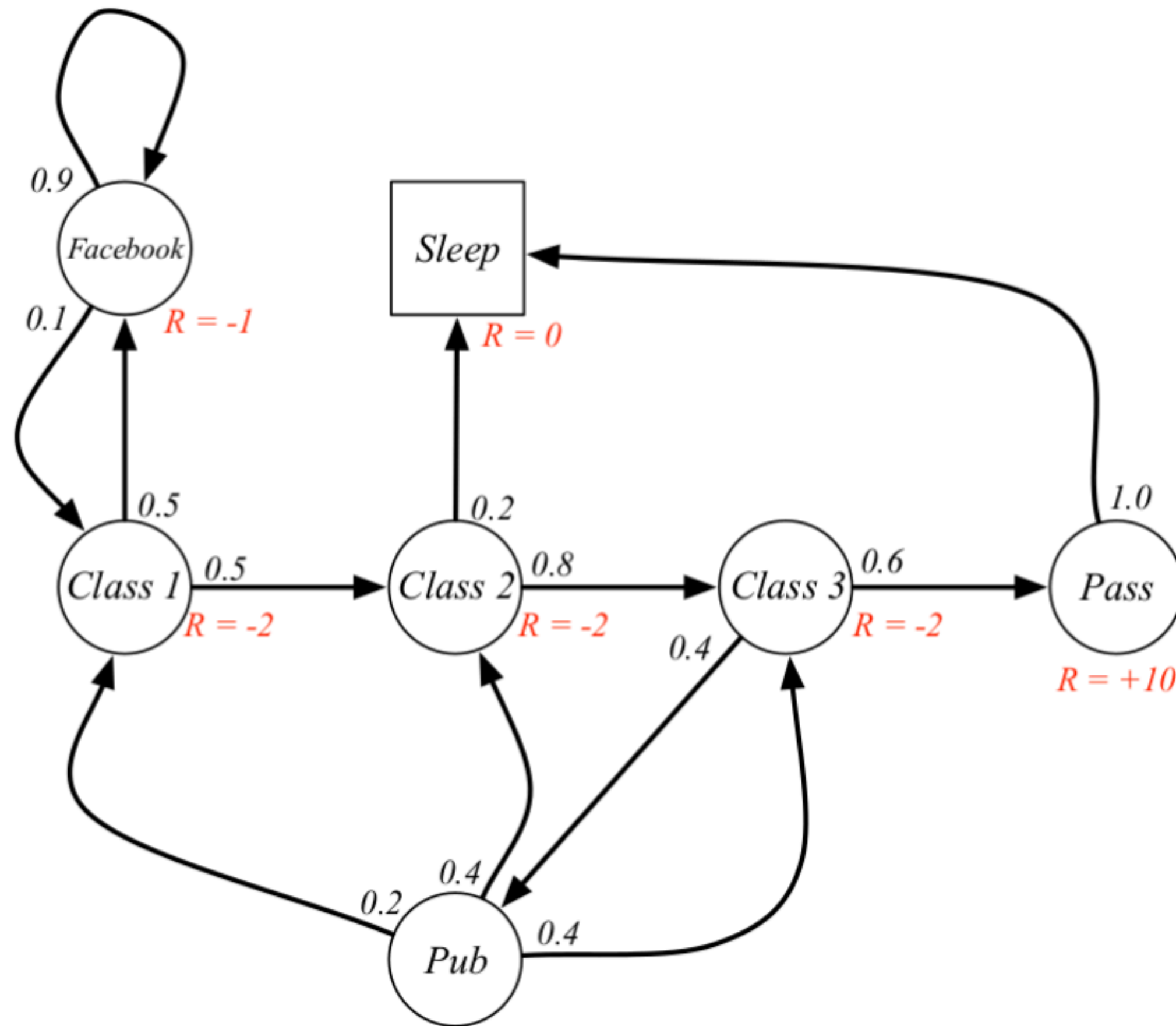# What is the best way to wander?

# Readings for today
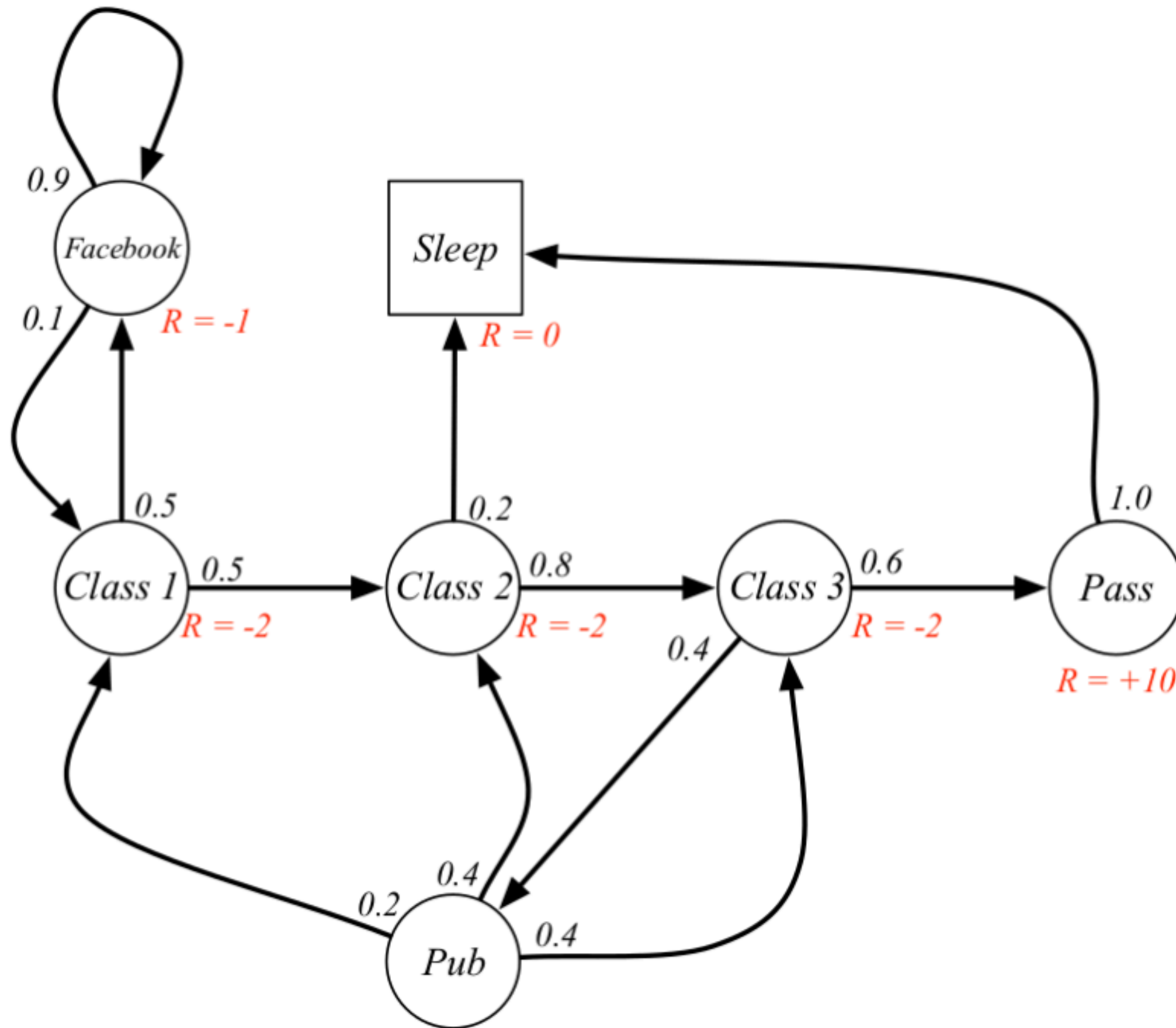
- Ashraf M. (2021). Reinforcement Learning Demystified: Markov Decision Processes (Part 1). Become Sentient.

