# Overview of Reinforcement Learning-Part I
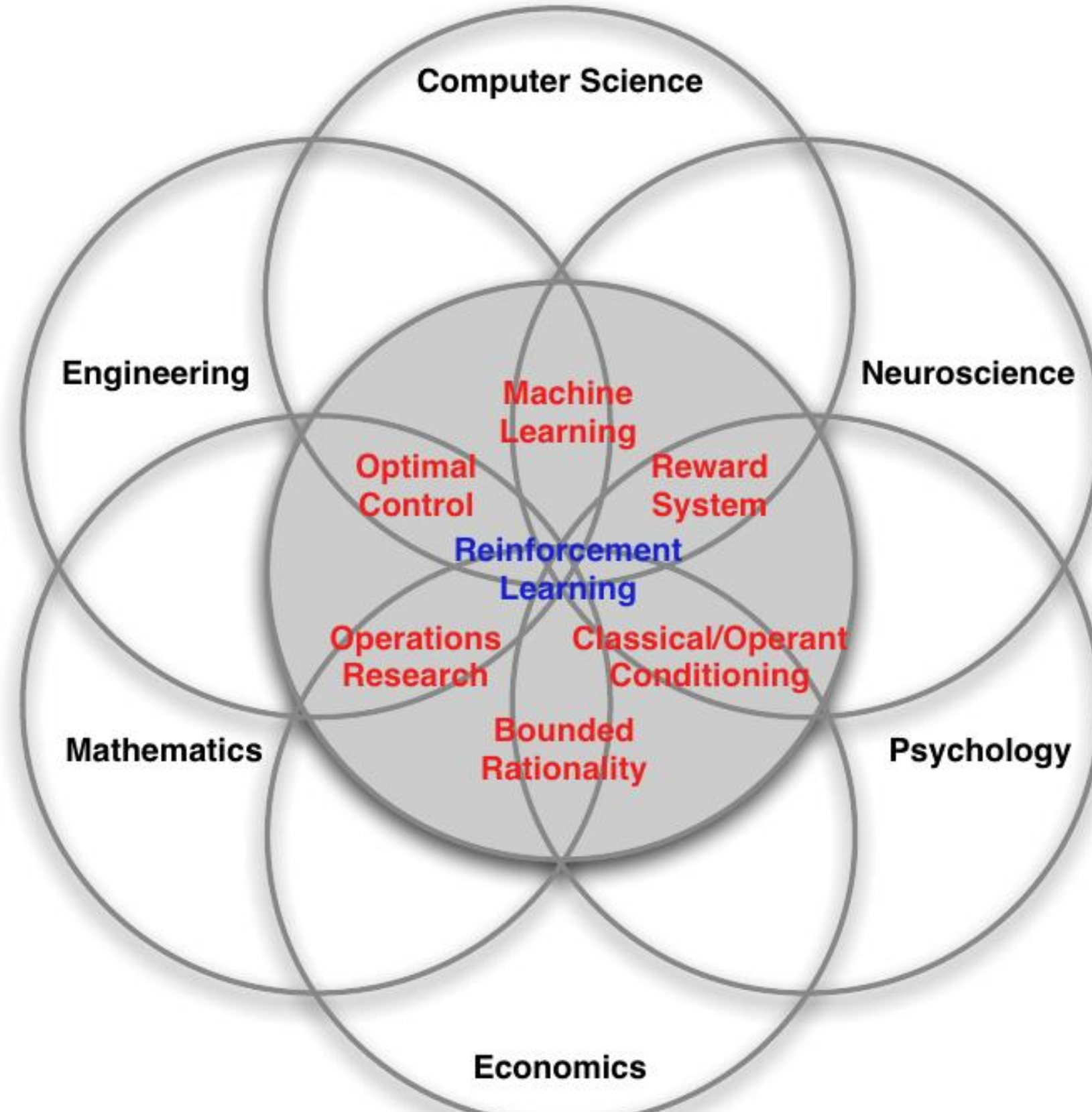
Prof. Jae Young Choi

Pattern Recognition and Machine Intelligence Lab. (PMI)
