

# Reinforcement learning

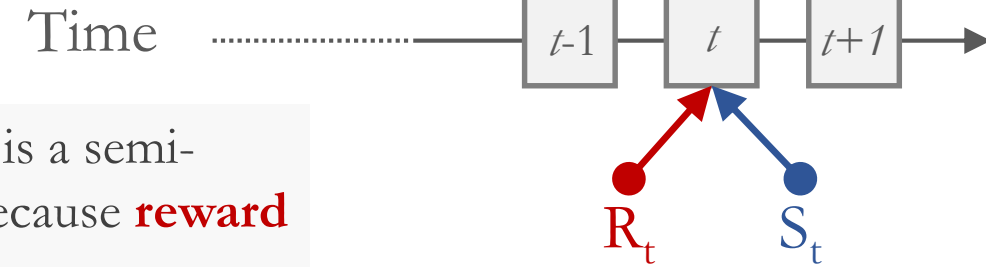
Introduction to reinforcement learning and deep reinforcement learning

Markov Decision Process (MDP)

# Reinforcement learning

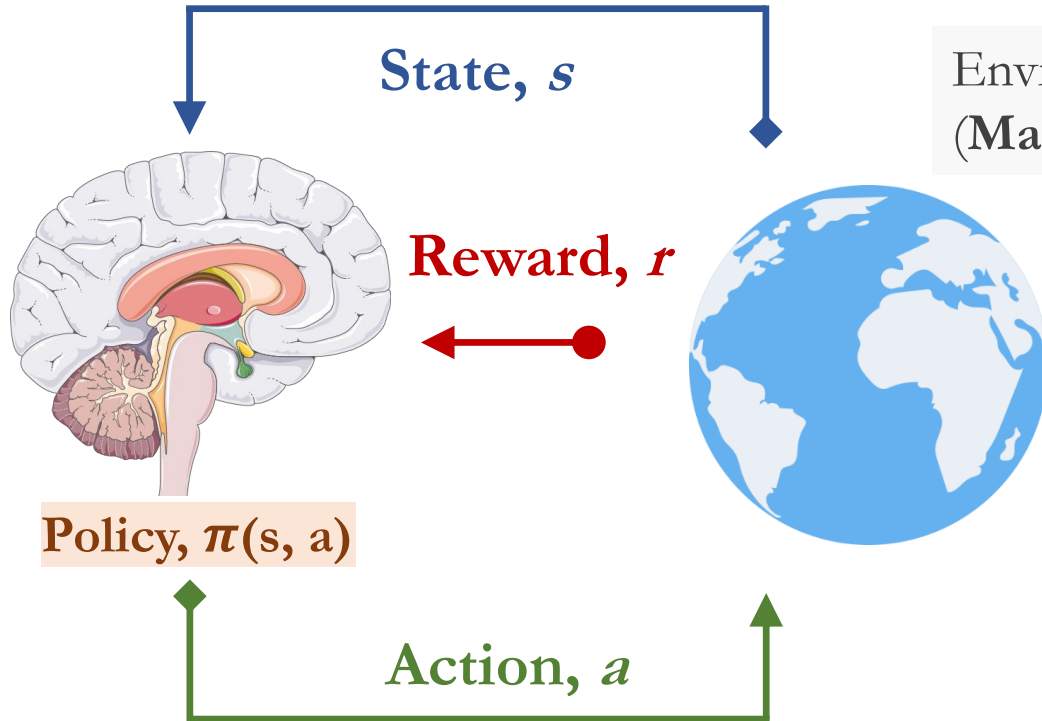
Reinforcement learning is a framework for learning how to interact with the environment from experience.

