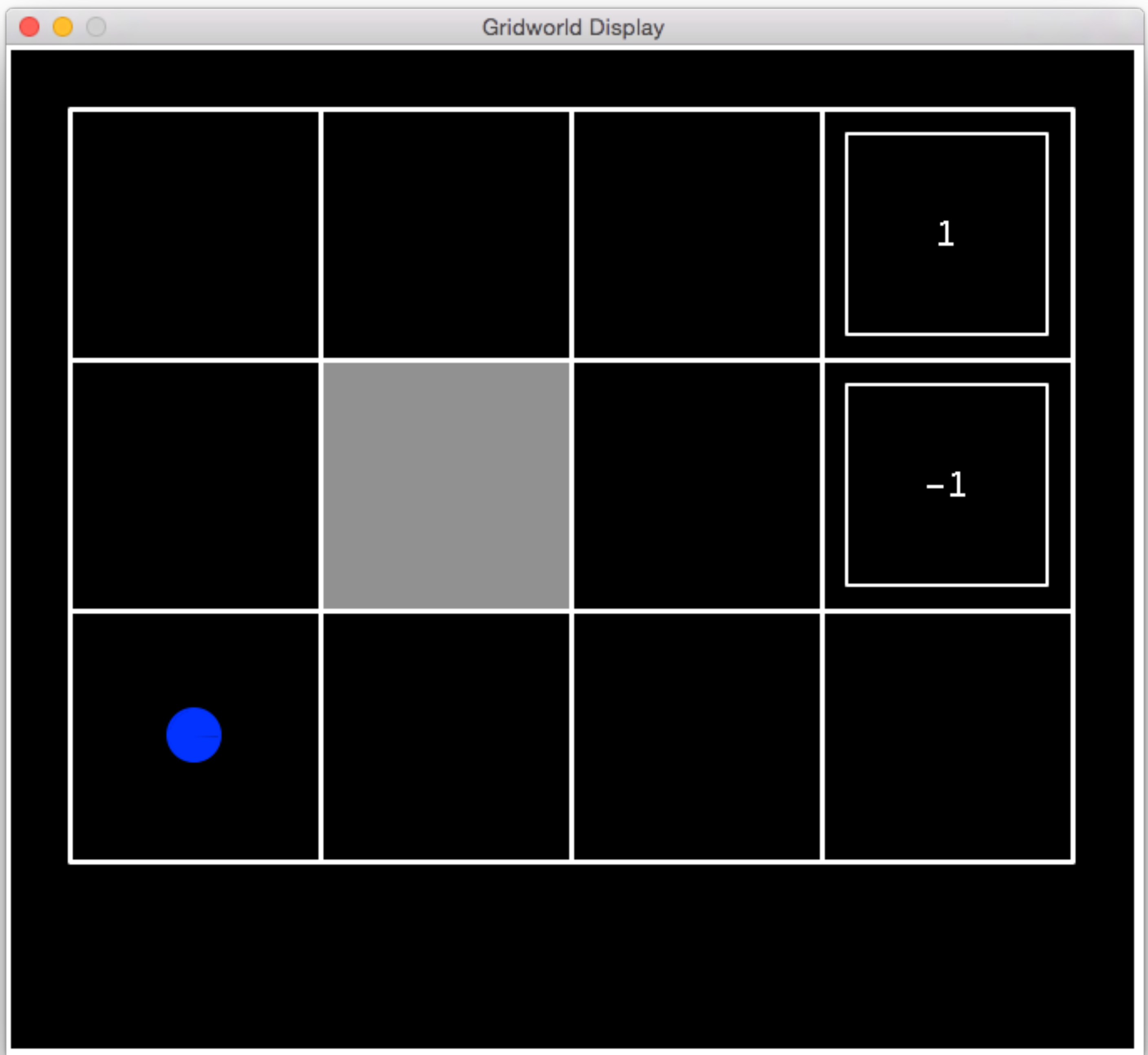


# Markov Decision Processes

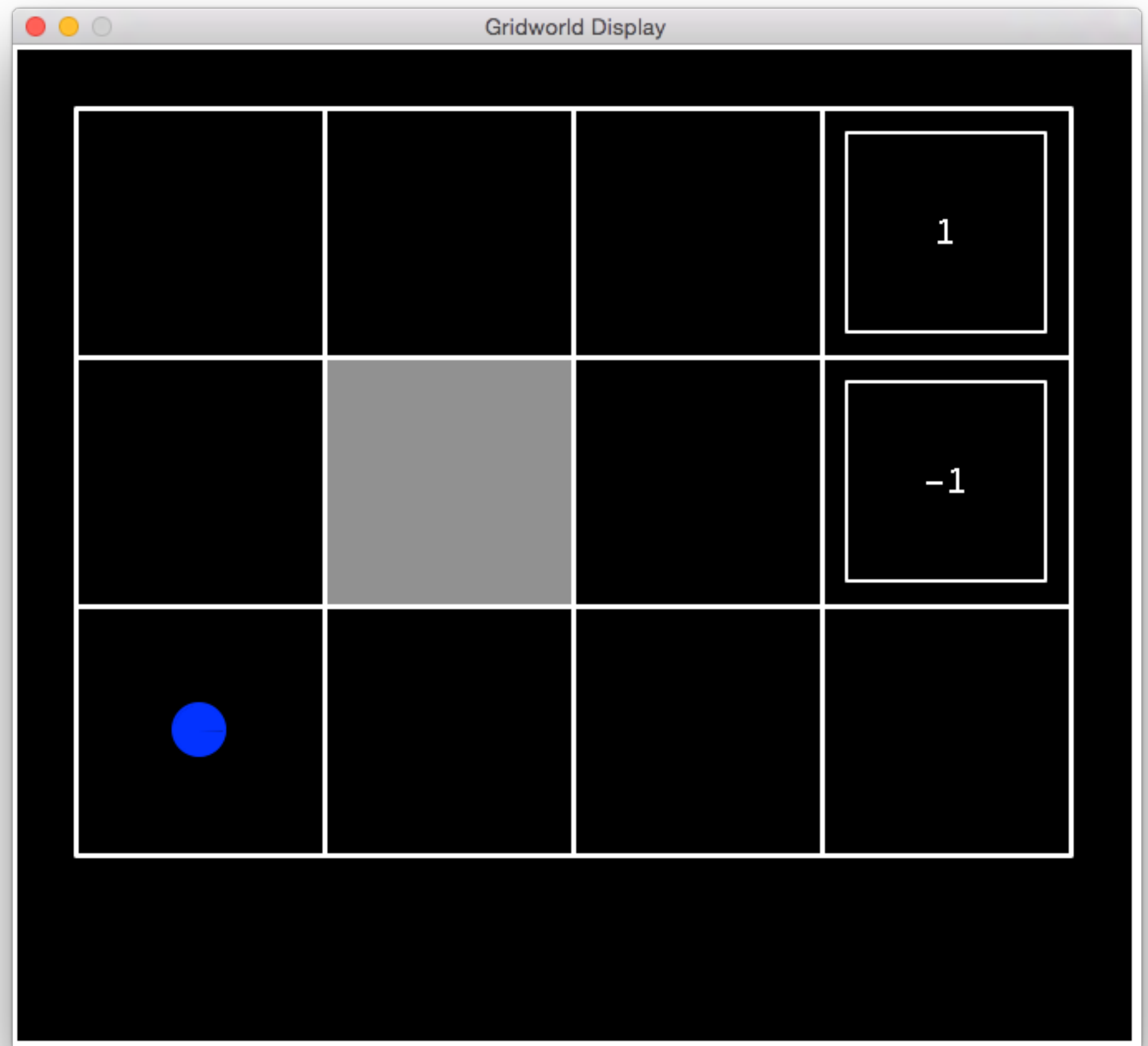
Virginia Tech CS5804  
Spring 2015

# Outline

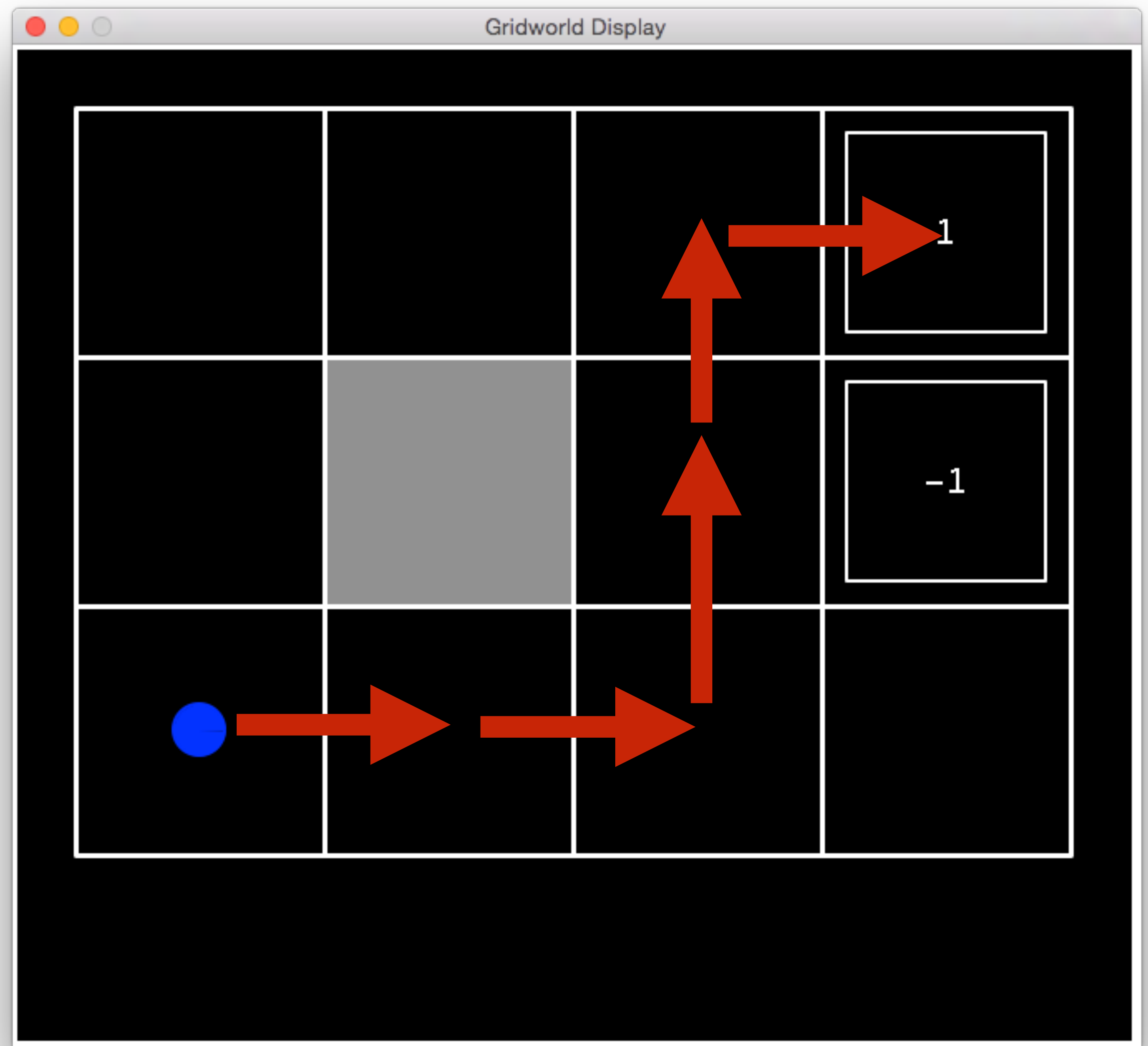
- Markov decision process: richer environment representation
- Reward functions
- Optimizing policies via value iteration