What is the best way to strategically shift from one state to another?
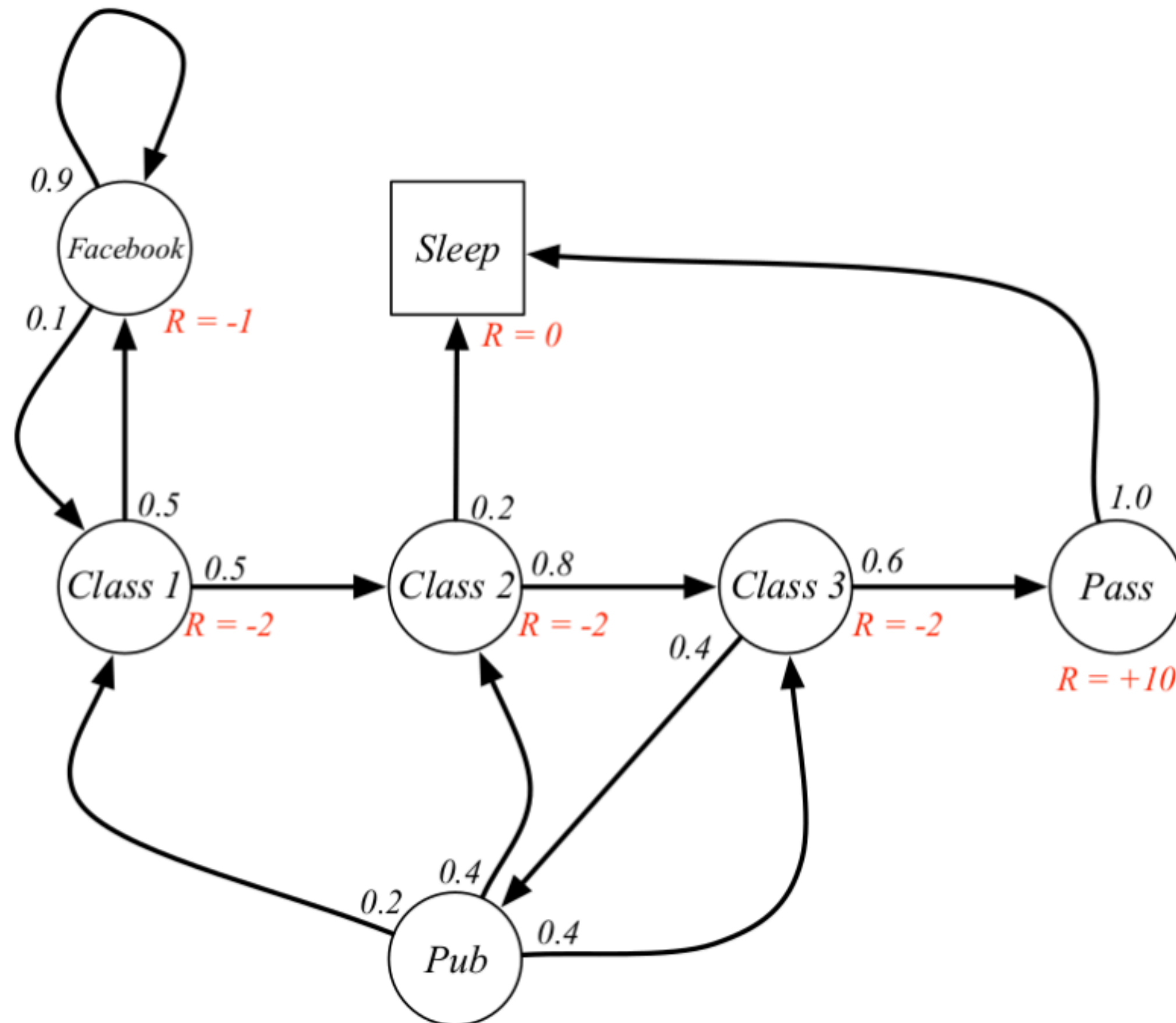
*"The future is independent of the past given the present."*