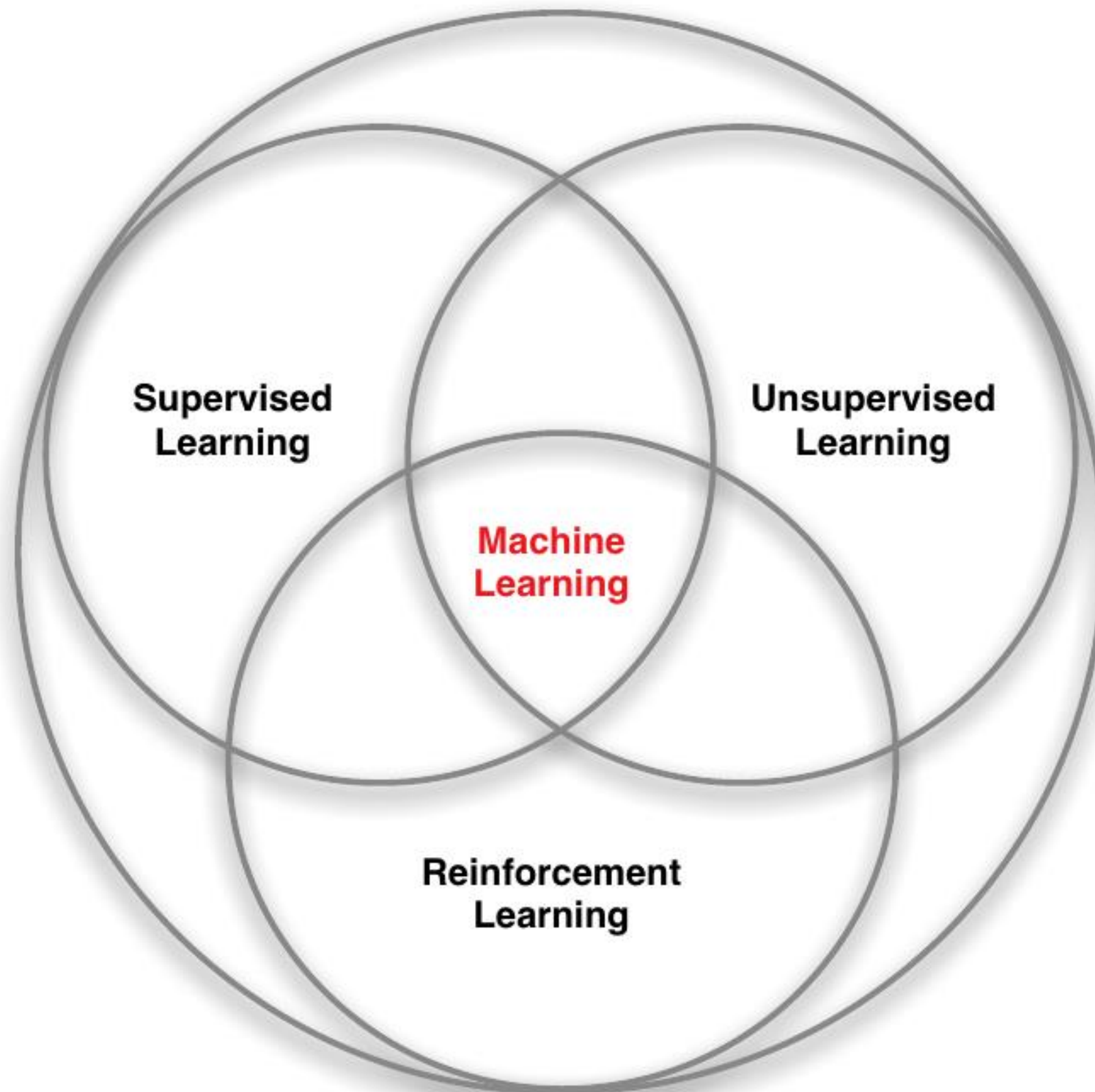
Division of Computer Engineering
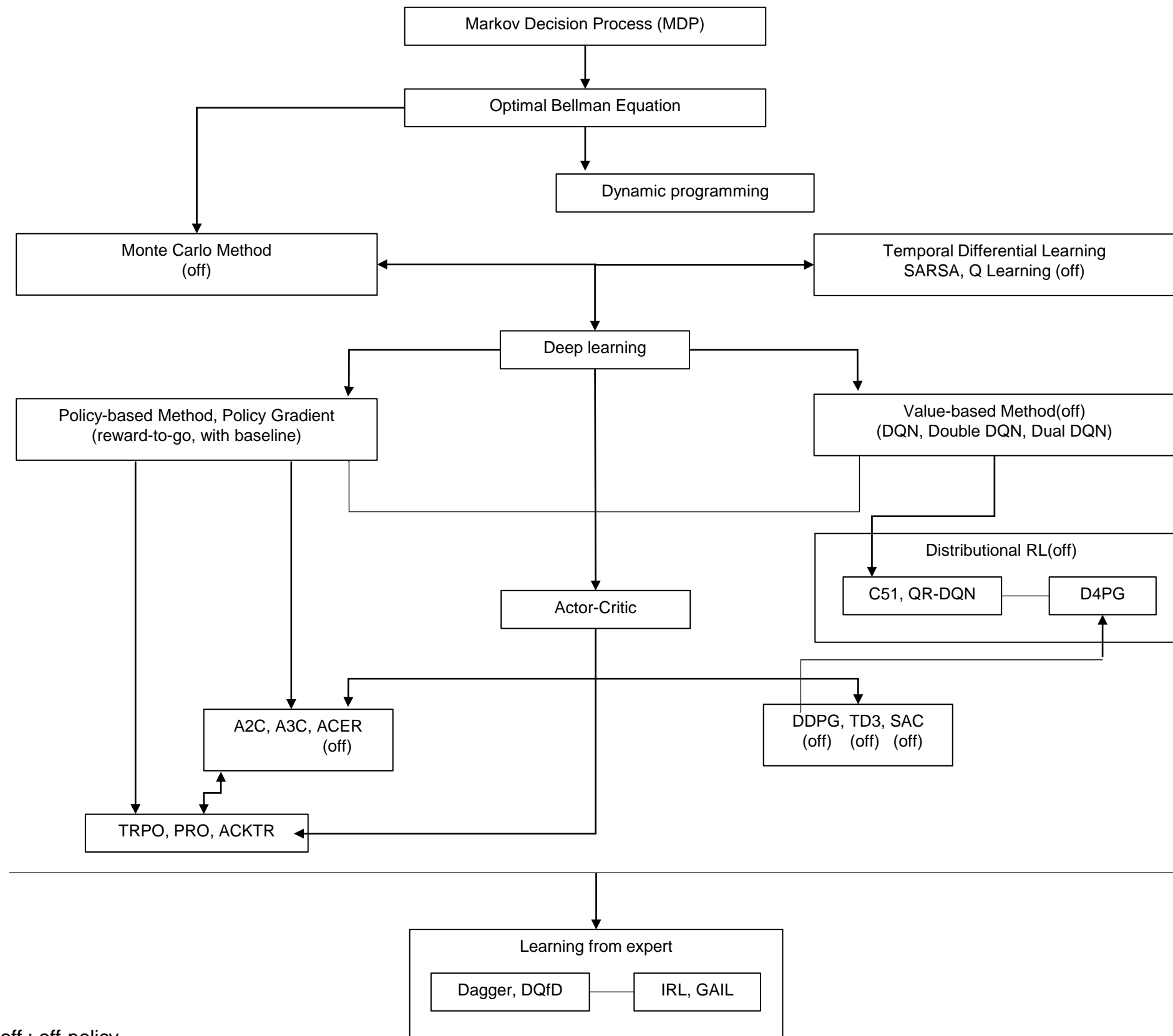
Hankuk University of Foreign Studies