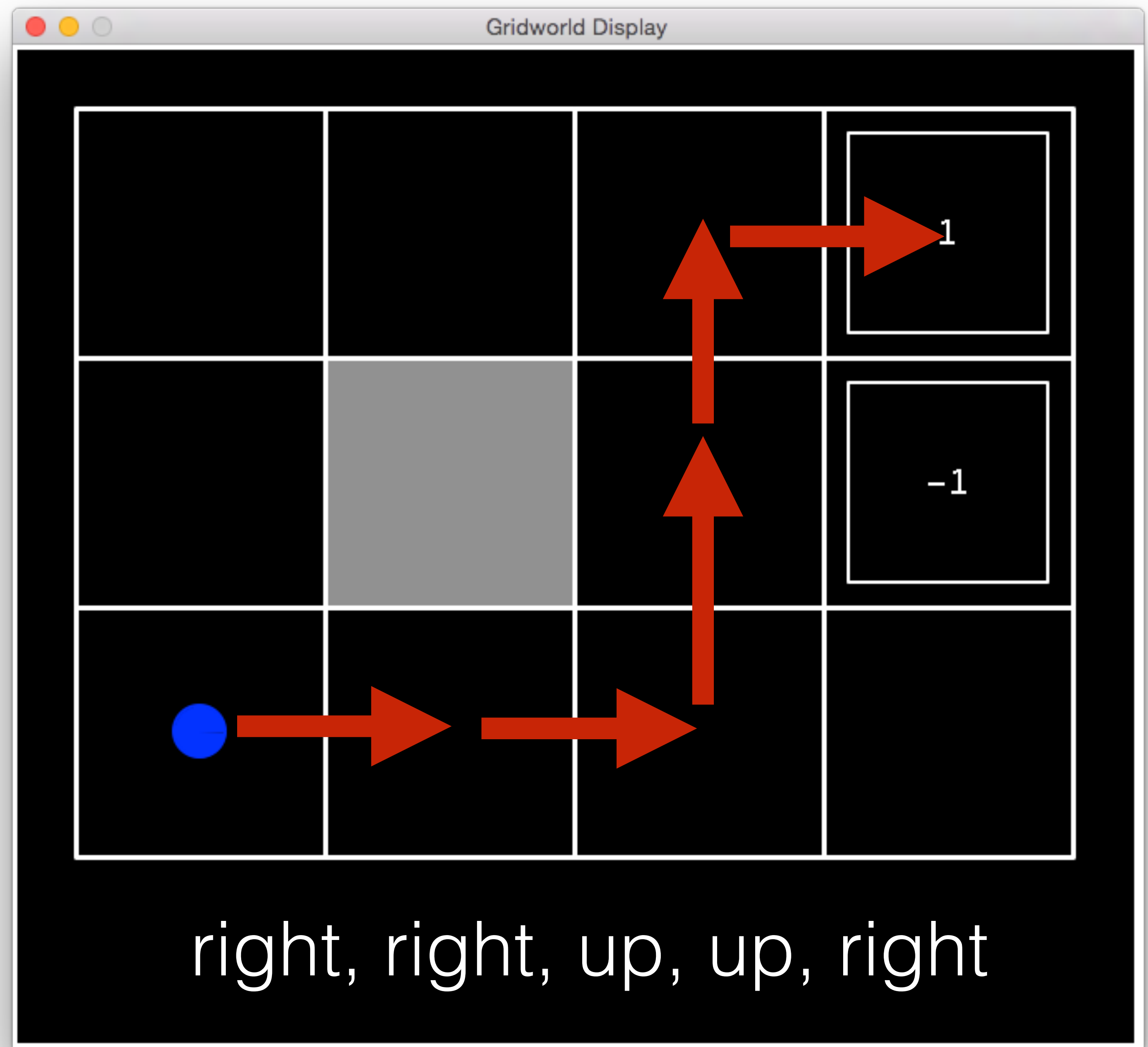
collect reward



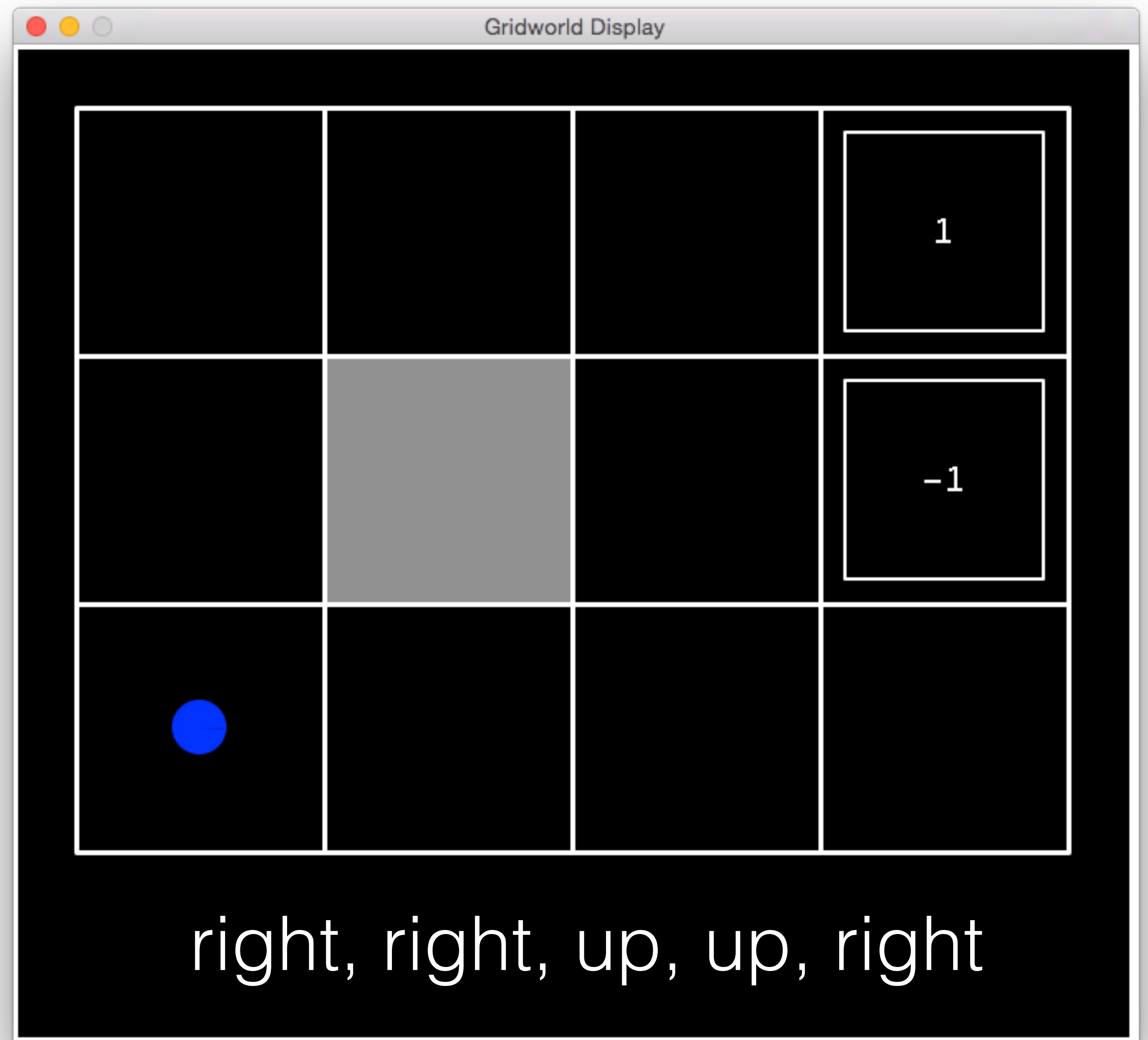
collect reward



collect reward

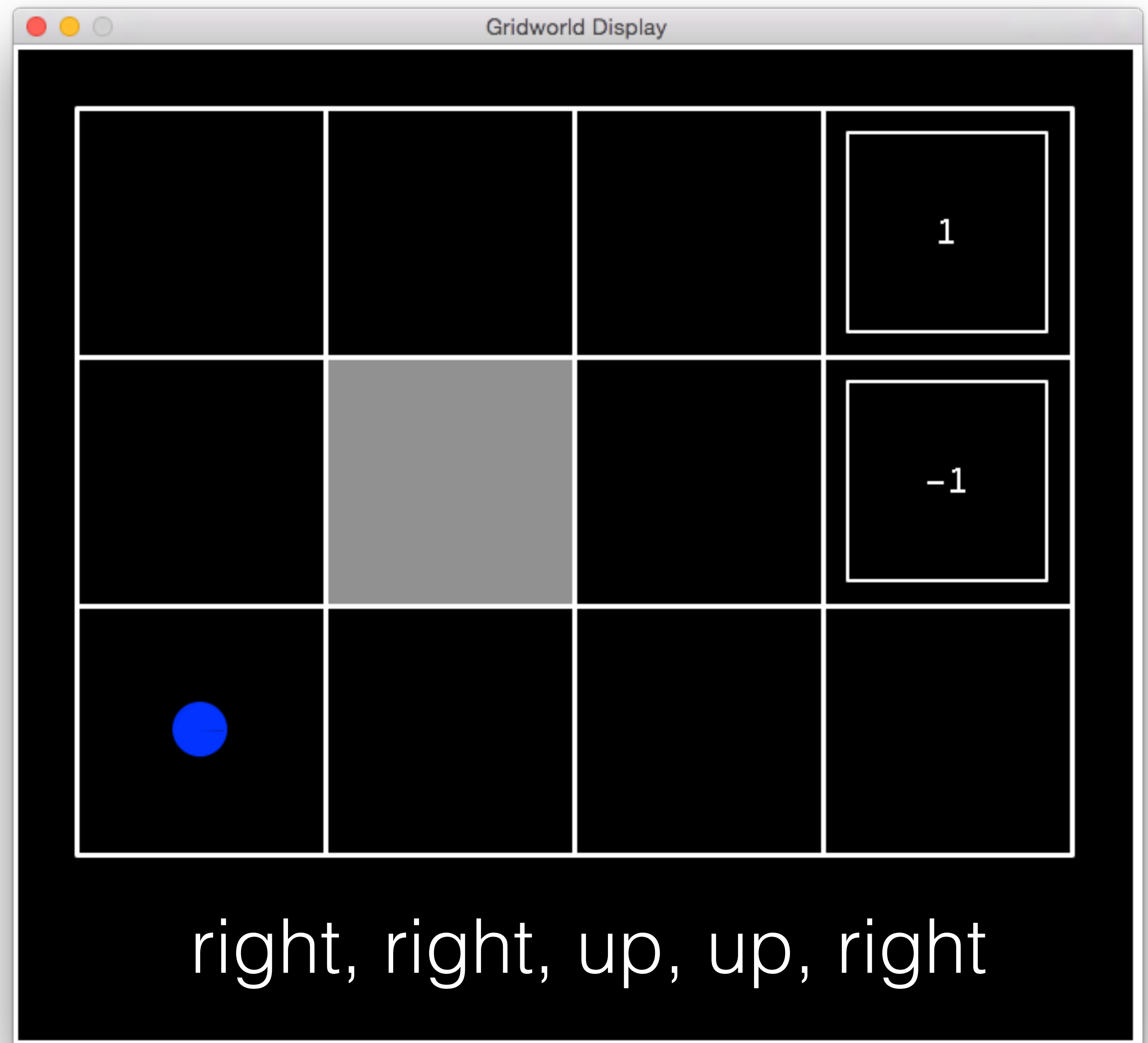


collect reward



collect reward

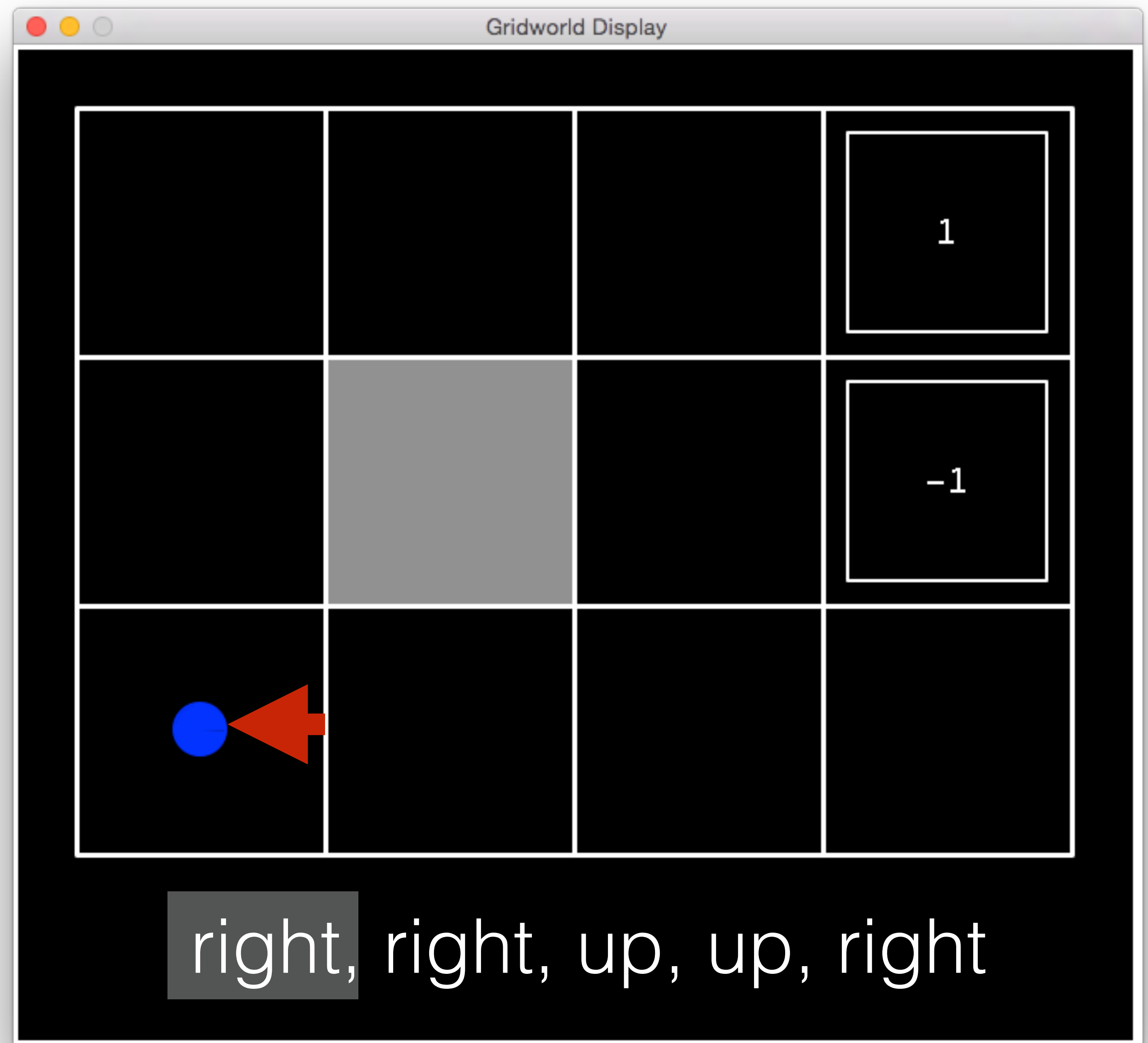
stochastic transitions





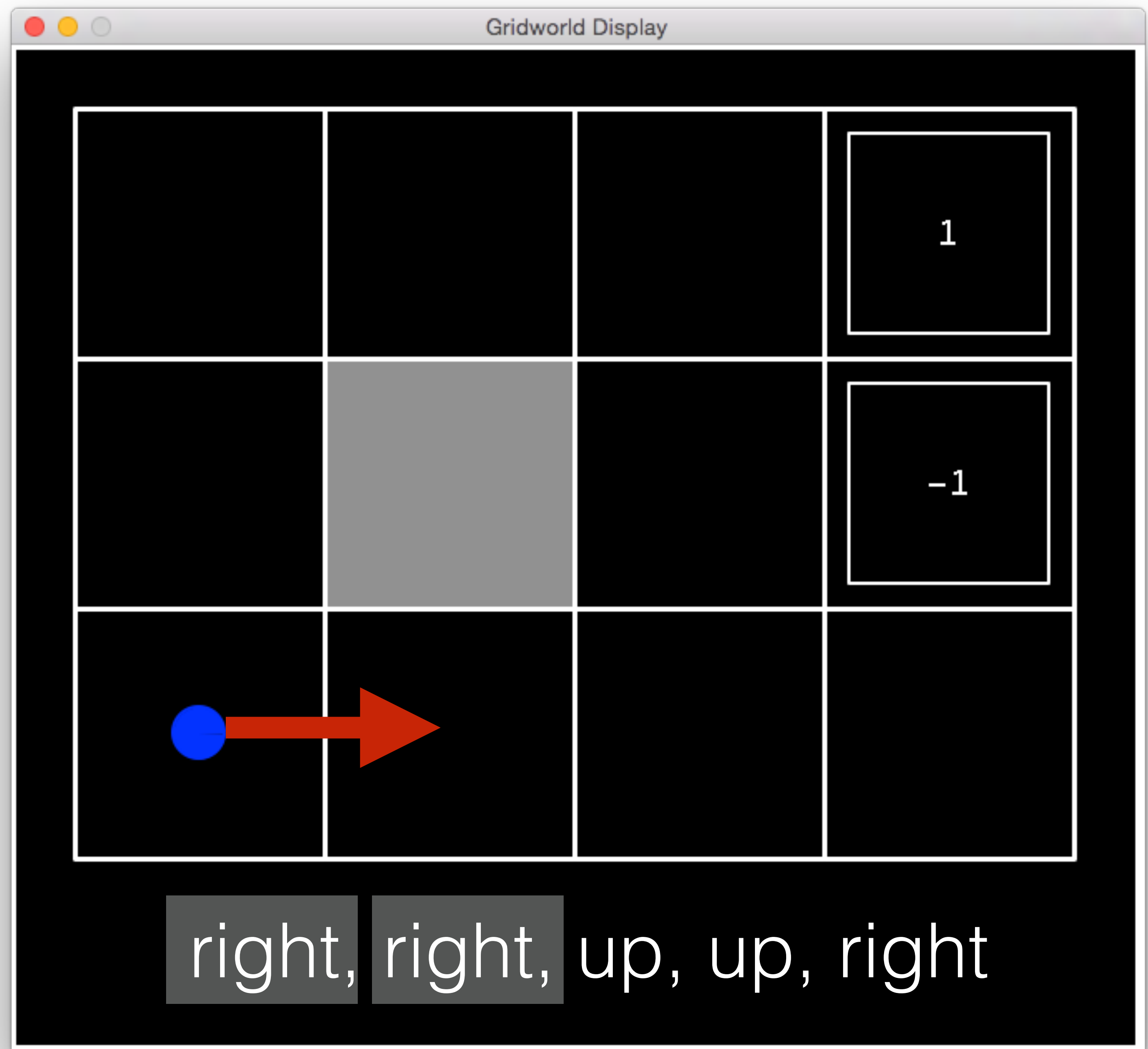
collect reward

stochastic transitions



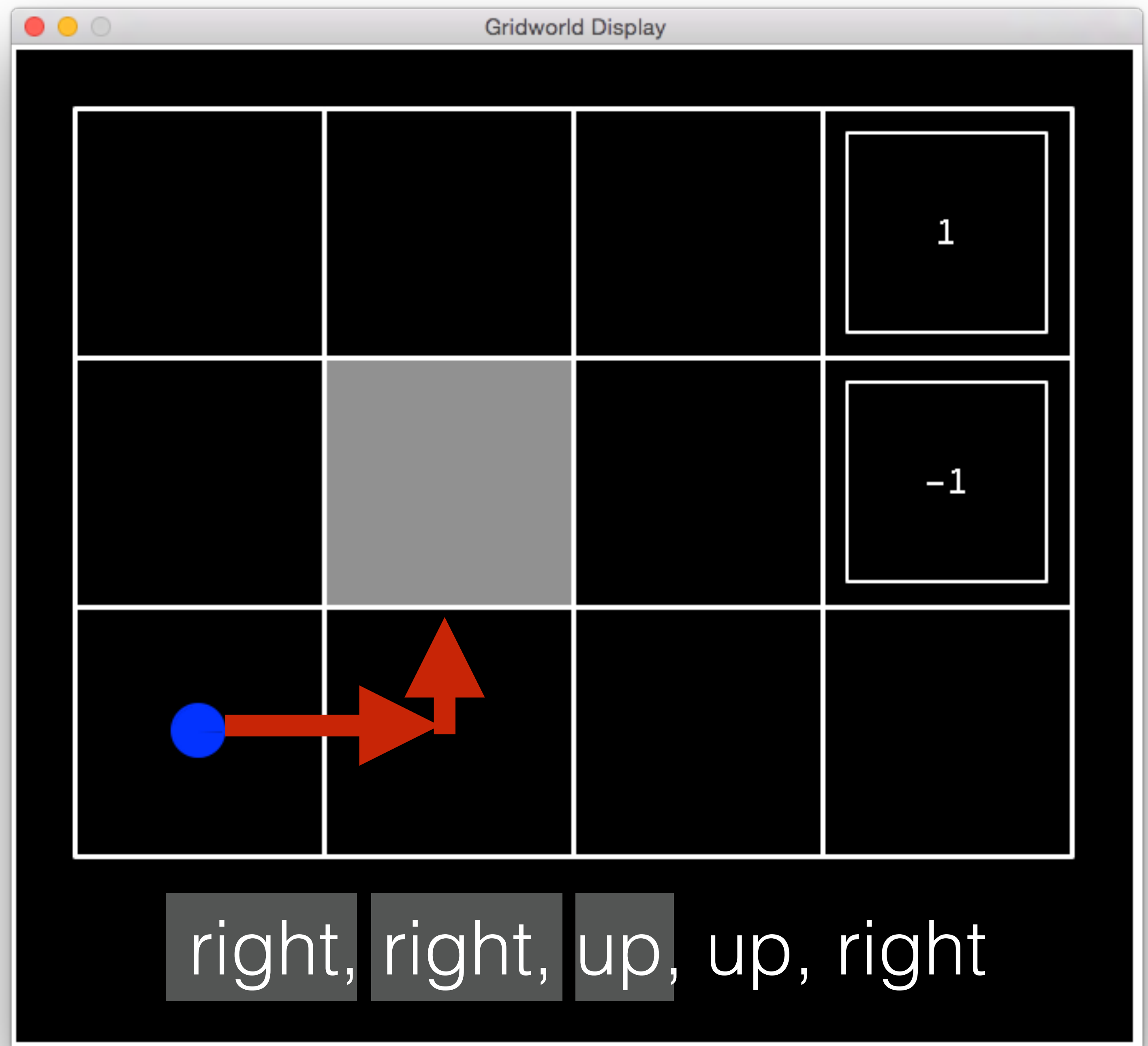
collect reward

stochastic transitions



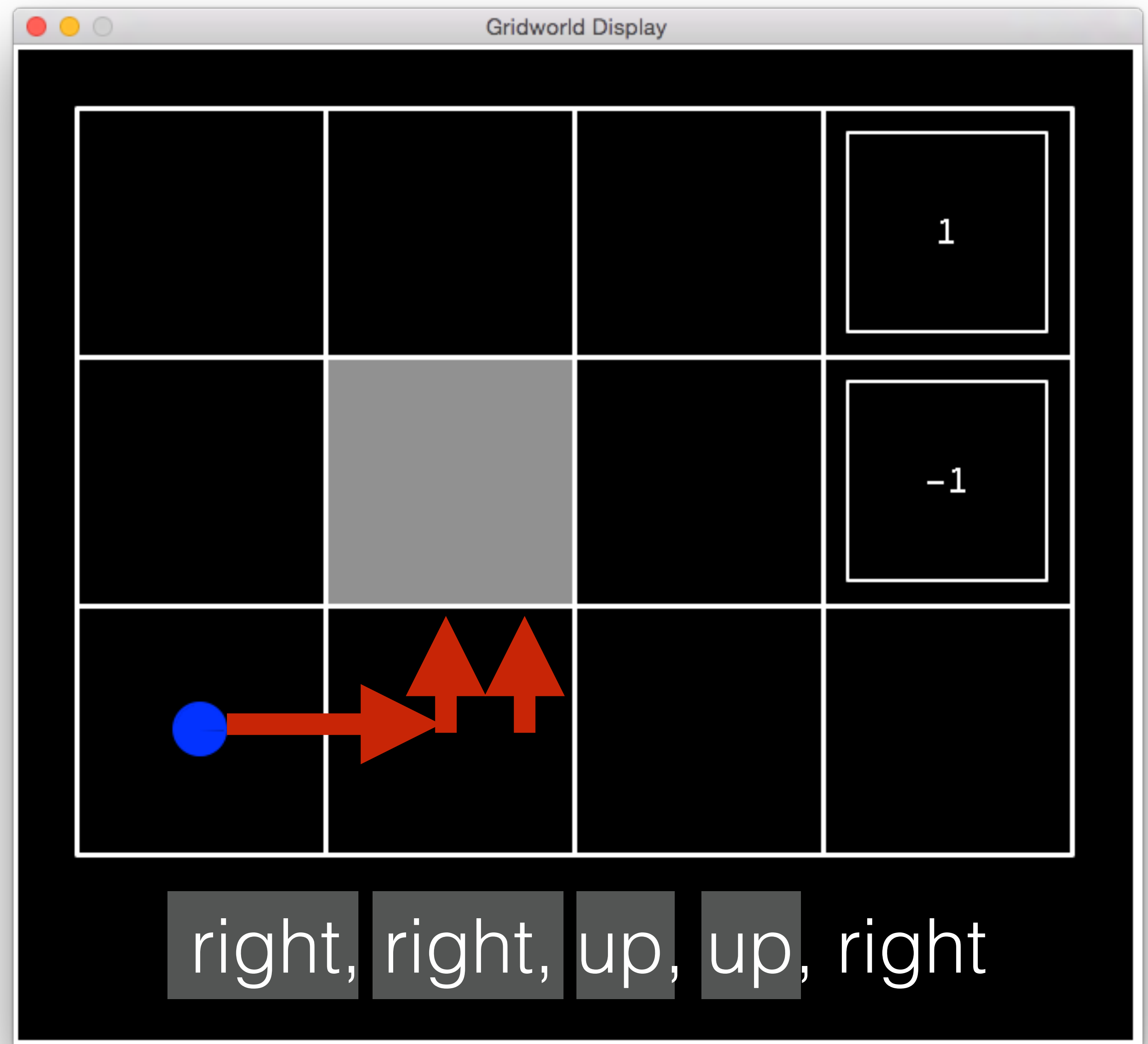
collect reward

stochastic transitions



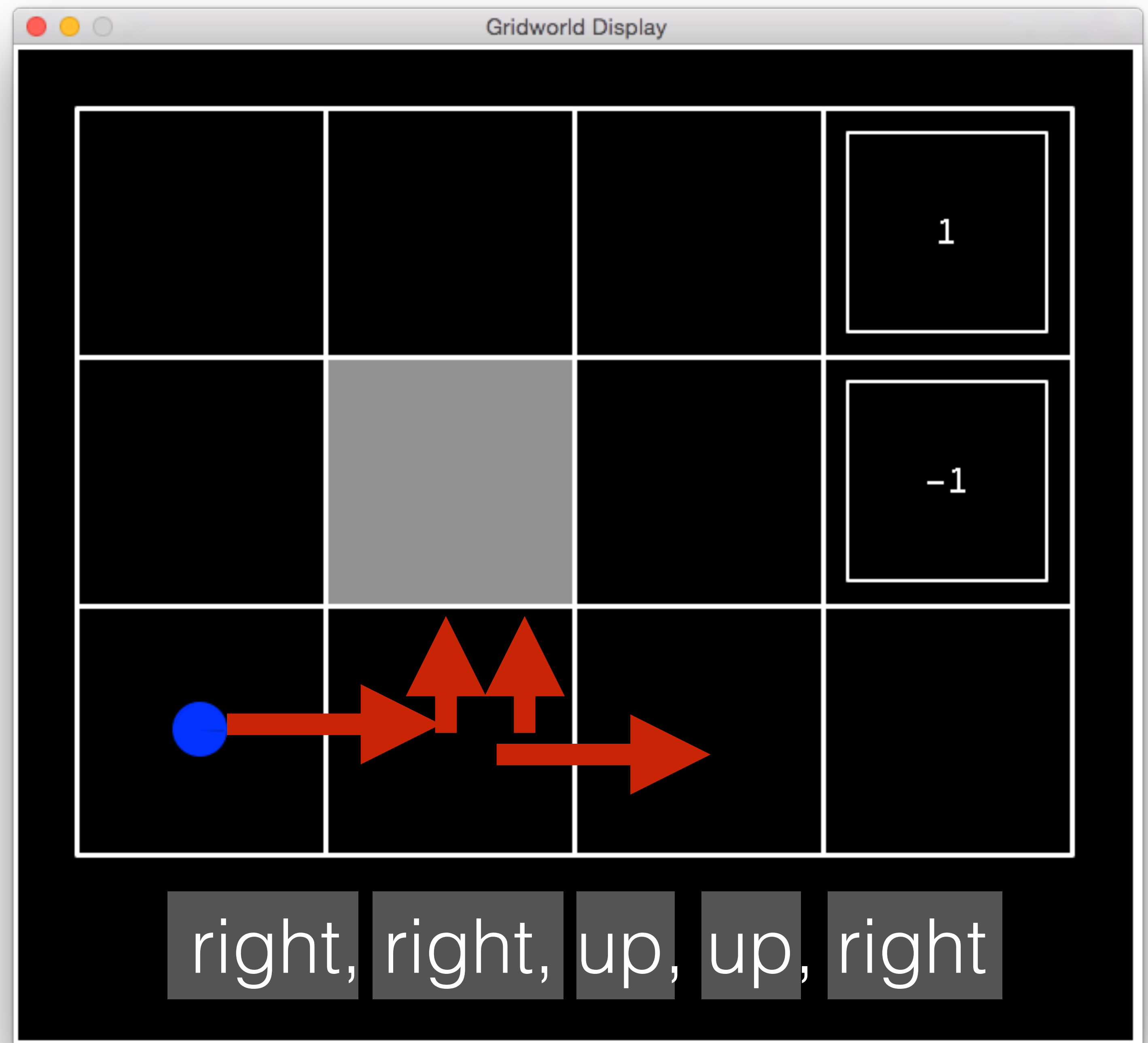
collect reward

