# Overview of Reinforcement Learning-Part II

Prof. Jae Young Choi

Pattern Recognition and Machine Intelligence Lab. (PMI)

Division of Computer Engineering

Hankuk University of Foreign Studies

# 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: $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$

# Value Function (1)

- Value function is a prediction of future reward
- Used to evaluate the goodness/badness of states
- And therefore to select between actions, e.g.

$$v_\pi(s) = \mathbb{E}_\pi \left[ R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots \mid S_t = s \right]$$

- A value function is a prediction of future reward
  - "How much reward will I get from action $a$ in state $s$?"
- $Q$-value function gives expected total reward
  - from state $s$ and action $a$
  - under policy $\pi$
  - with discount factor $\gamma$

$$Q^\pi(s, a) = \mathbb{E}\left[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \ldots \mid s, a\right]$$

- Value functions decompose into a Bellman equation

$$Q^\pi(s, a) = \mathbb{E}_{s', a'}\left[r + \gamma Q^\pi(s', a') \mid s, a\right]$$

► An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_\pi Q^\pi(s, a) = Q^{\pi^*}(s, a)$$
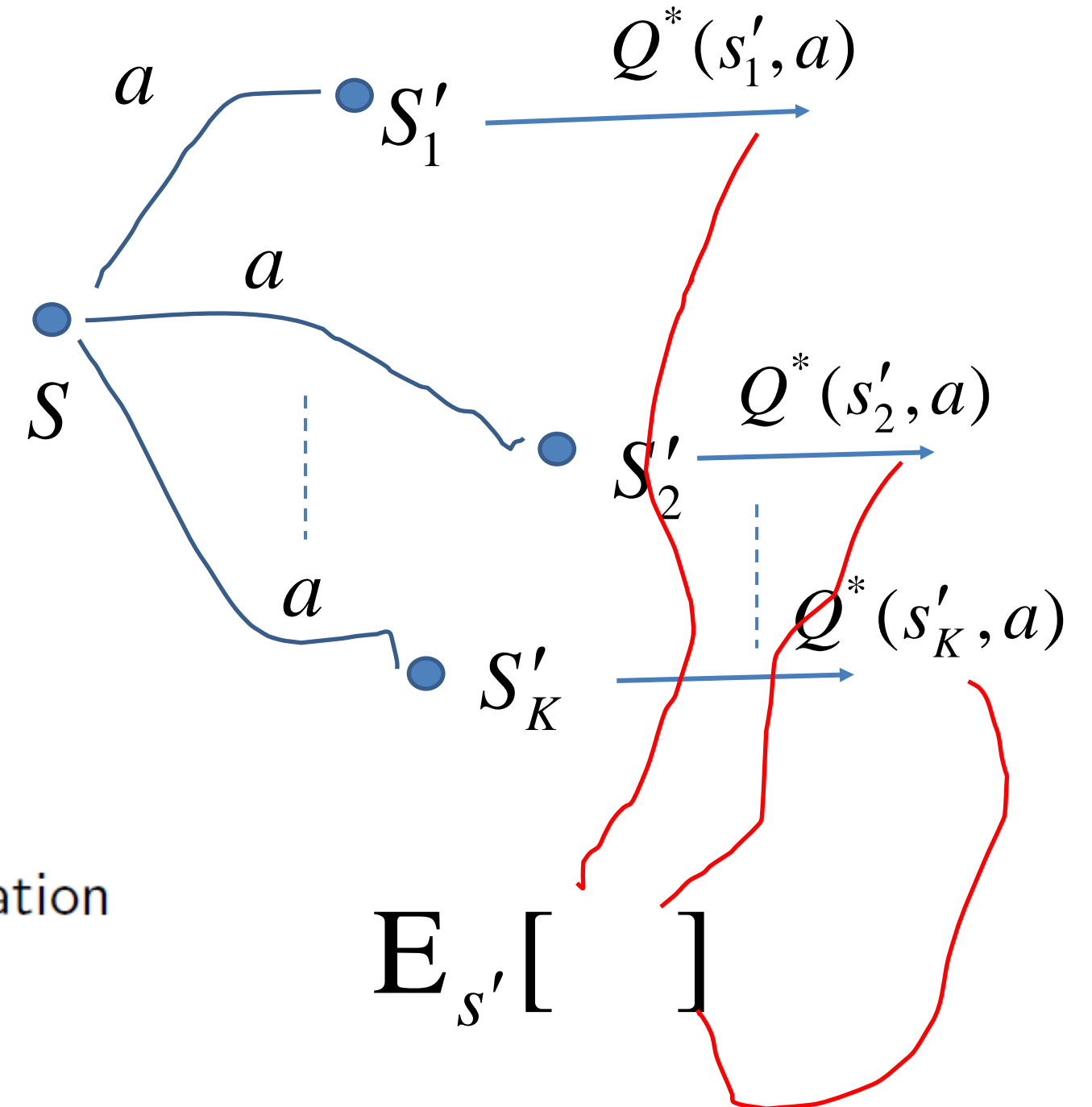
► Once we have $Q^*$ we can act optimally,

$$\pi^*(s) = \operatorname*{argmax}_a Q^*(s, a)$$

► Optimal value maximises over all decisions. Informally:

$$Q^*(s, a) = r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \ldots$$

$$= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$$

► Formally, optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'}\left[ r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