State $\mathbf{S}_t$ has the Markov property if and only if
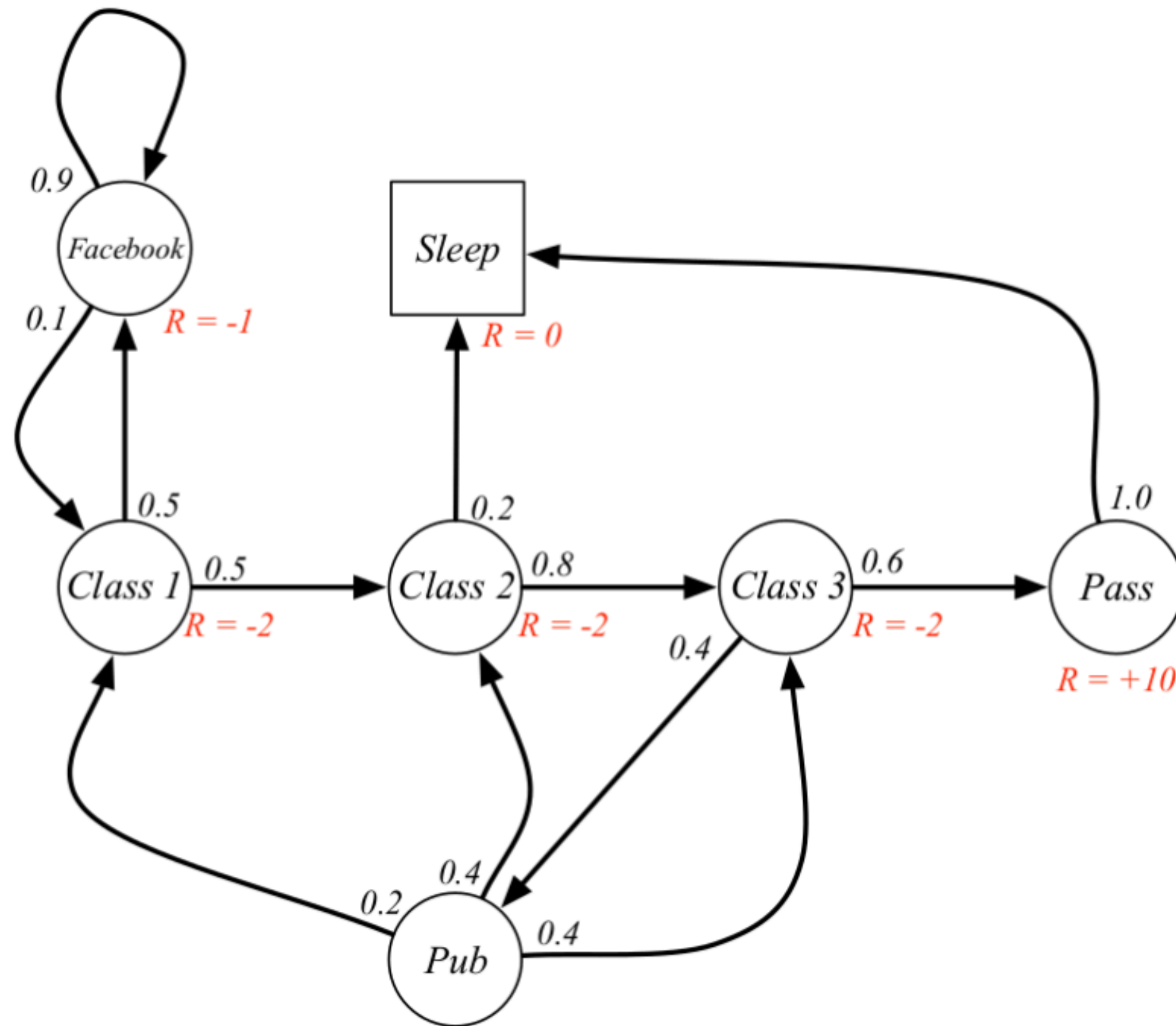
$$P(S_{t+1} \mid S_t) = P(S_{t+1} \mid S_1, \ldots, S_t)$$

**State transition**

$$\mathcal{P} = \text{from} \begin{bmatrix} \mathcal{P}_{11} & \ldots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \ldots & \mathcal{P}_{nn} \end{bmatrix} \overset{to}{}$$
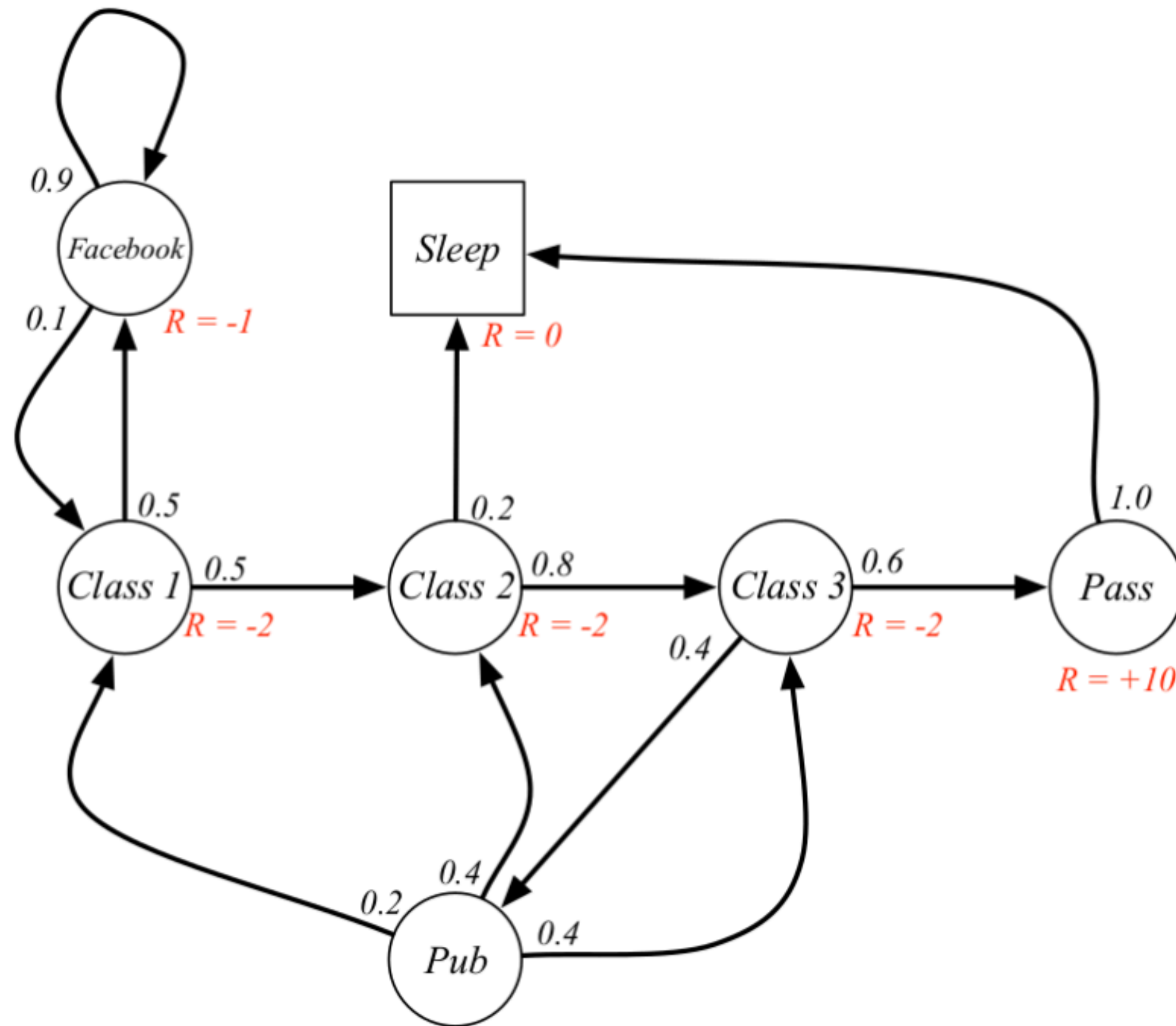
A *Markov process* is a memory-less random process, i.e. a sequence of random states $S_1, S_2, \ldots$ with the Markov property

