

CPSC 340: **Machine Learning and Data Mining**

Reinforcement Learning -- Bonus Lecture

Ben Norman
(slides adapted from Daniele Reda, and Helen Zhang)
Spring 2022 (2021W2)

Talk Content

- Part 1: What, Why, When
 - What is Reinforcement Learning, RL?
 - Why RL? Why not RL? When to do RL?
- Part 2: How Reinforcement Learning Works
- Part 3: Applications of RL!

Part 1

- What is RL?
- Why RL?
- When RL?

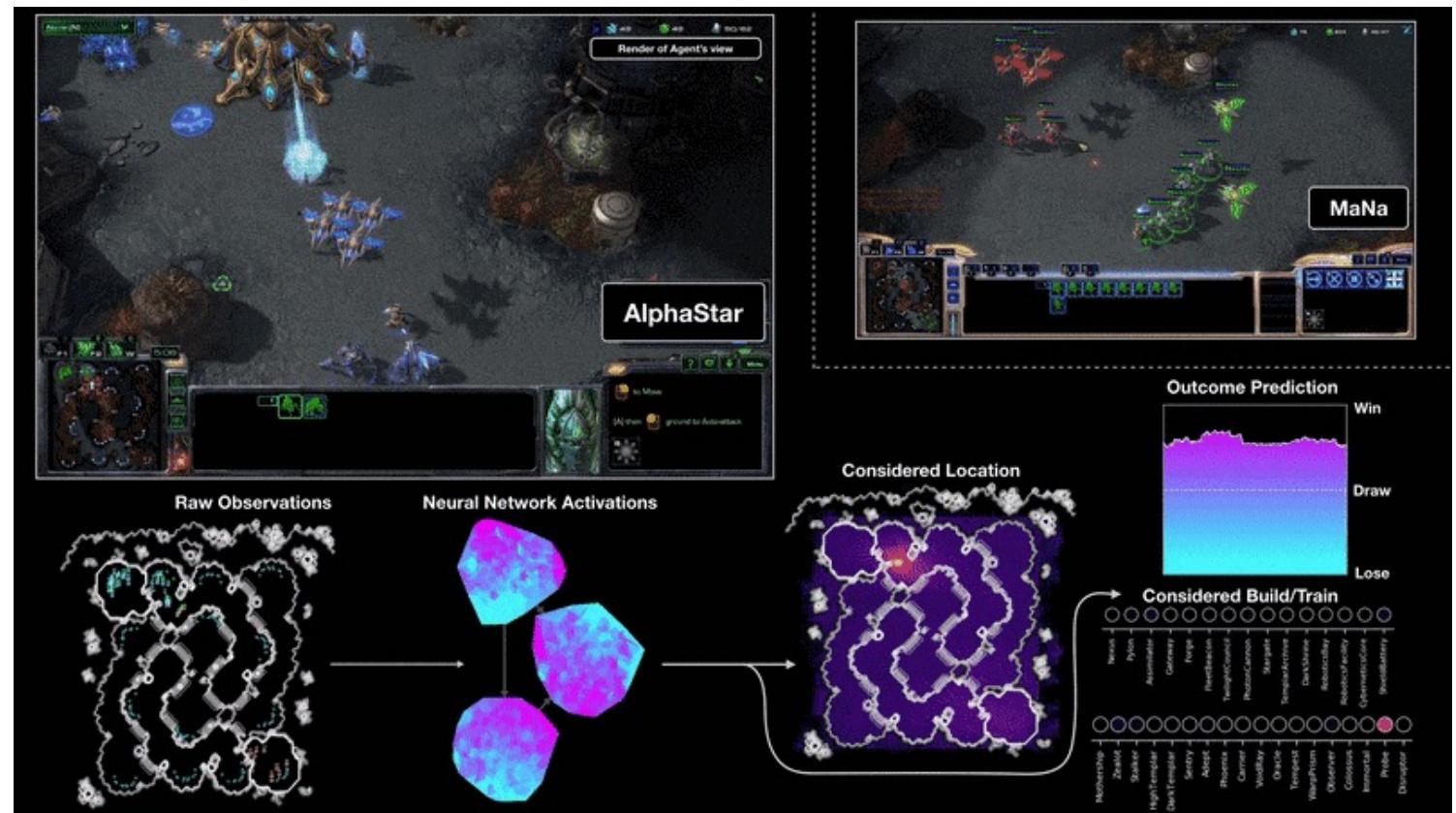
What are the problems we want solve?