Most of the time, RL is a semi-supervised learning because **reward** is time-delayed

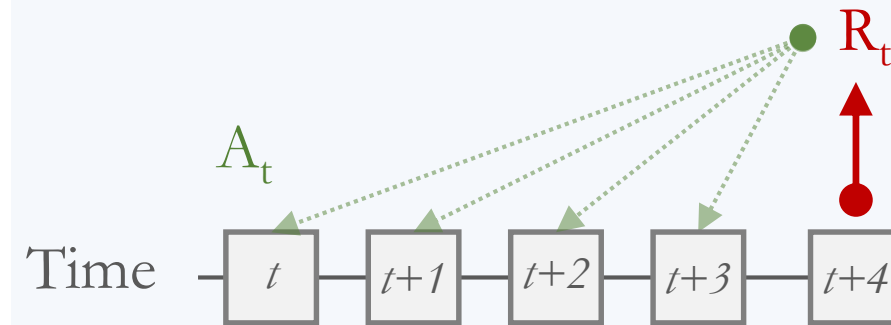


Exploration | Exploitation

Environment is modelled as probabilistic  
(**Markov Decision Process, MDP**)



Credit Assignment Problem



Source: <https://www.youtube.com/watch?v=0MNVhXEX9to>



Advanced cognitive modeling • Spring 2021

• Nicolas Legrand • Postdoctoral fellow • Embodied Computation Group



ECG  
embodied  
computation  
group

AARHUS UNIVERSITY

# Key concepts

**Model:** predict what the environment will do next.

$$p(s', r | s, a) = P(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

**Value function:** prediction of expected rewards.

$$v_{\pi}(s) = \mathbb{E}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s]$$

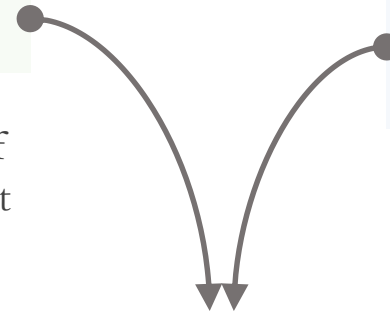
*Discount rate*

The value of a state **s** given a policy  $\pi$  is my expectation of how much reward I will get in the future if I start in that state and enact that policy.

**Policy:** how the agent pick its actions.

**Deterministic**       $\alpha = \pi(s)$

**Stochastic**       $\alpha \sim \pi(a|s)$



**Q-learning**

$Q^{\pi}(s, a)$  = quality of state/action pair

$$Q(s, a) = Q^{old}(s_t, a_t) + \alpha(r_t + \max_a Q(S_{t+1}, a) - Q^{old}(s_t, a_t))$$

Given a state **s** and an action **a**, and assuming that I will do the best thing I can in the future, what is the quality of being in that state and taking that action.

**Policy learning | Value learning**

Source: <https://www.youtube.com/watch?v=K67RJH3V7Yw&list=PLMsTLcO6ettgmyLVrcPvFLYi2Rs-R4JOE&index=4>



