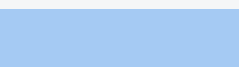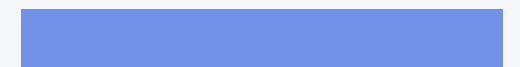# Lecture-2 Markov Decision Process
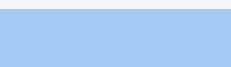
# Agenda

# 01.

# Markov Process

# Markov Property

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

A state $S_t$ is *Markov* if and only if

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

- The state captures all relevant information from the history

- Once the state is known, the history may be thrown away

- The state is a sufficient statistic of the future

# Markov Chain

- A Markov chain is a memoryless random process, i.e. a sequence of random states S1, S2,... with the Markov property.

## Definition

A *Markov Process* (or *Markov Chain*) is a tuple $\langle \mathcal{S}, \mathcal{P} \rangle$

- $\mathcal{S}$ is a (finite) set of states

- $\mathcal{P}$ is a state transition probability matrix,
  $$\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$$

# State Transition Matrix

- State Transition Probability

$$\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$$

- State Transition Matrix

$$\mathcal{P} = \text{from} \begin{bmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{bmatrix} \quad to$$

# State Transition Matrix



$$
\mathcal{P} = \begin{array}{c} \\ C1 \\ C2 \\ C3 \\ Pass \\ Pub \\ FB \\ Sleep \end{array}
\begin{array}{ccccccc}
C1 & C2 & C3 & Pass & Pub & FB & Sleep \\
 & 0.5 & & & & 0.5 & \\
 & & 0.8 & & & & 0.2 \\
 & & & 0.6 & 0.4 & & \\
 & & & & & & 1.0 \\
0.2 & 0.4 & 0.4 & & & & \\
0.1 & & & & & 0.9 & \\
 & & & & & & 1 \\
\end{array}
$$

# Episode Sampling



- C1 C2 C3 Pass Sleep

- C1 FB FB C1 C2 Sleep

- C1 C2 C3 Pub C2 C3 Pass Sleep

- C1 FB FB C1 C2 C3 Pub C1 FB FB FB C1 C2 C3 Pub C2 Sleep

# Markov Reward Process

- A Markov reward process is a Markov chain with values.

## Definition

A *Markov Reward Process* is a tuple $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $\mathcal{S}$ is a finite set of states

- $\mathcal{P}$ is a state transition probability matrix,
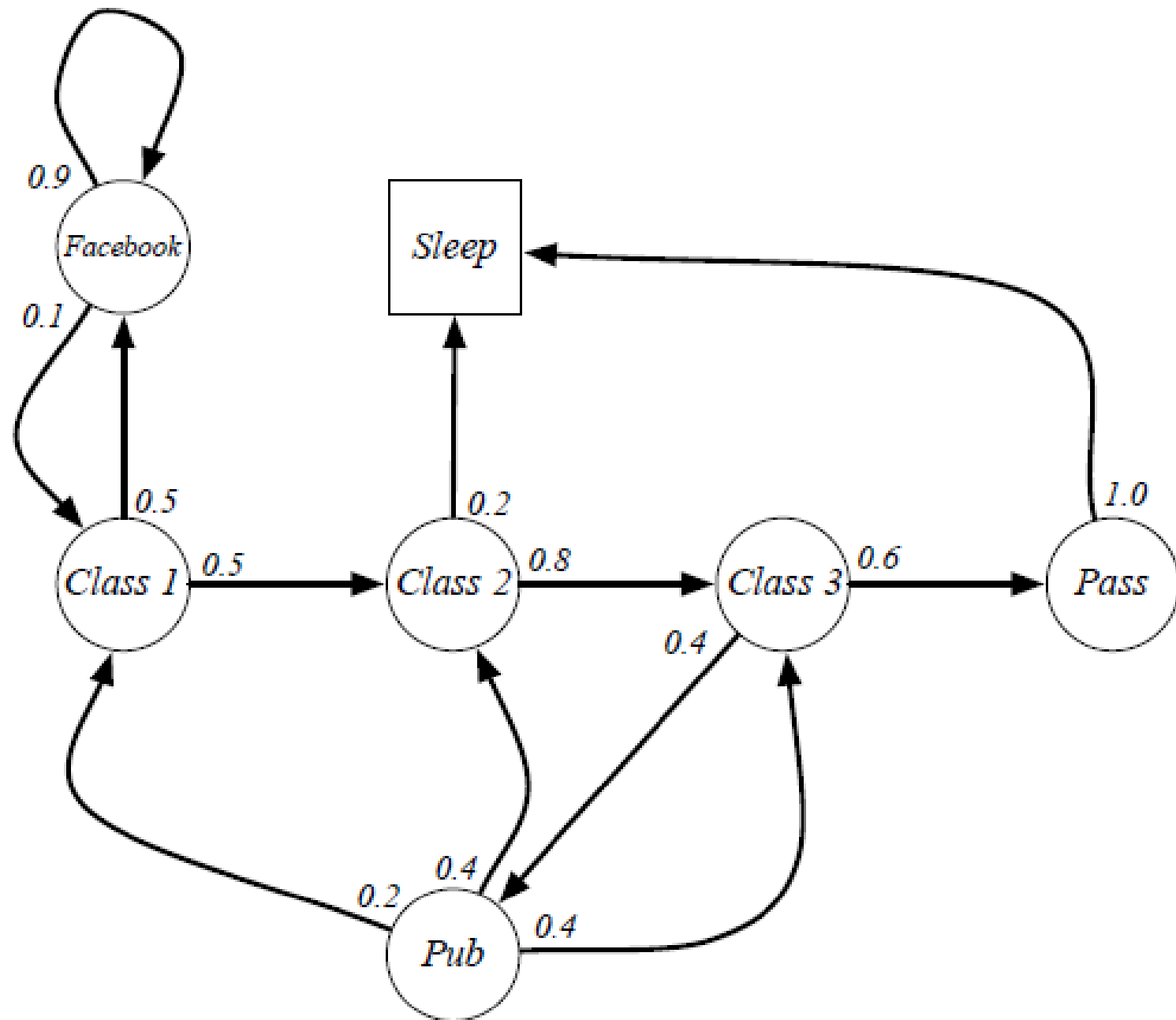$\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$