# Overview of Reinforcement Learning Techniques



Markov Decision Process (MDP)

Optimal Bellman Equation

Dynamic programming

Monte Carlo Method (off)

Temporal Differential Learning SARSA, Q Learning (off)

Deep learning

Policy-based Method, Policy Gradient (reward-to-go, with baseline)

Value-based Method(off) (DQN, Double DQN, Dual DQN)

Distributional RL(off)

C51, QR-DQN — D4PG

Actor-Critic

A2C, A3C, ACER (off)

DDPG, TD3, SAC (off) (off) (off)

TRPO, PRO, ACKTR

Learning from expert

Dagger, DQfD — IRL, GAIL

*off : off-policy

4

# Characteristics of Reinforcement Learning

What makes reinforcement learning different from other machine learning paradigms?

- There is no supervisor, only a *reward* signal
- Feedback is delayed, not instantaneous
- Time really matters (sequential, non i.i.d data)
- Agent's actions affect the subsequent data it receives

# Examples of Reinforcement Learning

- Fly stunt manoeuvres in a helicopter

- Defeat the world champion at Backgammon

- Manage an investment portfolio

- Control a power station

- Make a humanoid robot walk

- Play many different Atari games better than humans

# Rewards

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

Reinforcement learning is based on the reward hypothesis

## Definition (Reward Hypothesis)

*All* goals can be described by the maximisation of expected cumulative reward
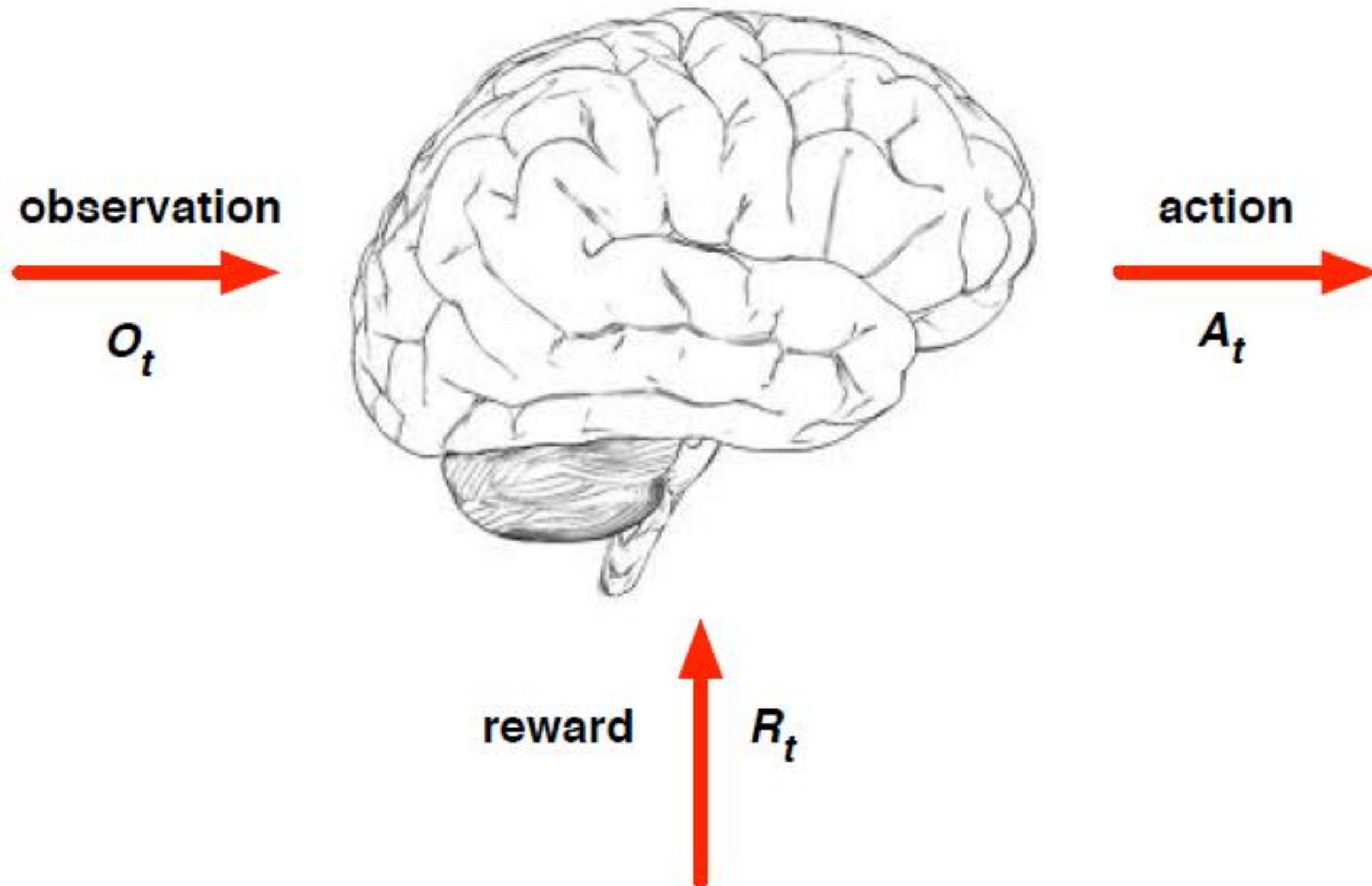
Do you agree with this statement?

# Example of Rewards

- Fly stunt manoeuvres in a helicopter
    - +ve reward for following desired trajectory
    - −ve reward for crashing
- Defeat the world champion at Backgammon
    - +/−ve reward for winning/losing a game
- Manage an investment portfolio
    - +ve reward for each $ in bank
- Control a power station
    - +ve reward for producing power
    - −ve reward for exceeding safety thresholds
- Make a humanoid robot walk
    - +ve reward for forward motion
    - −ve reward for falling over
- Play many different Atari games better than humans
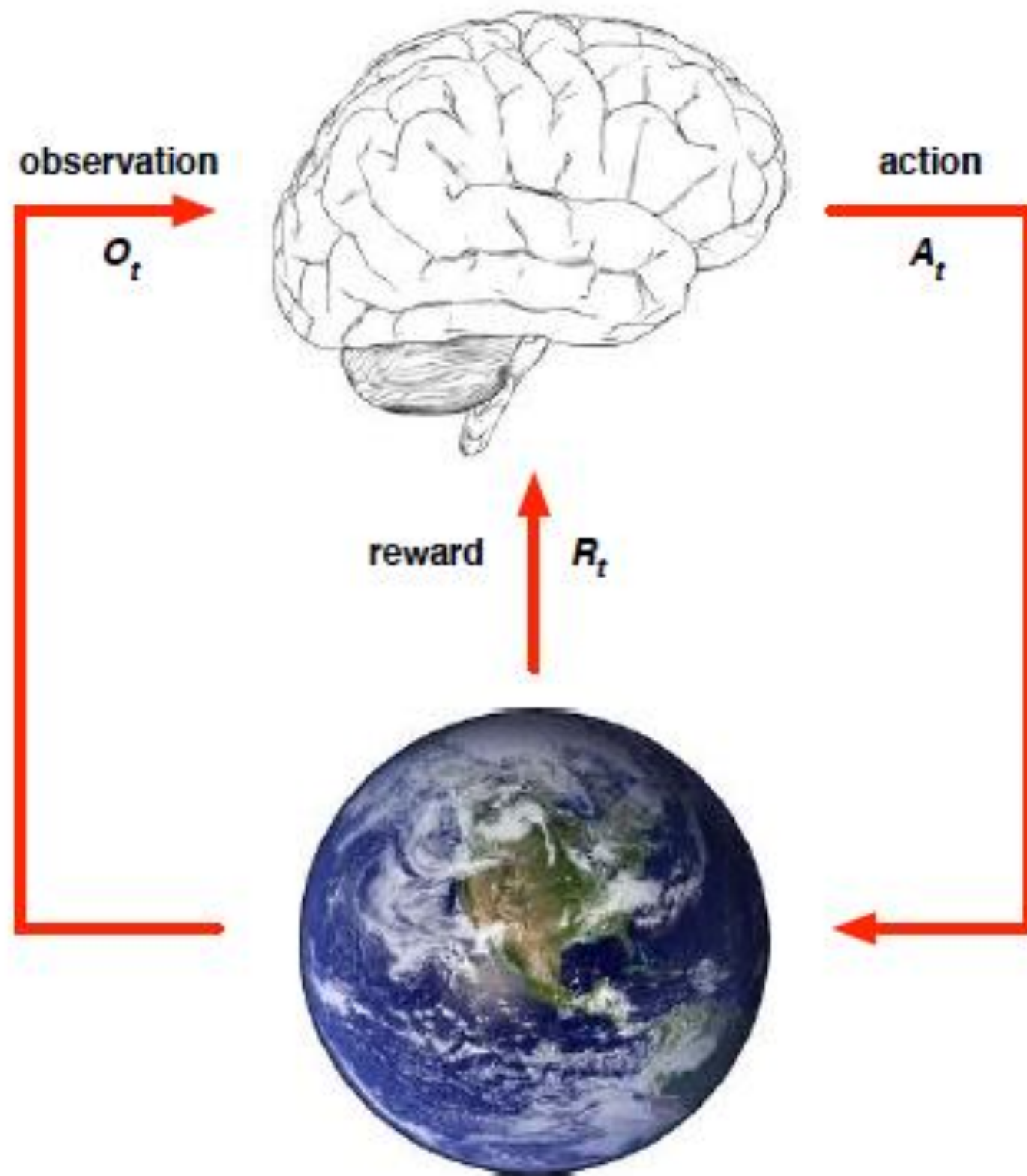    - +/−ve reward for increasing/decreasing score