stochastic transitions



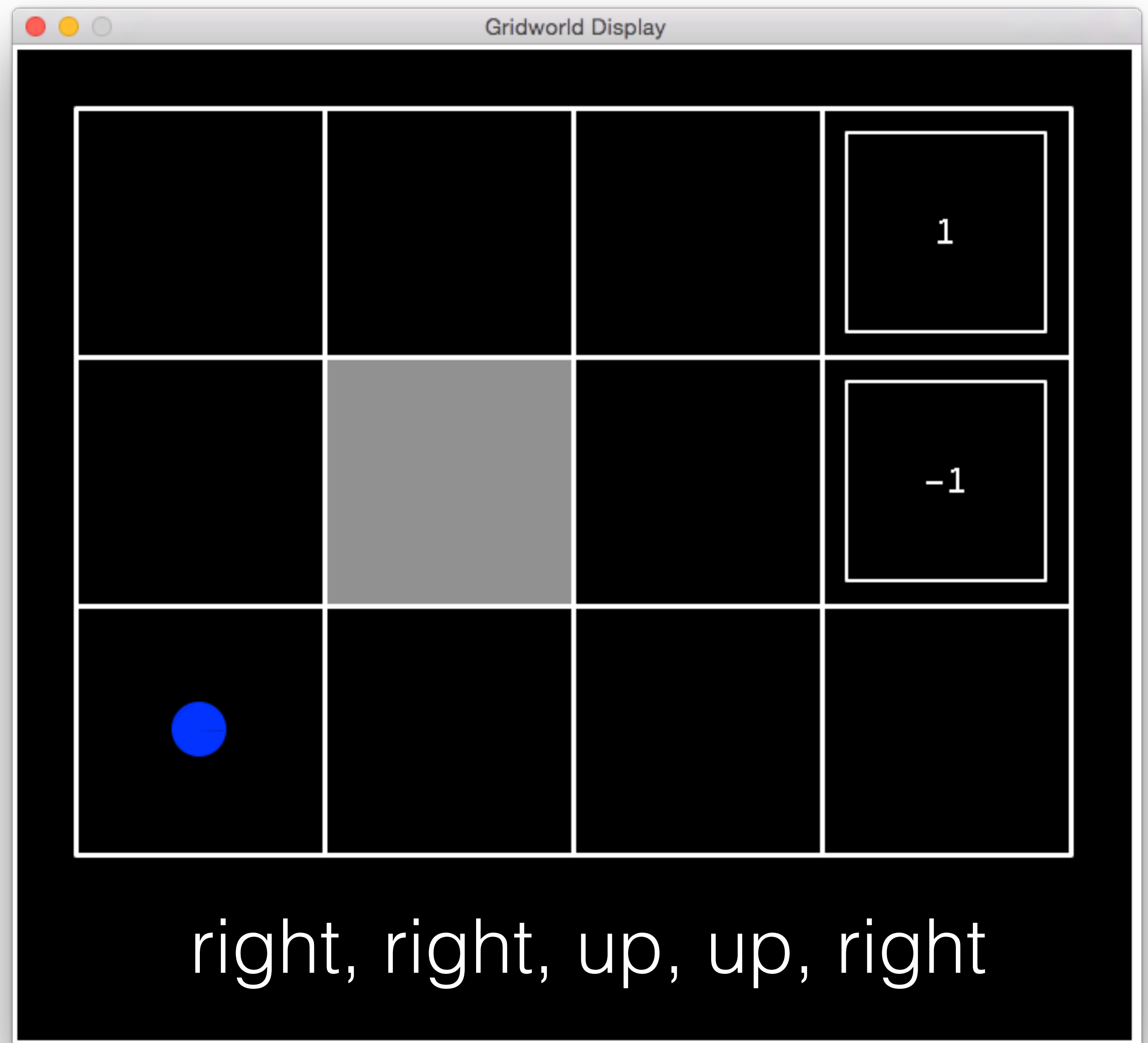
collect reward

stochastic transitions



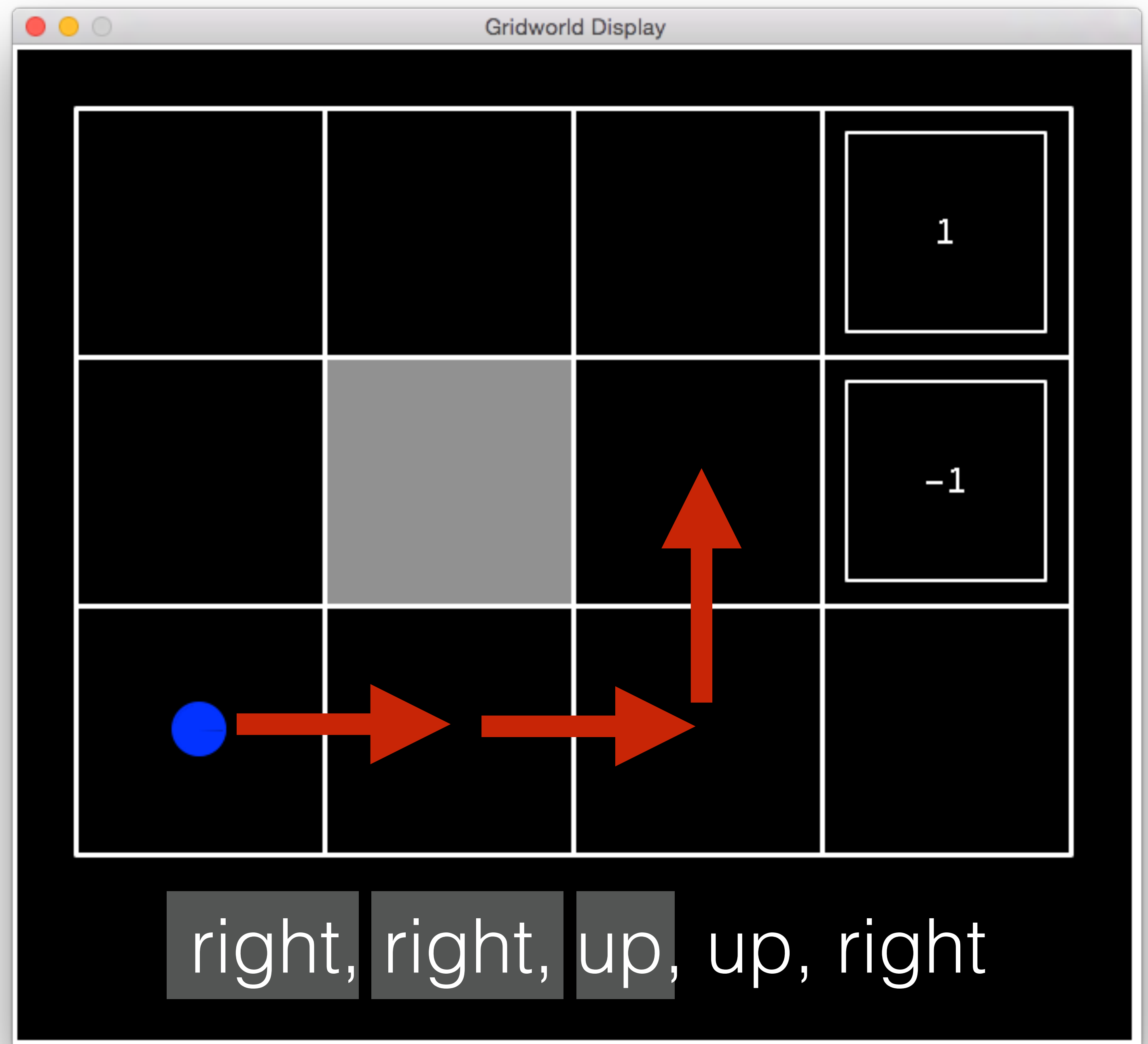
collect reward

stochastic transitions



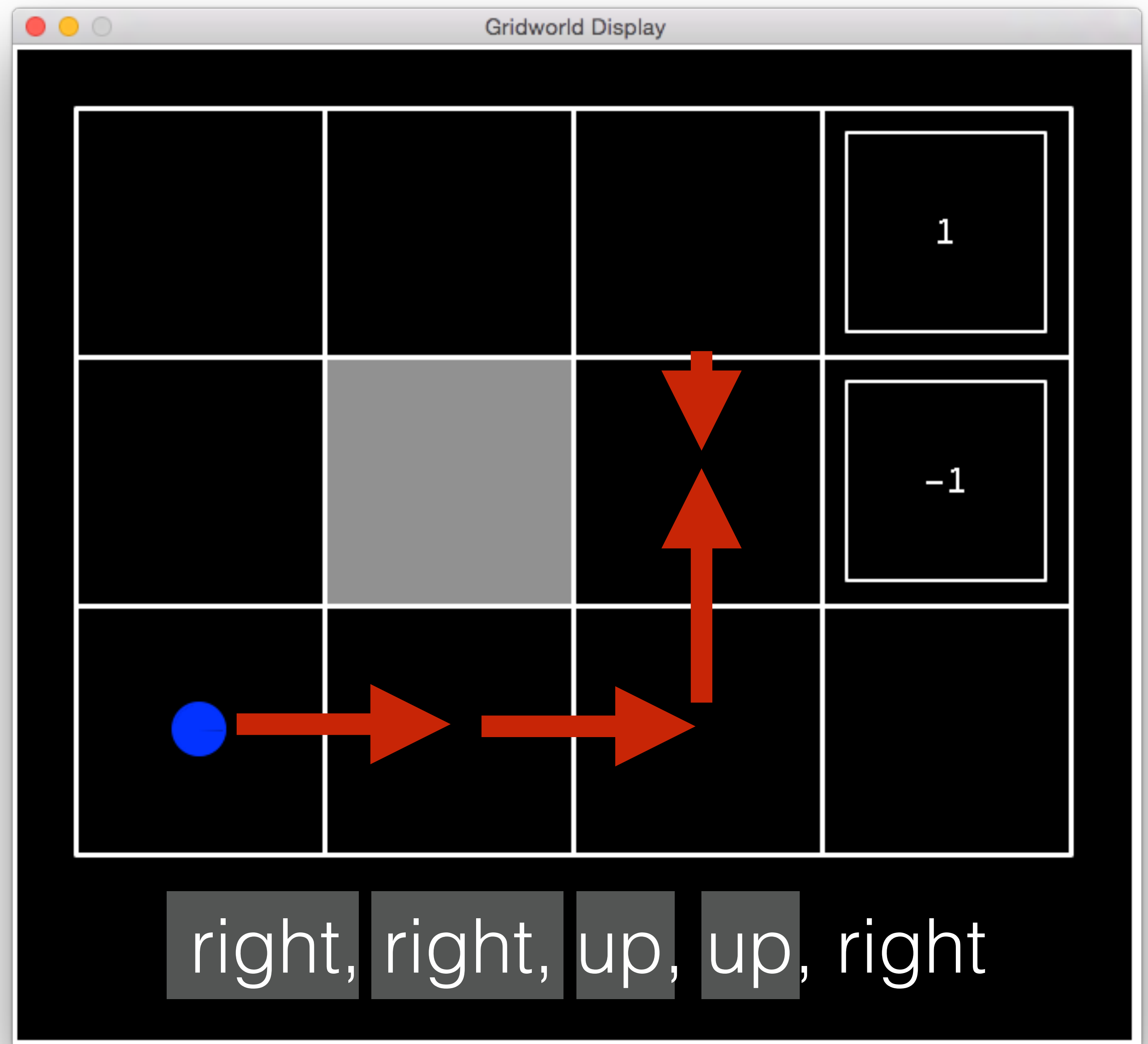
collect reward

stochastic transitions



collect reward

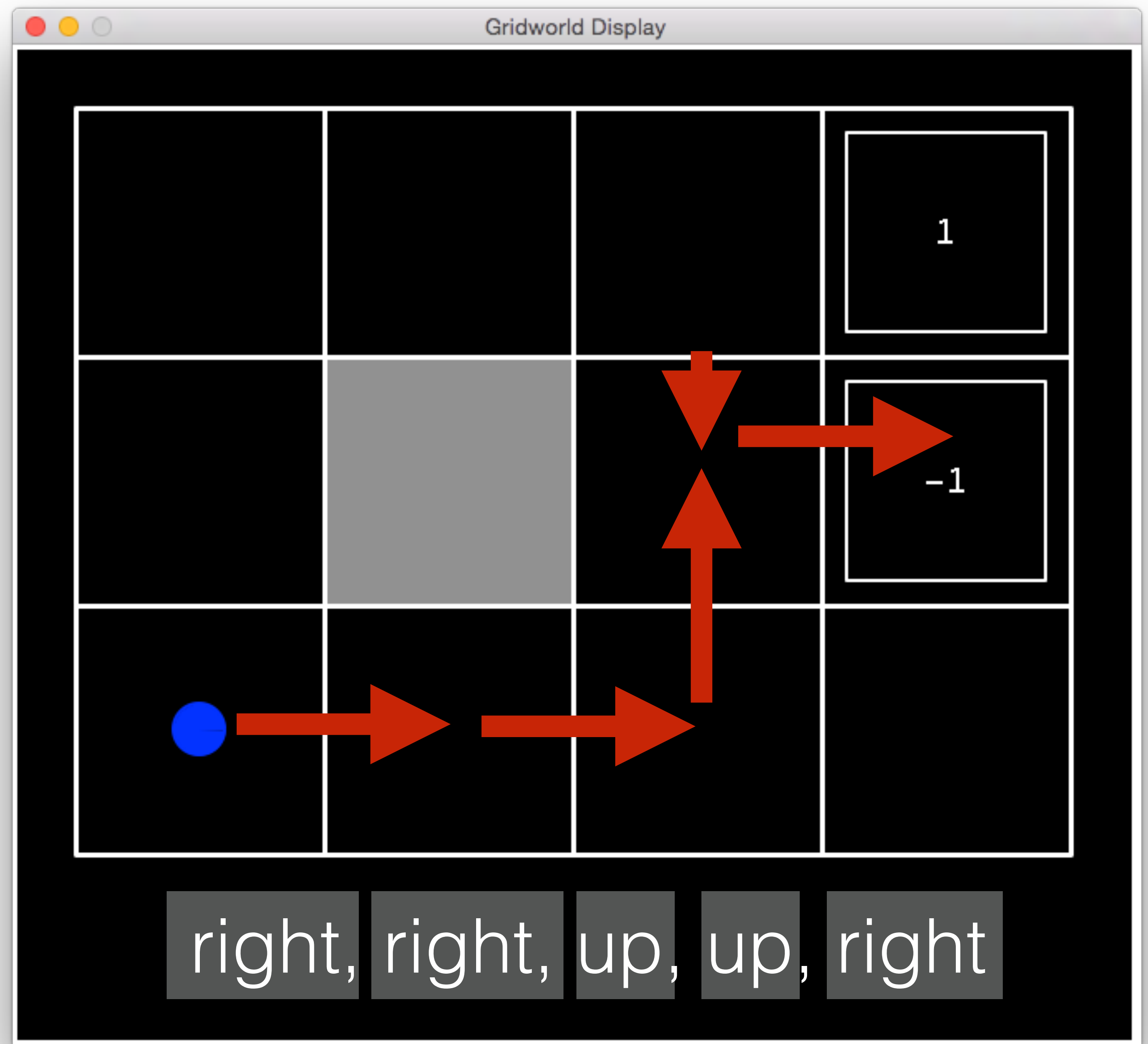
stochastic transitions





collect reward

stochastic transitions

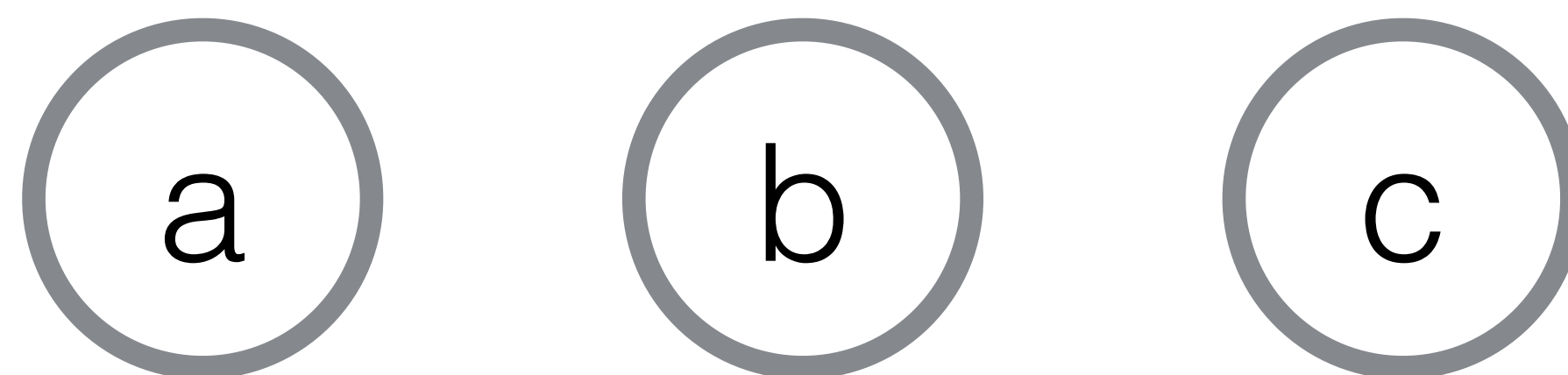


# Actions and Transitions

- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$

# Actions and Transitions

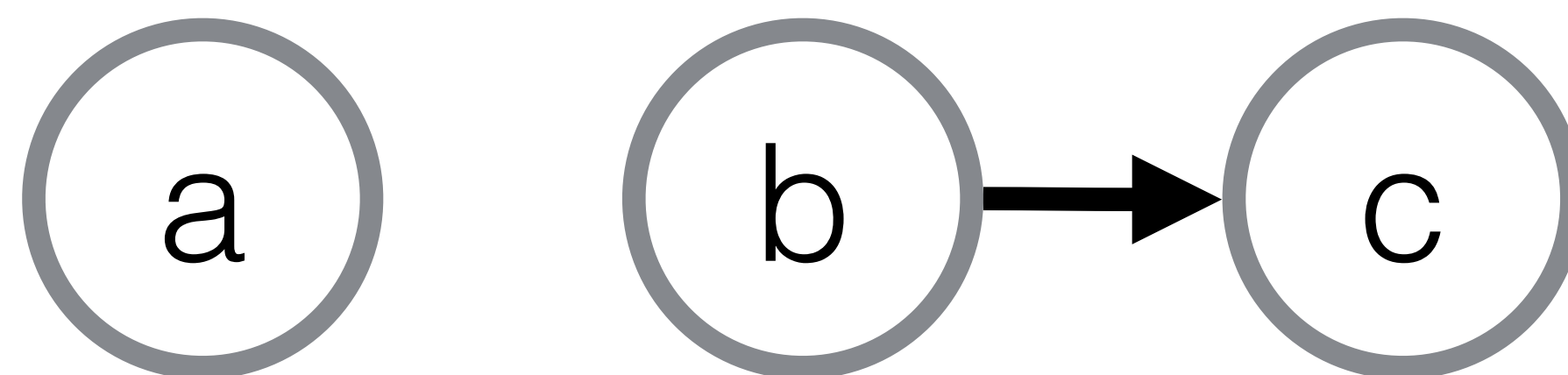