$$(\mathbf{S}, \mathbf{P})$$

state $\quad$ transition function

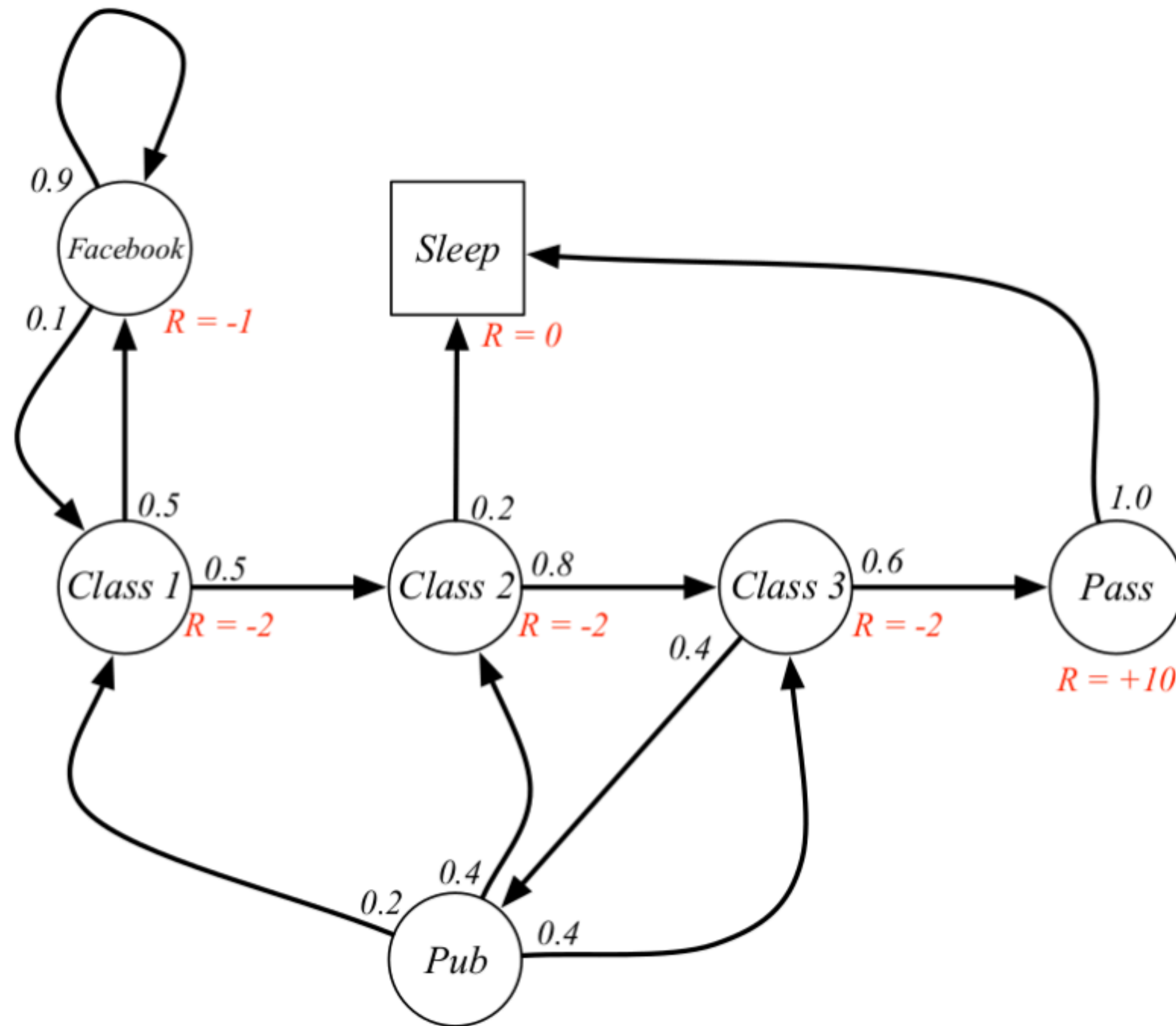A *Markov reward* process is a Markov process with a value judgement.