# Example of Value Function

- Each square represents a state
- **The return value is -1** everywhere on each transition
- **4 actions** for each state: north, south, east, west
- **Goal states** are the upper left corner and lower right corner



**Value function for randomly policy**



**Optimal Value function**



**Optimal Policy**

# Model

- A <span style="color:red">model</span> predicts what the environment will do next

- $\mathcal{P}$ predicts the next state

- $\mathcal{R}$ predicts the next (immediate) reward, e.g.

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = a]$$
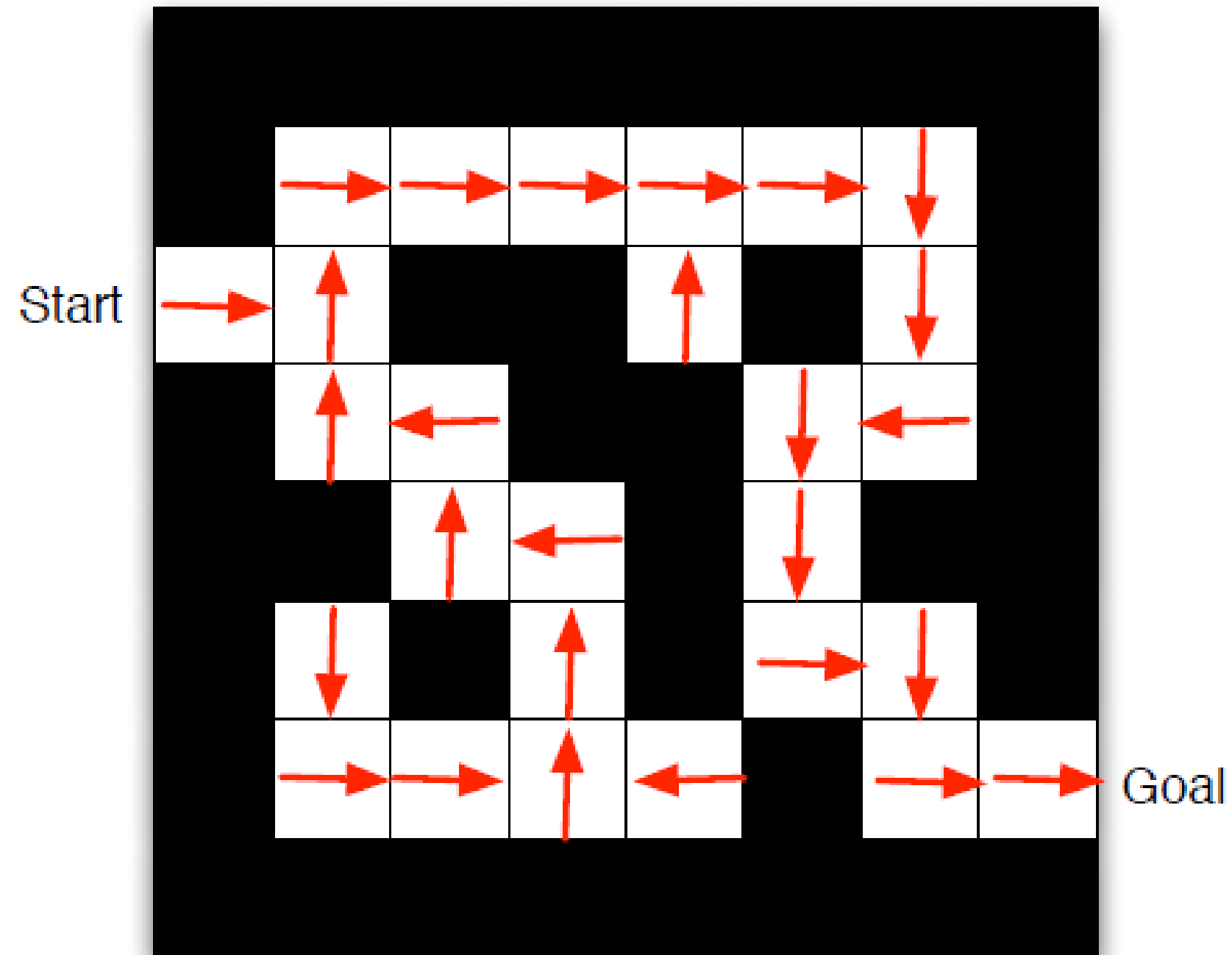$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$
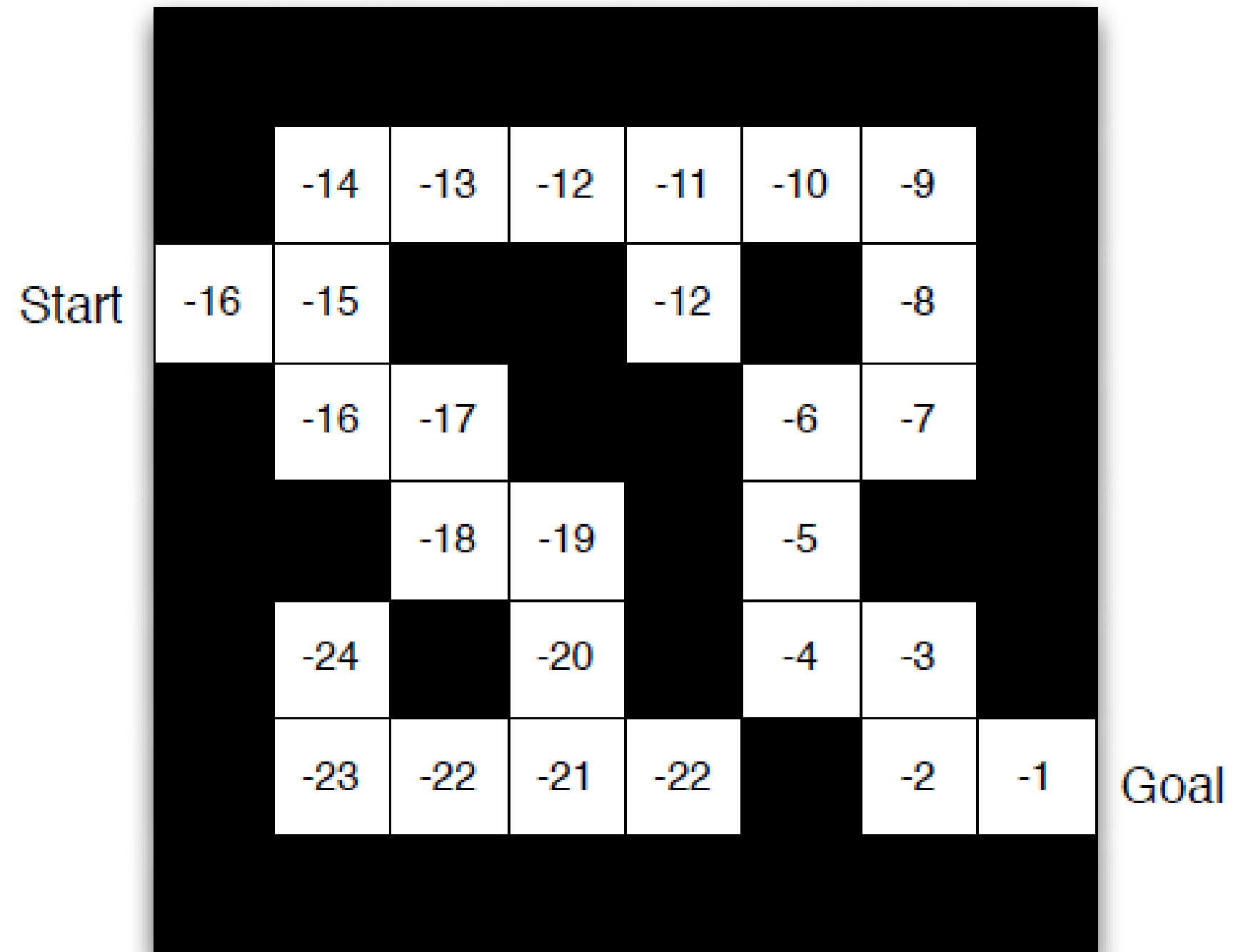
# Maze(미로) Example