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



**a** = right

# Actions and Transitions

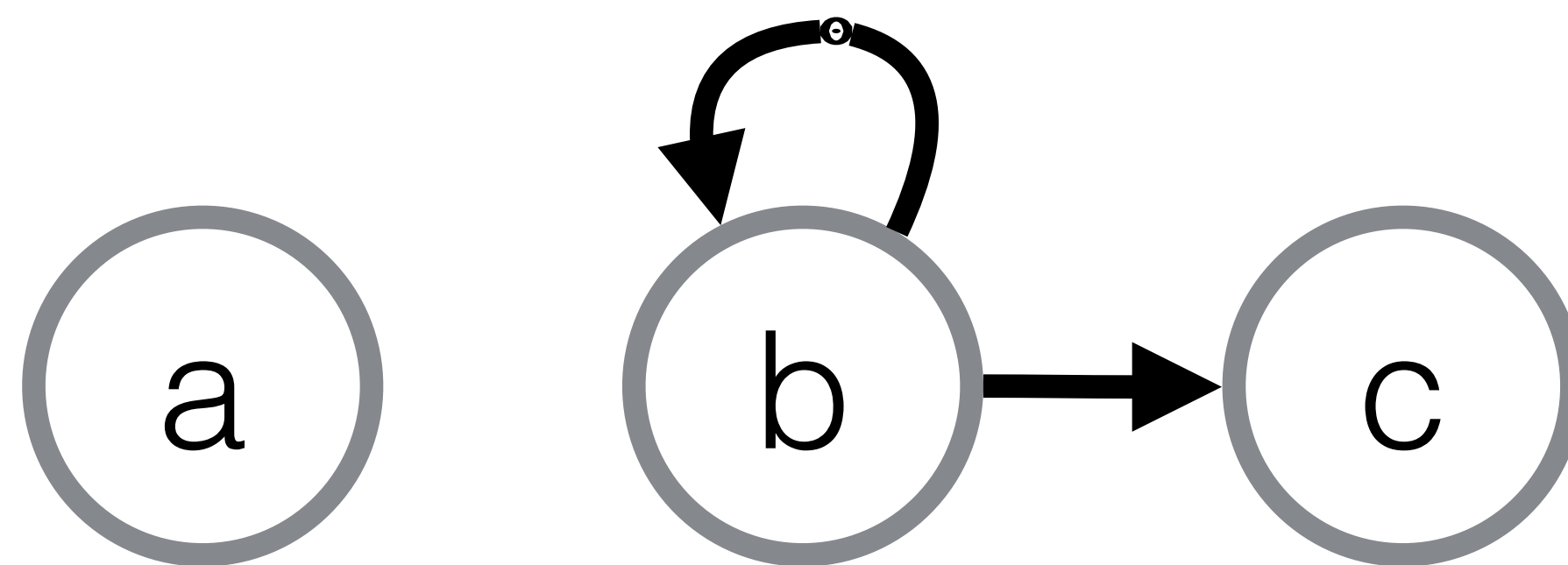
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

# Actions and Transitions

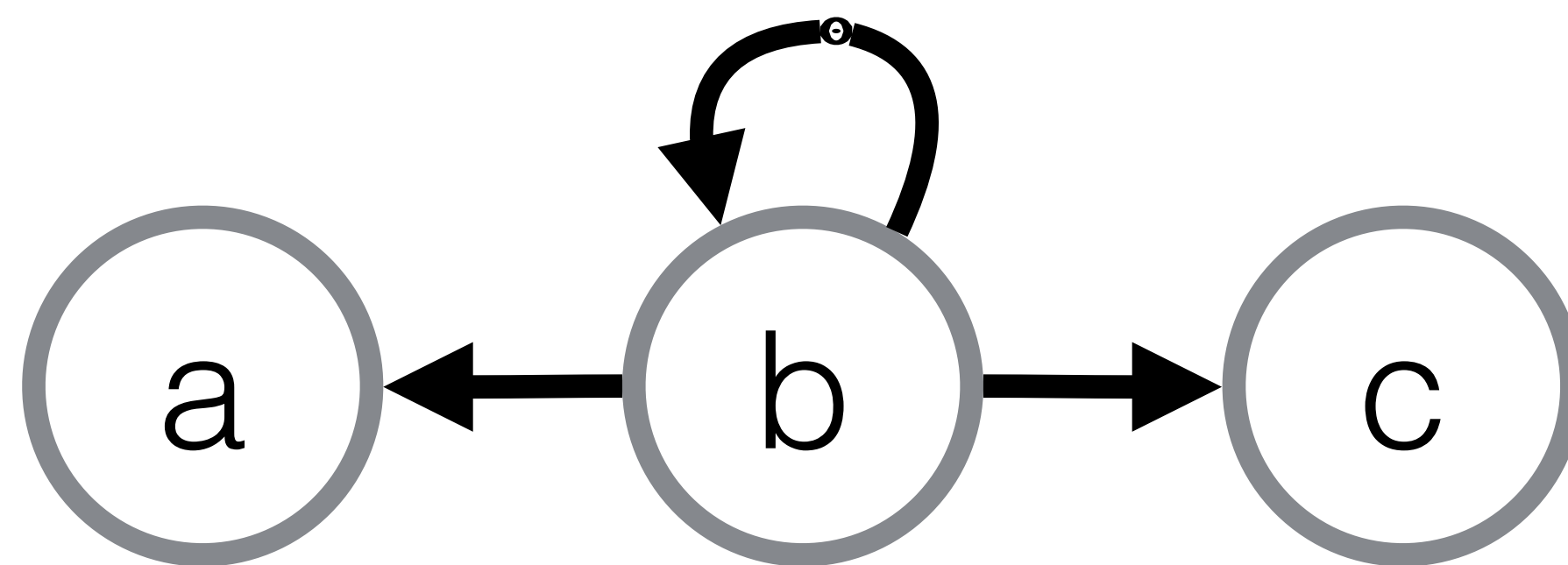
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

# Actions and Transitions

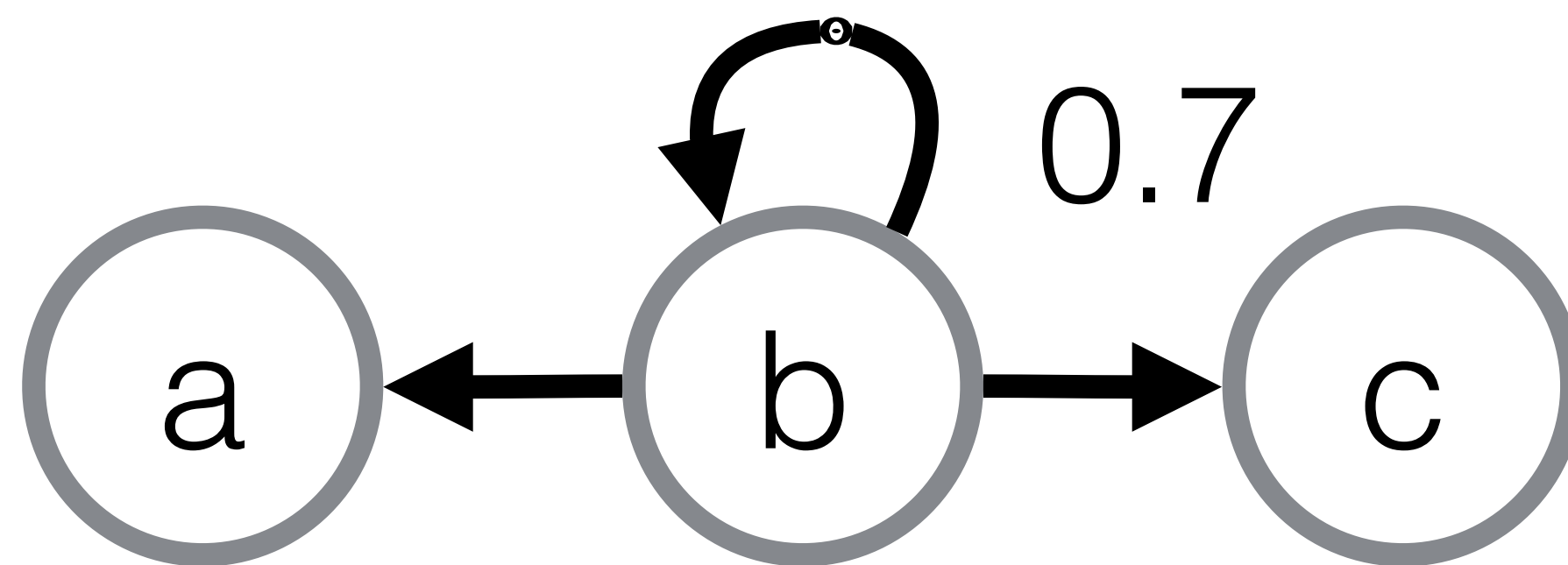
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

# Actions and Transitions

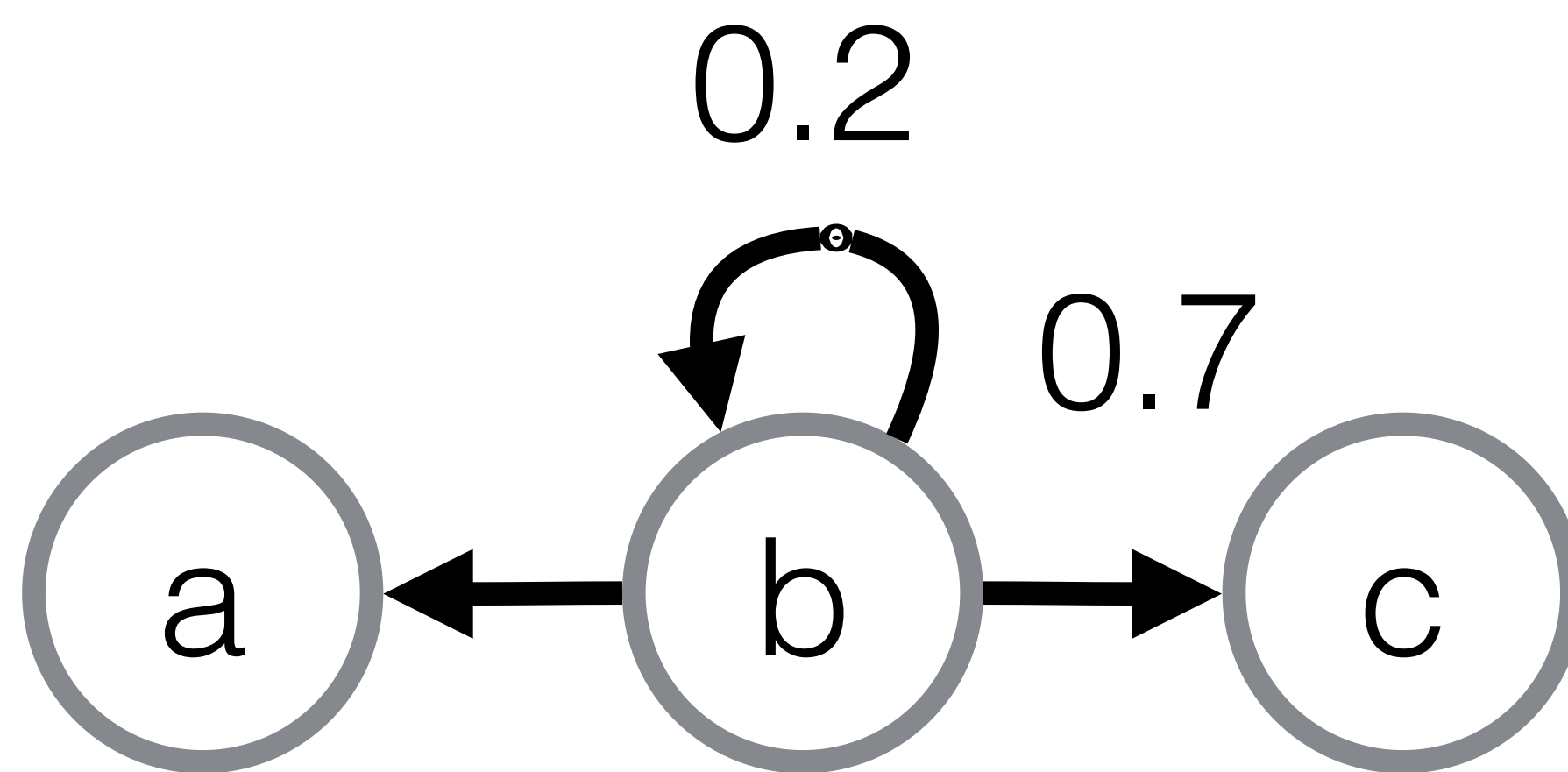
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