$$R_S = \mathbb{E}(R_{t+1}, S_t = S)$$

reward function — discount factor

$$(\mathbf{S}, \mathbf{P}, \dot{\mathbf{R}}, \gamma)$$

state — transition function

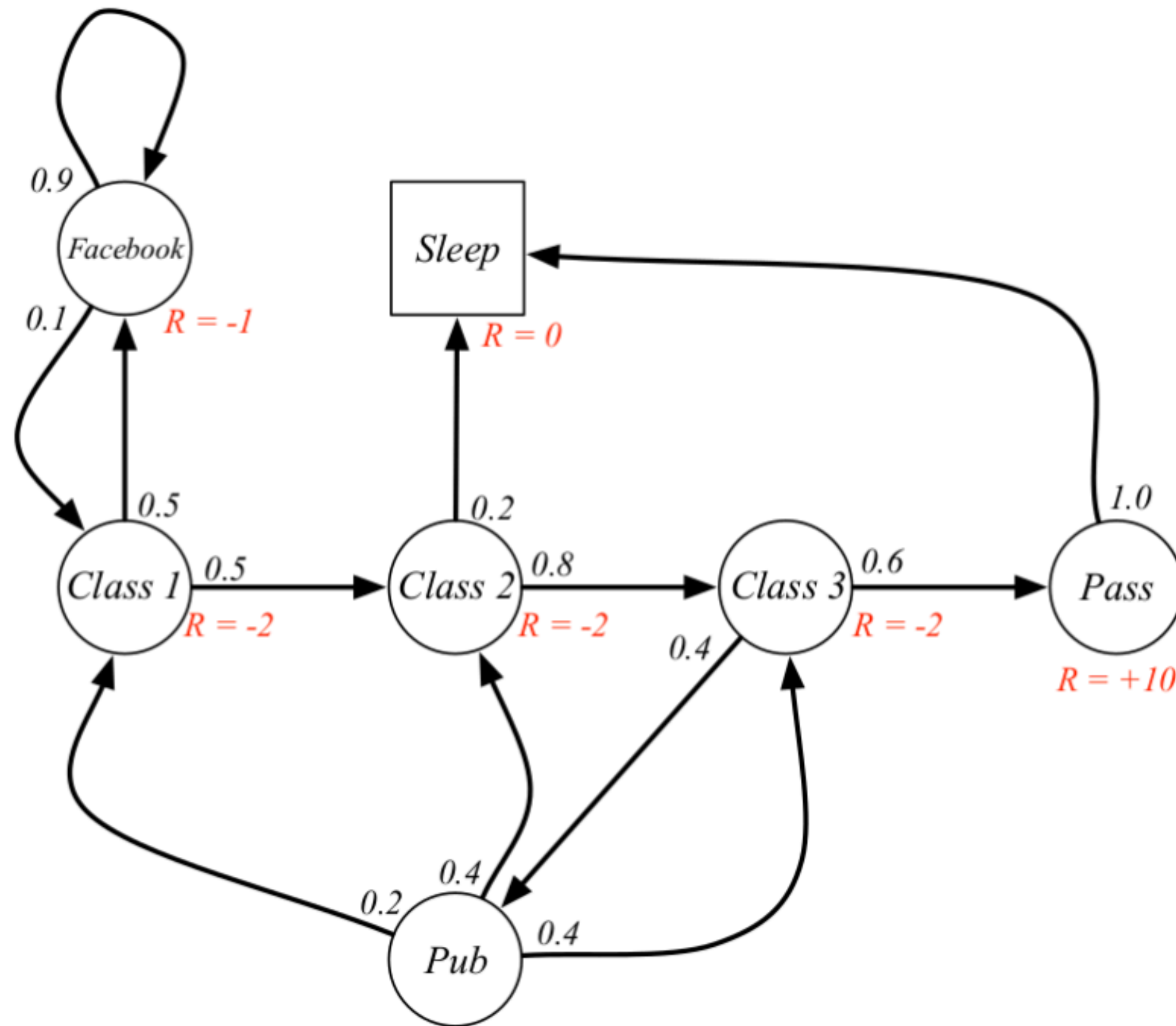# Markov *reward* process



**Gain function:** Far future awards are less valuable than more immediate awards.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

# Markov *reward* process



**State-value function:** Expected return starting from state $S$

$$v(s) = \mathbb{E}(G_t \mid S_t = s)$$

Return for the path of [Class 1 →Class 2 → Class 3 → Pass →Sleep] is:

$$v = -2 - 2 \times \frac{1}{2} - 2 \times \frac{1}{4} + 10 \times \frac{1}{8} = -2.25$$

What is the *optimal* path through potential states that has the highest value?
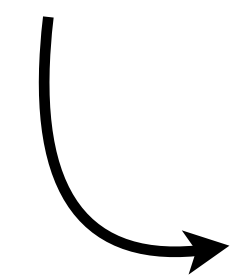
Bellman equation

$$v(s) = \mathbb{E}(R_{t+1} + \gamma v(S_{t+1}) \,|\, S_t = s)$$

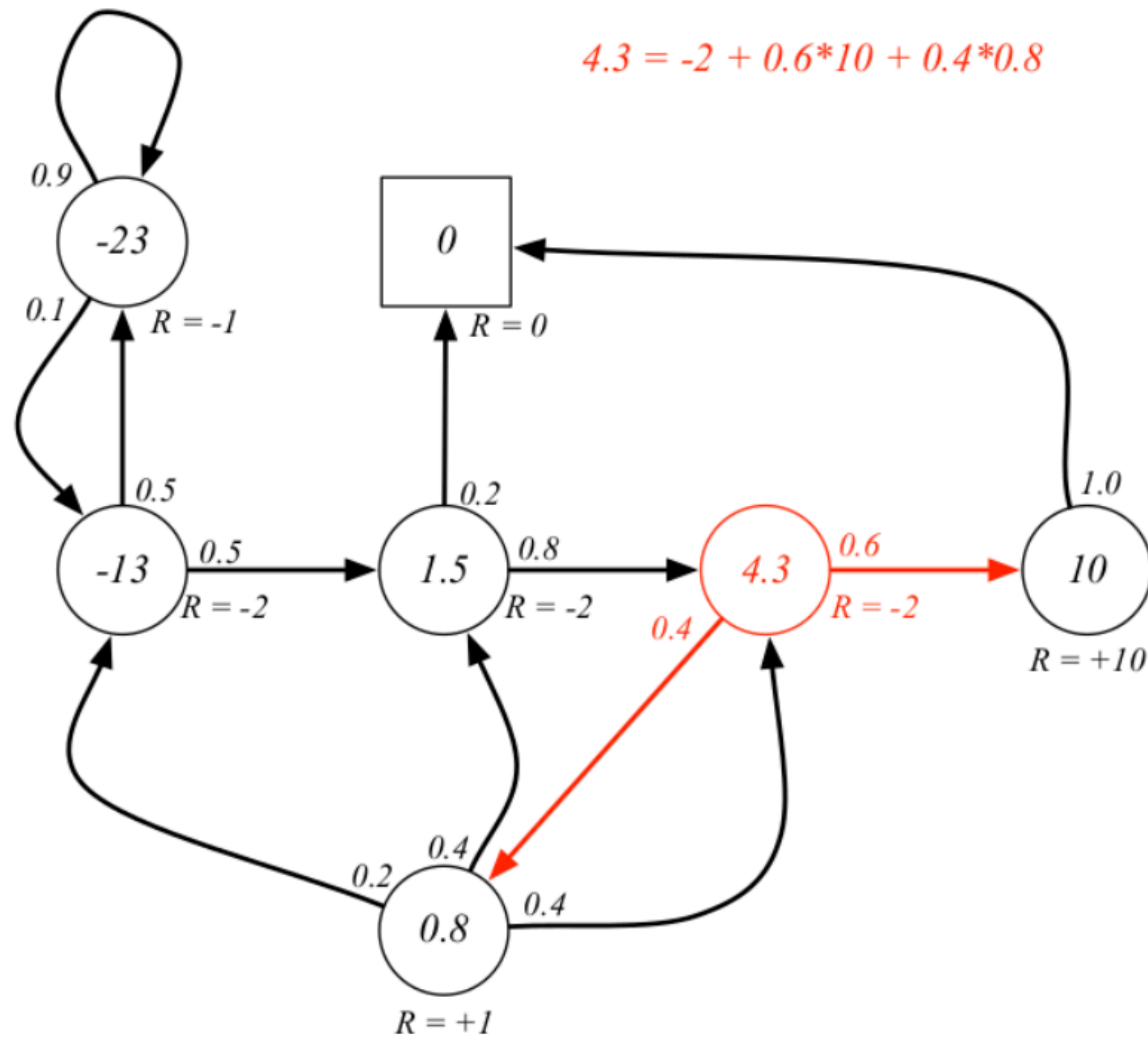$$v(s) = \mathbb{E}(G_t \,|\, S_t = s)$$

$$= \mathbb{E}(R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots \,|\, S_t = s)$$