- Rewards: -1 per time-step
- Actions: N, E, S, W
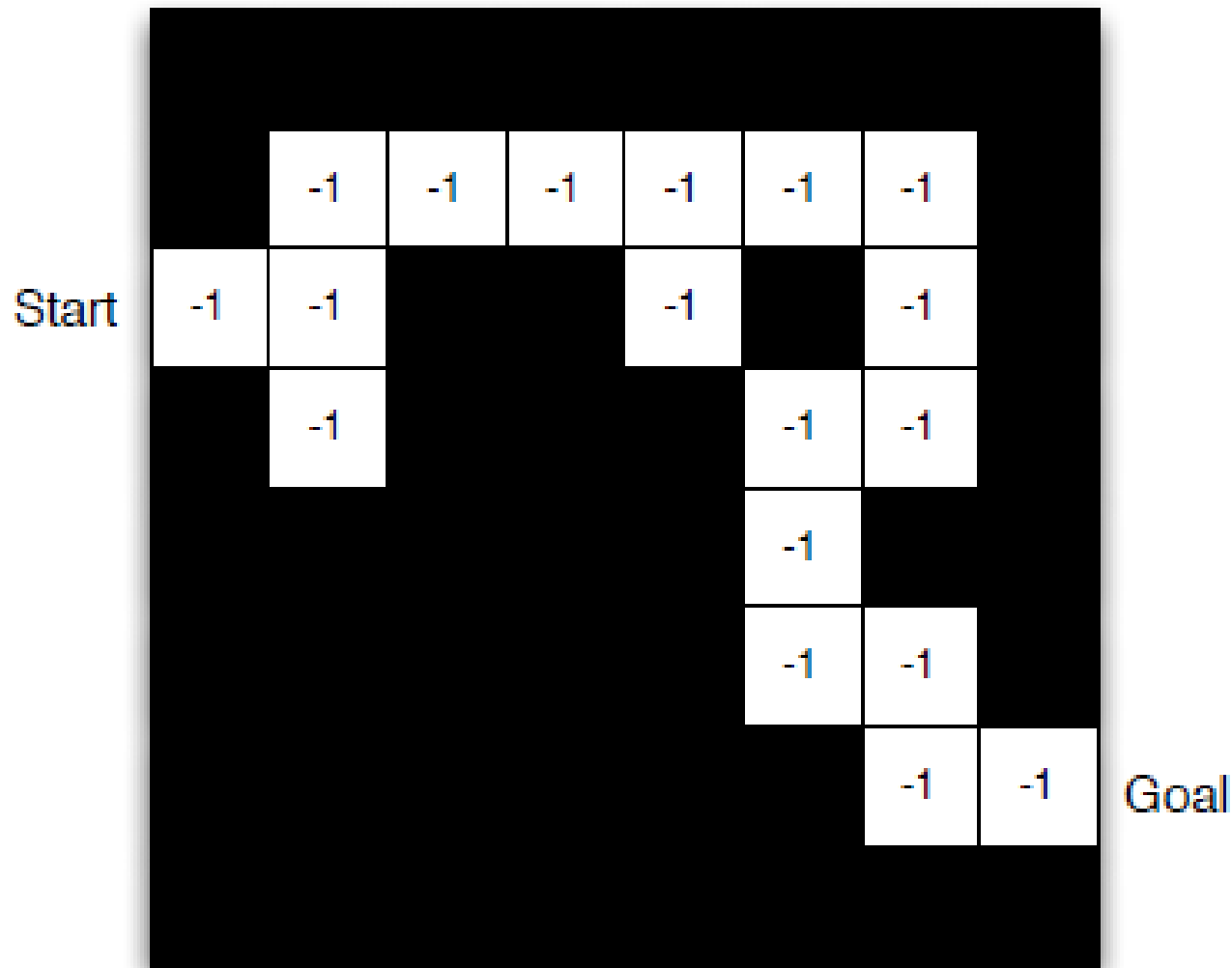- States: Agent's location

# Maze(미로) Example : Policy



■ Arrows represent policy $\pi(s)$ for each state $s$

# Maze(미로) Example : Value Function

- Numbers represent value $v_\pi(s)$ of each state $s$

# Maze(미로) Example : Model

- Agent may have an internal model of the environment
- Dynamics: how actions change the state
- Rewards: how much reward from each state
- The model may be imperfect

- Grid layout represents transition model $\mathcal{P}_{ss'}^{a} = P(S' \mid S, a)$
- Numbers represent immediate reward $\mathcal{R}_{s}^{a}$ from each state $s$ (same for all $a$)

Value-based RL

- ▶ Estimate the optimal value function $Q^*(s, a)$
- ▶ This is the maximum value achievable under any policy

Policy-based RL

- ▶ Search directly for the optimal policy $\pi^*$
- ▶ This is the policy achieving maximum future reward

