# Lecture 1: Introduction to Reinforcement Learning

## The RL Problem

### Rewards

1. A **reward** $R_t$ is a scalar feedback signal
2. Indicates how well agent is doing at step $t$
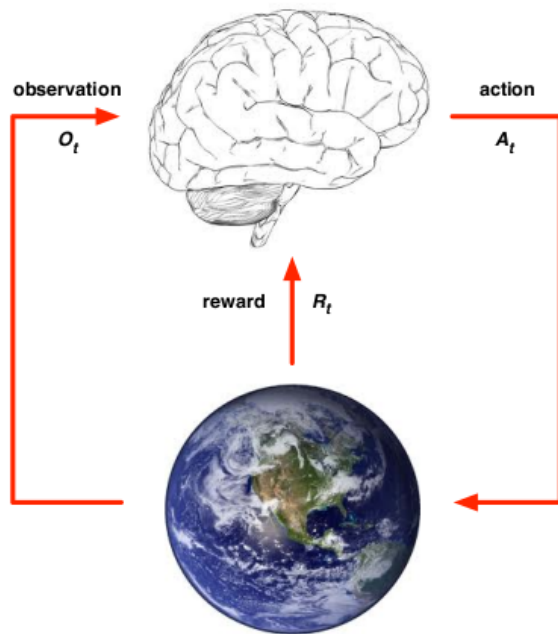3. The agent's job is to maximise cumulative reward

**Reward hypothesis** : All goals can be described by the maximisation of the expected cumulative reward.

### Sequential Decision Making

- Goal: select actions to miximise total future reward

- Actions may have long term consequences

- Reward may be delayed

- It may be better to sacrifice immediate reward to gain more long-term reward

- Examples:

    - A financial investment (may take months to mature)
    - Refuelling a helicopter (might prevent a crash in several hours)
    - Blocking opponent moves (might help winning chances many moves from now)
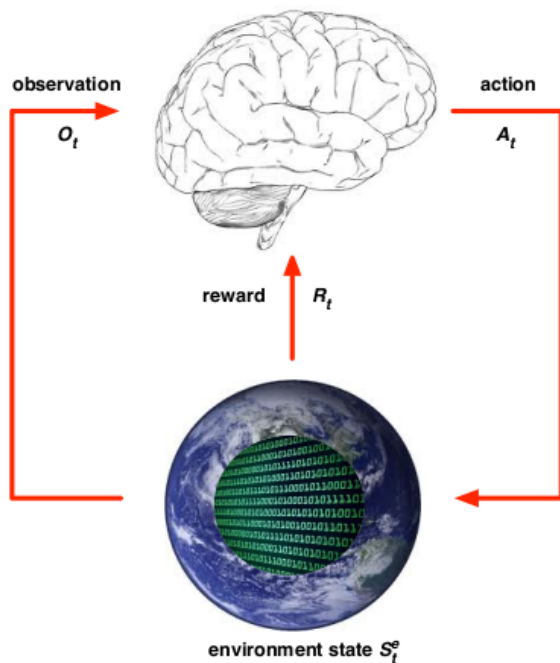
## Environments

# Agent and Environment



- At each step $t$ the agent:
    - Executes action $A_t$
    - Receives observation $O_t$
    - Receives scalar reward $R_t$
- The environment:
    - Receives action $A_t$
    - Emits observation $O_{t+1}$
    - Emits scalar reward $R_{t+1}$
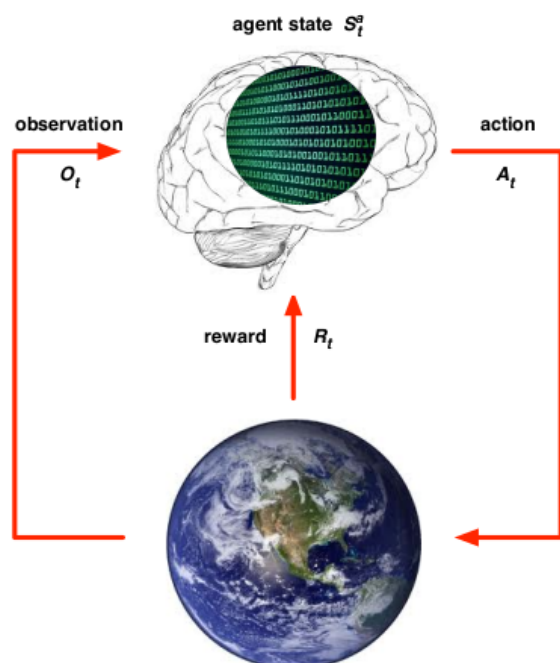- $t$ increments at env. step

- The **history** is the sequence of observations, actions, rewards
    - $H_t = O_1, R_1, A_1, \ldots, A_{t-1}, O_t, R_t$
- i.e. all observable variables up to time $t$
- i.e. the sensorimotor stream of a robot or embodied agent
- What happens next depends on the history:
    - The agent selects actions
    - The environment selects observations/rewards
- **State** is the information used to determine what happens next
- Formally, state is a function of the history:
    - $S_t = f(H_t)$