- $\mathcal{R}$ is a reward function, $\mathcal{R}_s = \mathbb{E}\left[R_{t+1} \mid S_t = s\right]$

- $\gamma$ is a discount factor, $\gamma \in [0, 1]$

# Markov Reward Process

# Markov Decision Process

- A Markov decision process (MDP) is a Markov reward process with decisions.

- It is an environment in which all states are Markov.

## Definition

A *Markov Decision Process* is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
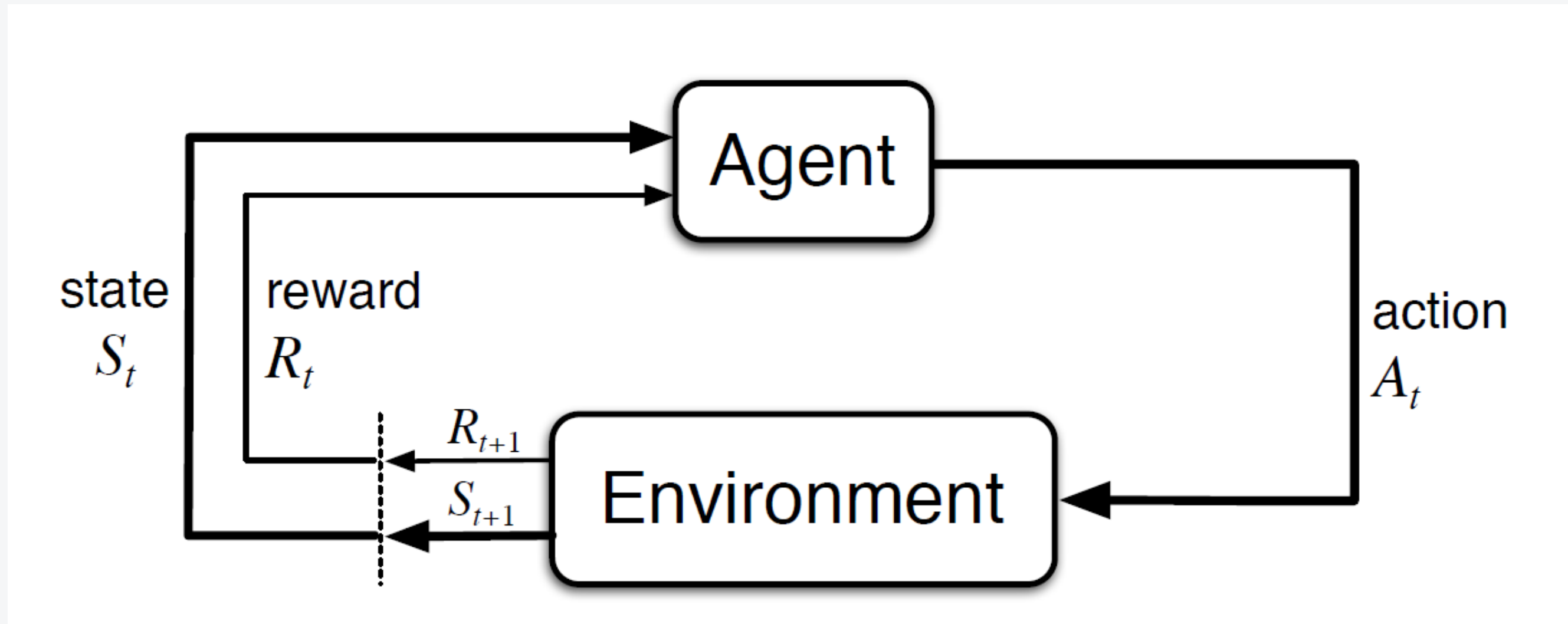
- $\mathcal{S}$ is a finite set of states

- $\mathcal{A}$ is a finite set of actions

- $\mathcal{P}$ is a state transition probability matrix,
  $\mathcal{P}_{ss'}^{a} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s, A_t = a\right]$

- $\mathcal{R}$ is a reward function, $\mathcal{R}_s^a = \mathbb{E}\left[R_{t+1} \mid S_t = s, A_t = a\right]$

- $\gamma$ is a discount factor $\gamma \in [0, 1]$.