Having a Robot Walk

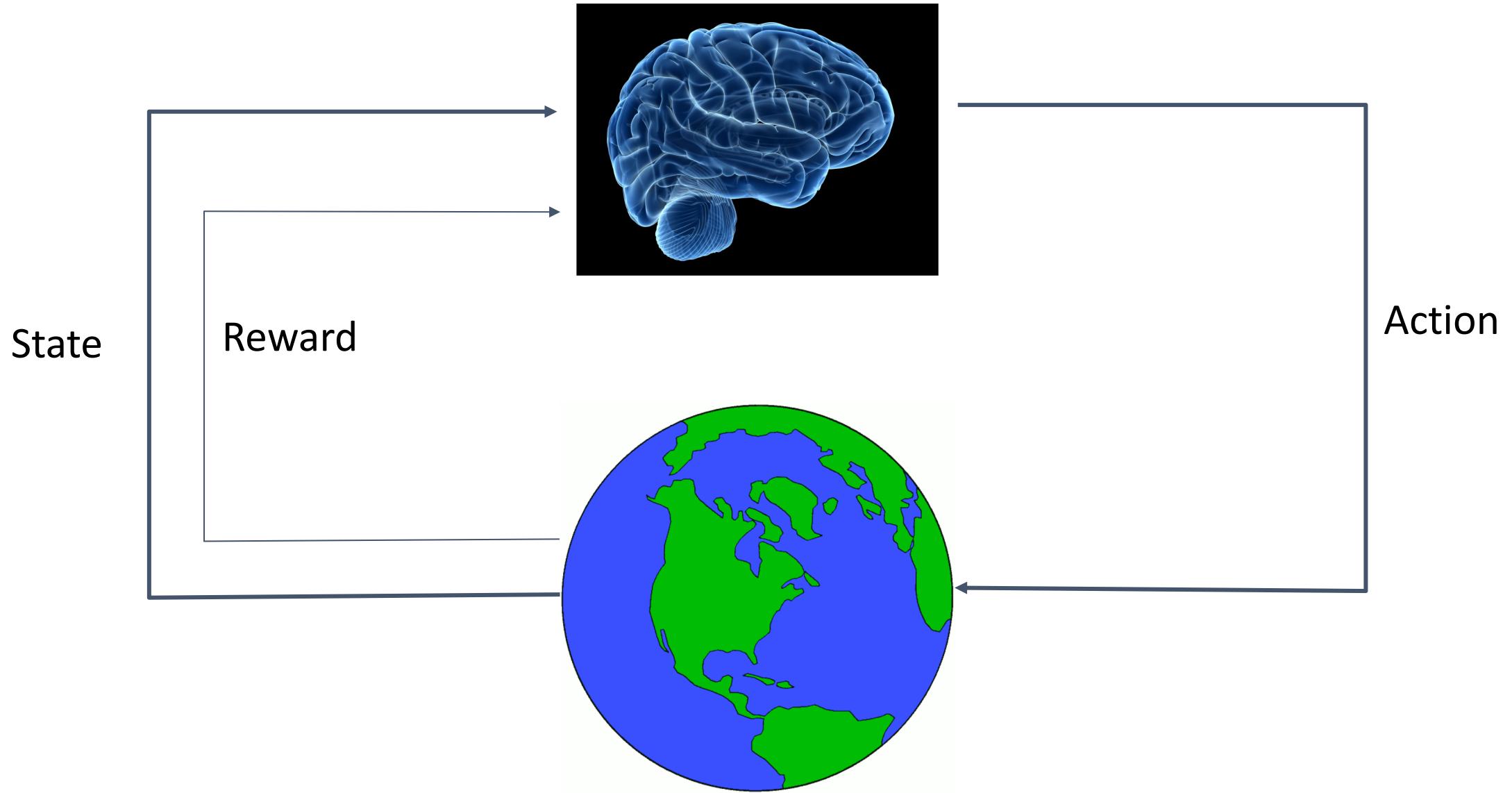
Investing

Chess

Starcraft



Problem Setting



What is the nature of these problems?

The agent gets information, and must take actions

The info:

State: Your limb positions, Your stock portfolio, The chess board, The Starcraft Game State

Output:

The next **action:** Muscle Activations, Buying and Selling Stock, A single chess move, Mouse clicks and Key presses

Note: we can have Complex Time and Action Dependencies

How can we solve these problems?

Two main approaches:

- We can do a form of **Supervised Learning** called **Imitation Learning** and **imitate** and **distill** expert play. (State is your x, Action is your y)