# Actions and Transitions

- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$

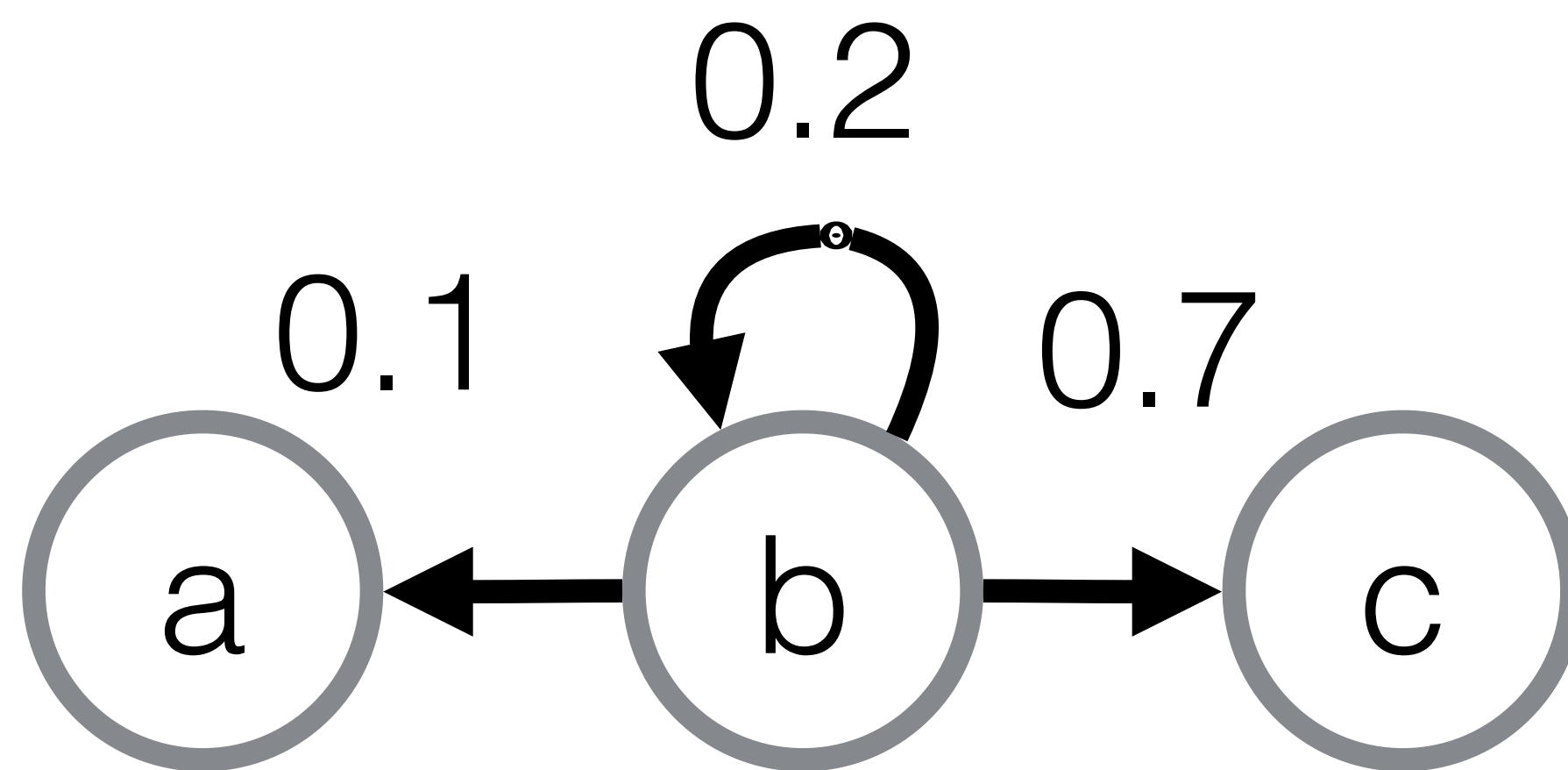


$\mathbf{a}$  = right



# Actions and Transitions

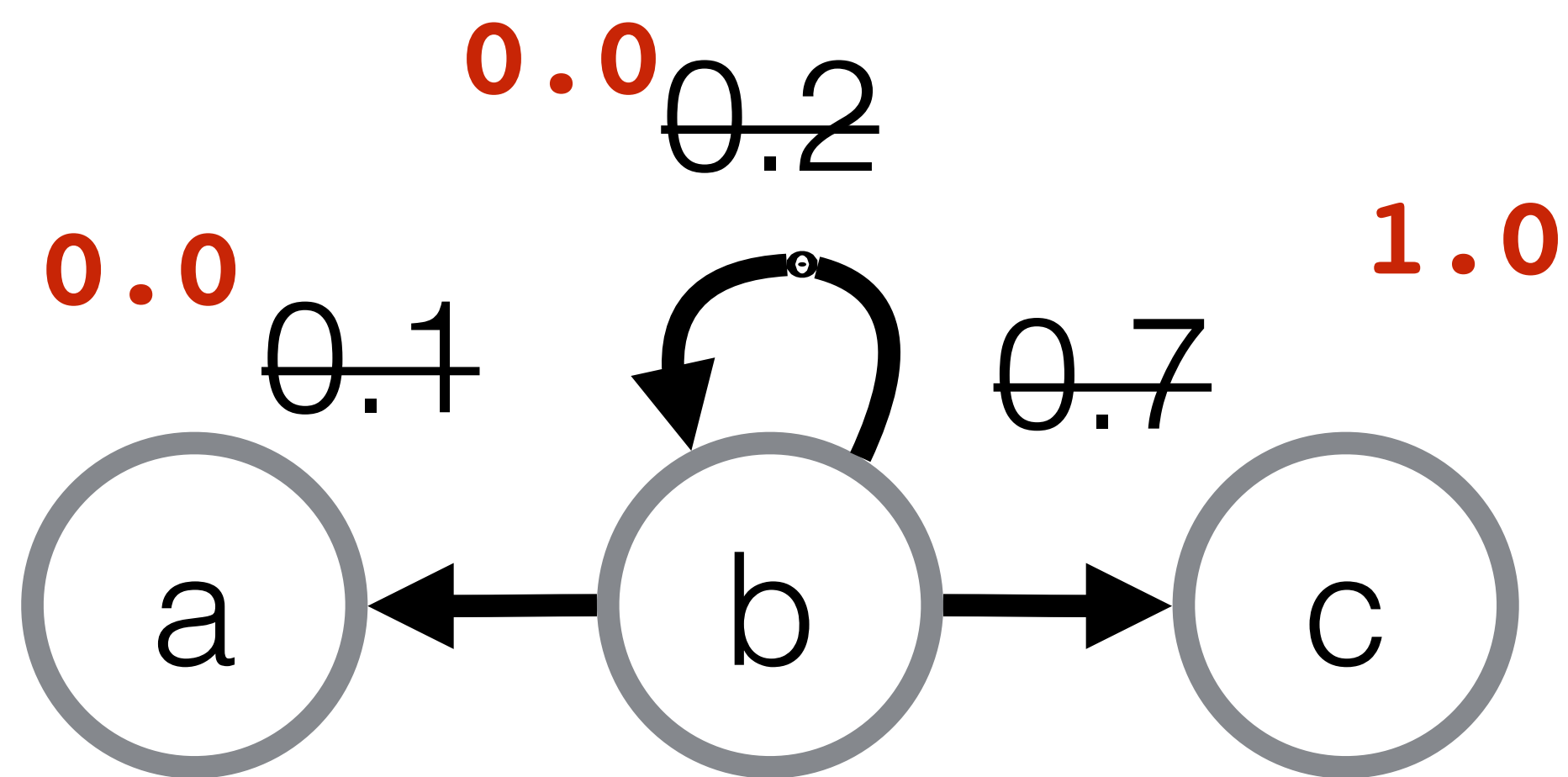
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

# Actions and Transitions

- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = right

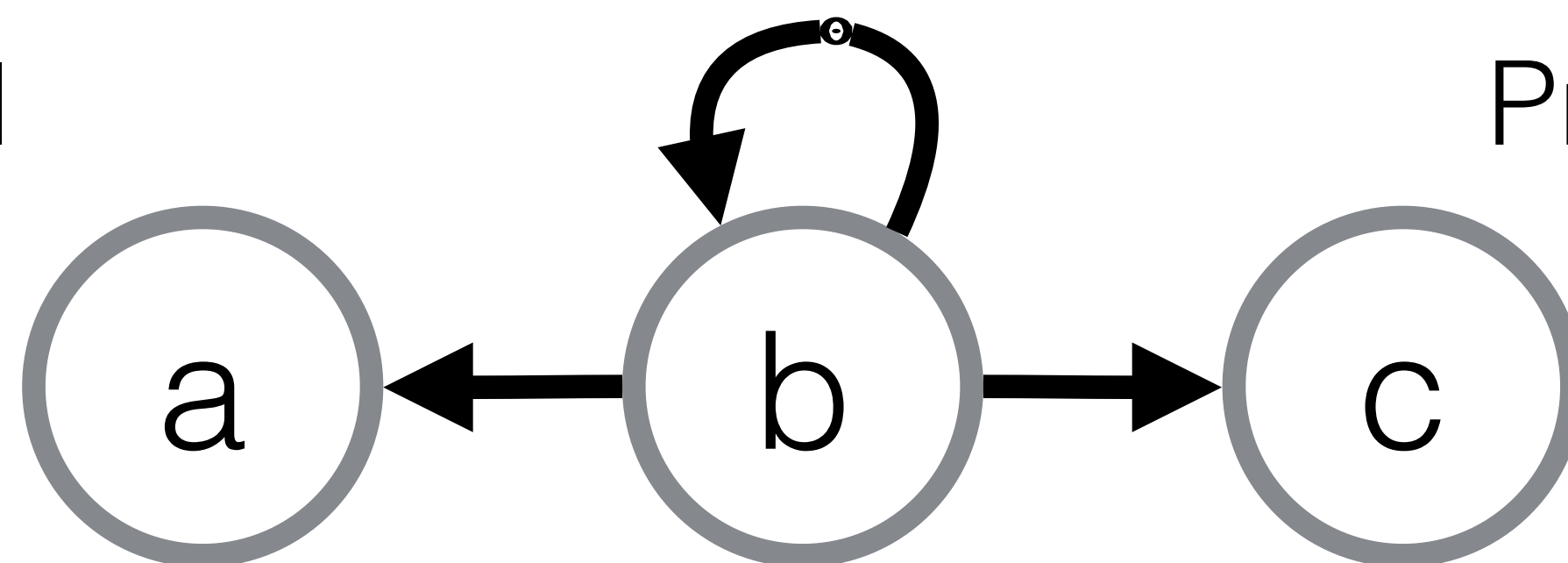
# Actions and Transitions

- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$

$$\Pr(b \mid b, \text{right}) = 0.2$$

$$\Pr(a \mid b, \text{right}) = 0.1$$

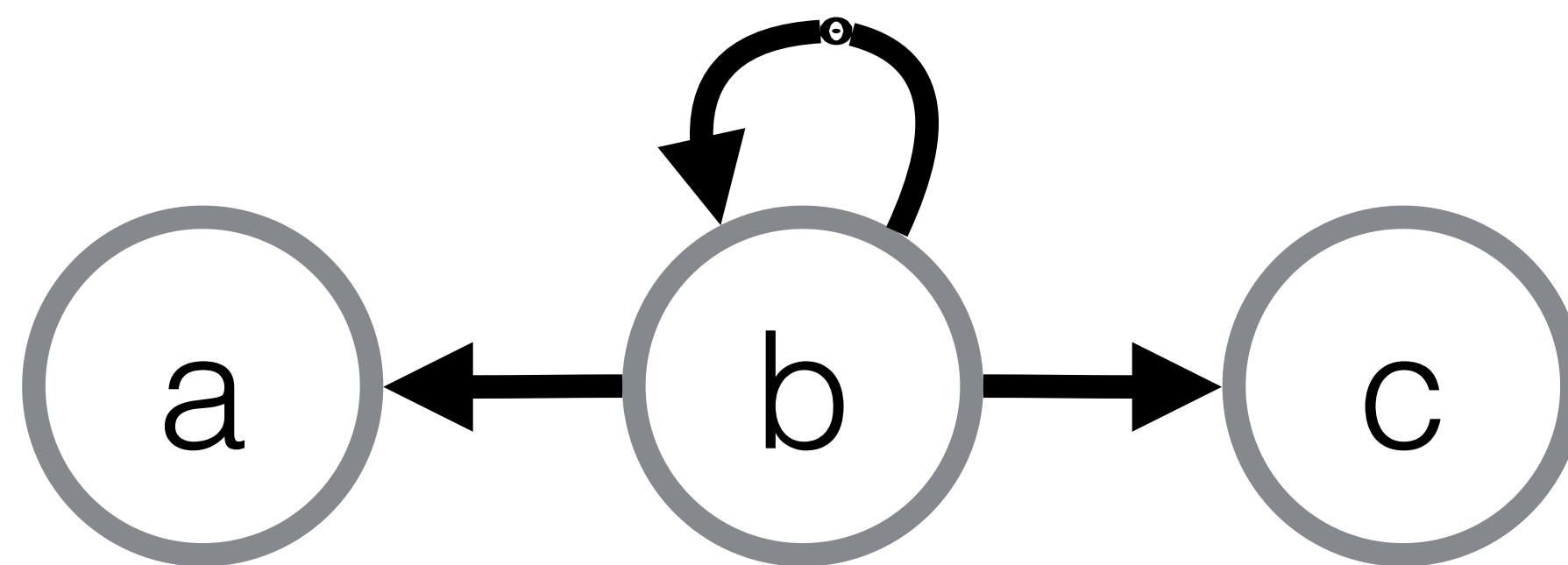
$$\Pr(c \mid b, \text{right}) = 0.7$$



$\mathbf{a} = \text{right}$

# Actions and Transitions

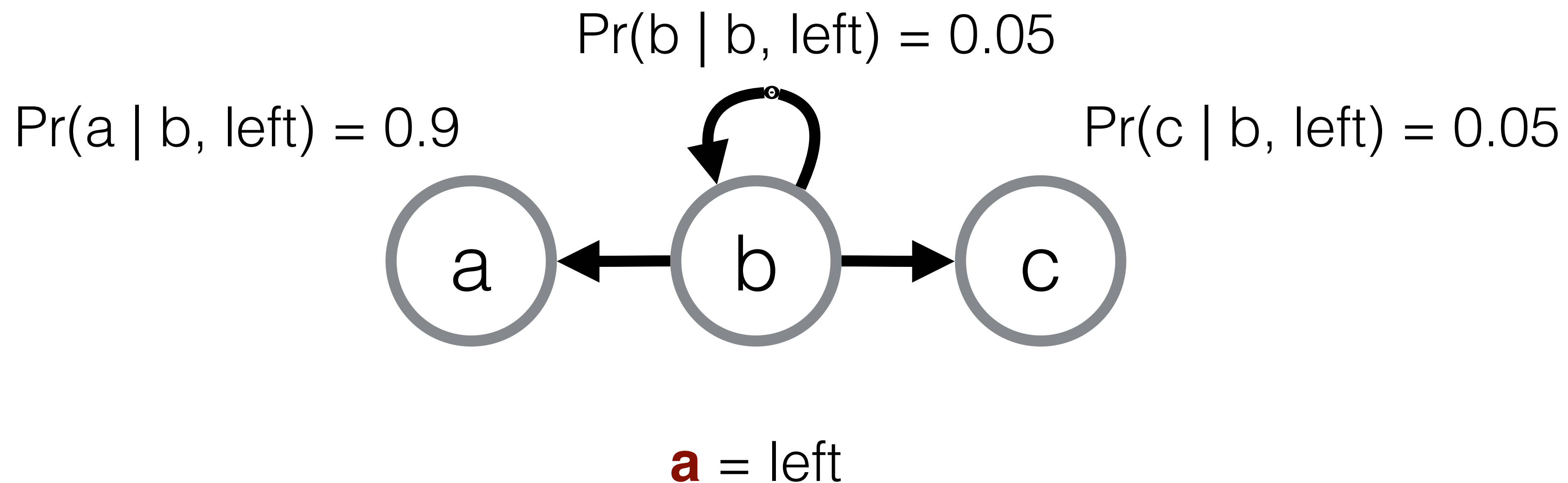
- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$



$\mathbf{a}$  = left

# Actions and Transitions

- $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$ 
  - Probability we **transition** to  $\mathbf{s}'$  if we choose **action**  $\mathbf{a}$  in state  $\mathbf{s}$

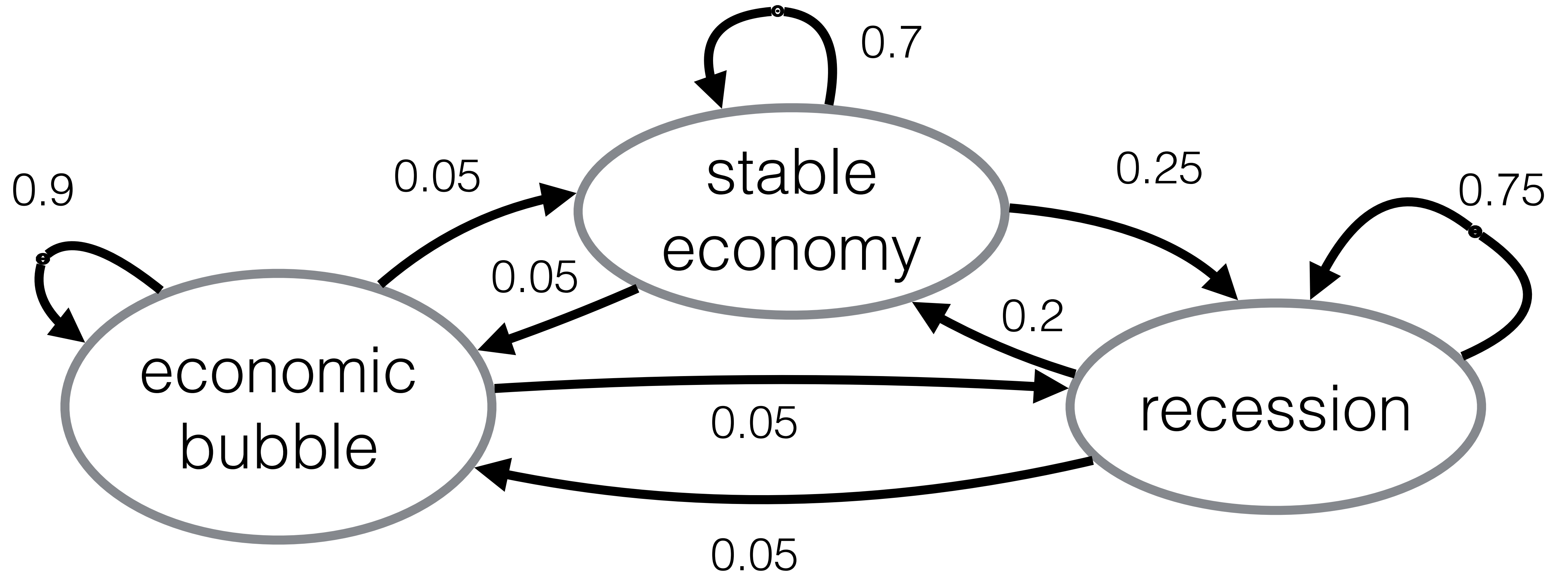


# Preview: Markov Models

Markov Decision Process:  $\Pr(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$