# Markov Decision Process

$$p(s',r\,|\,s,a) \;\doteq\; \Pr\{S_t\!=\!s', R_t\!=\!r \mid S_{t-1}\!=\!s, A_{t-1}\!=\!a\},$$
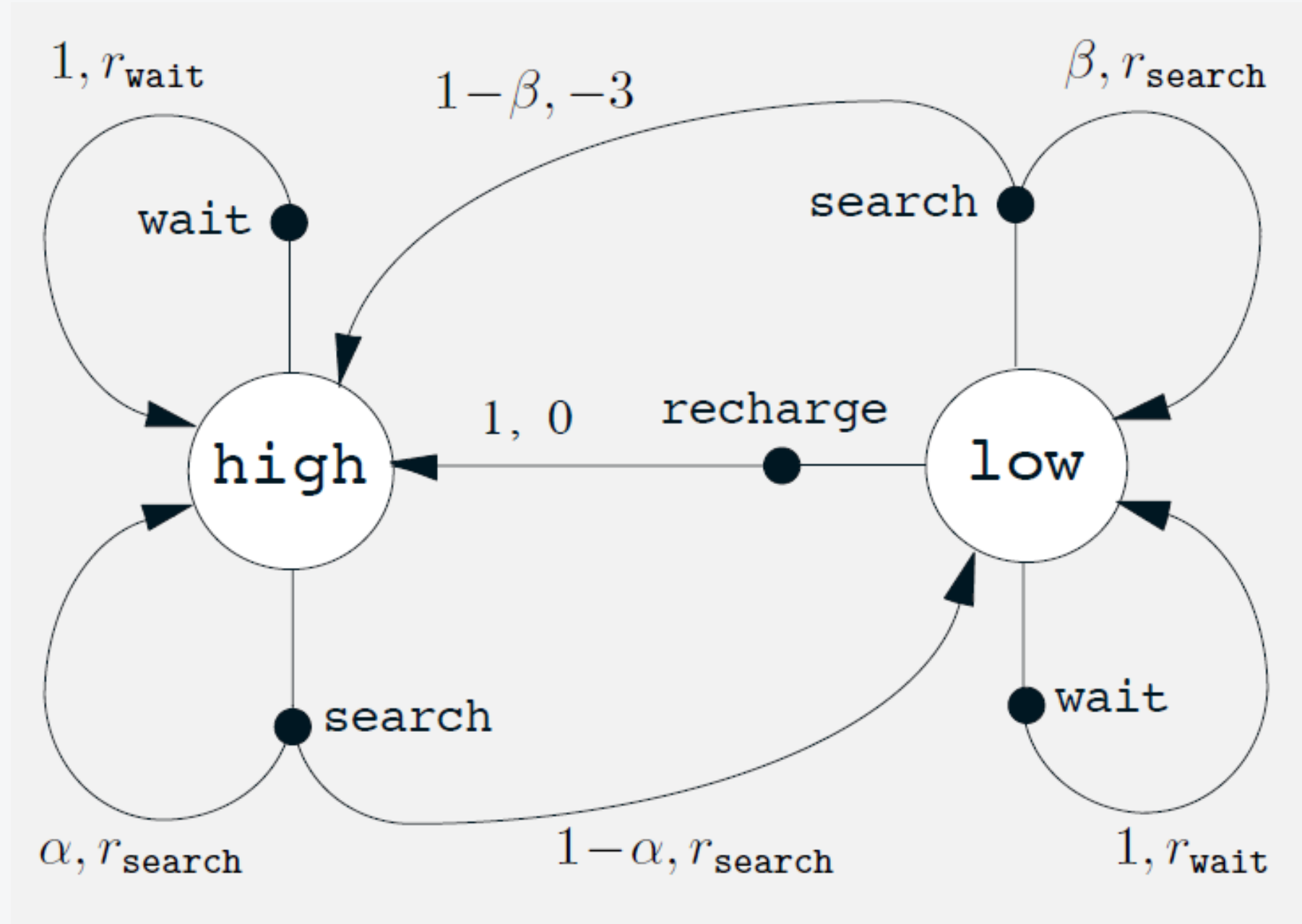
$$\sum_{s'\in\mathcal{S}}\sum_{r\in\mathcal{R}} p(s',r\,|\,s,a) = 1, \ \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}(s).$$