Effective and Popular

- We can do **Reinforcement Learning** and **explore** and **discover** play

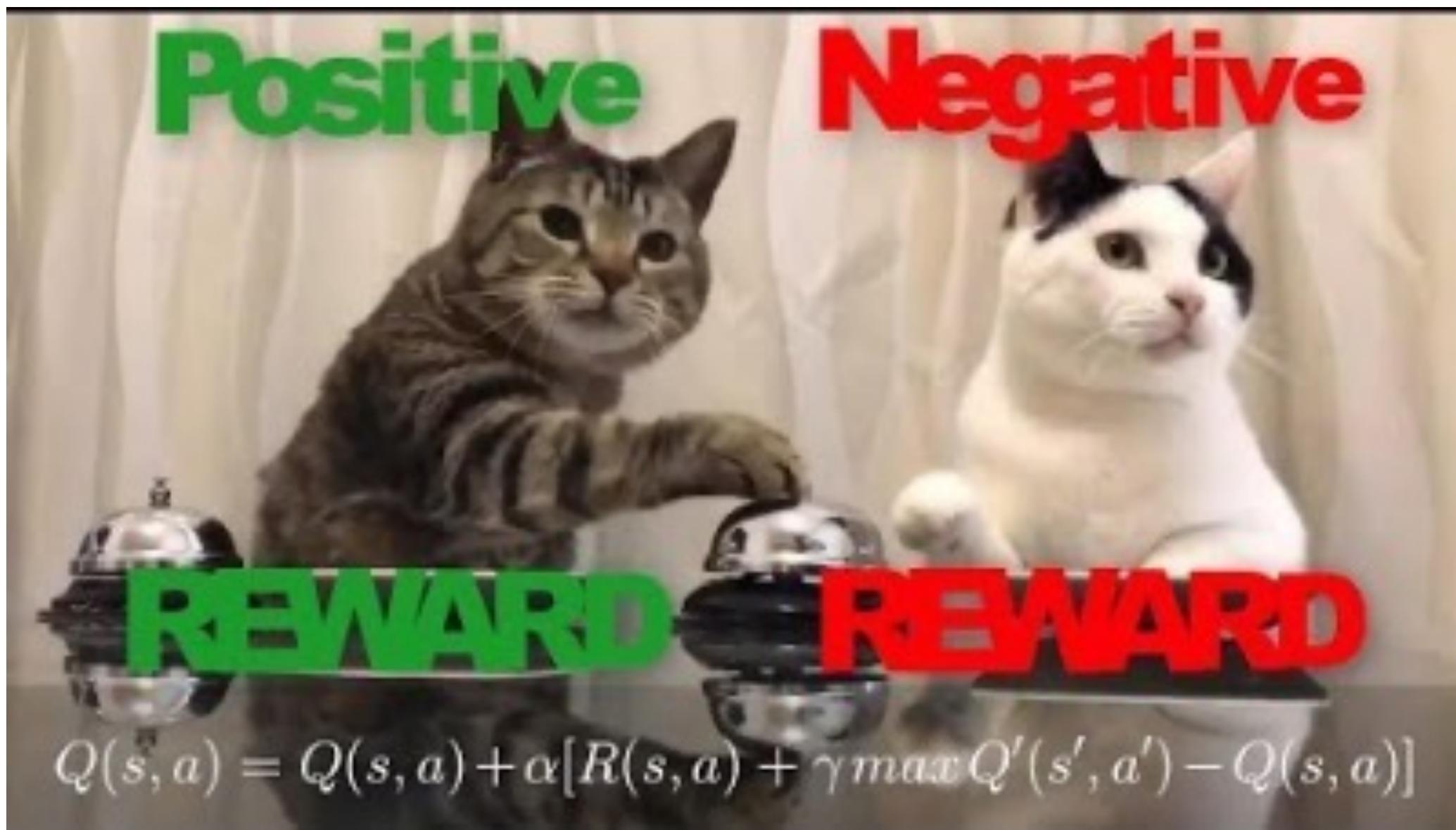
Effective and Popular

Types of Feedback

Imagine you take a mock test for your mid terms:

- **Model Answers Feedback**, “question 1 part a) you should make the point that ...” : **Supervised Learning**
- **Score Feedback**, “you got 200”: **Reinforcement Learning**
- **No Feedback**, “ ” : **Unsupervised**, or Self-Supervised Learning

Learning from scores is far harder than learning from model answers.



$$Q(s, a) = Q(s, a) + \alpha[R(s, a) + \gamma \max Q'(s', a') - Q(s, a)]$$

Summary

- How is Reinforcement Learning, RL, different from Supervised Learning?

In Supervised Learning you have Labels

In Reinforcement Learning you have Score, or Reward

Imitation Learning copies and distils the expert plays

Reinforcement Learning explores, and discovers strategies

Why Reinforcement Learning?

- No need for Labels

Score is cheap, and can often be **generated** by algorithm

- Beyond distilling, **discovery**. The potential for **superhuman** performance!
- Can be more **robust**

Note: imitation learning *can* be better than human too

(imagine 5 experts each amazing at one specialty, the best imitation is amazing at all 5 specialties)

Why not Reinforcement Learning?

- **Scoring** can be **difficult**, sometimes having labels is easier,

e.g.



- RL can have **low sample efficiency**, often run in simulations, but simulating can be hard



RL can give super-human performance

AlphaGo, and its successors, AlphaGo Zero, AlphaZero, MuZero... really are superhuman, in a *meaningful* way

- discovered new corner sequences (Joseki)

