

# Markov Decision Processes

CMPUT 261: Introduction to Artificial Intelligence

S&B §3.0-3.5

# Lecture Outline

1. Recap & Logistics
2. Markov Decision Processes
3. Returns & Episodes
4. Policies & Value Functions
5. Bellman Equations

*After this lecture, you should be able to:*

- define a Markov decision process
- represent a problem as a Markov decision process
- define a policy
- explain whether a task is episodic or continuing
- give expressions for the state-value function and the action-value function
- state the Bellman optimality equations
- give expressions for episodic and discounted continuing returns

# Logistics

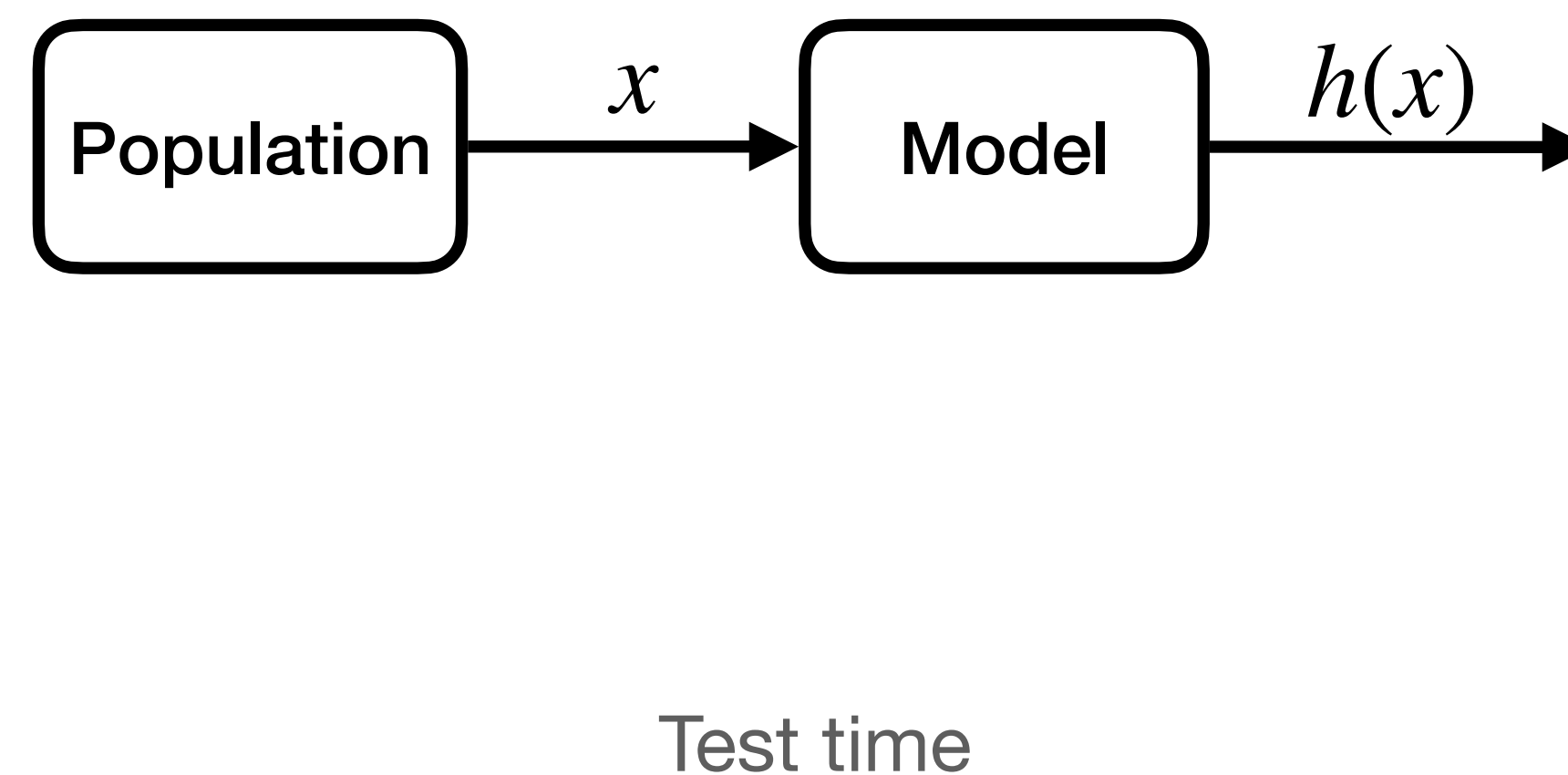
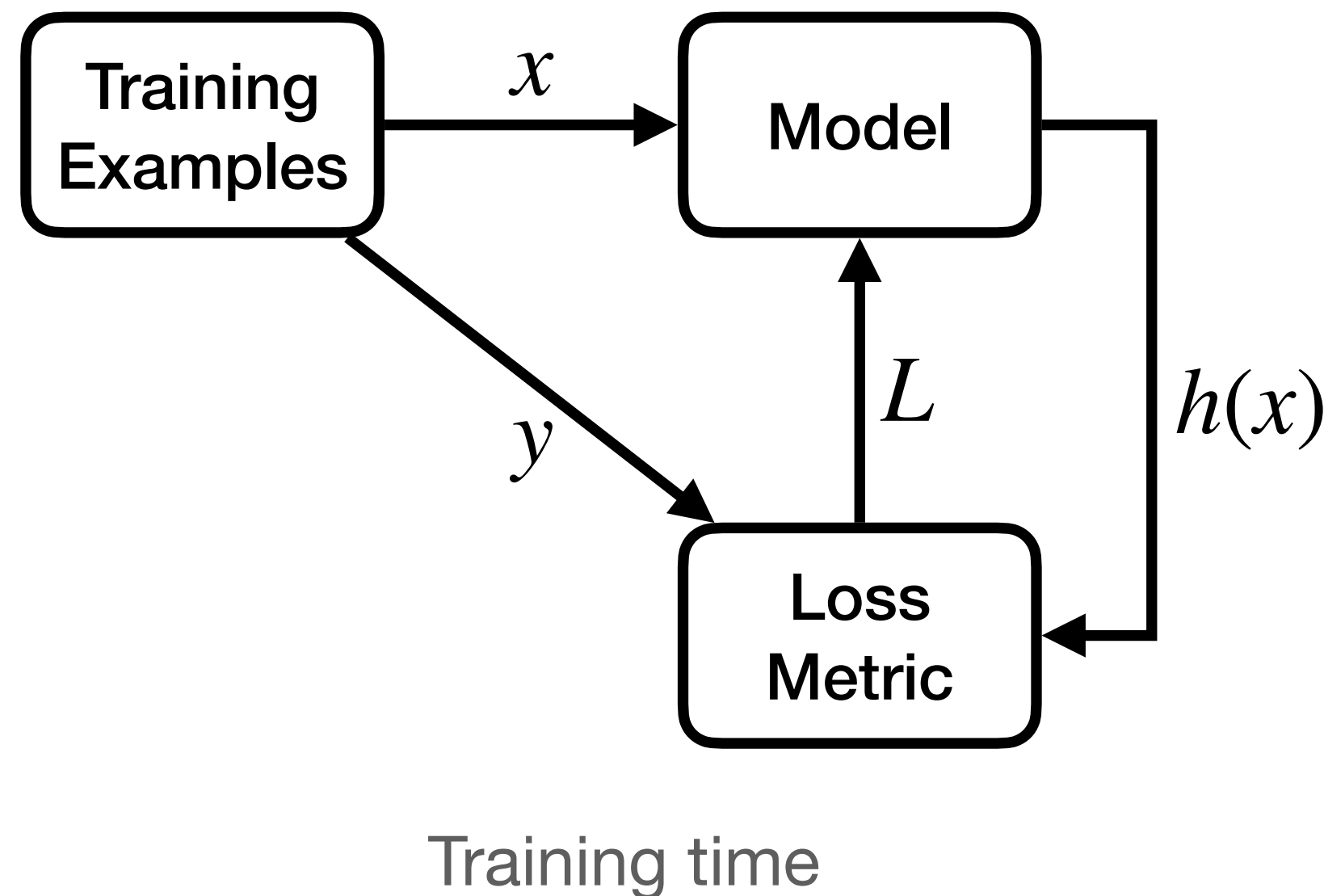
- **Assignment 3** is due **Tuesday Nov 21** at 11:59pm
  - Late submissions until the following Friday with 20% deduction
- Next week is **reading week**
  - No labs or lectures