# Markov Decision Process



$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \ldots$$

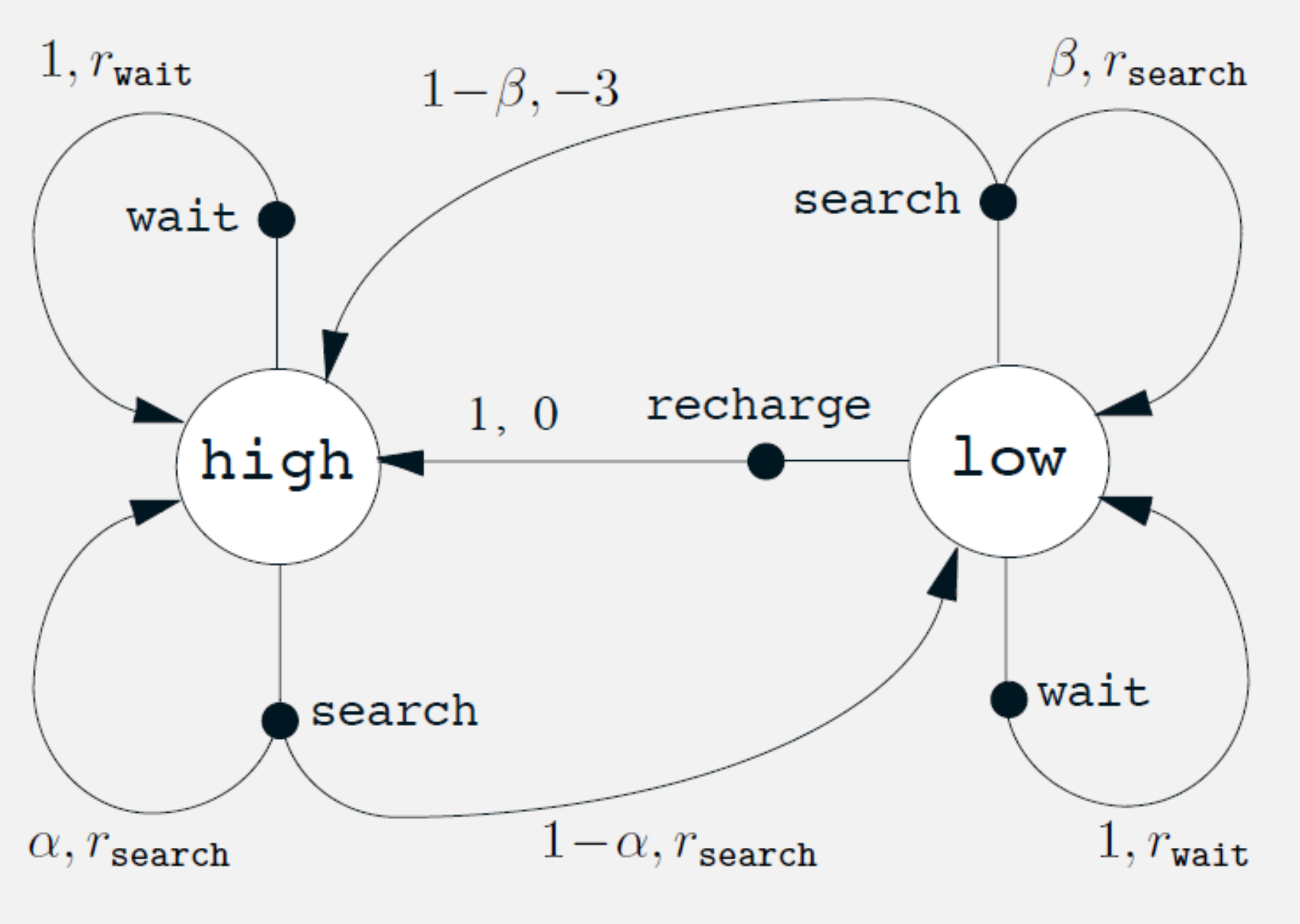# Markov Decision Process



$$\mathcal{S} = \{\text{high}, \text{low}\}.$$

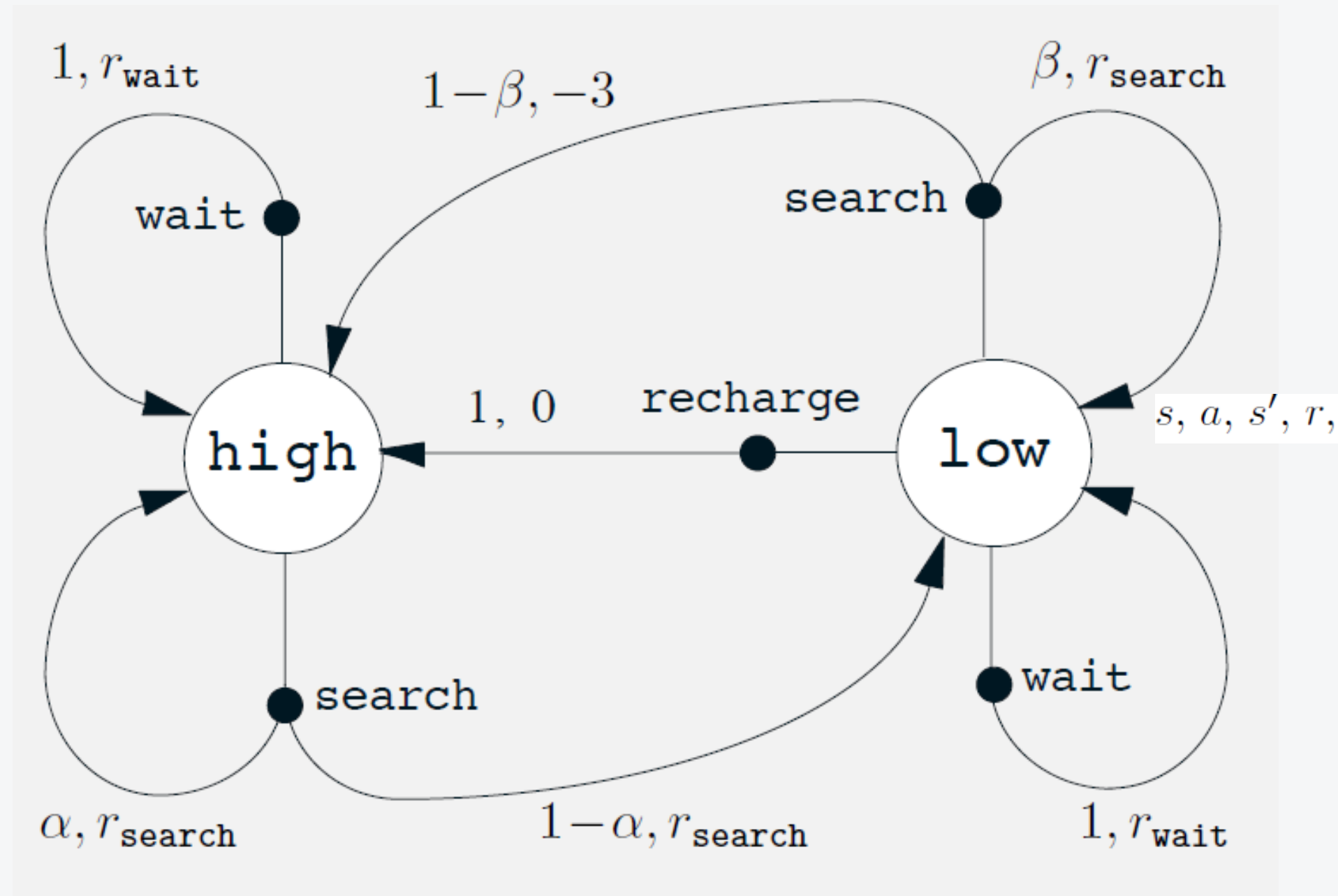$$\mathcal{A}(\text{high}) = \{\text{search}, \text{wait}\}$$