Not aware of any meaningful super-human imitation learning

Part 2

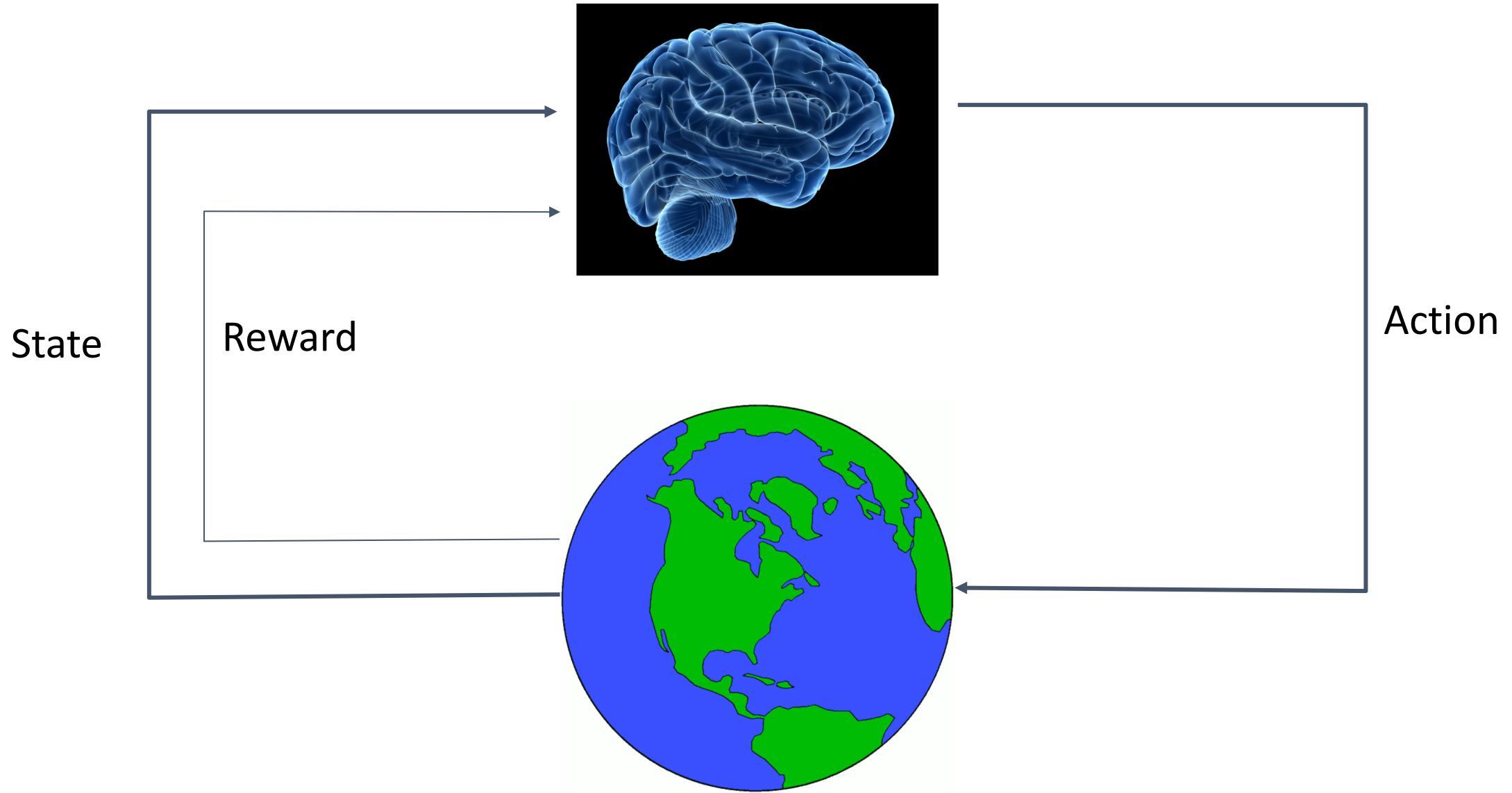
- How does Reinforcement Learning work?

How does Reinforcement Learning Work?

Two basic approaches:

- Value Function Estimation : estimate the **long-term** value of states (and actions)
- Policy Gradient : policy in state **s** do action **a**, update our policy to increase likelihood of high scoring policies

Problem Setting



Goal

- maximize total reward

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots + \gamma^{T-t-1} r_T$$

- T to infinity
- Why discount factor?
 - Convergence, ‘it works mathematically’
 - Estimates of Future Reward are more uncertain.

What is the value of a state?

Given a policy, π , we can define the Value function of a state

$$V_\pi(s) = \mathbb{E}_\pi(R_t | S_t = s)$$

We want the value under the optimal policy, π_*

$$V_*(s) = \mathbb{E}_*(R_t | S_t = s)$$

Satisfies a ‘Bellman Equation’:

$$V^*(s) = \max_{\alpha} \left[R(s, \alpha) + \gamma \sum_{s'} p(s' | s, \alpha) V^*(s') \right]$$

How to choose action though?

Idea: **value state action pairs**

Definition:

$$Q_\pi(s, a) = \mathbb{E}_\pi(R_t | S_t = s, A_t = a)$$

Bellman Equation:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} p(s' | s, a) \max_{\alpha} (Q^*(s', \alpha))$$

The Q-Function specifies a policy

	Actions						
	0	1	2	3	4	5	
State <i>s</i>	0	-1	-1	-1	-1	0	-1
	1	-1	-1	-1	0	1	100
	2	-1	-1	-1	0	-1	-1
	3	-1	0	0	-1	0	-1
	4	0	-1	-1	0	-1	100
	5	-1	0	-1	-1	0	100

How to learn it?

There is a theorem: Policy Improvement Theorem

Given any (tabular) Q-function, if we update our Q-function, through the Bellman Equations, the policy will improve

$$Q^{n+1}(s, a) = R(s, a) + \gamma \max_{\alpha} (Q^n(s', \alpha))$$

Note, that in expectation is gives the Bellman Equation

Problem? How do we ensure exploration?