# Sequential Decision Making

- Goal: *select actions to maximise 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)

# Agent and Environment

observation

$O_t$

action

$A_t$

reward $R_t$

# 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 $R_{t+1}$
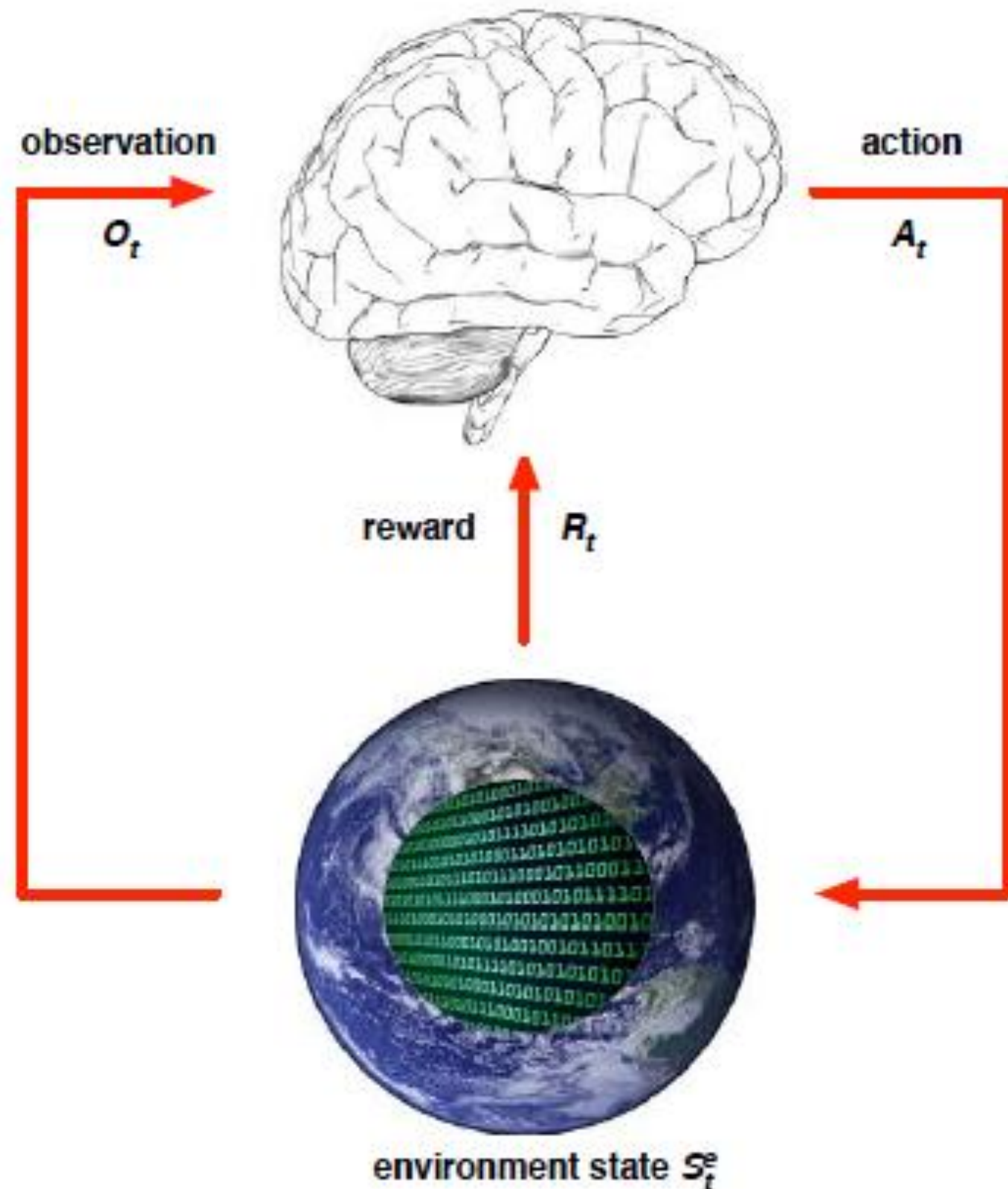- $t$ increments at env. step