$$\mathcal{A}(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}.$$

# Markov Decision Process



| $s$ | $a$ | $s'$ | $p(s'\,|\,s,a)$ | $r(s,a,s')$ |
|------|--------|------|-----------------|-------------|
| high | search | high | $\alpha$ | $r_{\text{search}}$ |
| high | search | low | $1 - \alpha$ | $r_{\text{search}}$ |
| low | search | high | $1 - \beta$ | $-3$ |
| low | search | low | $\beta$ | $r_{\text{search}}$ |
| high | wait | high | $1$ | $r_{\text{wait}}$ |
| high | wait | low | $0$ | - |
| low | wait | high | $0$ | - |
| low | wait | low | $1$ | $r_{\text{wait}}$ |
| low | recharge | high | $1$ | $0$ |
| low | recharge | low | $0$ | - |

# Task



Give a table analogous to previous table, but for $p(s', r \mid s, a)$

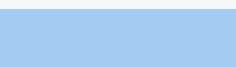It should have columns

$$s, \; a, \; s', \; r,$$
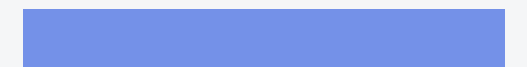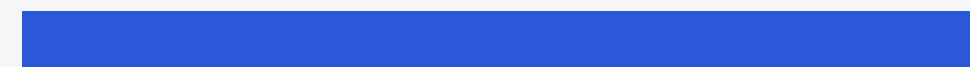
$$p(s', r \mid s, a)$$

and a row for every 4-tuple for
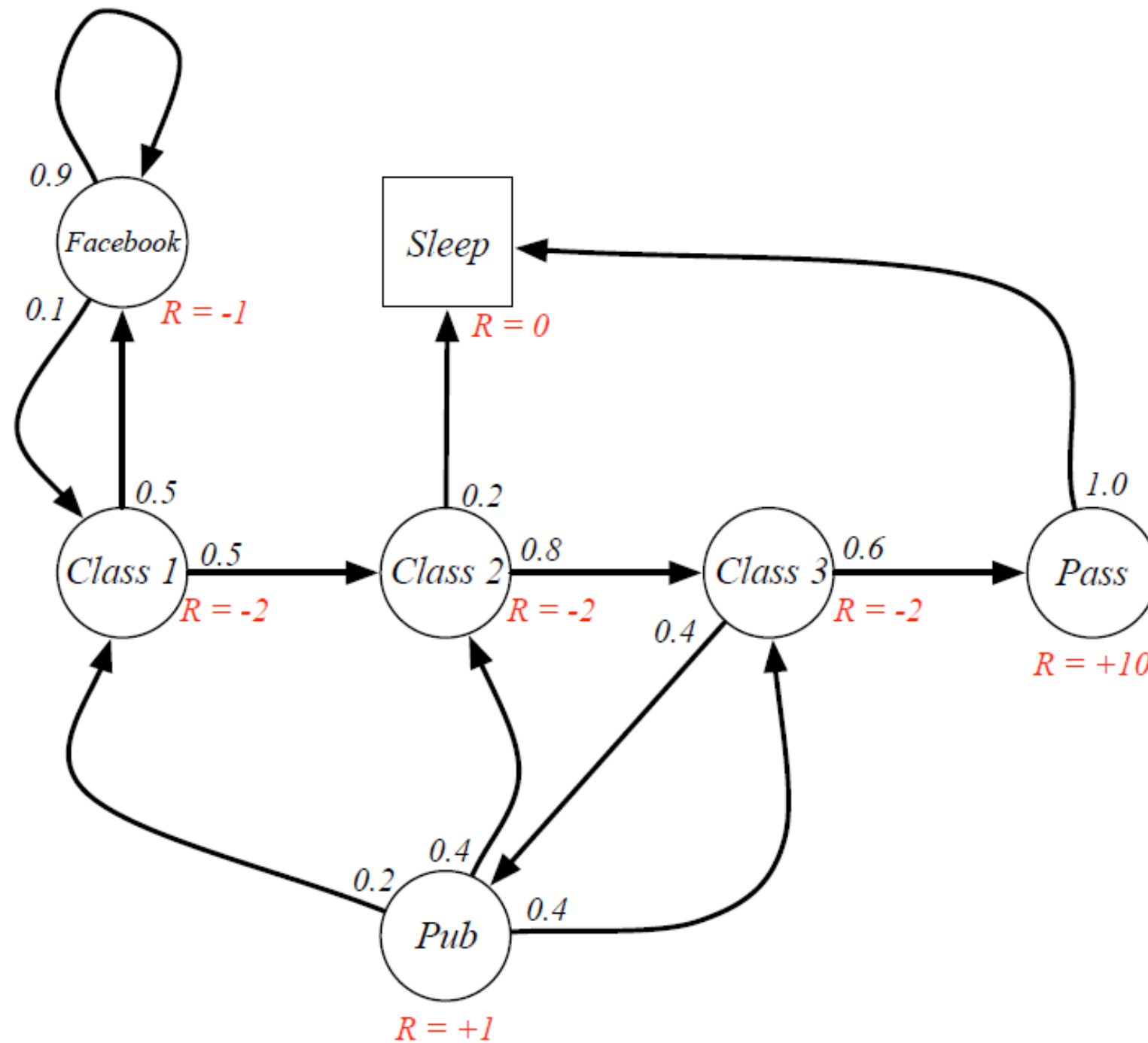
which $p(s', r \mid s, a) > 0$

# 02.

# Reward and Return

# **Return**

- The agent's goal is to maximize the cumulative reward in the long run.