# Recap: Deep Learning

- **Feedforward neural networks** are **extremely flexible** parametric models that can be trained by gradient descent
- **Convolutional neural networks** add **pooling** and **convolution** operations
  - Vastly more efficient to train on **vision** tasks, due to **fewer parameters** and domain-appropriate **invariances**
- **Recurrent neural networks** process elements of a **sequence one at a time**, while maintaining **state**
  - Same function with **same parameters** applied to each (element + state)
- **Transformers** process elements of a **sequence** in **parallel**
  - Each output element depends on **weighed sum** of transformed input elements, using same parameters
  - Weights are dot product of input element's **key** and output element's **query**
  - Keys and queries are computed using the **same parameters** for all elements

# Recap: Supervised Learning

Neural networks are typically used to solve **supervised learning** tasks: Selecting a **hypothesis**  $h : X \rightarrow Y$  that maps from **input** features to **target** features.



# Example: CanBot

- CanBot's job is to find and recycle empty cans
- At any given time, its battery charge is either **high** or **low**
- It can do three actions: **search** for cans, **wait**, or **recharge**
- *Goal:* Find cans efficiently without running out of battery charge

## Questions:

1. Is this an instance of a **supervised learning** problem?
2. Is this an instance of a **search** problem?

# Reinforcement Learning

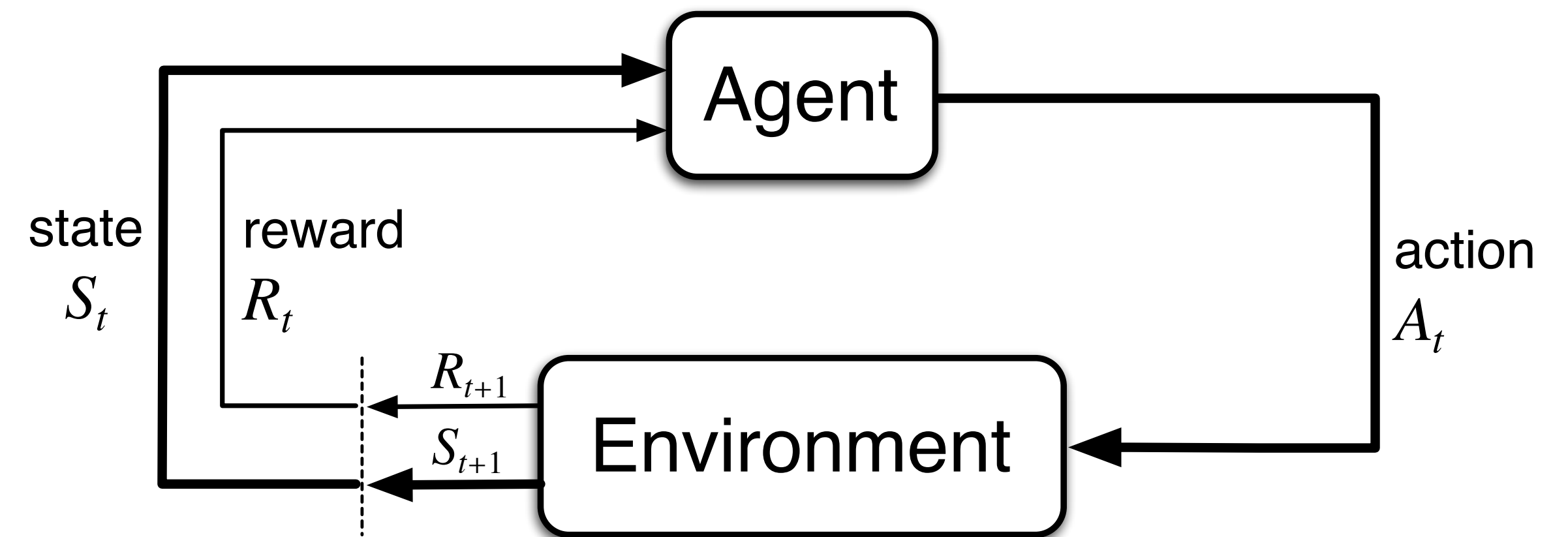
In a **reinforcement learning** task, an agent learns how to **act** based on feedback from the **environment**.

- The agent's actions may change the environment
  - Actions now can impact effects of actions in later timesteps
- The "right answer" is not known
  - Goal is to maximize total reward collected
- The task may be either **episodic** or **continuing**
- The agent makes decisions **online**: determines how to act while interacting with the environment

# Interacting with the Environment

At each time  $t = 1, 2, 3, \dots$

1. Agent receives input denoting **current state**  $S_t$
2. Agent chooses **action**  $A_t$
3. Next time step, agent receives **reward**  $R_{t+1}$  and **new state**  $S_{t+1}$ , chosen according to a distribution  $p(s', r \mid s, a)$



This interaction between agent and environment produces a **trajectory**:  
 $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$



# Markov Decision Process

## Definition:

A **Markov decision process** is a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p)$ , where

- $\mathcal{S}$  is a set of **states**,
- $\mathcal{A}$  is a set of **actions**,
- $\mathcal{R} \in \mathbb{R}$  is a set of **rewards**,
- $p(s', r \mid s, a) \in [0, 1]$  defines the **dynamics** of the process, and
- the probabilities from  $p$  **completely** characterize the environment's dynamics

# Dynamics

The four-argument dynamics function returns the probability of every **state transition**:

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

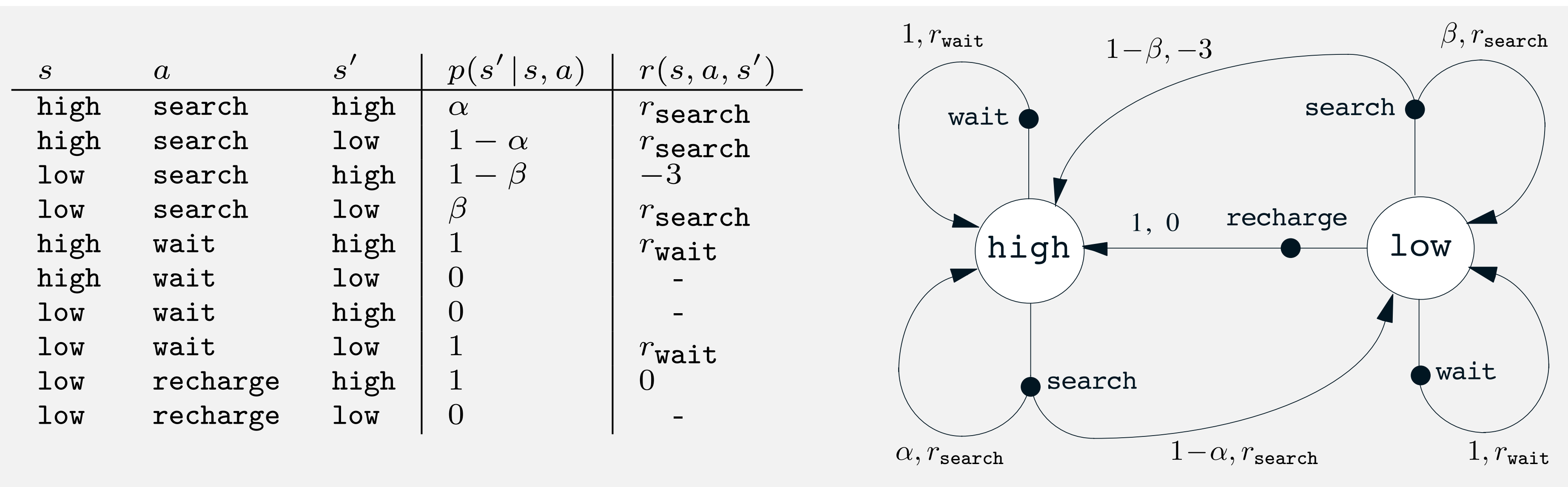
It is often convenient to use **shorthand notation** rather than the full four-argument dynamics function:

$$\begin{aligned} p(s' | s, a) &\doteq \Pr(S_t = s' | S_{t-1} = s, A_{t-1} = a) &&= \sum_{r \in \mathcal{R}} p(s', r | s, a) \\ r(s, a) &\doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] &&= \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a) \\ r(s, a, s') &\doteq \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] &&= \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)} \end{aligned}$$

# CanBot as a Reinforcement Learning Agent

**Question:** How can we represent CanBot as a **reinforcement learning** agent?

- Need to define **states**, **actions**, **rewards**, and **dynamics**



# Reward Hypothesis

**Definition:** *Reward hypothesis*

An agent's goals and purposes can be entirely represented as the maximization of the **expected value** of the **cumulative sum** of a **scalar signal**.

# Returns for Episodic Tasks

## Question:

What does "maximize the expected value of the cumulative sum of rewards" *mean*?

**Definition:** A task is **episodic** if it ends after some **finite number**  $T$  of time steps in a special **terminal state**  $S_T$ .

**Definition:** The **return**  $G_t$  after time  $t$  is the sum of rewards received after time  $t$ :  $G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$ .

**Answer:** The return  $G_t$  is a **random variable**. In an episodic task, we want to maximize its **expected value**  $\mathbb{E}[G_t]$ .

# Returns for Continuing Tasks

**Definition:** A task is **continuing** if it does not end (i.e.,  $T = \infty$ ).

- In a continuing task, we can't just maximize the sum of rewards (**why?**)
- Instead, we maximize the **discounted return**:

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \end{aligned}$$

$\gamma \leq 1$  is the **discount factor**

- Returns are **recursively** related to each other:

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

# Policies

**Question:** How should an agent in a Markov decision process choose its **actions**?

- **Markov assumption:** The state incorporates all of the necessary information about the history up until this point
  - i.e., Probabilities of future rewards & transitions are the same from state  $S_t$  **regardless of how you got there**
- So the agent can choose its actions based **only** on  $S_t$
- This is called a (memoryless) **policy**:  $\pi(a \mid s) \in [0,1]$  is the probability of taking **action**  $a$  given that the **current state** is  $s$