$$= \mathbb{E}(R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \ldots) \,|\, S_t = s)$$
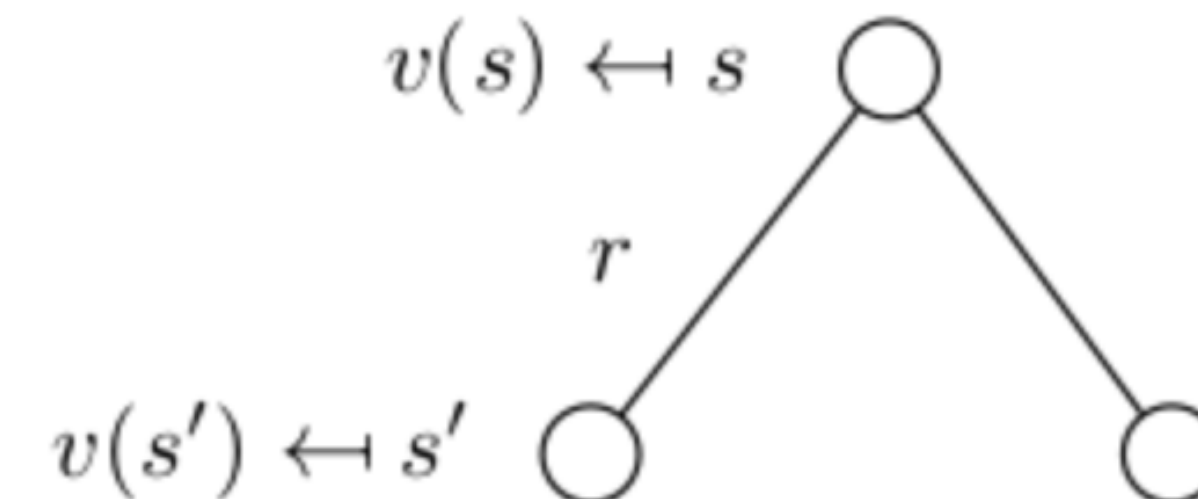
$$= \mathbb{E}(R_{t+1} + \gamma G_{t+1} \,|\, S_t = s)$$

$$G_{t+1} \rightarrow v(S_{t+1})$$
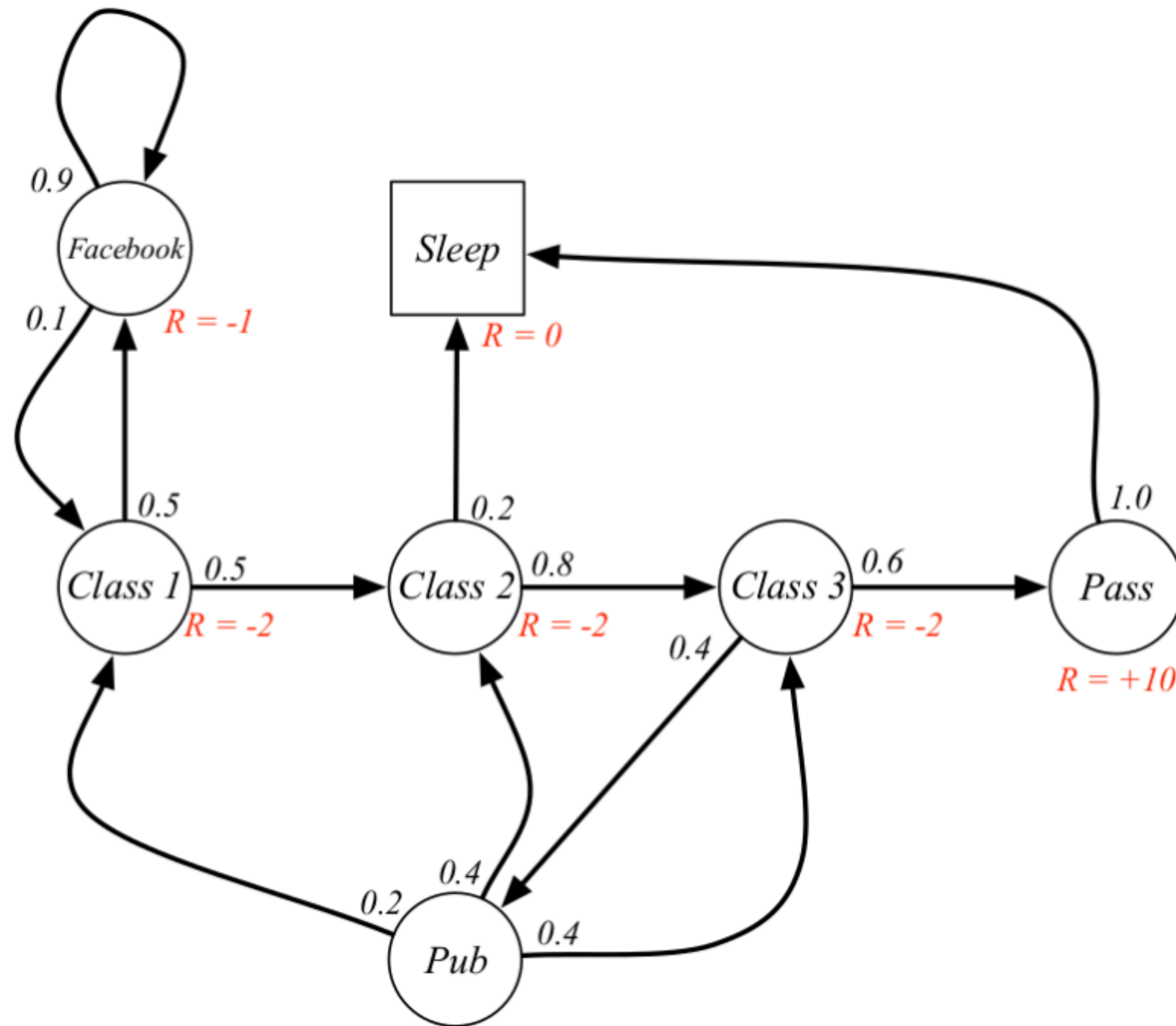
# The Bellman equation



$4.3 = -2 + 0.6*10 + 0.4*0.8$

The value both depends on the reward <u>and</u> the transition probability .