# History and State

- The history is the sequence of observations, actions, rewards

$$H_t = O_1, R_1, A_1, ..., A_{t-1}, O_t, R_t$$

- 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
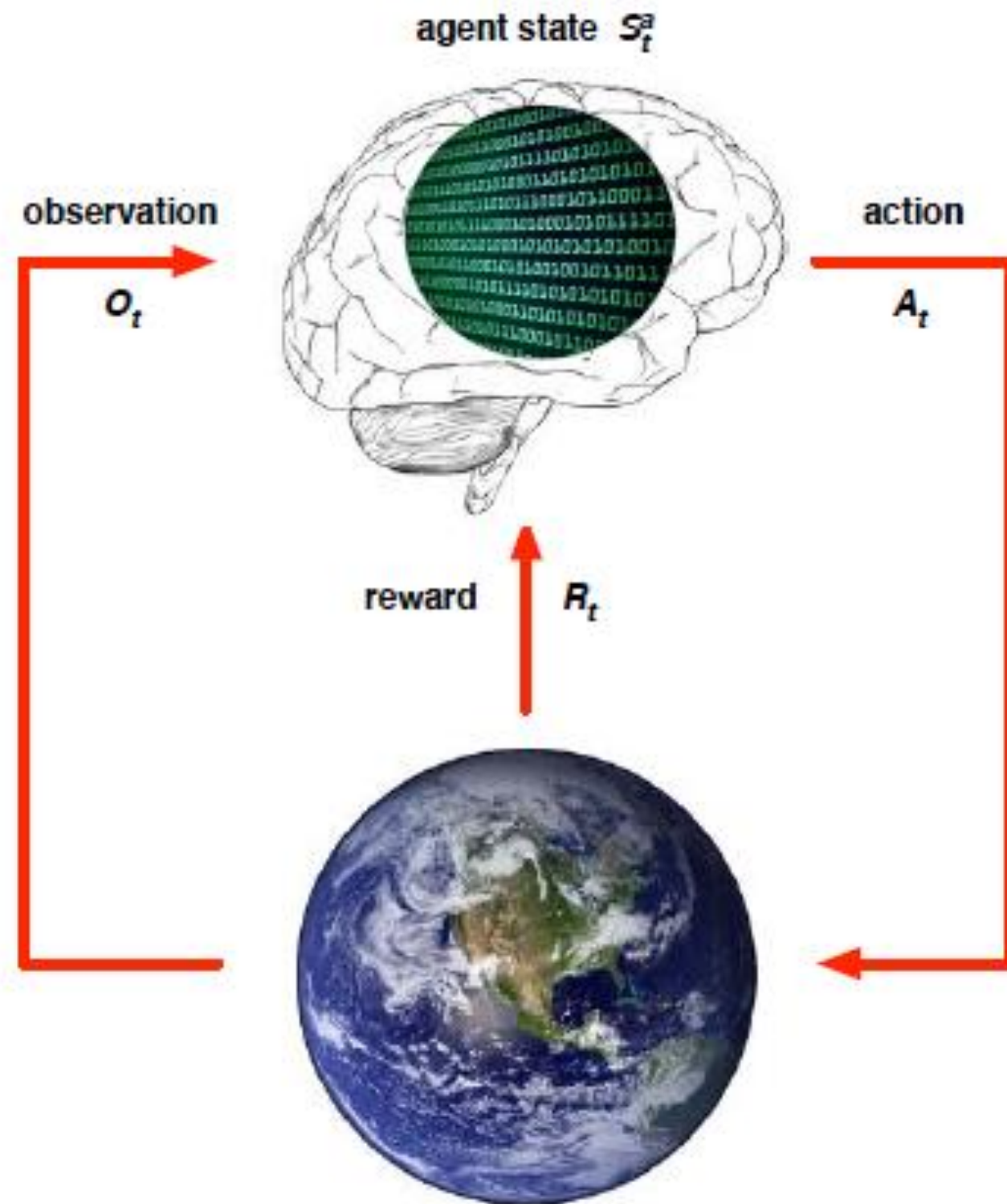
- Formally, state is a function of the history:

$$S_t = f(H_t)$$

# 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



agent state $S_t^a$

observation $O_t$

action $A_t$

reward $R_t$

- 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:

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

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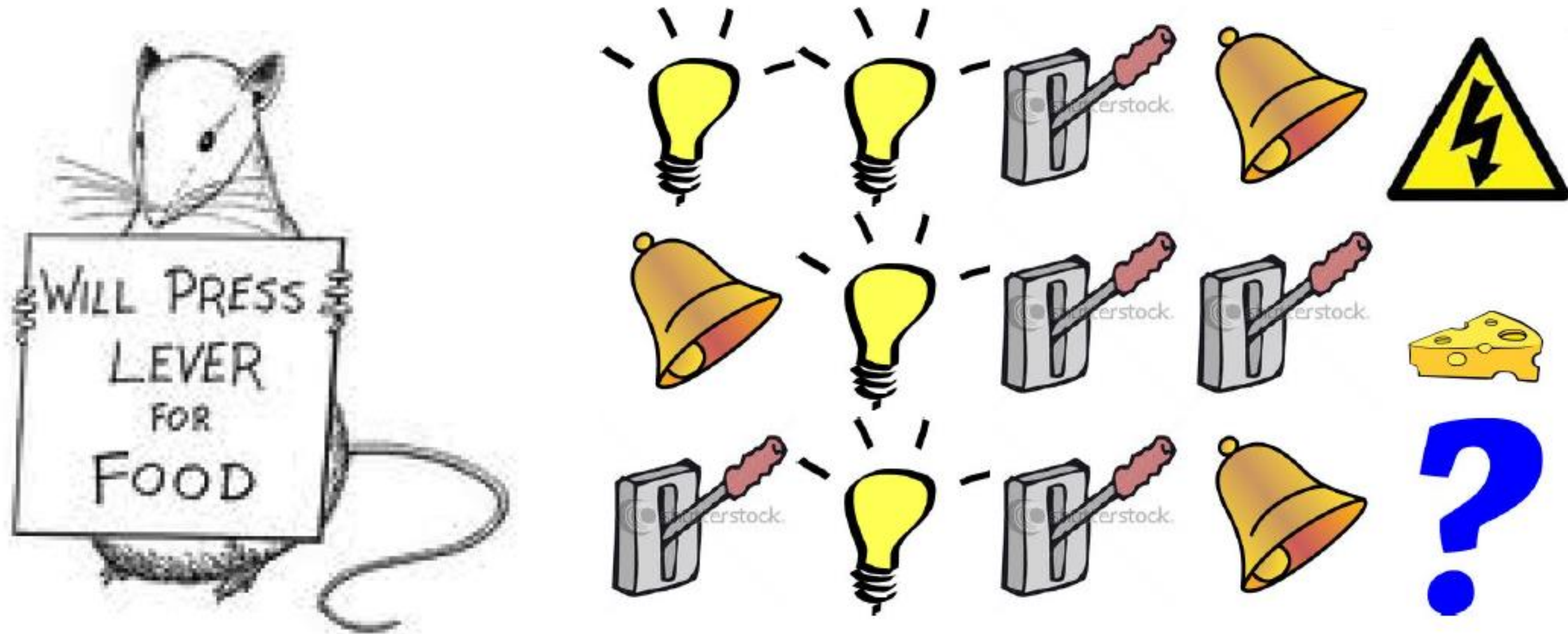
## Definition