The simplest possible idea:

Epsilon greedy: take the best action $1 - \varepsilon$ of the time, random action ε of the time

Q-learning

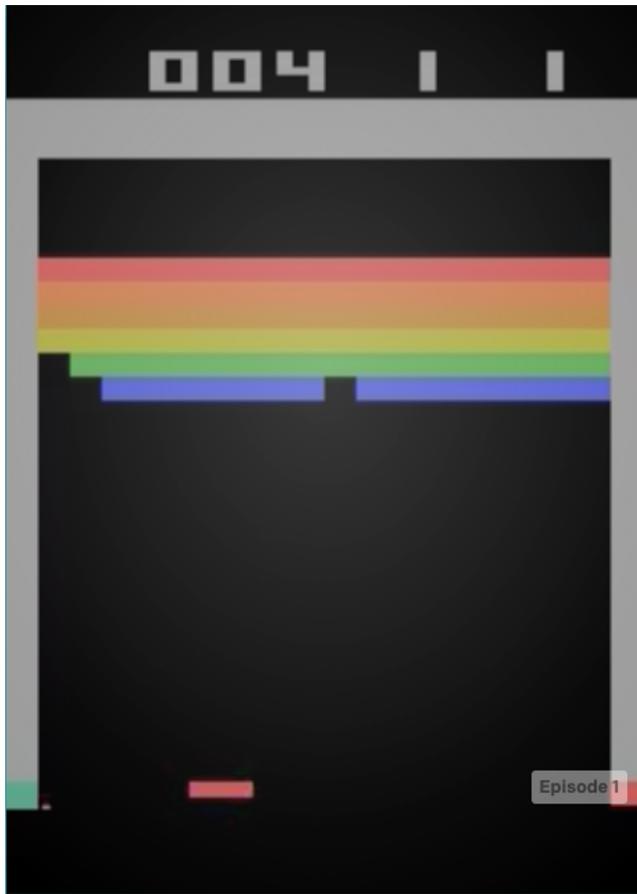
Initialize Q-table with random values.

1. Choose action a to perform in current state s . (ϵ -greedy)
2. Perform a and receive reward $R(s,a)$.
3. Observe new state $S(s,a)$.
4. Update Q-table.

$$Q' (s, a) \leftarrow R (s, a) + \gamma \max_{\alpha} \{ Q' (S(s, a), \alpha) \}$$

PROBLEM:
TABLE CAN EASILY EXPLODE IN DIMENSIONS

Let's look at an example: ATARI



State (the actual image)
84x84x4 pixels (gray-
scale)



2 Actions
Left-Right

Could we use a q-table?

Atari Breakout example:

State = raw pixels of last 4 frames (84x84) with 256 different possible values.