Model-based RL

- ▶ Build a model of the environment
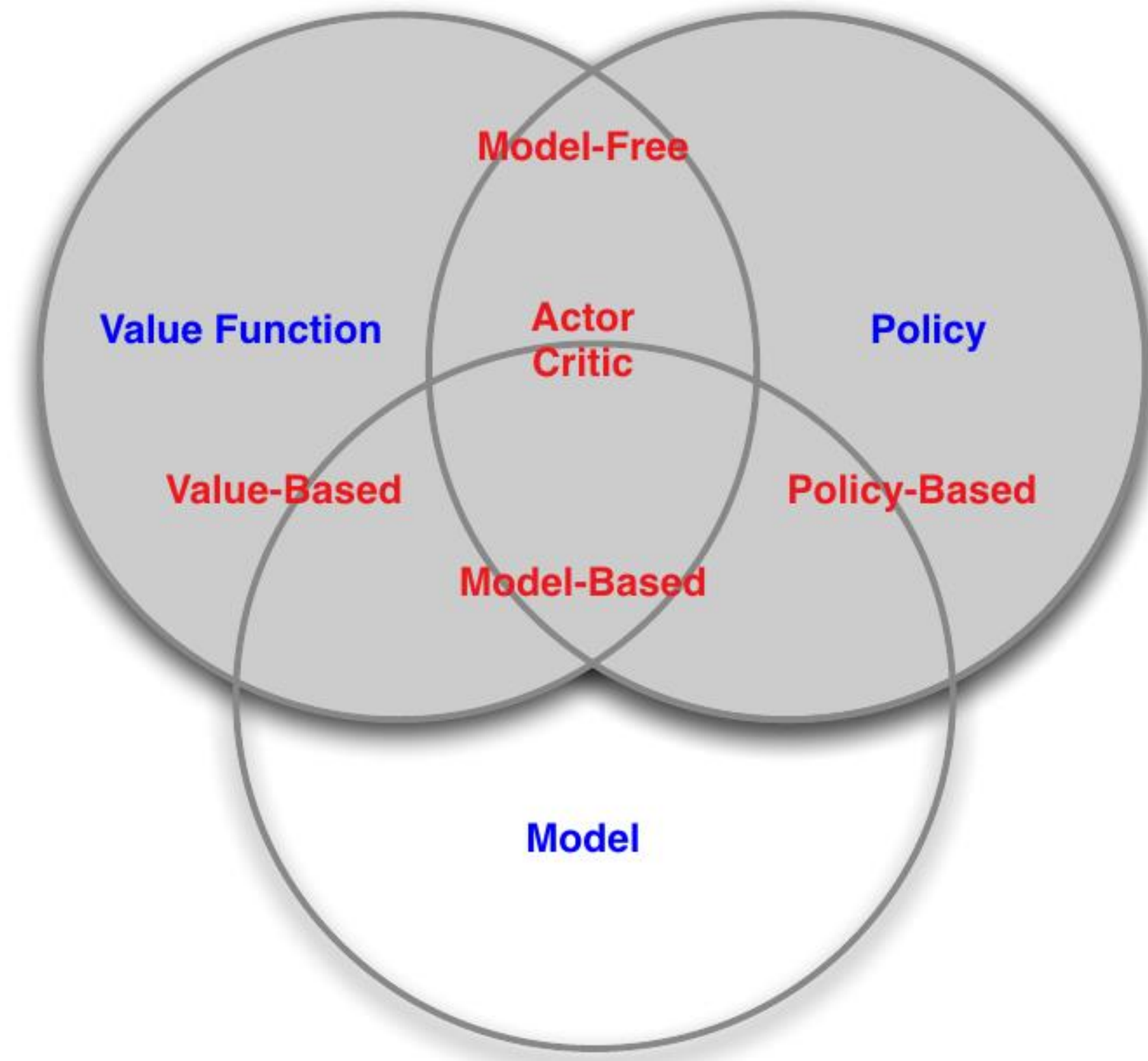- ▶ Plan (e.g. by lookahead) using model