Markov Process  $\Pr(\mathbf{s}' \mid \mathbf{s})$

# Preview: Markov Models



Markov Process  $\Pr(\mathbf{s}' | \mathbf{s})$

# Reward function $R(s)$





# Policy $\pi(s)$

	-100	-100	-100	-100	-100	
1						10
	-100	-100	-100	-100	-100	

# Policy $\pi(s)$

	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	
<b>right</b> 1	<b>right</b>	<b>right</b>	<b>right</b>	<b>right</b>	<b>right</b>	<b>stay</b> 10
	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	

# Policy $\pi(s)$

	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	<b>down</b> -100	
<b>stay</b> 1	<b>left</b>	<b>left</b>	<b>right</b>	<b>right</b>	<b>right</b>	<b>stay</b> 10
	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	<b>up</b> -100	

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T R(s_t)$$

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T R(s_t)$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T \gamma^t R(s_t)$$

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T R(s_t)$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^T \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$



# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0 = s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

# How Good is a Policy?

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0 = s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

$$\pi_s^* = \arg \max_{\pi} U^\pi(s)$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0 = s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

$$\pi_s^* = \arg \max_{\pi} U^\pi(s)$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0 = s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t) \right]$$

$$\pi_s^* = \arg \max_{\pi} U^\pi(s) = \pi_{s'}^* \text{ for any } s'$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0=s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

$$\pi_s^* = \arg \max_{\pi} U^\pi(s) = \pi_{s'}^* \text{ for any } s'$$

$$U(s) = U^{\pi^*}(s)$$

$$U([s_0, s_1, \dots, s_T]) = \sum_{t=0}^{\infty} \gamma^t R(s_t) \quad \gamma \in (0, 1]$$

$$U^\pi(s) = \mathbb{E}_{\text{Pr}([s_0, s_1, \dots] | s_0=s, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

$$\pi_s^* = \arg \max_{\pi} U^\pi(s) = \pi_{s'}^* \text{ for any } s'$$

$$U(s) = U^{\pi^*}(s)$$

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s' | s, a) U(s')$$

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$U(s')$  = expected utility given optimal play from  $s'$



$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$U(s')$  = expected utility given optimal play from  $s'$

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$U(s')$  = expected utility given optimal play from  $s'$