# State-Value Function

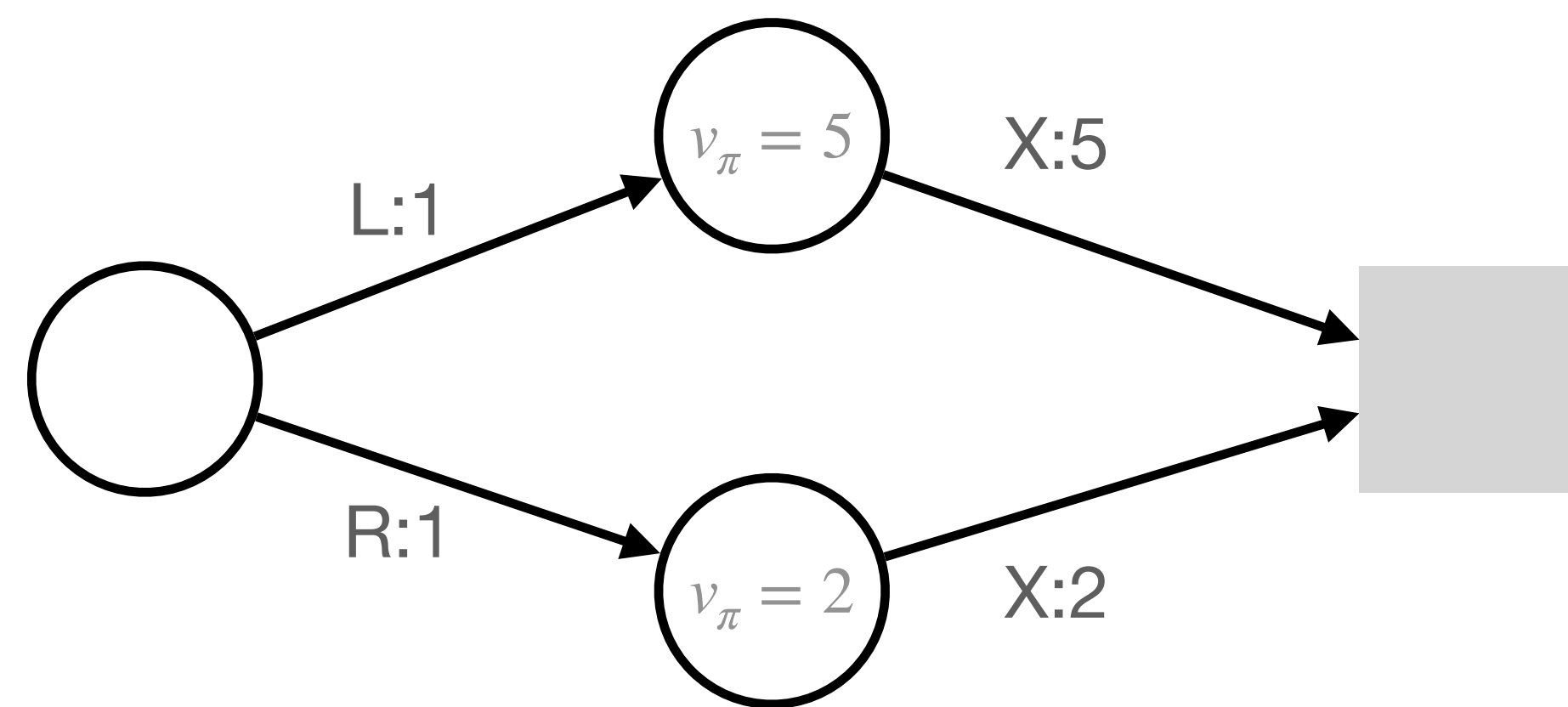
- Once you know the **policy**  $\pi$  and the **dynamics**  $p$ , you can compute the probability of every possible state transition starting from any given state
- It is often valuable to know the **expected return** starting from a given state  $s$  under a given policy  $\pi$  (**why?**)
- The **state-value function**  $v_\pi$  returns this quantity:

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t | S_t = s] \quad \forall t$$
$$= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$



# Using State-Value Function

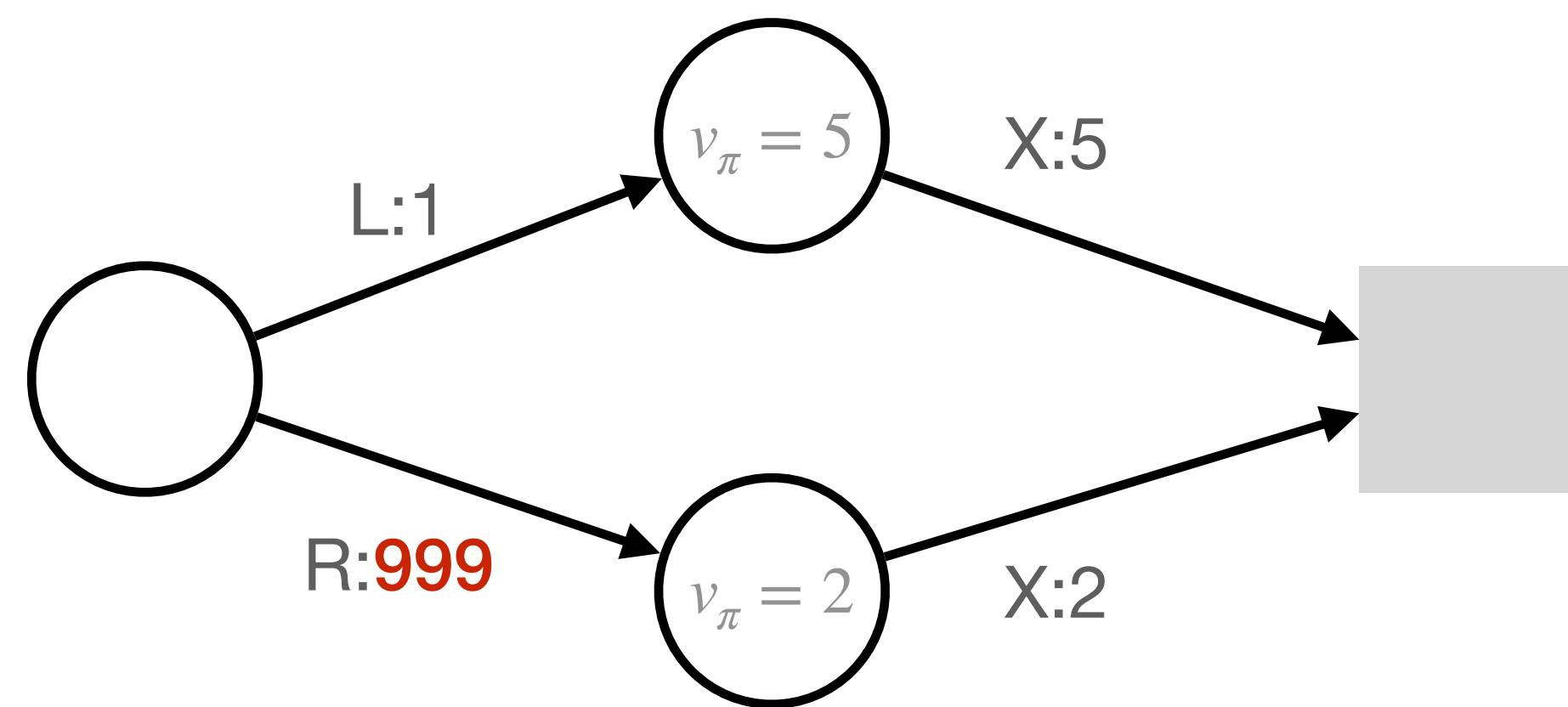
**Question:** Suppose state transitions are deterministic. Does it make sense to always choose the action that leads to the next state  $s'$  with the highest  $v_{\pi}(s)$ ?