**Advanced cognitive modeling • Spring 2021**

• Nicolas Legrand • Postdoctoral fellow • Embodied Computation Group



ECG  
embodied  
computation  
group

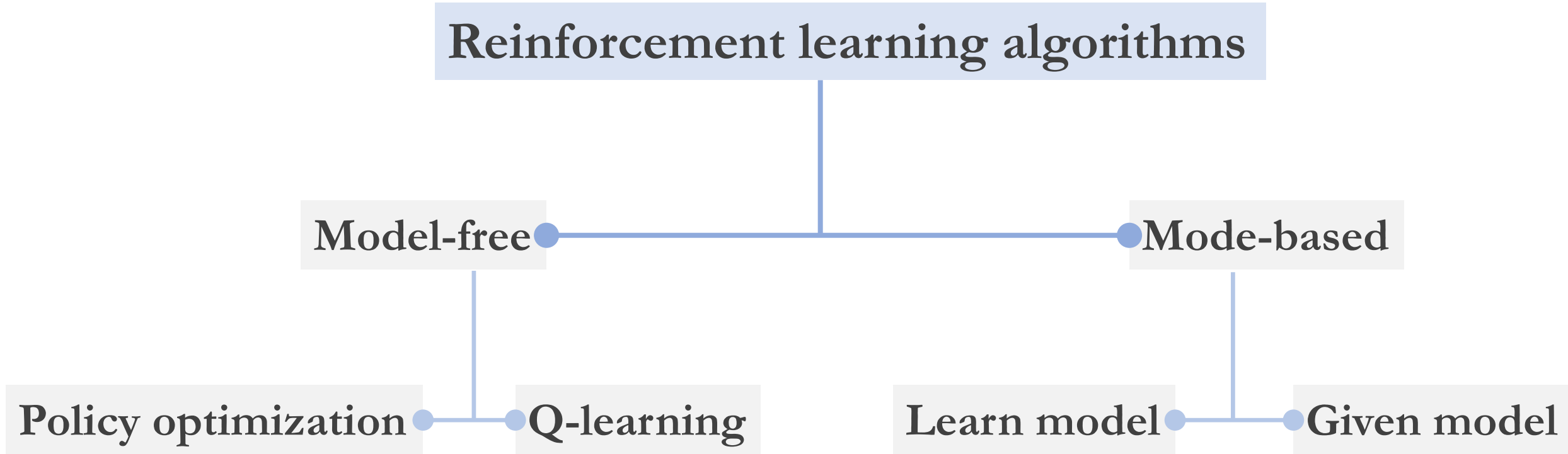
AARHUS UNIVERSITY

# RL Algorithms

Hindsight Experience Replay

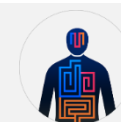
Save all behaviors and code reward for different goal.

[https://www.youtube.com/watch?v=0Ey02HT\\_1Ho](https://www.youtube.com/watch?v=0Ey02HT_1Ho)



Advanced cognitive modeling • Spring 2021

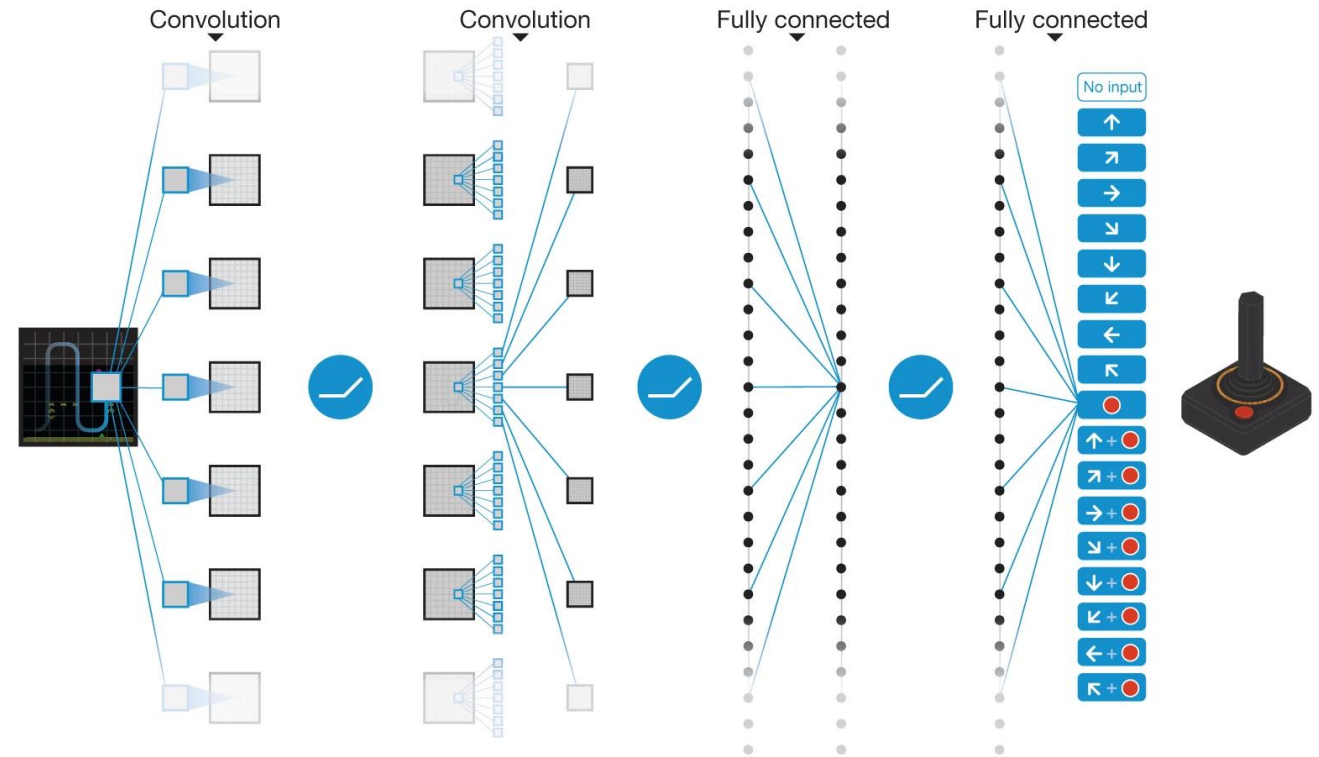
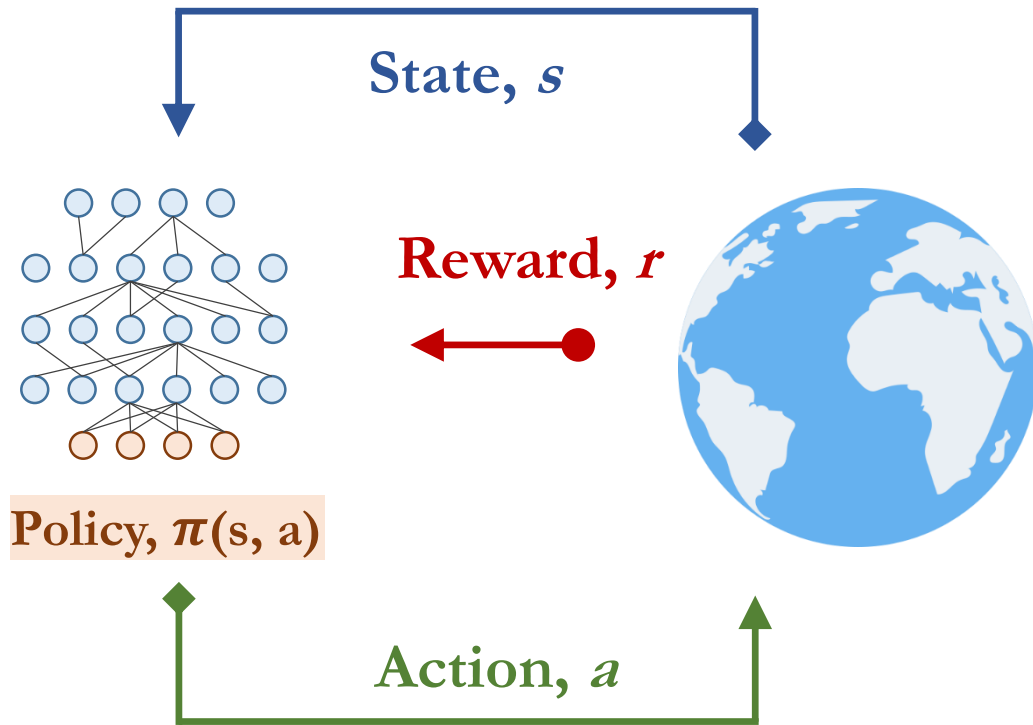
• Nicolas Legrand • Postdoctoral fellow • Embodied Computation Group



ECG  
embodied  
computation  
group

AARHUS UNIVERSITY

# Deep reinforcement learning



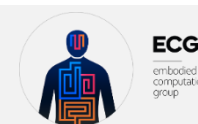
Mnih et al. (2015)

Source: <https://www.youtube.com/watch?v=IUiKAD6cuTA>



Advanced cognitive modeling • Spring 2021

• Nicolas Legrand • Postdoctoral fellow • Embodied Computation Group



ECG  
embodied  
computation  
group

AARHUS UNIVERSITY

# Examples

## Hide and seek

<https://www.youtube.com/watch?v=Lu56xVlZ40M>

## Flexible muscle-based locomotion for bipedal creatures

<https://vimeo.com/79098420>

## Atari video games

<https://www.youtube.com/watch?v=TmPfTpjtdgg&t=43s>

## AlphaGo Move 37

<https://www.youtube.com/watch?v=JNrXgpSEEIE>

## Cart-Pole

<https://www.youtube.com/watch?v=XiigTGKZfks>



# Markov Decision Processes



Advanced cognitive modeling • Spring 2021

• Nicolas Legrand • Postdoctoral fellow • Embodied Computation Group



ECG  
embodied  
computation  
group

AARHUS UNIVERSITY