# Categorizing RL agents (1)

- Value Based
  - No Policy (Implicit)
  - Value Function

- Policy Based
  - Policy
  - No Value Function

- Actor Critic
  - Policy
  - Value Function

- ■ Model Free
  - ■ Policy and/or Value Function
  - ■ No Model

- ■ Model Based
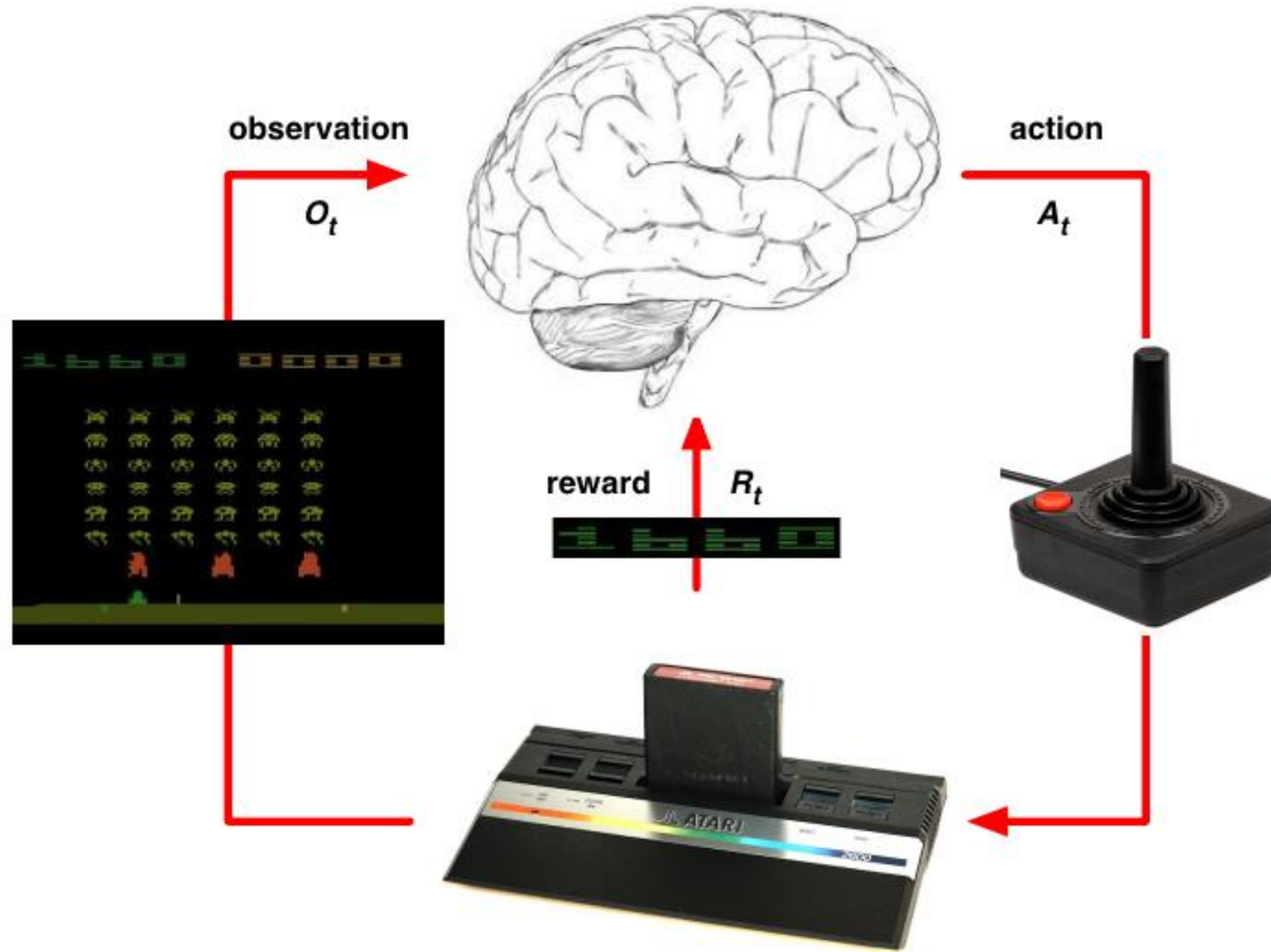  - ■ Policy and/or Value Function
  - ■ Model