# Using State-Value Function

**Question:** Suppose state transitions are deterministic. Does it make sense to always choose the action that leads to the next state  $s'$  with the highest  $v_{\pi}(s)$ ?

*Not always;* the reward for the **transition itself** is also important!



# Action-Value Function

The **action-value function**  $q_{\pi}(s, a)$  estimates the expected return  $G_t$  starting from state  $s$  if we

1. Take action  $a$  in state  $S_t = s$ , and then
2. Follow policy  $\pi$  for every state  $S_{t+1}$  afterward

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \end{aligned}$$

**Question:**  
How is this any different from the **state-value** function  $v_{\pi}(s)$ ?

# Bellman Equations

Value functions satisfy a **recursive consistency condition** called the **Bellman equation**:

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t | S_t = s]$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \mathbb{E}_{\pi}[R_{t+1} | S_t = s] + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_t = s]$$

$$= \sum_a \sum_{s'} \sum_r \overbrace{\Pr[S_{t+1} = s', R_{t+1} = r, A_t = a | S_t = s]} \left[ r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right]$$

$$= \sum_a \sum_{s'} \sum_r \underbrace{\Pr[S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a]} \underbrace{\Pr[A_t = a | S_t = s]} \left[ r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right]$$

$$= \sum_a \underbrace{\pi(a | s)} \sum_{s'} \sum_r \underbrace{p(s', r | s, a)} \left[ r + \gamma \underbrace{\mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s']} \right]$$

$$= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) \left[ r + \gamma v_{\pi}(s') \right] \quad \blacksquare$$

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

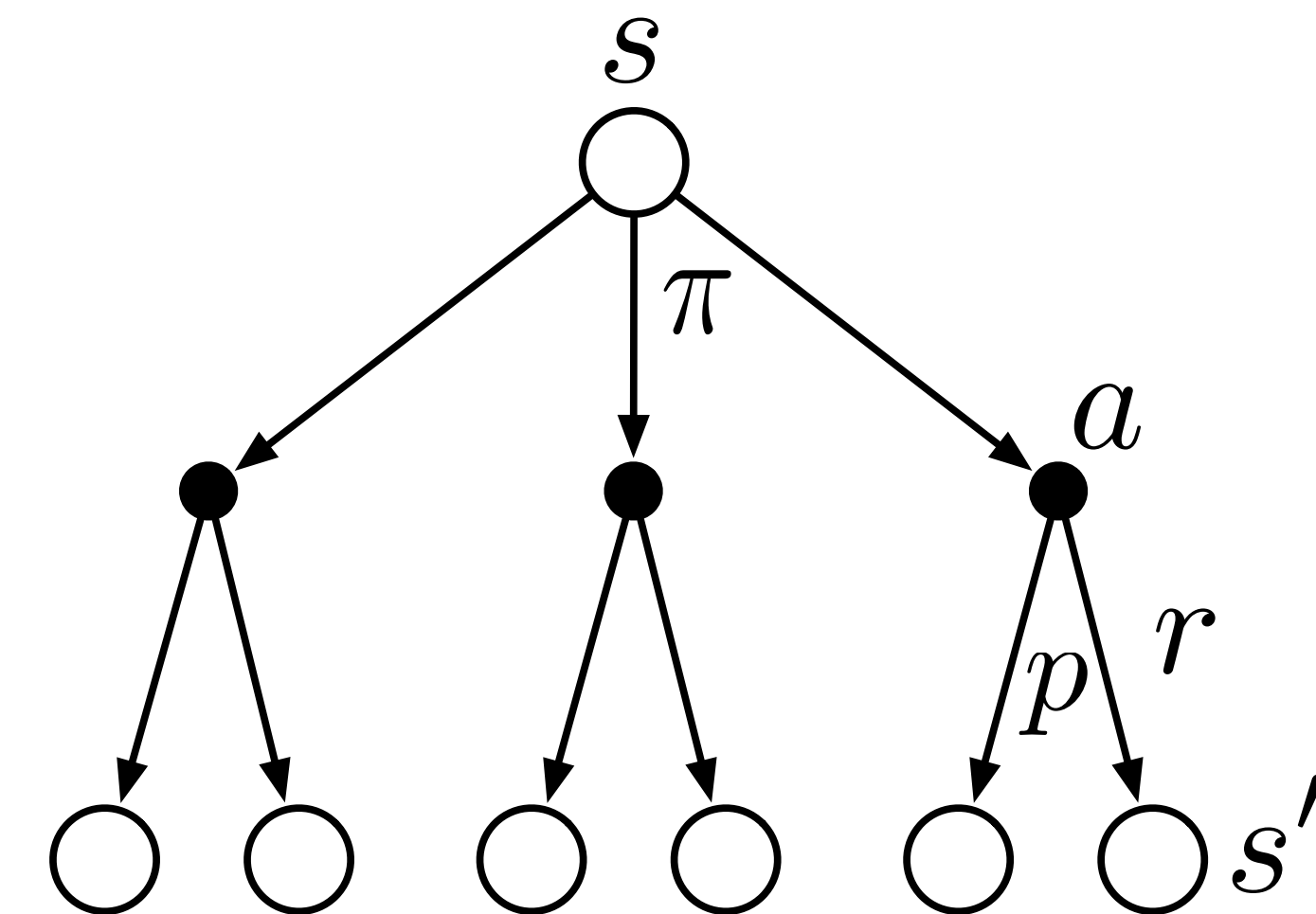
$$\mathbb{E}[A + cB] = \mathbb{E}[A] + c\mathbb{E}[B]$$