- Cumulative reward is called Return.

- Return, denoted Gt, is defined as some specific function of the reward sequence.
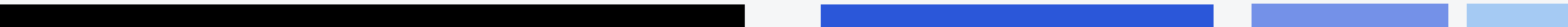
# Episodic Tasks



- C1 C2 C3 Pass Sleep

- C1 FB FB C1 C2 Sleep

- C1 C2 C3 Pub C2 C3 Pass Sleep

- C1 FB FB C1 C2 C3 Pub C1 FB FB FB C1 C2 C3 Pub C2 Sleep

# Episodic Tasks

- Episodic tasks: interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze.

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T,$$

- where T is a final time step at which a terminal state is reached, ending an episode.

# Continuing Tasks

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

$$\begin{aligned}
G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \\
&= R_{t+1} + \gamma\left(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \cdots\right) \\
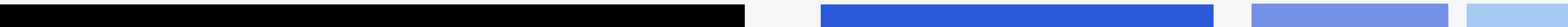&= R_{t+1} + \gamma G_{t+1}
\end{aligned}$$

# Discount Factor

- The *discount* $\gamma \in [0, 1]$ is the present value of future rewards
- The value of receiving reward $R$ after $k + 1$ time-steps is $\gamma^k R$.
- This values immediate reward above delayed reward.
    - $\gamma$ close to 0 leads to "myopic" evaluation
    - $\gamma$ close to 1 leads to "far-sighted" evaluation

# **Task**

*In a Markov decision process, a large discount factor γ means that short term rewards are much more influential than long term rewards.*
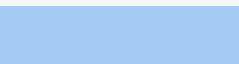
- True

- False

# Task

*Exercise 3.8* Suppose $\gamma = 0.5$ and the following sequence of rewards is received $R_1 = -1$, $R_2 = 2$, $R_3 = 6$, $R_4 = 3$, and $R_5 = 2$, with $T = 5$. What are $G_0$, $G_1$, ..., $G_5$? Hint: Work backwards. □

**Task**

*Exercise 3.9* Suppose $\gamma = 0.9$ and the reward sequence is $R_1 = 2$ followed by an infinite sequence of 7s. What are $G_1$ and $G_0$? □

# 03.

# Value Function

# Value Function

- Functions of states (or of state–action pairs) that estimate how good it is for the agent to be in a given state (or how good it is to perform a given action in a given state).

- The notion of "how good" is defined in terms of expected return.

- Value functions are defined with respect to particular policies.
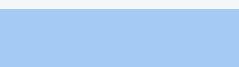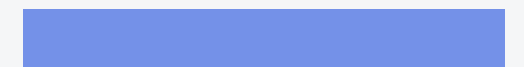
# Value Function

## State Value

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \,\middle|\, S_t = s\right], \text{ for all } s \in \mathcal{S},$$

## Action Value

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \,\middle|\, S_t = s, A_t = a\right].$$
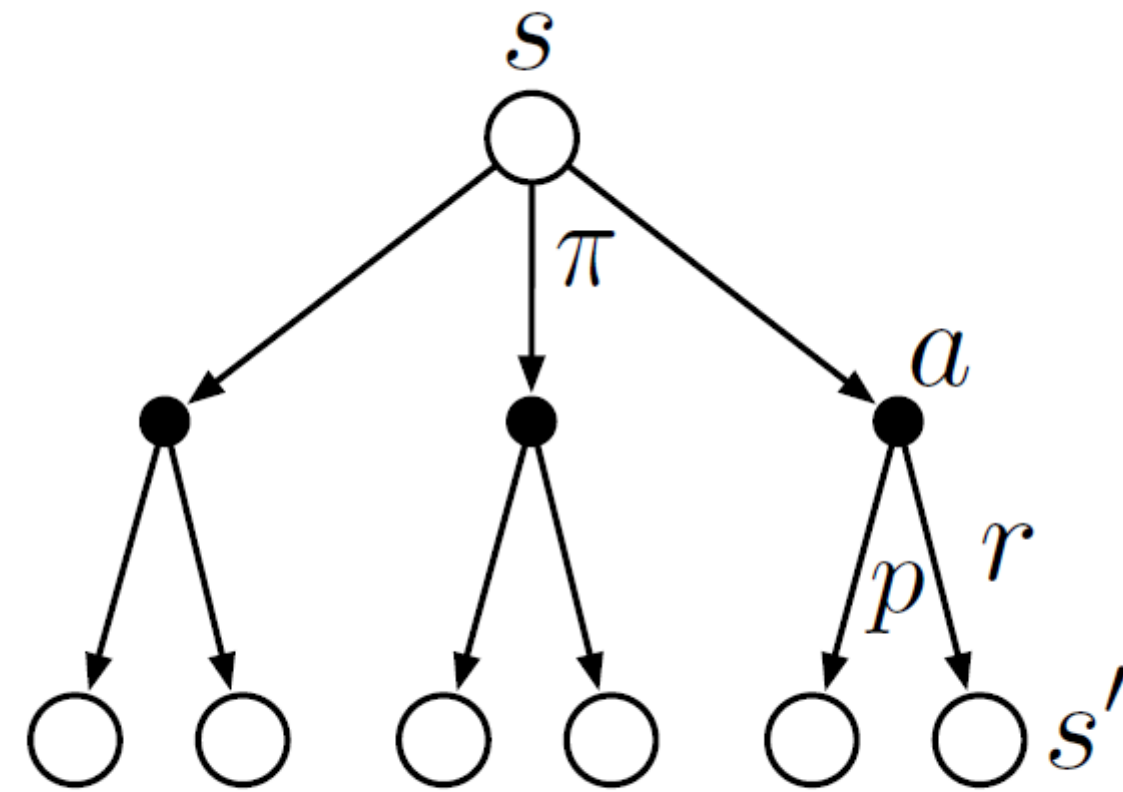
04.