Actions = 2 actions available

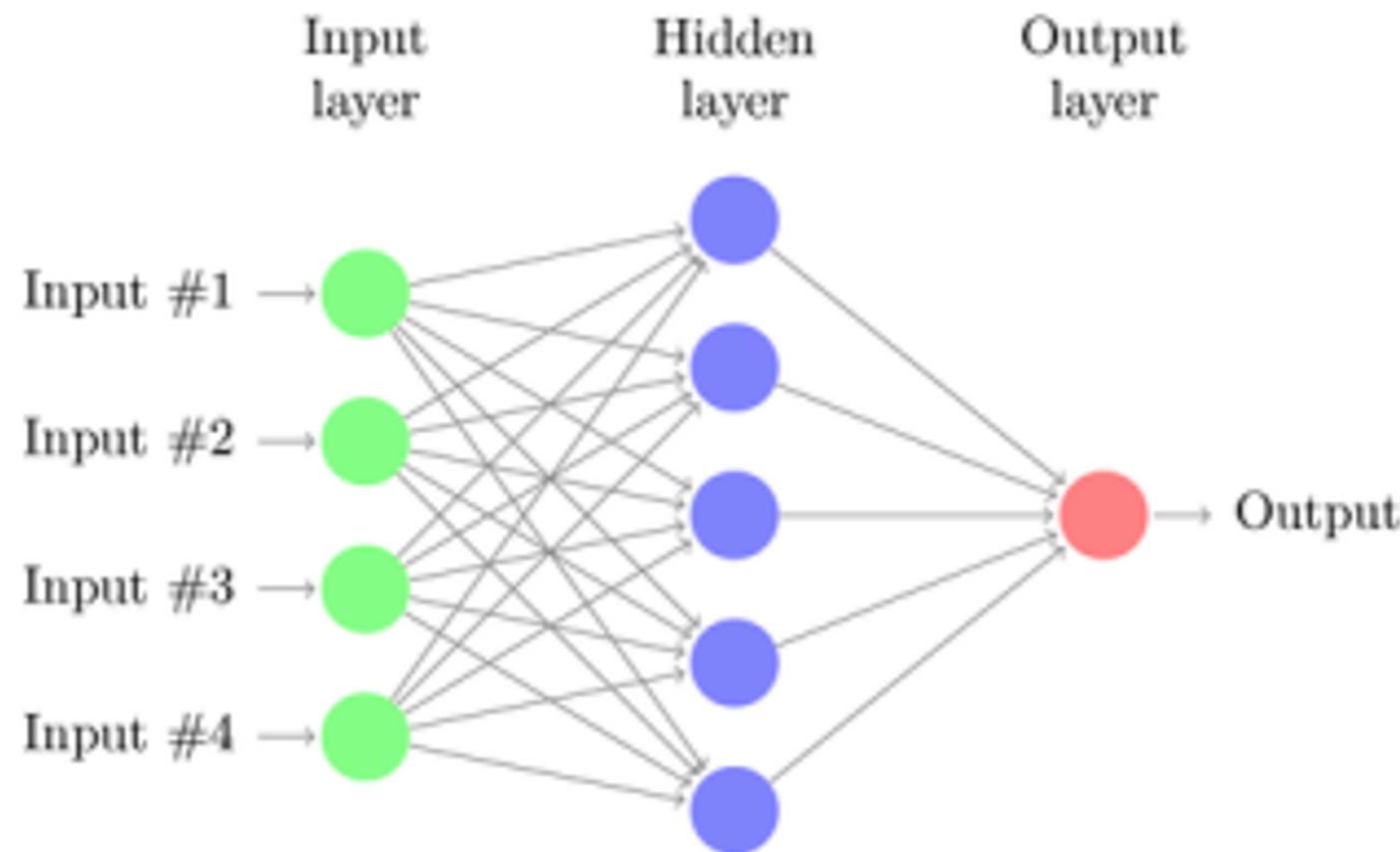


$$(4 \times 84 \times 84 \times 256) \times 2 = 14450688 \text{ different values}$$

Solution : Functional Approximation

Neural networks!

Neural Networks



Neural Networks

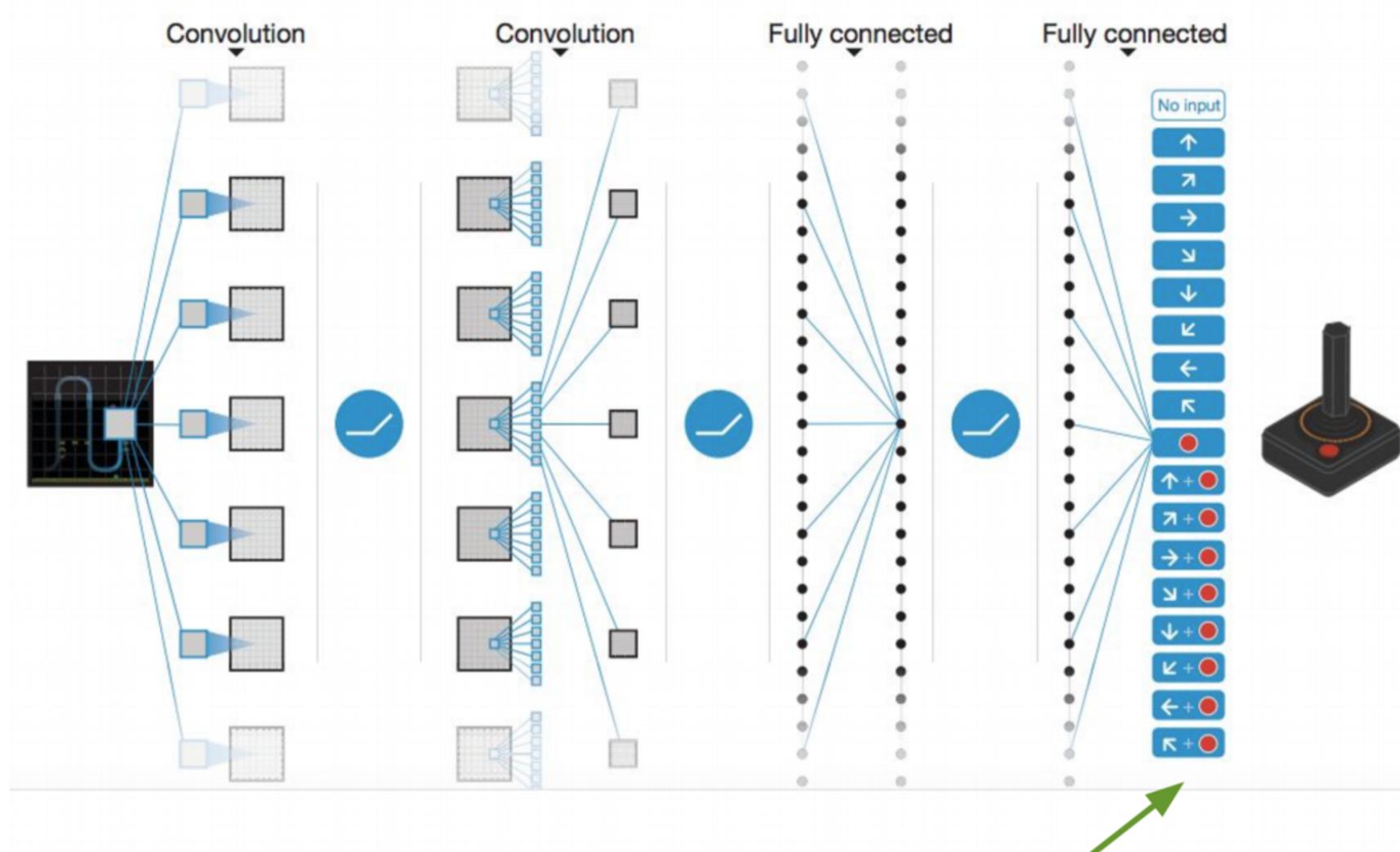
$$loss = \left(\frac{r + \gamma \max_{a'} \hat{Q}(s, a') - Q(s, a)}{\frac{\text{Target}}{\text{Prediction}}} \right)^2$$

Reward Decay Rate

Target Prediction

$$Q' (s, a) \leftarrow \mathcal{R} (s, a) + \gamma \max_{\alpha} \{ Q' (\mathcal{S}(s, a), \alpha) \}$$