- $v_{\pi}$  is the **unique solution** to  $\pi$ 's Bellman equation
- There is also a Bellman equation for  $\pi$ 's **action-value function**

# Backup Diagrams

Backup diagrams help to visualize the flow of **information back to a state** from its successor states or action-state pairs:

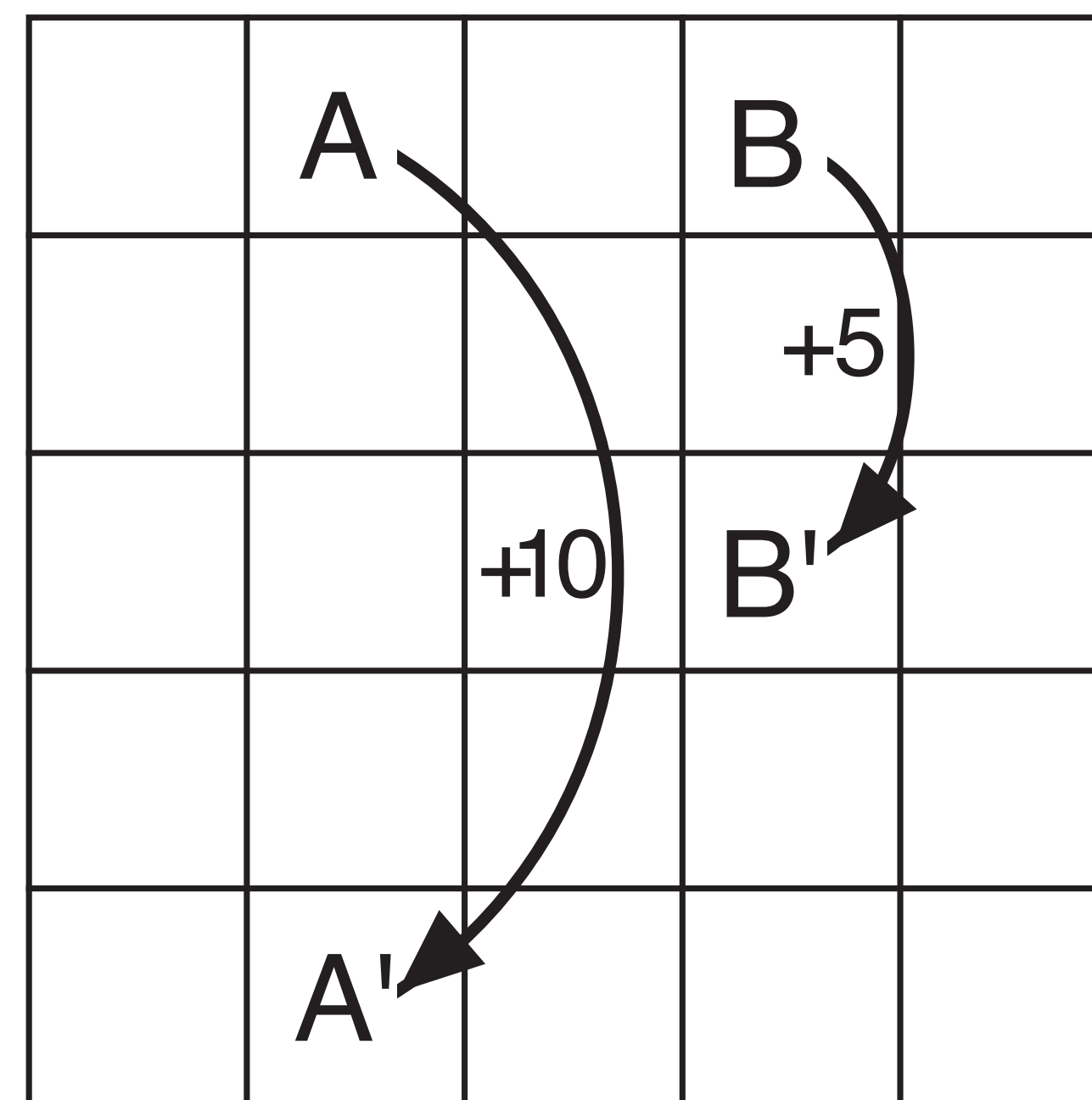
$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t | S_t = s] \\ &= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \end{aligned}$$