$$v(s) = \mathbb{E}(R_{t+1} + \gamma v(S_{t+1}) \,|\, S_t = s)$$

$$= \mathbf{R}_S + \gamma \sum_{s' \in S} \mathbf{P}_{ss'} v(s')$$

$v(s) \leftarrowtail s$
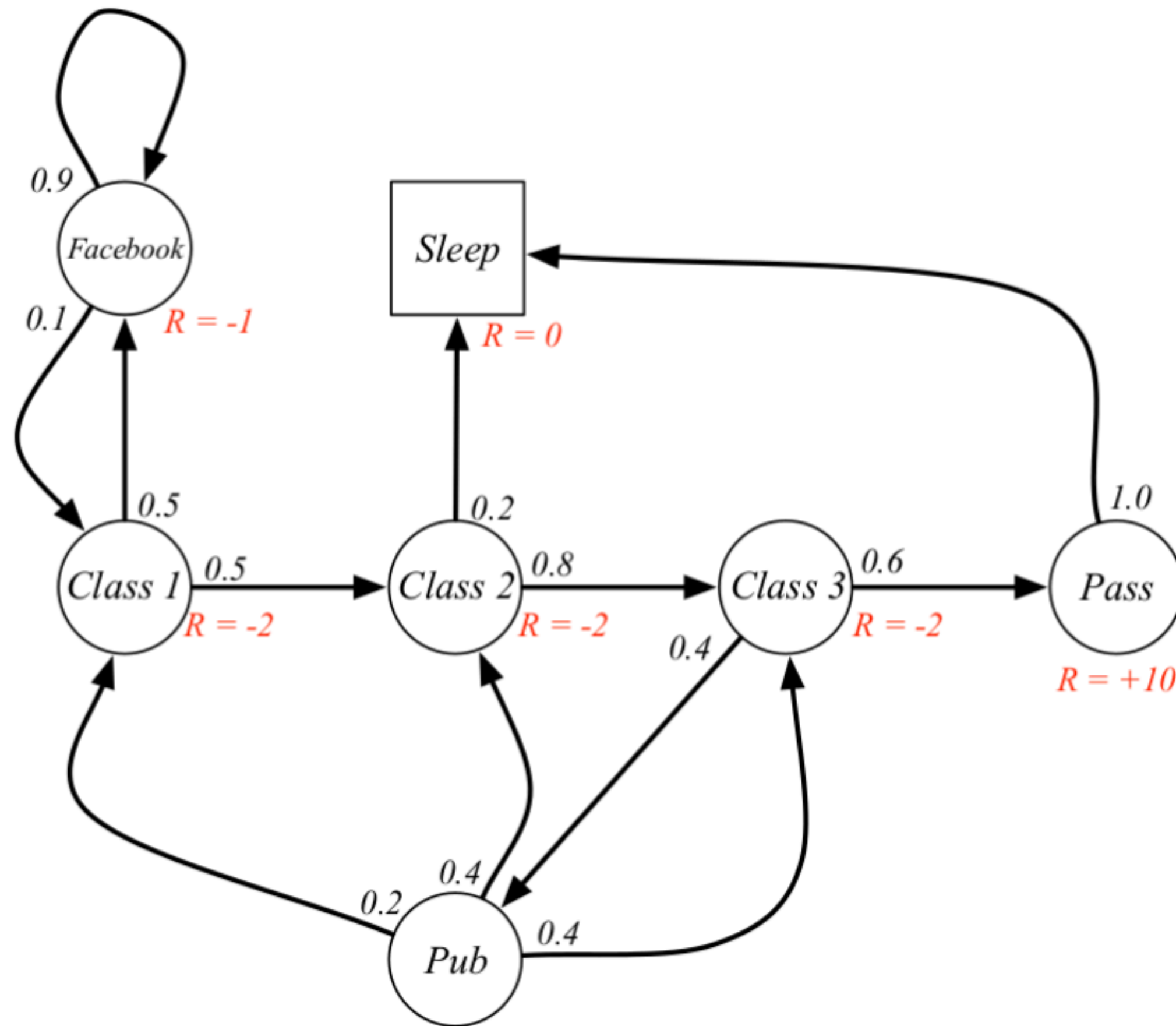
$r$

$v(s') \leftarrowtail s'$

A Markov reward process, with a decision policy $\pi$.