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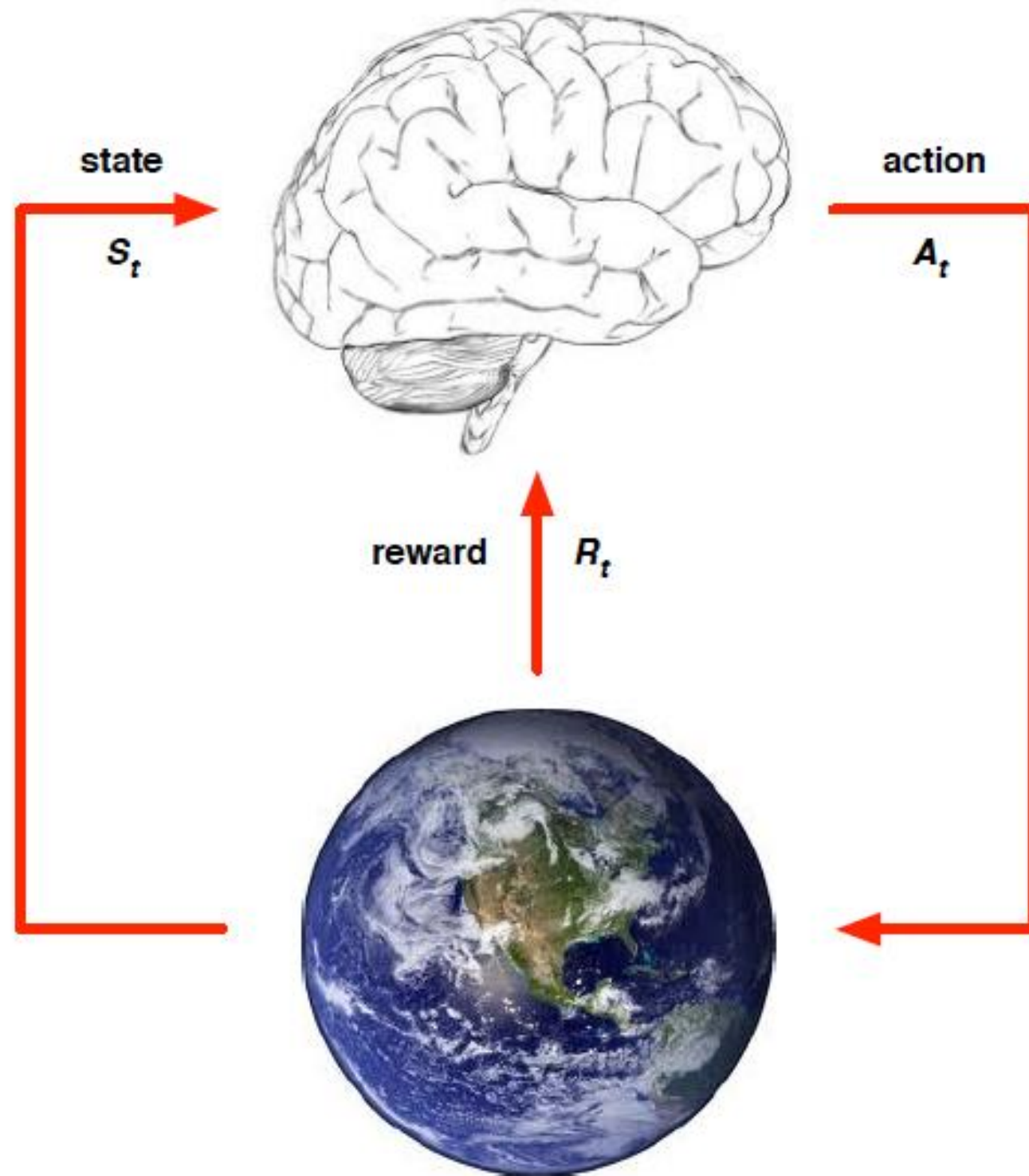
$$H_{1:t} \rightarrow S_t \rightarrow H_{t+1:\infty}$$

- 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

- What if agent state = last 3 items in sequence?
- What if agent state = counts for lights, bells and levers?
- What if agent state = complete sequence?

# Agent State



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)

# State



▶ Experience is a sequence of observations, actions, rewards

$$o_1, r_1, a_1, ..., a_{t-1}, o_t, r_t$$

▶ The state is a summary of experience

$$s_t = f(o_1, r_1, a_1, ..., a_{t-1}, o_t, r_t)$$

▶ In a fully observed environment

$$s_t = f(o_t)$$

# Episodes

❖ During agent's lifetime, its experience is presented as episodes

❖ Every episode is a sequence of observations (states), actions, rewards

| Episode 1 | $o_1,a_1,r_1$ | $o_2,a_2,r_2$ | $o_3,a_3,r_3$ | $o_4,a_4,r_4$ | $o_5,a_5,r_5$ | $o_6,a_6,r_6$ | $R = r_1 + r_2 + \ldots + r_6$ |
|---|---|---|---|---|---|---|---|
| Episode 2 | $o_1,a_1,r_1$ | $o_2,a_2,r_2$ | $o_3,a_3,r_3$ | $o_4,a_4,r_4$ | | | $R = r_1 + r_2 + r_3 + r_4$ |
| Episode 3 | $o_1,a_1,r_1$ | $o_2,a_2,r_2$ | $o_3,a_3,r_3$ | $o_4,a_4,r_4$ | $o_5,a_5,r_5$ | | $R = r_1 + r_2 + \ldots + r_5$ |
| Episode 4 | $o_1,a_1,r_1$ | $o_2,a_2,r_2$ | $o_3,a_3,r_3$ | | | | $R = r_1 + r_2 + r_3$ |

**Sample episodes with observations, actions, and rewards**