# Bellman Equation

# Bellman Equation

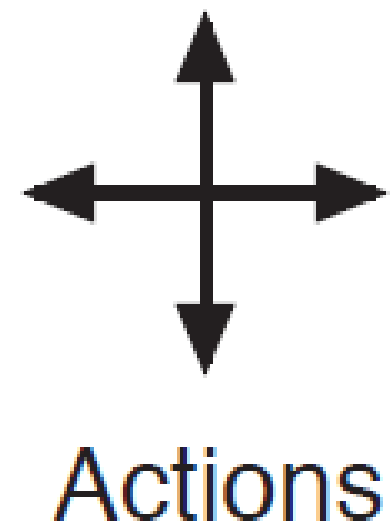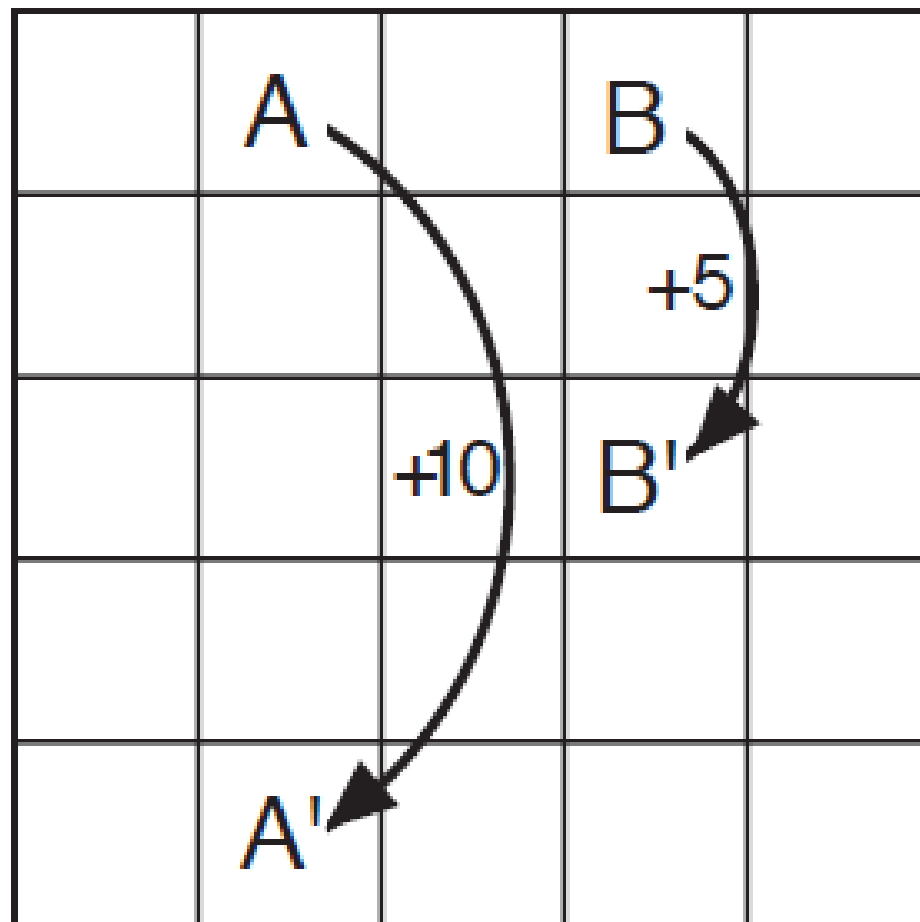$$
\begin{aligned}
v_\pi(s) &\doteq \mathbb{E}_\pi[G_t \mid S_t = s] \\
&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a)\left[r + \gamma \mathbb{E}_\pi[G_{t+1}|S_{t+1} = s']\right] \\
&= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)\left[r + \gamma v_\pi(s')\right], \quad \text{for all } s \in \mathcal{S},
\end{aligned}
$$

# Backup Diagram
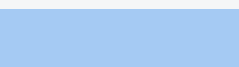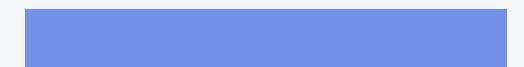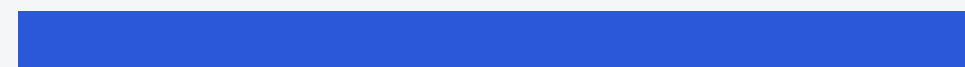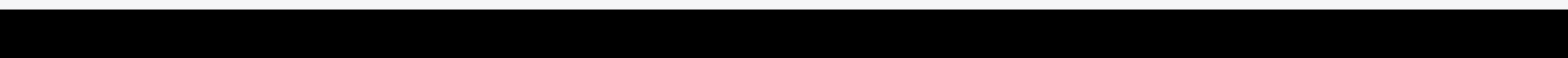


Backup diagram for $v_\pi$

# Task



- A gridworld representation of a simple finite MDP.

