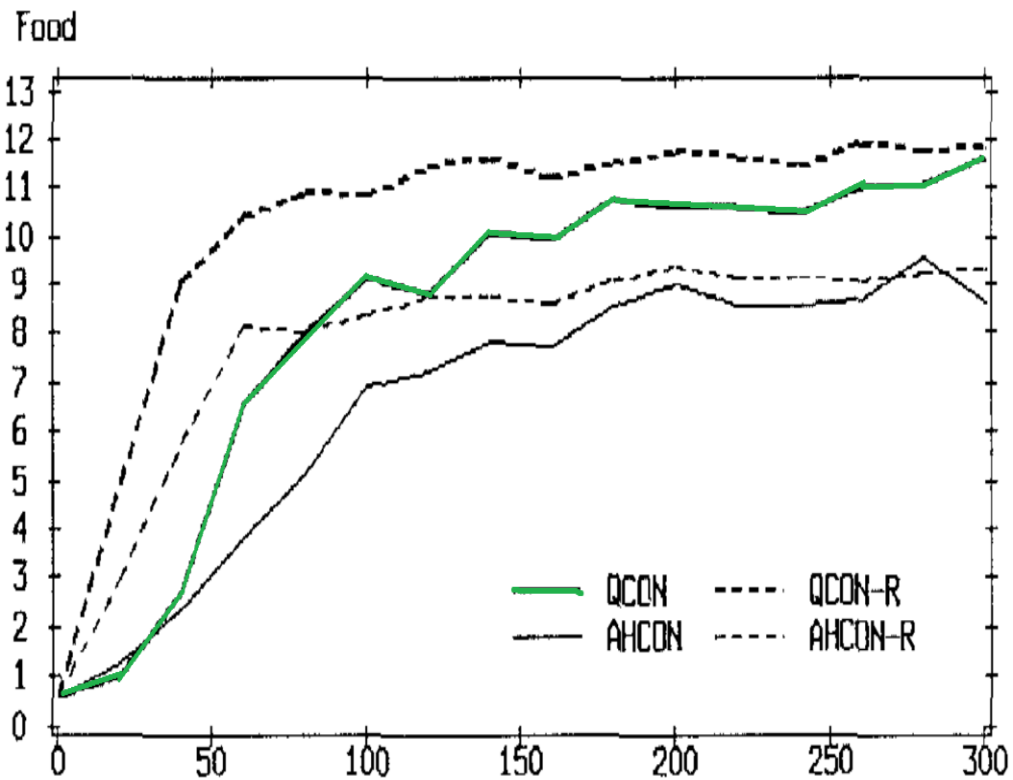
Résultats

Résultats

4



- ❖ *“Programming robots using reinforcement learning and teaching”, Long-Ji LIN, 1991*
- ❖ Cinq rewards :
 - + 0.4 lorsque l’agent collecte de la nourriture
 - - 1.0 lorsque l’agent meurt
 - - 0.5 lorsque l’agent rentre dans un mur
 - - 0.5 lorsque l’agent revient à son ancienne position
 - 0.0 sinon



**Merci de votre attention !
Des questions ?**