$$\pi(a \,|\, s) = P(A_t = a \,|\, S_t = s)$$

discount factor

reward function

action

$$(\mathbf{S}, \mathbf{P}, \mathbf{R}, \gamma, \mathbf{A})$$

state space

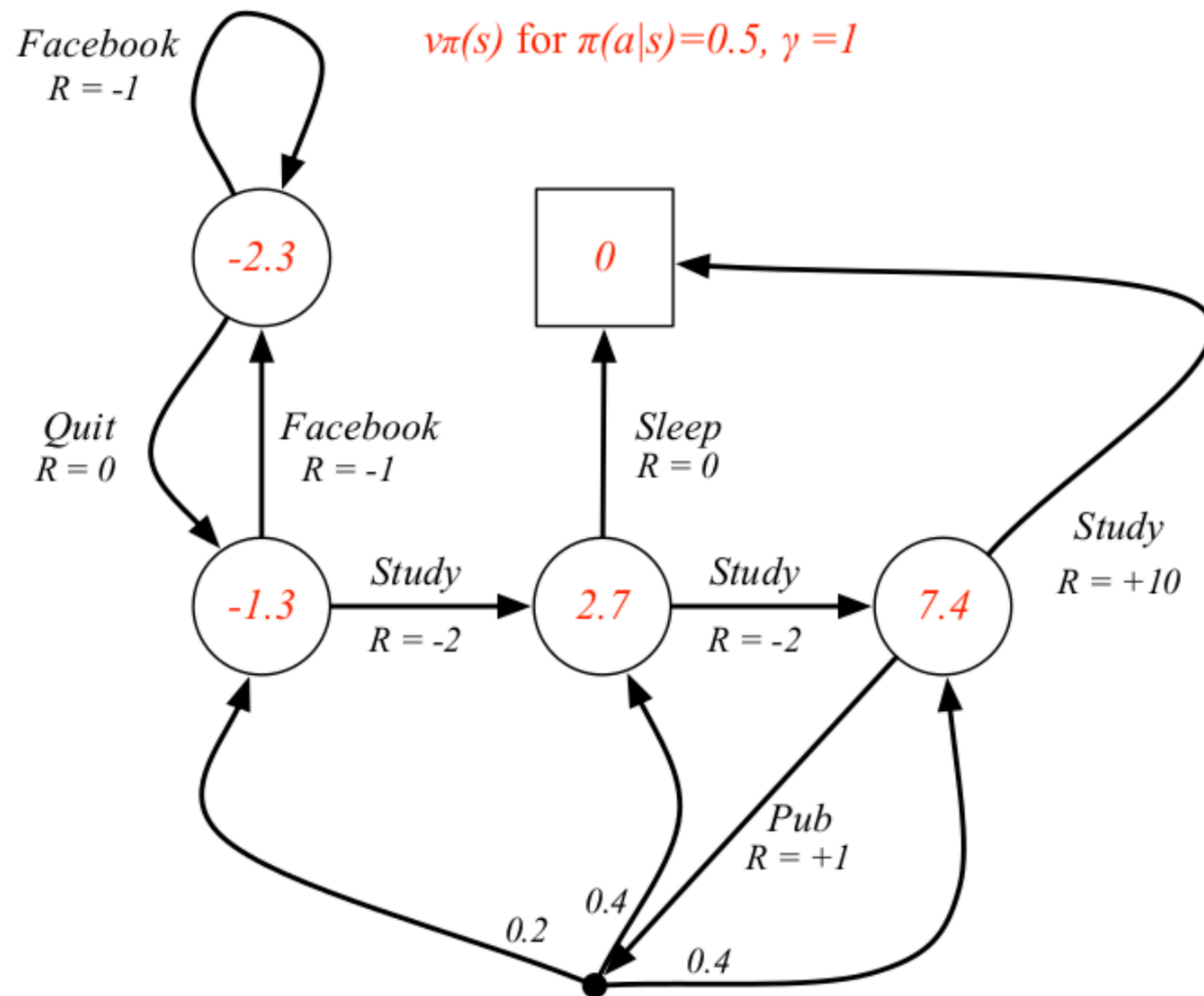transition function

# Markov *decision* process



**State transition function**

$$P_{s,s'}^{\pi} = \sum_{a \in A} \pi(a \mid s) P_{s,s'}^{a}$$

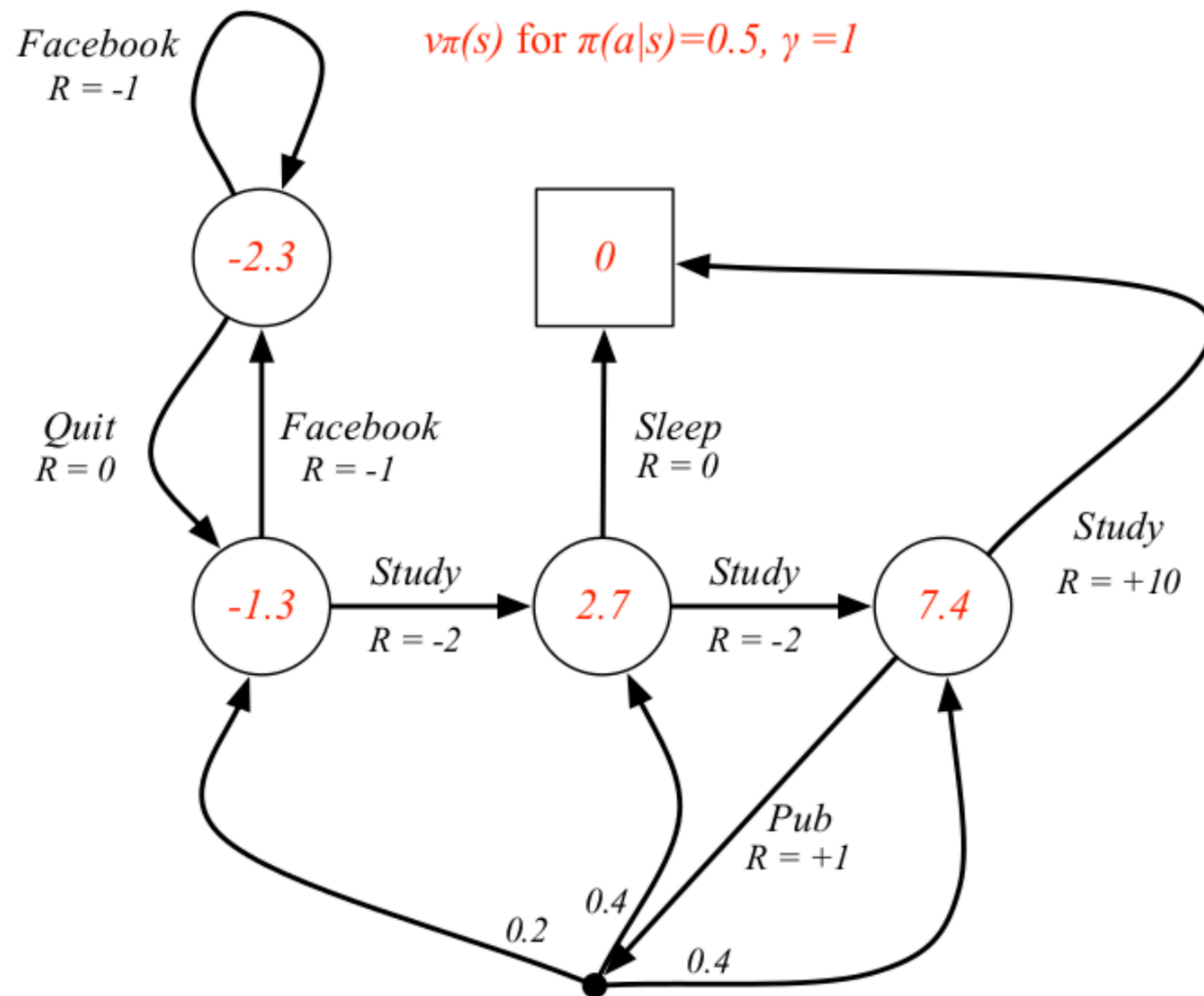**Reward function**

$$R_{s}^{\pi} = \sum_{a \in A} \pi(a \mid s) R_{s}^{a}$$

Facebook
R = -1

vπ(s) for π(a|s)=0.5, γ =1

-2.3

0

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

-1.3

Study
R = -2

2.7

Study
R = -2

7.4

Pub
R = +1

0.4

0.2

0.4

The **state-value** function $V_\pi$(s) of an MDP is the expected return starting from state S, and then following policy $\pi$.

$$v_\pi(s) = \mathbb{E}_\pi(G_t \mid S_t = s)$$

$$= \mathbb{E}_\pi(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s)$$

# Markov *decision* process



The **state action-value** function $q_\pi(s, a)$ is the expected return starting from state s, taking action a, and following policy $\pi$.

$$q_\pi(s, a) = \mathbb{E}_\pi(G_t \mid S_t = s, A_t = a)$$

$$= \mathbb{E}_\pi\left(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a\right)$$

# Take home message

- Markov decision processes (MDPs) show how the future can be independent of the past conditioned on the present.

- MDPs determine the ideal sequence of decisions to maximize *future* rewards when transition probabilities and rewards are known.