# 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
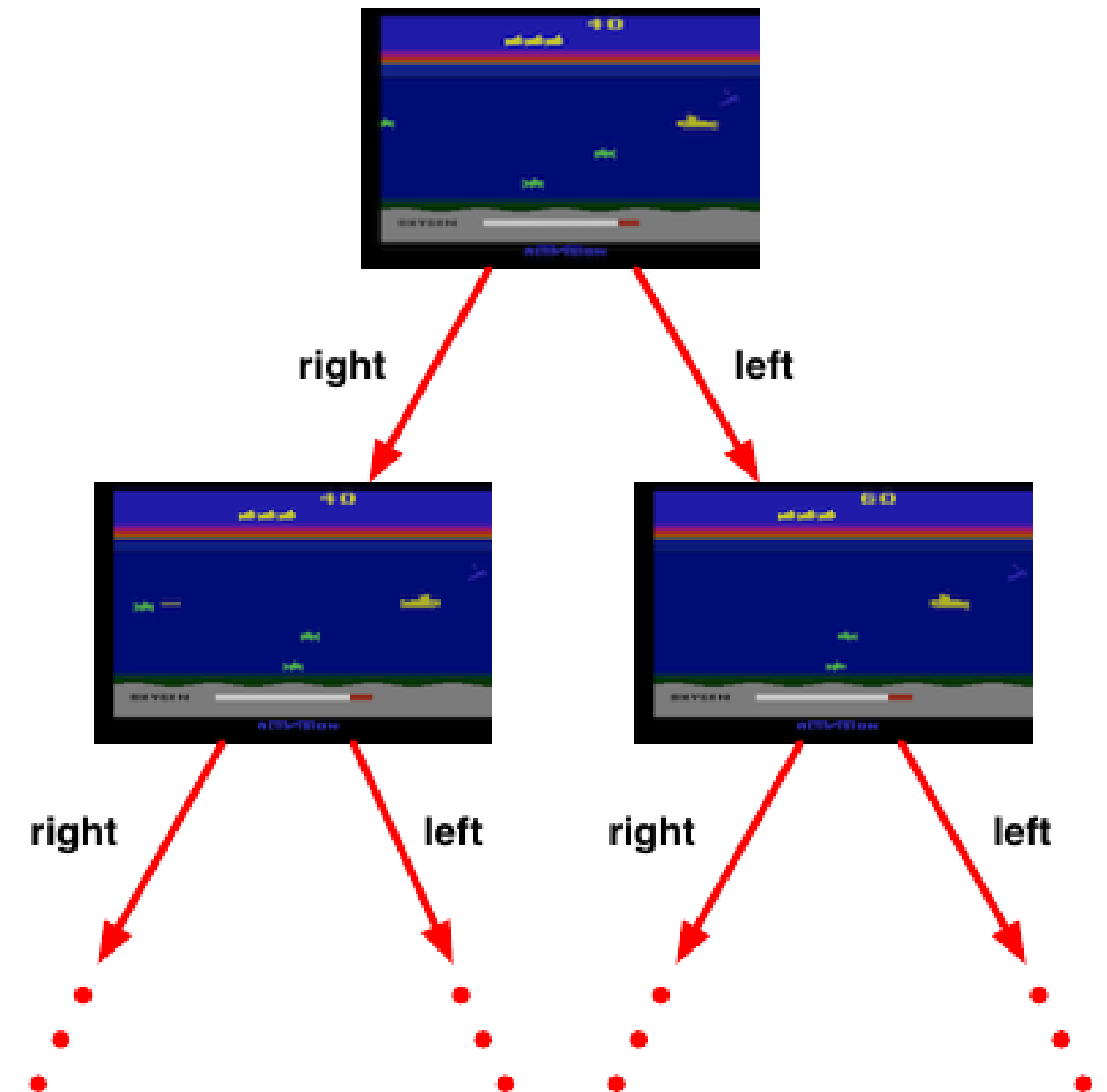
# Atari Example: Reinforcement Learning

observation

$O_t$

action

$A_t$

reward $R_t$

- Rules of the game are unknown
- Learn directly from interactive game-play
- Pick actions on joystick, see pixels and scores

- Rules of the game are known
- Can query emulator
  - perfect model inside agent's brain
- If I take action $a$ from state $s$:
  - what would the next state be?
  - what would the score be?
- Plan ahead to find optimal policy
  - e.g. tree search

# Exploration and Exploitation (1)

- 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 and Exploitation (2)

- *Exploration* finds more information about the environment
- *Exploitation* exploits known information to maximise reward
- It is usually important to explore as well as exploit

# Examples

- Restaurant Selection

  Exploitation  Go to your favourite restaurant
  Exploration  Try a new restaurant

- Online Banner Advertisements

  Exploitation  Show the most successful advert
  Exploration  Show a different advert

- Oil Drilling

  Exploitation  Drill at the best known location
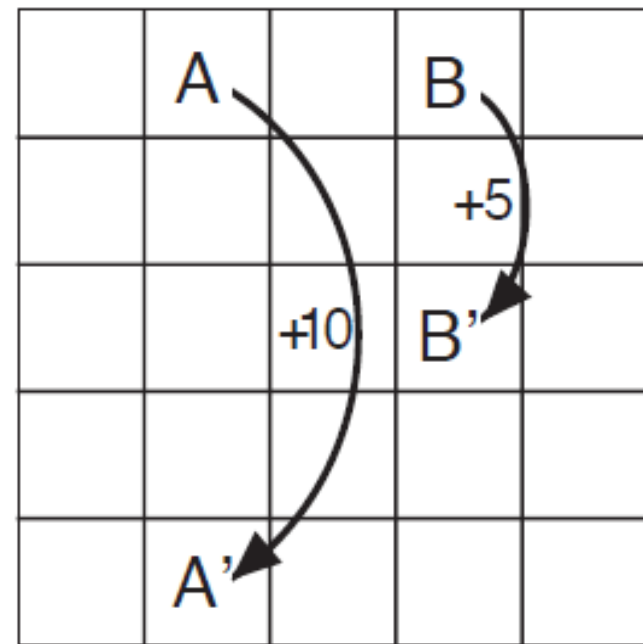  Exploration  Drill at a new location