# Environment State



- The environment state $S_t^e$ is the environment's private representation
- i.e. whatever data the environment uses to pick the next observation/reward
- The environment state is not usually visible to the agent
- Even if $S_t^e$ is visible, it may contain irrelevant information

# Agent State



- The agent state $S_t^a$ is the agent's internal representation
- i.e. whatever information the agent uses to pick the next action
- i.e. it is the information used by reinforcement learning algorithms
- It can be any function of history:

$$S_t^a = f(H_t)$$

# Information State

An information state (a.k.a. Markov state) contains all useful information from the history.
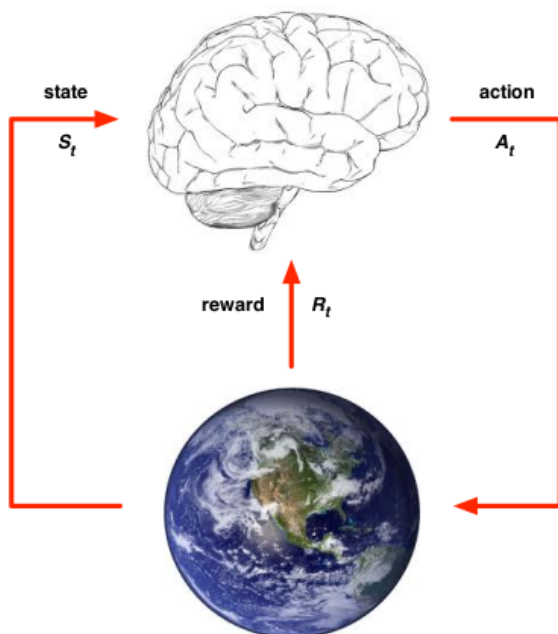
### Definition

A state $S_t$ is Markov if and only if

$$\mathbb{P}[S_{t+1} \mid S_t] = \mathbb{P}[S_{t+1} \mid S_1, ..., S_t]$$

- "The future is independent of the past given the present"

$$H_{1:t} \rightarrow S_t \rightarrow H_{t+1:\infty}$$

- Once the state is known, the history may be thrown away
- i.e. The state is a sufficient statistic of the future
- The environment state $S_t^e$ is Markov
- The history $H_t$ is Markov

# Fully Observable Environments



Full observability: agent directly observes environment state

$$O_t = S_t^a = S_t^e$$

- Agent state = environment state = information state
- Formally, this is a Markov decision process (MDP)
- (Next lecture and the majority of this course)

## Partially Observable Environments

- Partial observability: agent indirectly observes environment:
  - A robot with camera vision isn't told its absolute location
  - A trading agent only observes current prices
  - A poker playing agent only observes public cards
- Now agent state $\neq$ environment state
- Formally this is a partially observable Markov decision process (POMDP)
- Agent must construct its own state representation $S_t^a$, e.g.
  - Complete history: $S_t^a = H_t$
  - Beliefs of environment state: $S_t^a = (\mathbb{P}[S_t^e = s^1], ..., \mathbb{P}[S_t^e = s^n])$
  - Recurrent neural network: $S_t^a = \sigma(S_{t-1}^a W_s + O_t W_o)$

## Inside An RL Agent

### Major Components of an RL Agent

- An RL agent may include one or more of these components:
  - **Policy**: agent's behaviour function
  - **Value function**: how good is each state and/or action
  - **Model**: agent's representation of the environment

### Policy

- A policy is the agent's behaviour
- It is a map from state to action, e.g.
- Deterministic policy: $a = \pi(s)$
- Stochastic policy: $\pi(a|s) = P[A_t = a|S_t = s]$

### Value Function

- Value function is a prediction of future reward
- Used to evaluate the goodness/badness of states
- And therefore to select between actions, e.g.
  - $v_\pi(s) = E_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots |S_t = s]$

### Model

- A model predicts what the environment will do next
- $P$ predicts the next state
- $R$ predicts the next (immediate) reward, e.g.

$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$
$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

### Categorizing RL agents

- Value Based
  - No Policy (Implicit)
  - Value Function
- Policy Based
  - Policy
  - No Value Function
- Actor Critic
  - Policy
  - Value Funciton
- Model Free
  - Policy and/or Value Function
  - No Model
- Model Based
  - Policy and/or Value Function
  - Model

## Learning and Planning

### Two fundamental Problems in sequential decision making

- Reinforcement Learning:
  - The environment is initially unknown
  - The agent interacts with the environment
  - The agent improves its policy
- Planning:
  - A model of the environment is known
  - The agent performs computations with its model (without any external interaction)
  - The agent improves its policy
  - a.k.a deliberation, reasoning, introspection, pondering, thought, search

## Exploration and Exploitation

- Reinforcement learning is like trial-and-error learning
- The agent should discover a good policy
- From its experiences of the environment
- Without losing too much reward along the way
- Exploration finds more information about the environment
- Exploitation exploits known information to maximise reward
- It is usually important to explore as well as exploit

## Prediction and Control

- Prediction: evaluate the future
  - Given a policy
- Control: optimise the future
  - Find the best policy