- The cells of the grid correspond to the states of the environment.

- At each cell, four actions are possible: north, south, east, and west.

# Task

- Actions that would take the agent off the grid leave its location unchanged but also result in a reward of –1.

- Other actions result in a reward of 0, except those that move the agent out of the special states A and B.

- From state A, all four actions yield a reward of +10 and take the agent to A'. From state B, all actions yield a reward of +5 and take the agent to B'.

# Task

*Exercise 3.14* The Bellman equation (3.14) must hold for each state for the value function $v_\pi$ shown in Figure 3.2 (right) of Example 3.5. Show numerically that this equation holds for the center state, valued at +0.7, with respect to its four neighboring states, valued at +2.3, +0.4, −0.4, and +0.7. (These numbers are accurate only to one decimal place.) □

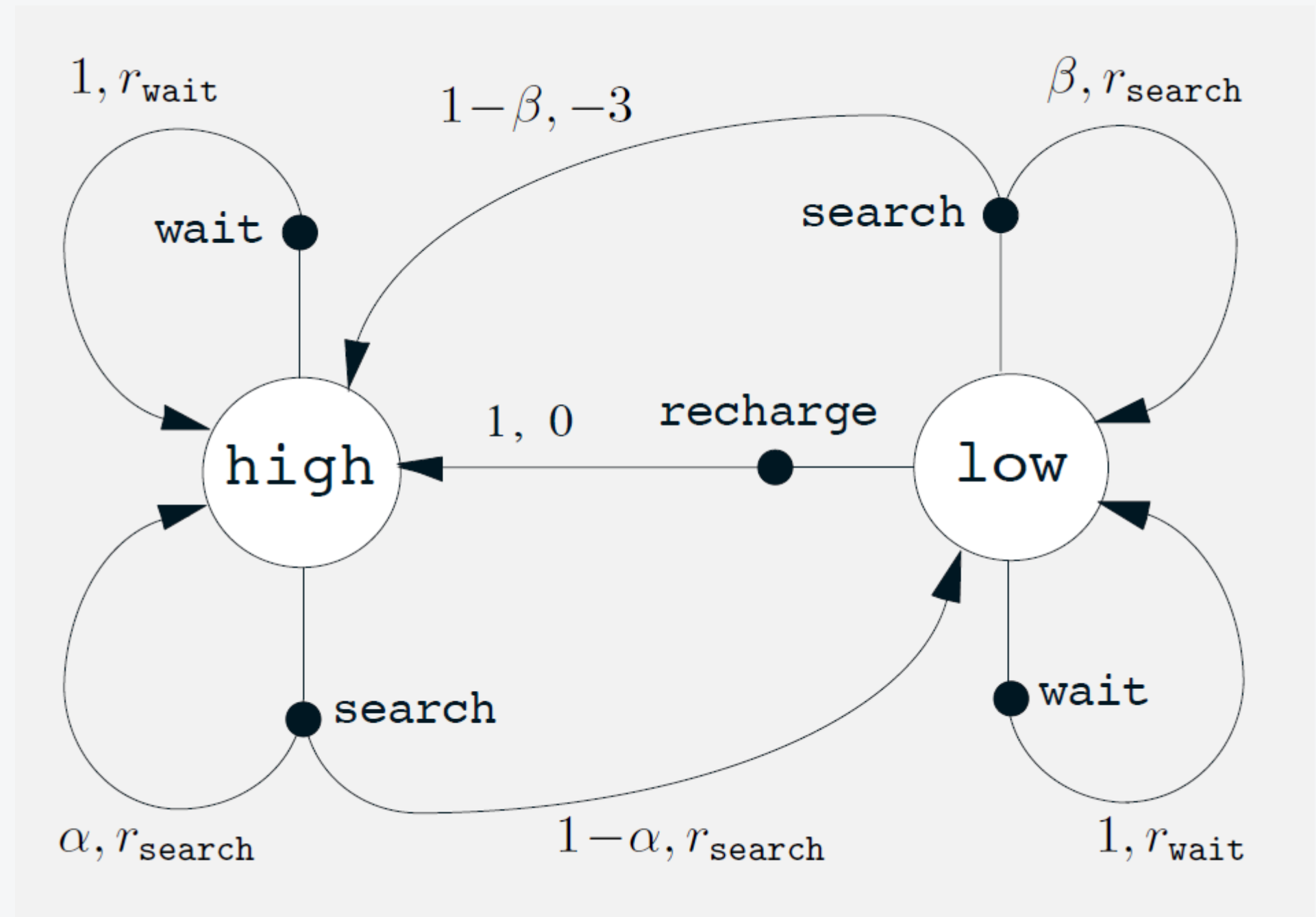| 3.3 | 8.8 | 4.4 | 5.3 | 1.5 |
|-----|-----|-----|-----|-----|
| 1.5 | 3.0 | 2.3 | 1.9 | 0.5 |
| 0.1 | 0.7 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.4 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

# Optimal Policy and Bellman Equation

$$v_*(s) \doteq \max_\pi v_\pi(s),$$

$$
\begin{aligned}
v_*(s) &= \max_{a \in \mathcal{A}(s)} q_{\pi_*}(s, a) \\
&= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \\
&= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
&= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\
&= \max_a \sum_{s', r} p(s', r \mid s, a)\big[r + \gamma v_*(s')\big].
\end{aligned}
$$

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a].$$

# Task

Give the Bellman optimality equations for this recycling robot.

# Thank You!