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

**Bellman equation**

# Value Iteration

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

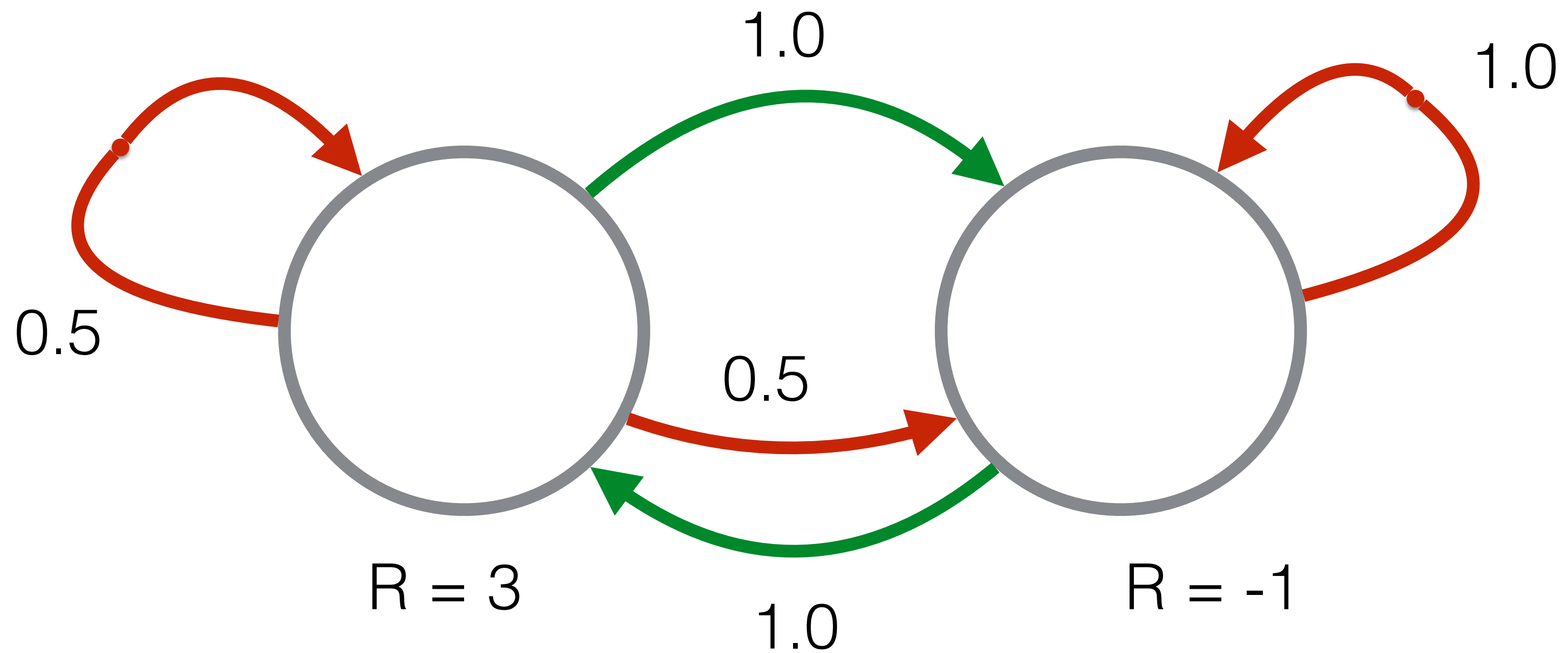
# Value Iteration

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

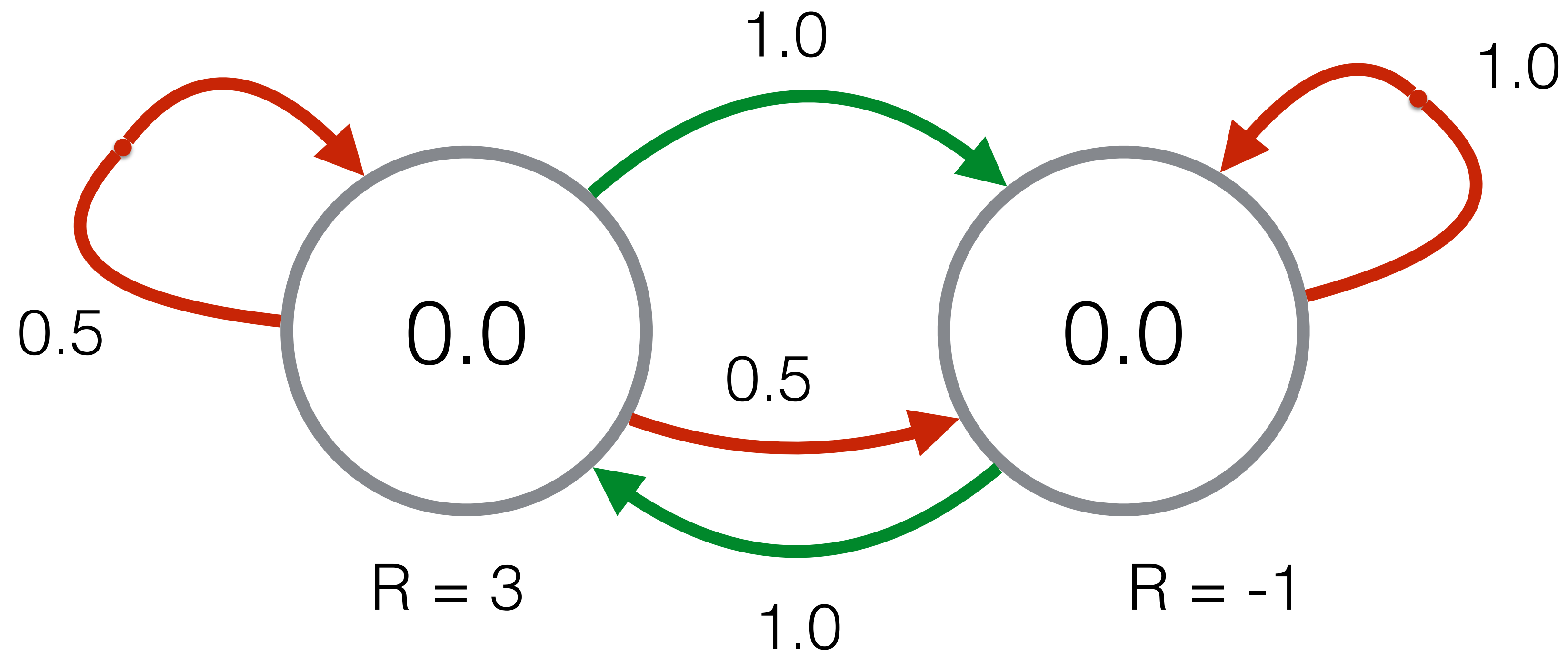
**Bellman equation**

# Value Iteration Example

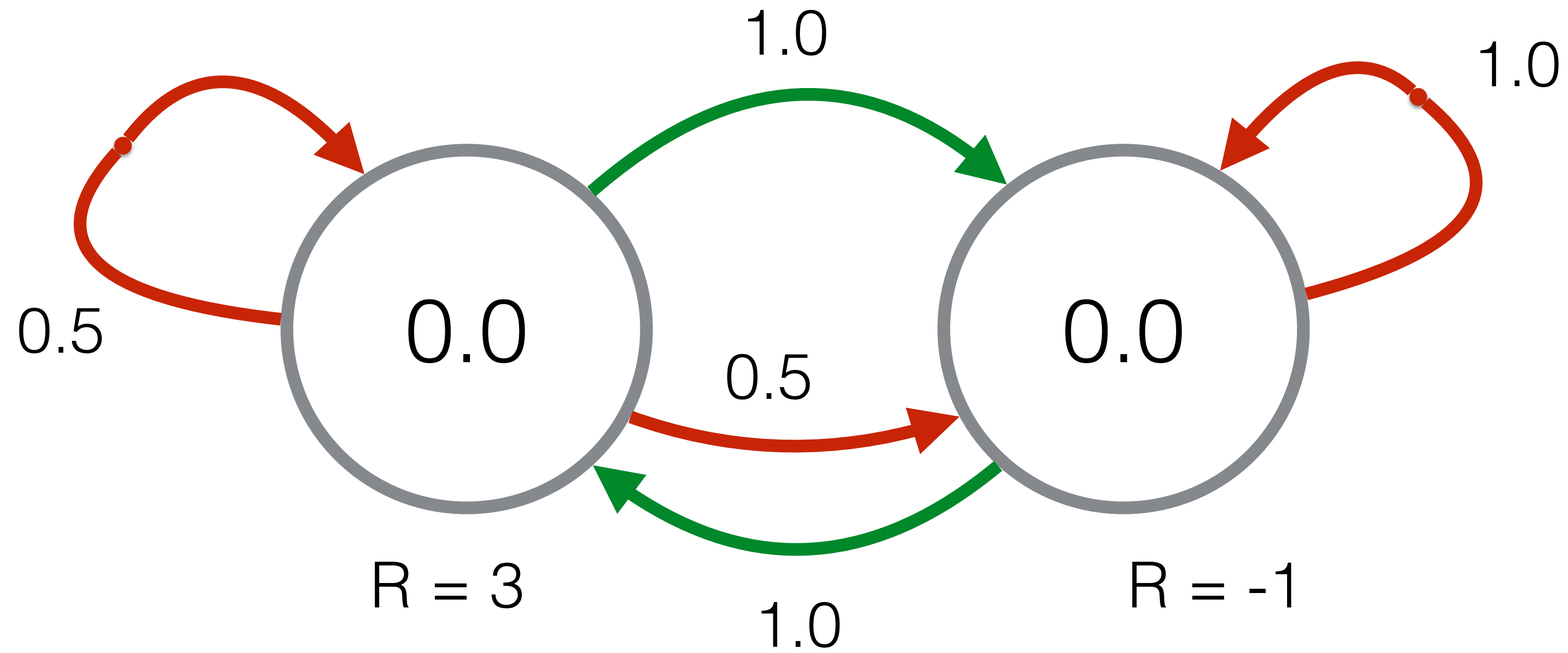


# Value Iteration Example

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$

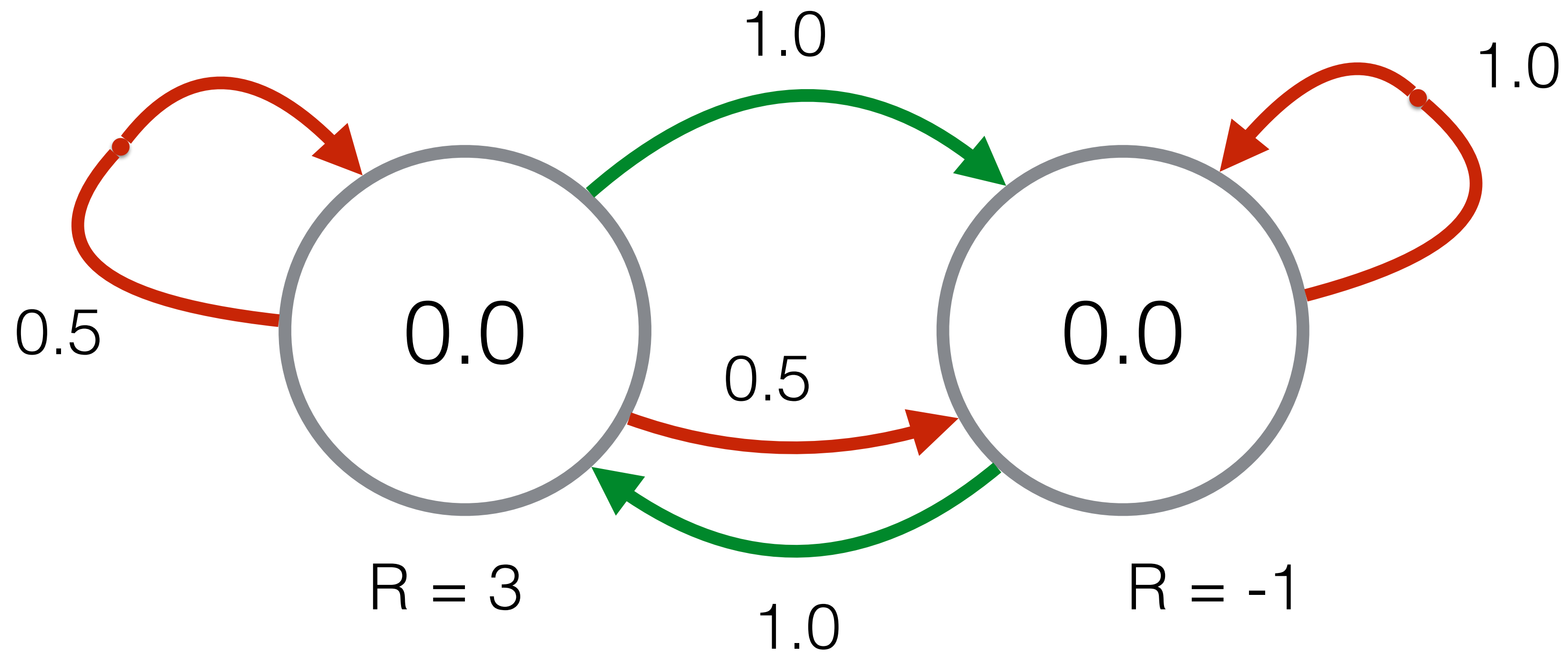


$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

$$\gamma = 0.5$$

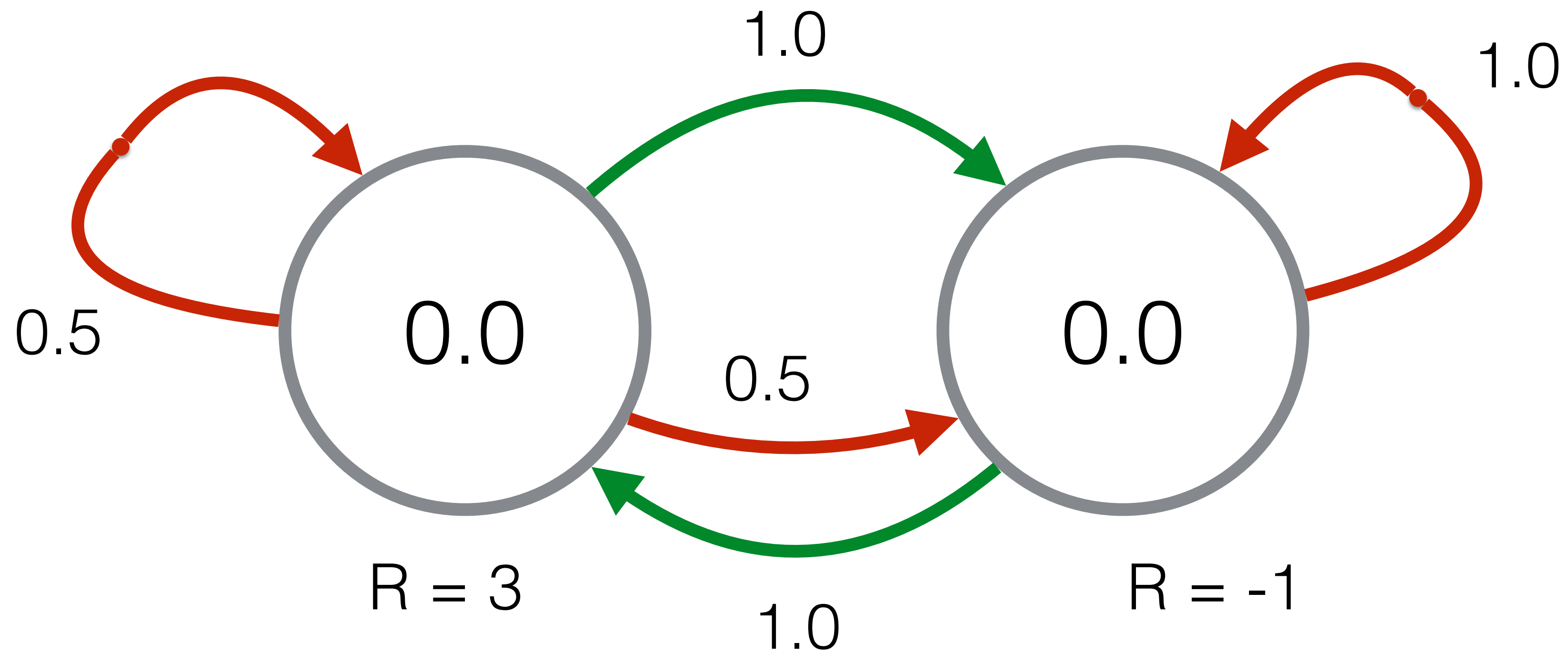


$$3 + 0.5 \max\{ 1.0 * 0.0, 0.5 * 0.0 + 0.5 * 0.0 \} = 3$$



$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

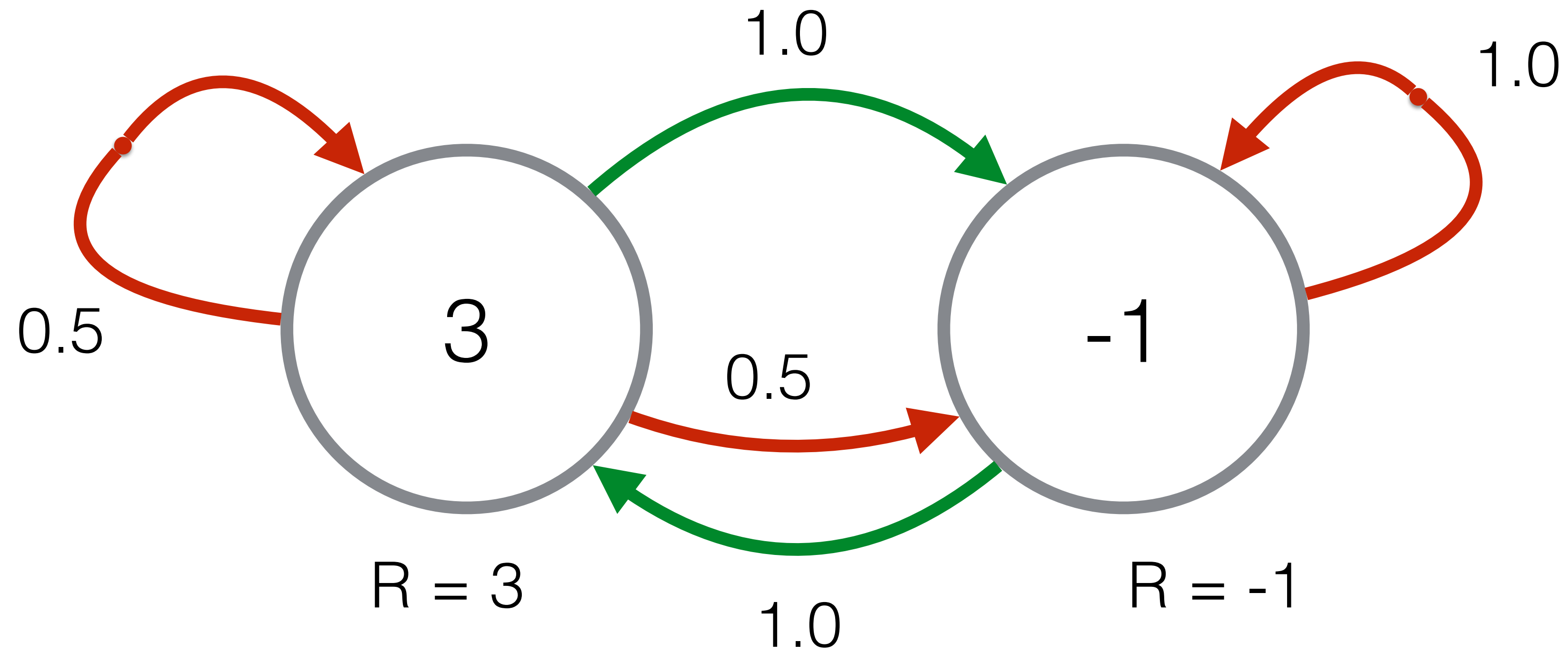
$$\gamma = 0.5$$



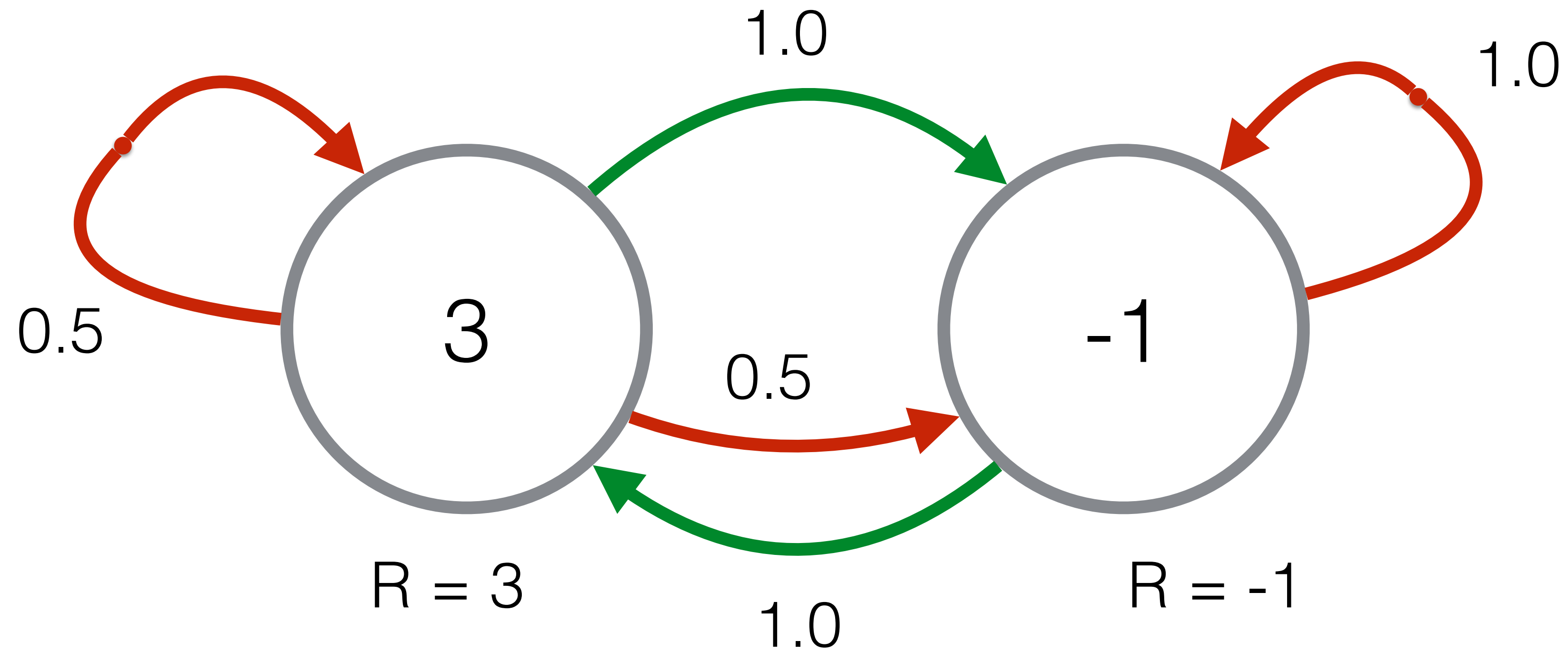
$$3 + 0.5 \max\{ 1.0 * 0.0, 0.5 * 0.0 + 0.5 * 0.0 \} = 3$$

$$-1 + 0.5 \max\{ 1.0 * 0.0, 1.0 * 0.0 \} = -1$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$

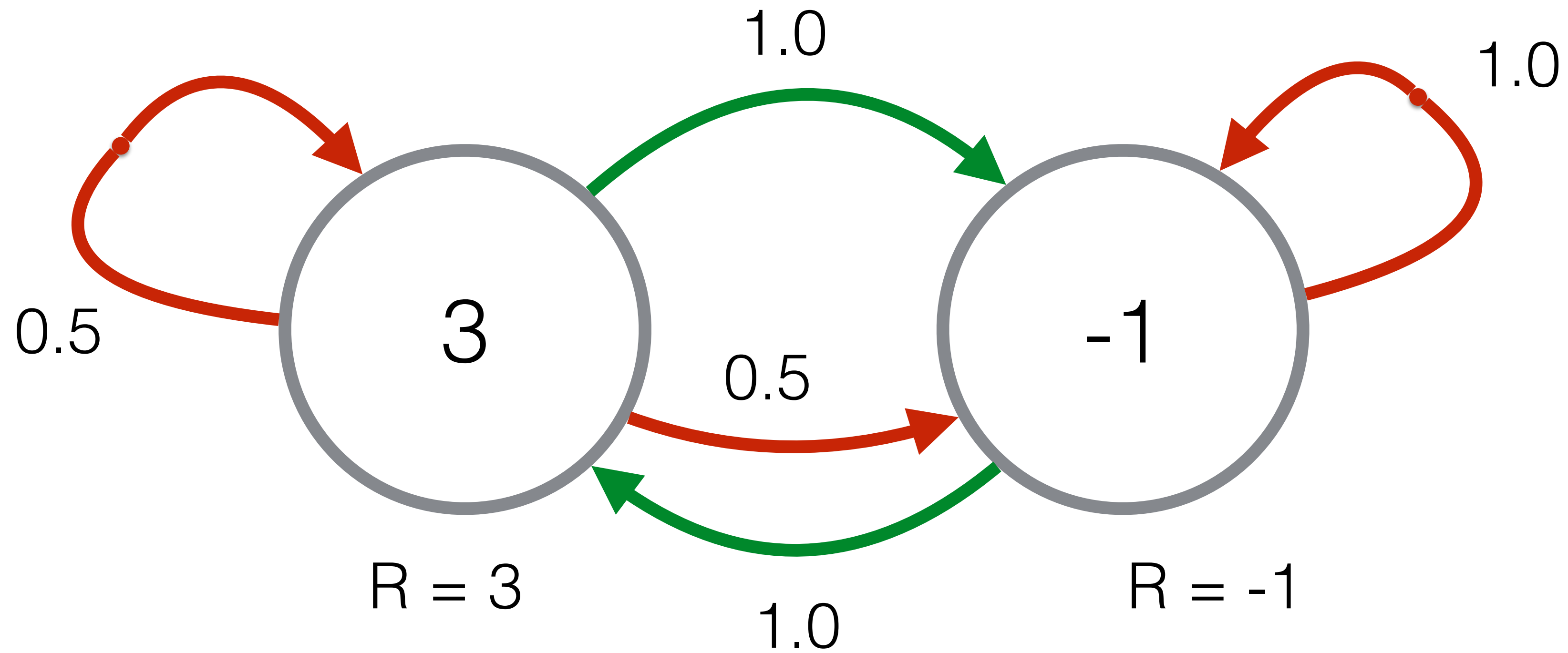


$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$3 + 0.5 \max\{ 1.0 * (-1), \quad 0.5 * 3 + 0.5 * (-1) \} = 3 + 0.5 \max\{ -1, 1 \} = 3.5$$

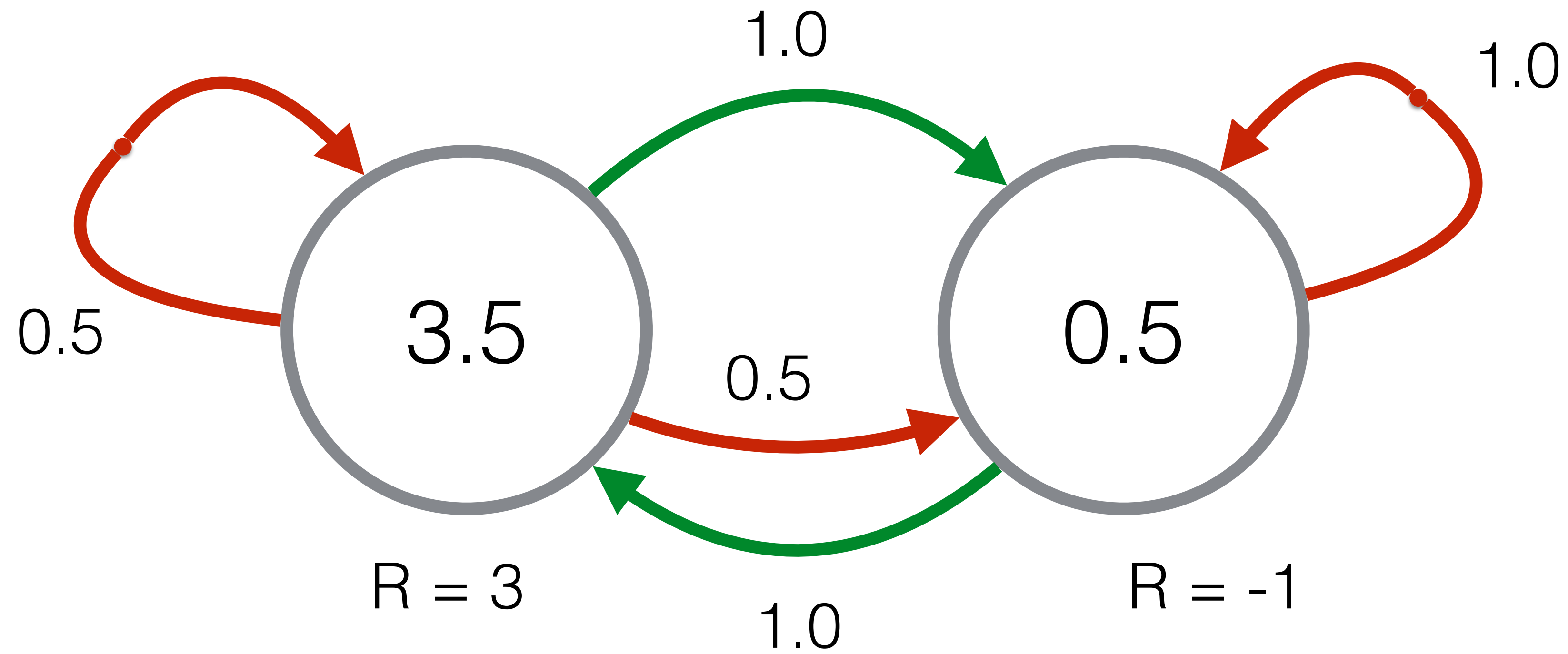
$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



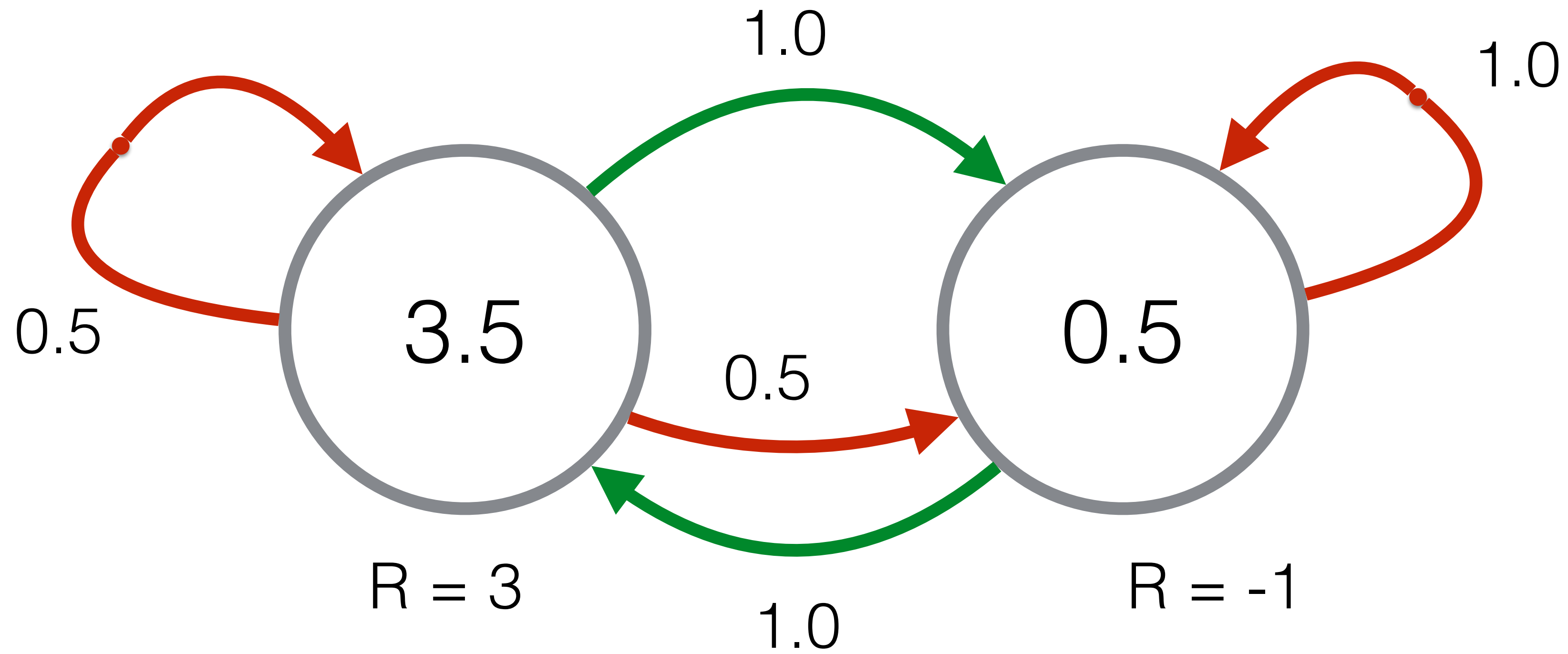
$$3 + 0.5 \max\{ 1.0 * (-1), \quad 0.5 * 3 + 0.5 * (-1) \} = 3 + 0.5 \max\{ -1, 1 \} = 3.5$$

$$-1 + 0.5 \max\{ 1.0 * 3, 1.0 * (-1) \} = 0.5$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$