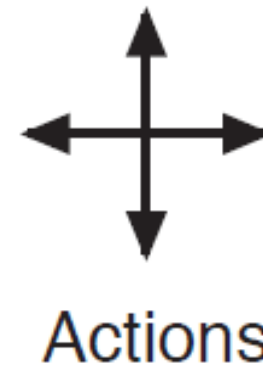- Game Playing

  Exploitation  Play the move you believe is best
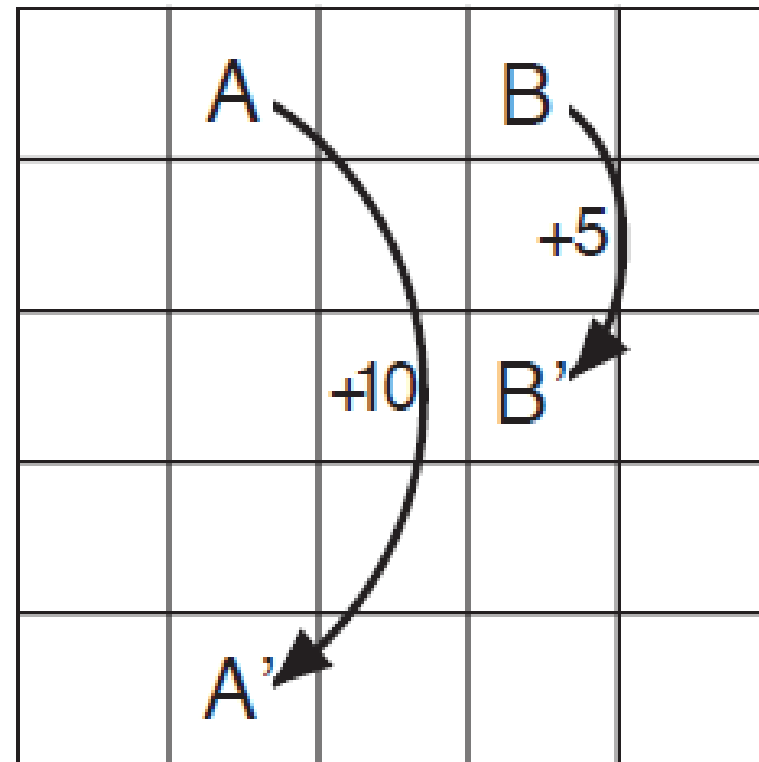  Exploration  Play an experimental move

(a)   Actions   (b)

What is the value function for the uniform random policy?

- At each grid cell, **four actions** are possible: north, south, east and west
- Actions taking agent off grid leaves its location unchanged but reward of '-1', other actions produces a reward of '0' except for state **A** and **B**
- From state A, all four actions yield a reward '+10' and take agent to A'
- From state B, all four actions yield a reward '+5' and take agent to B'
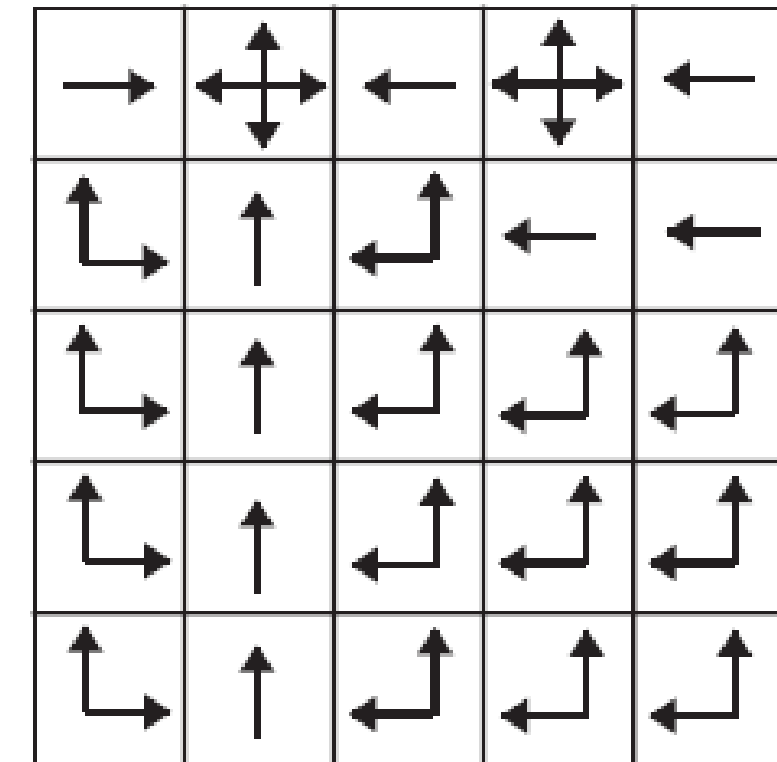- Random policy: agent selects all four actions with equal probability in all states

a) gridworld  b) $v_*$  c) $\pi_*$

What is the optimal value function over all possible policies?
What is the optimal policy?