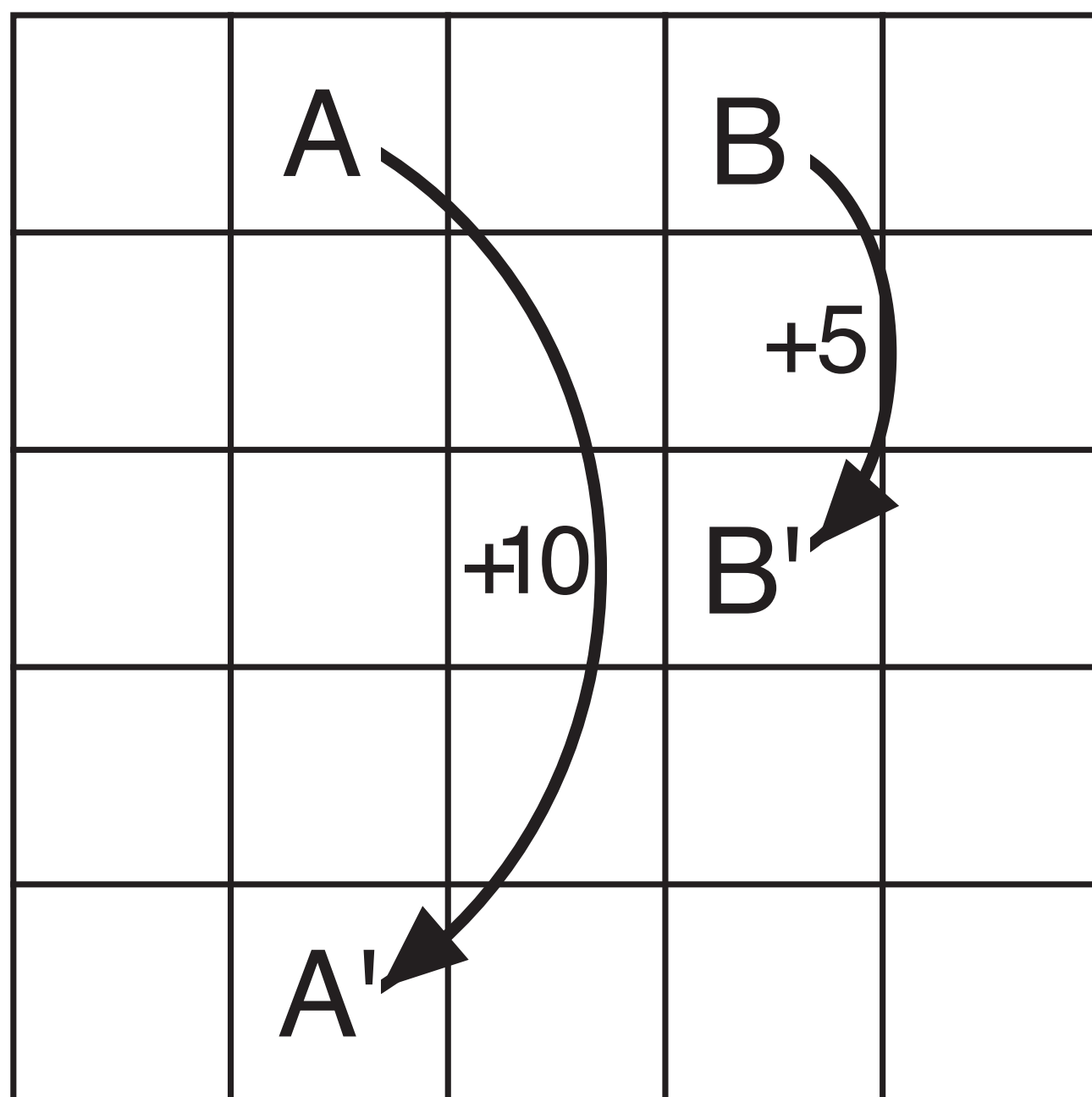
Backup diagram for  $v_{\pi}$

# GridWorld

- At each cell, can go north, south, east, west
- Try to go off the **edge**: reward of **-1**
- Leaving state **A**: takes you to state **A'**, reward of **+10**
- Leaving state **B**: takes you to state **B'**, reward of **+5**



# GridWorld



Reward dynamics

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

State-value function  $v_\pi$  for random policy

$$\pi(a | s) = 0.25$$

# Summary

- **Supervised learning models** are trained **offline** using **labelled training examples**, and then make **predictions**
- **Reinforcement learning agents** choose their **actions online**, and update their behaviour based on **rewards** from the **environment**
- We can formally represent reinforcement learning environments using **Markov decision processes**, for both **episodic** and **continuing** tasks
- Reinforcement learning agents maximize **expected returns**
- **Policies** map **states** to (distribution over) **actions**
- Given a **policy**  $\pi$ , every state  $s$  has an **expected value**  $v_{\pi}(s)$
- State-value and action-value functions satisfy the **Bellman equations**