DQN Framework



1 network, outputs Q value for each action

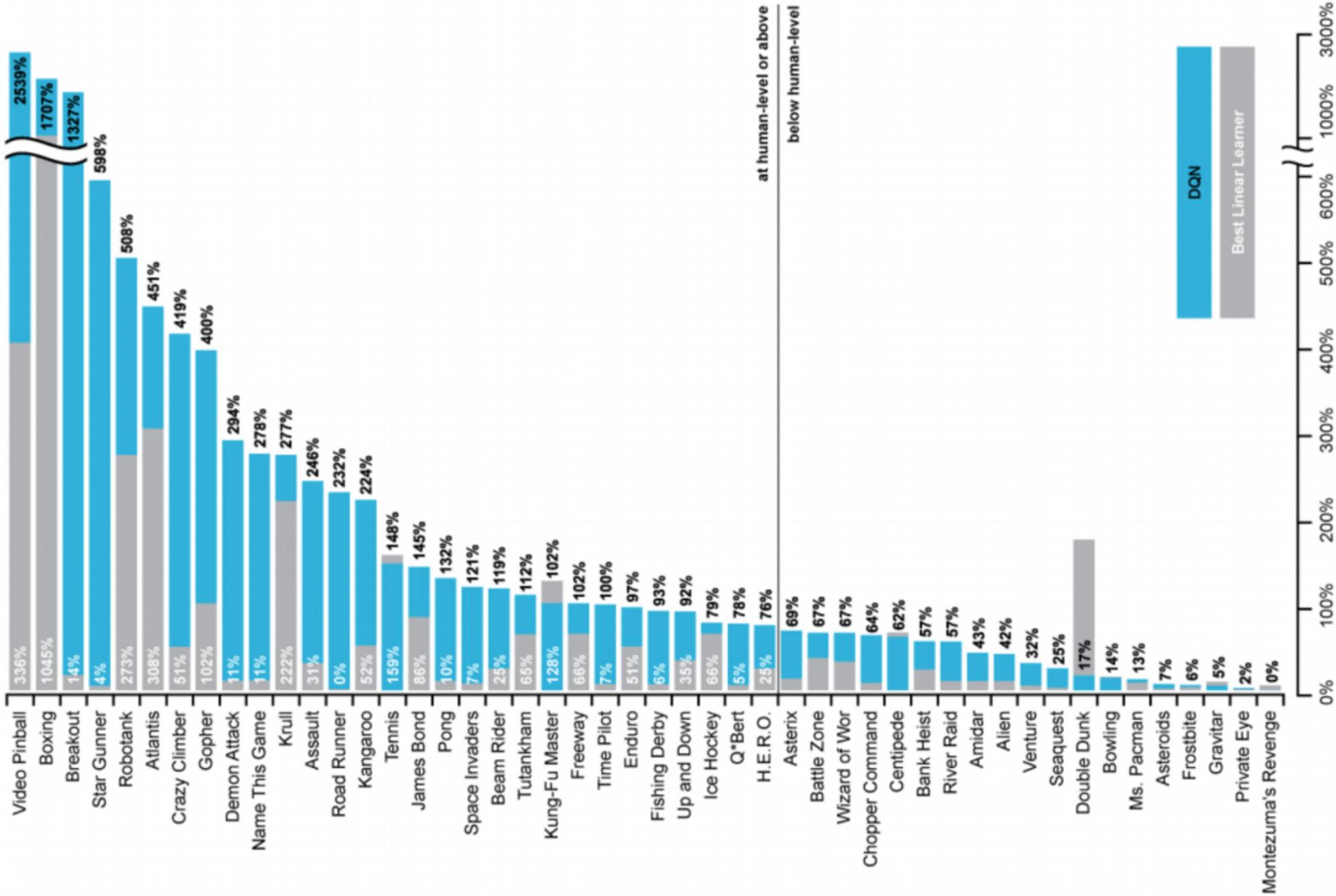
DQN Framework

Algorithm 1: DQN Pesudocode

```
1 Randomly initialize neural network  $NN$ 
2 Get initial state  $s_0$ 
3  $t = 0$ 
4 while 1 do
    //  $\epsilon$ -greedy strategy
    5  $r = \text{get random value from } (0, 1)$ 
    6 if  $r > \epsilon$  then
        7    $a_t = NN(s_0)$ 
    8 else
        9    $a_t = \text{random action between the ones available}$ 
    10 end
    11  $s_{t+1}, r_t, done = environment(a_t)$ 
    12 if  $done = True$  then
        13    $y = r_t$ 
    14 else
        15    $y = r + \gamma * \max_{a^i} NN(s_{t+1})$ 
    16 end
    17 Do gradient descent on  $y - NN(s_t, a_t)$  to update weights of  $NN$ 
    18  $t = t + 1$ 
19 end
```

Famous successes of RL: Atari Breakout

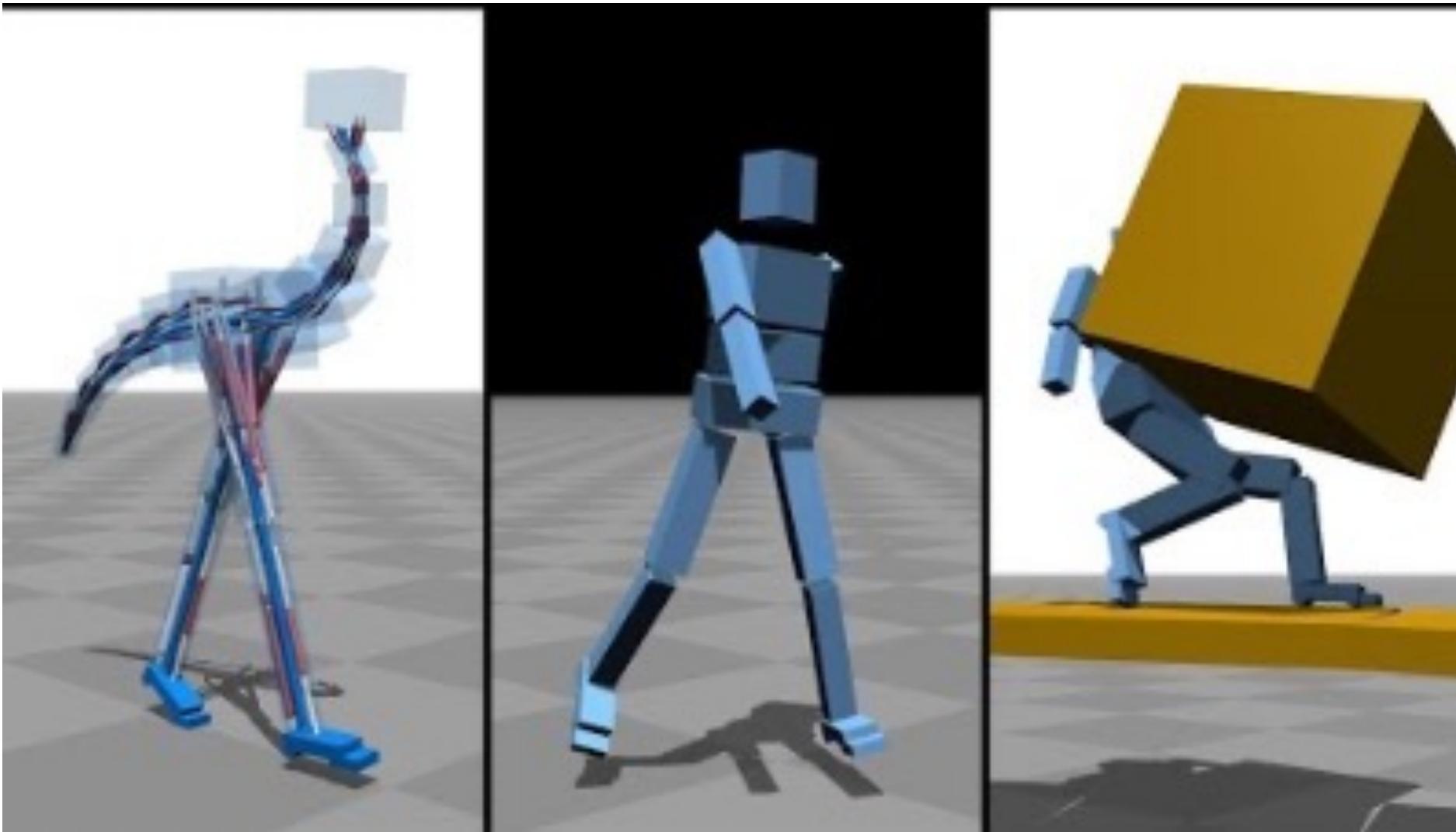




Part 3

- Applications of RL!

Flexible Muscle-Based Locomotion for Bipedal Creatures



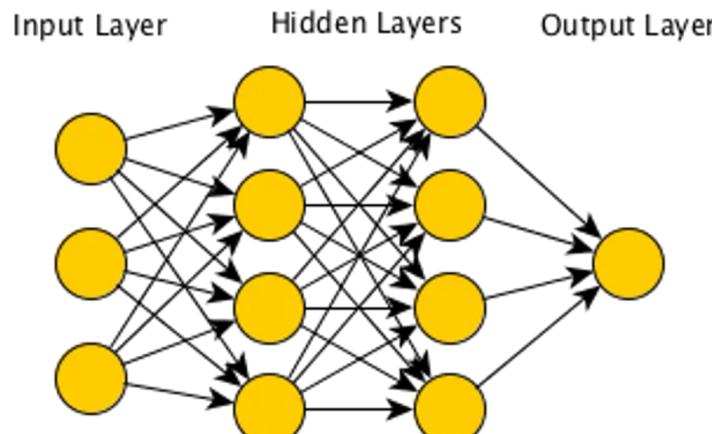
A real world example: Learning to drive with Reinforcement Learning



Learning to drive with RL



**Steering & Speed
Measurement**



**Steering & Speed
Command**

Reward: forward distance

Terminate when it goes out of
the lane

Kendall, Alex, et al. "Learning to drive in a day."



Where to go from here?

- [openAI Spinning Up RL](#)
- [Sutton book](#)
- [David Silver UCL Course](#)

Reinforcement Learning

An Introduction
second edition

Richard S. Sutton and Andrew G. Barto