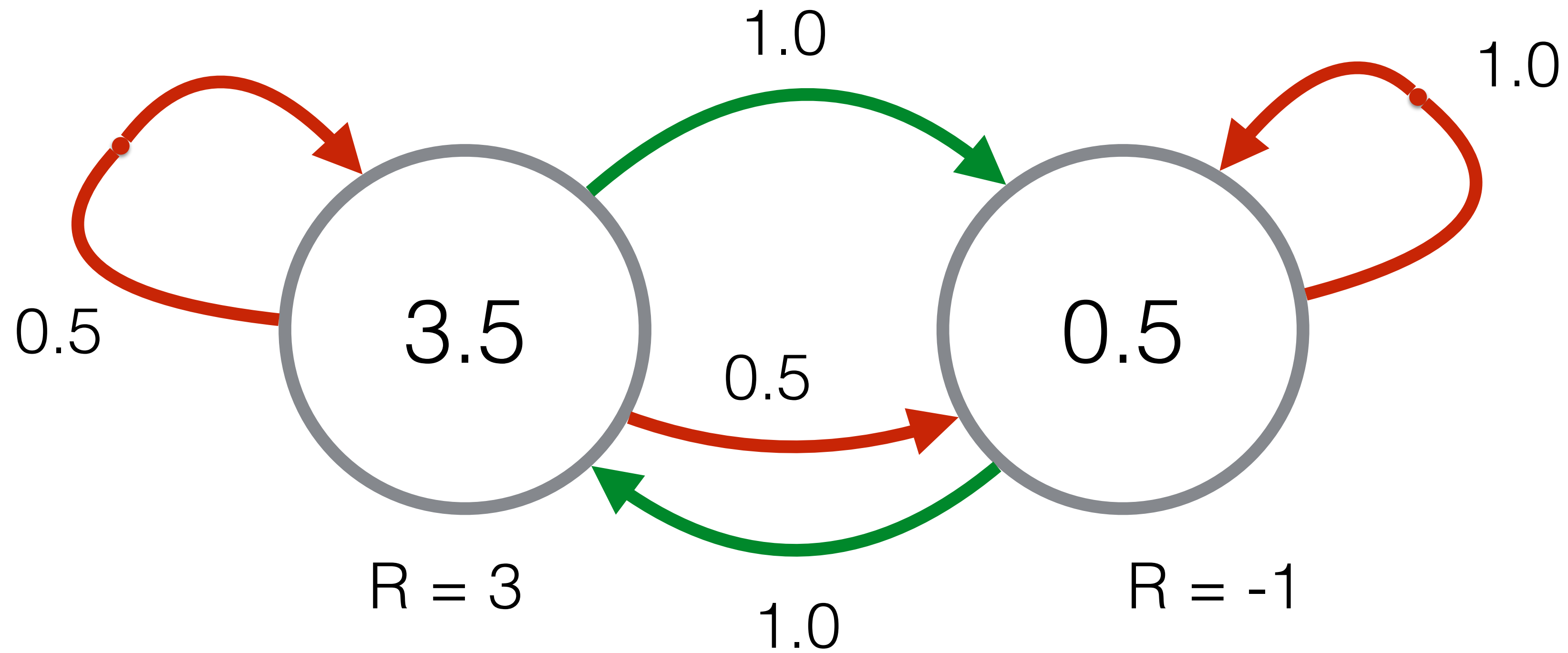


$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$3 + 0.5 \max\{ 1.0 * 0.5, \quad 0.5 * 3.5 + 0.5 * 0.5 \} = 3 + 0.5 \max\{ 0.5, 2 \} = 4$$

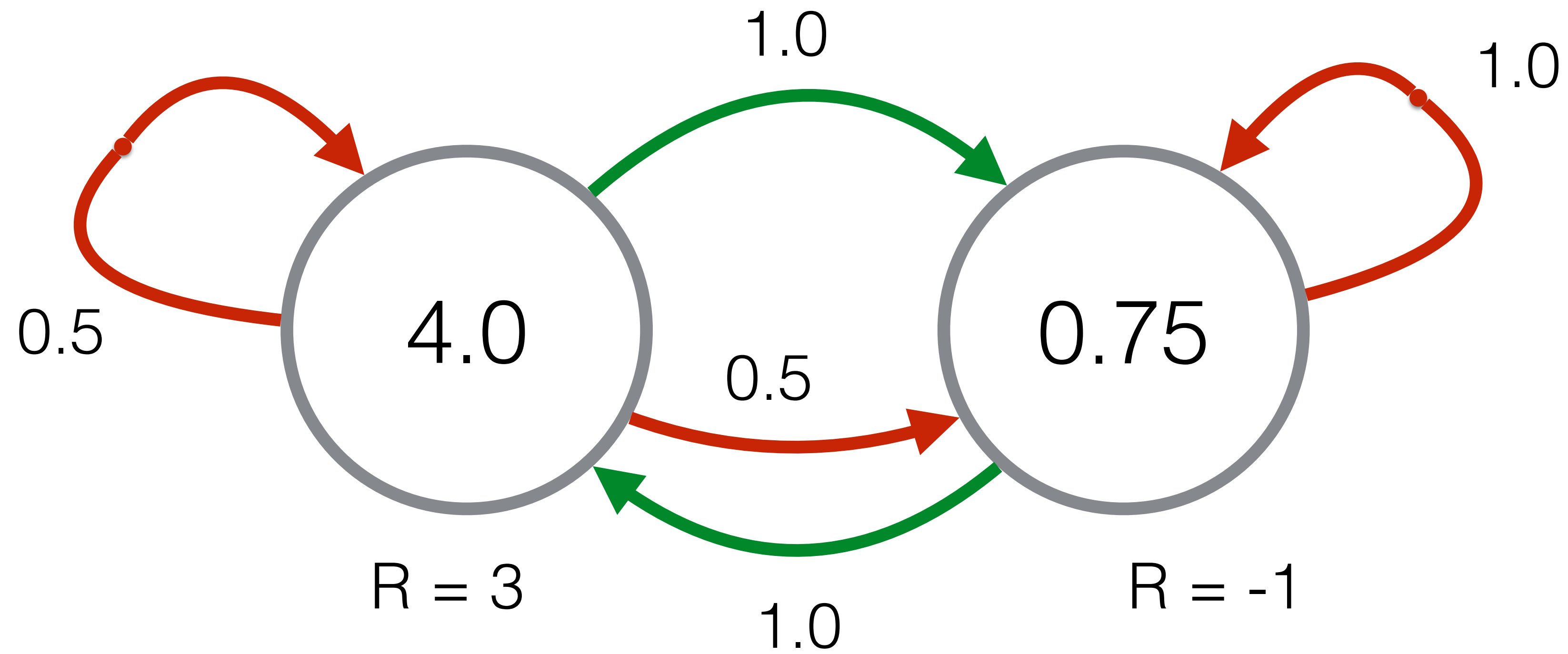
$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$3 + 0.5 \max\{ 1.0 * 0.5, \quad 0.5 * 3.5 + 0.5 * 0.5 \} = 3 + 0.5 \max\{ 0.5, 2 \} = 4$$

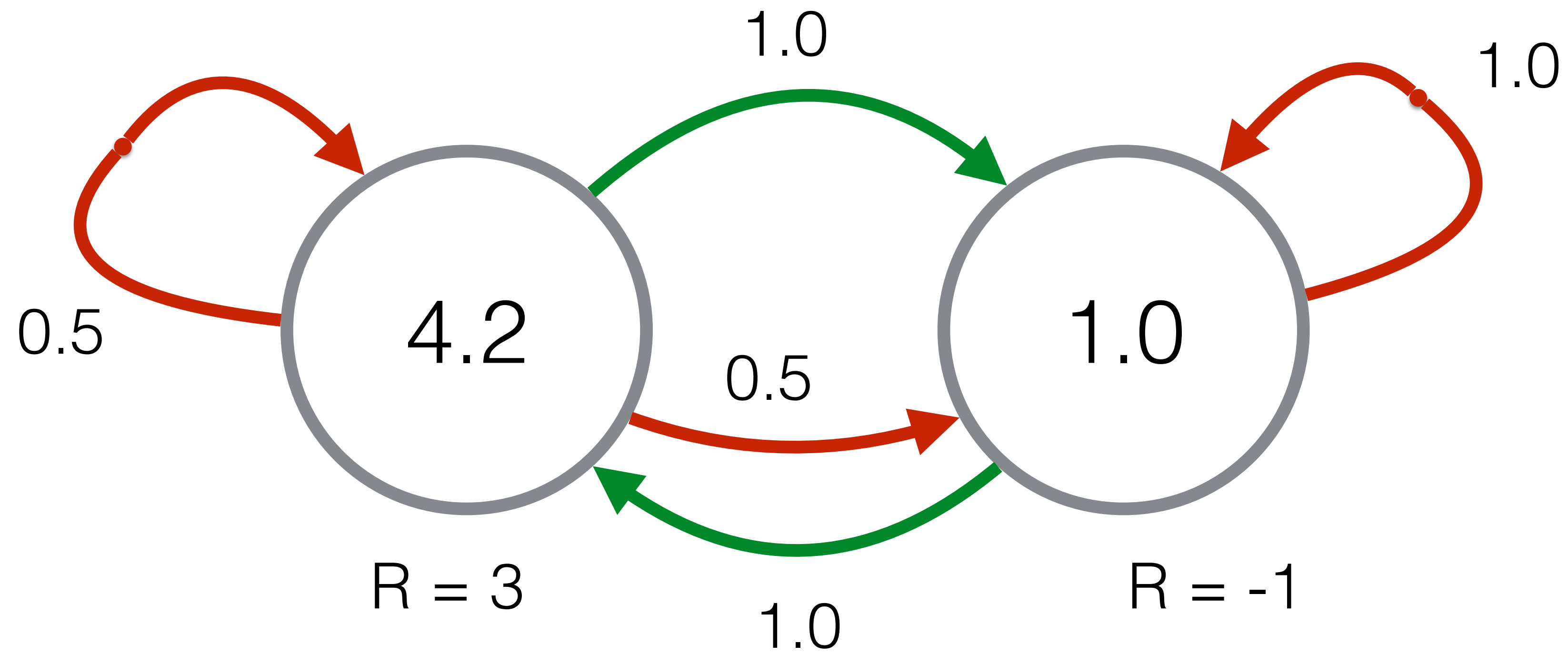
$$-1 + 0.5 \max\{ 1.0 * 3.5, 1.0 * 0.5 \} = 0.75$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$

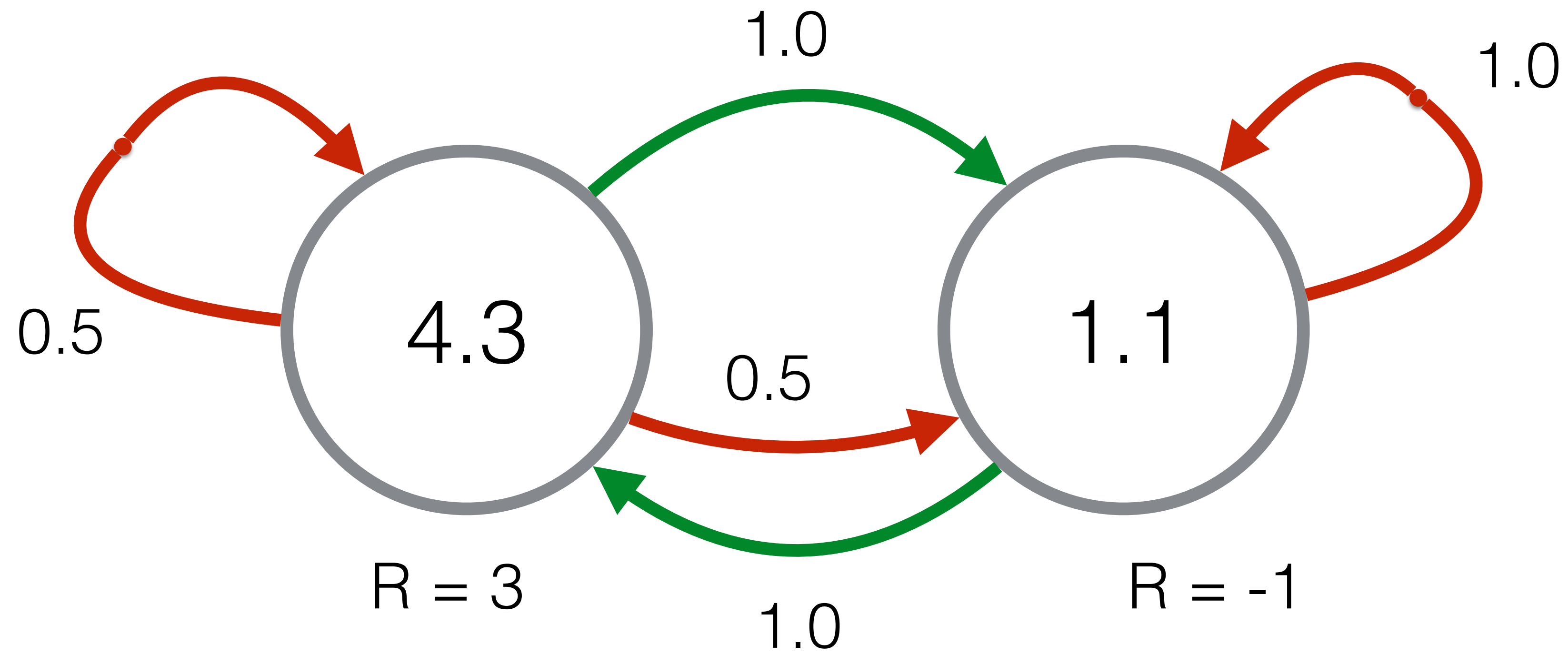




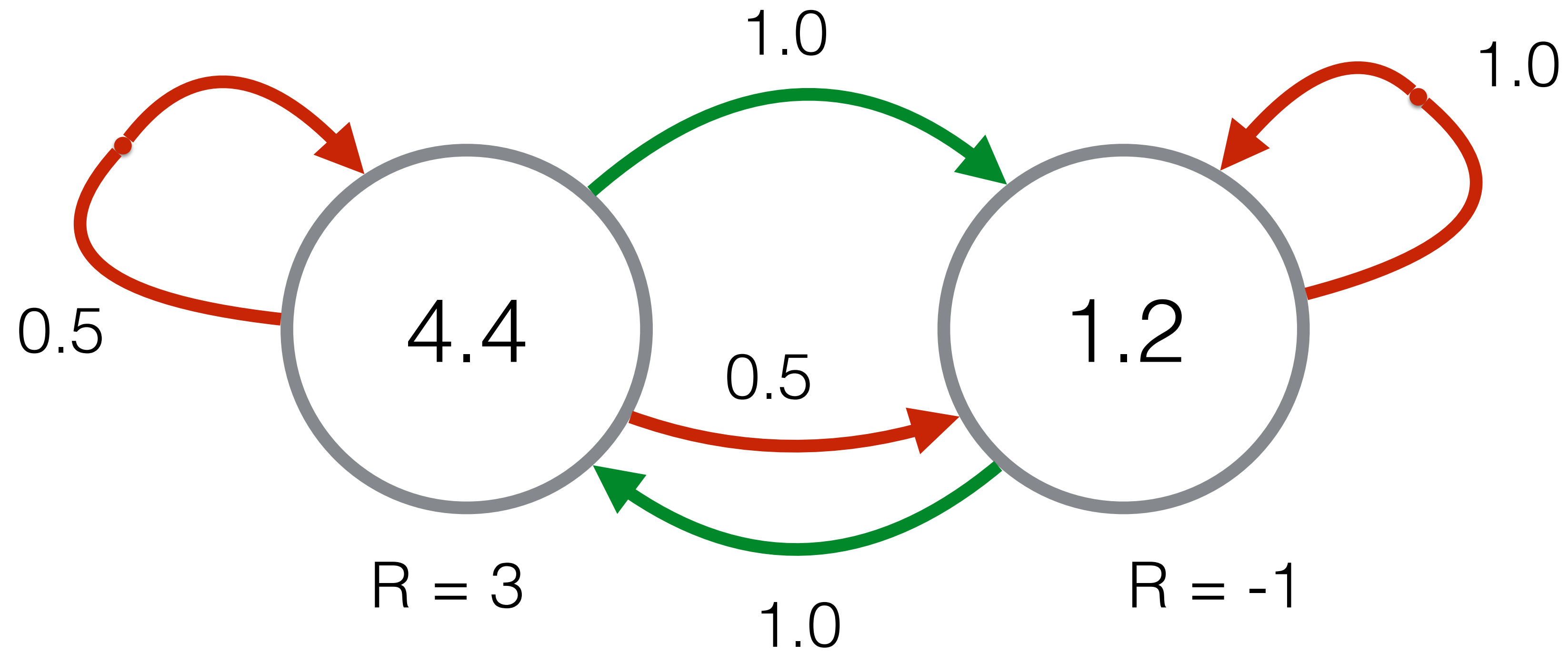
$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$

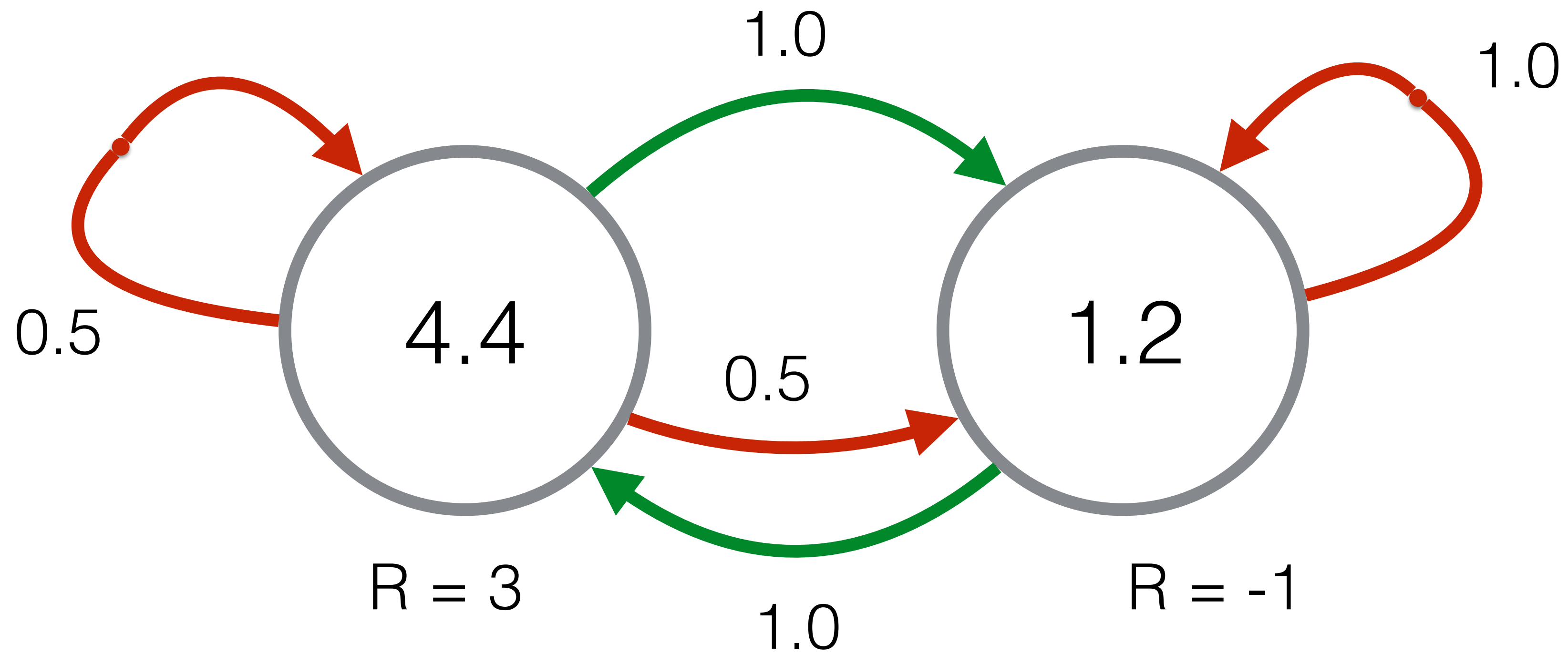


$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

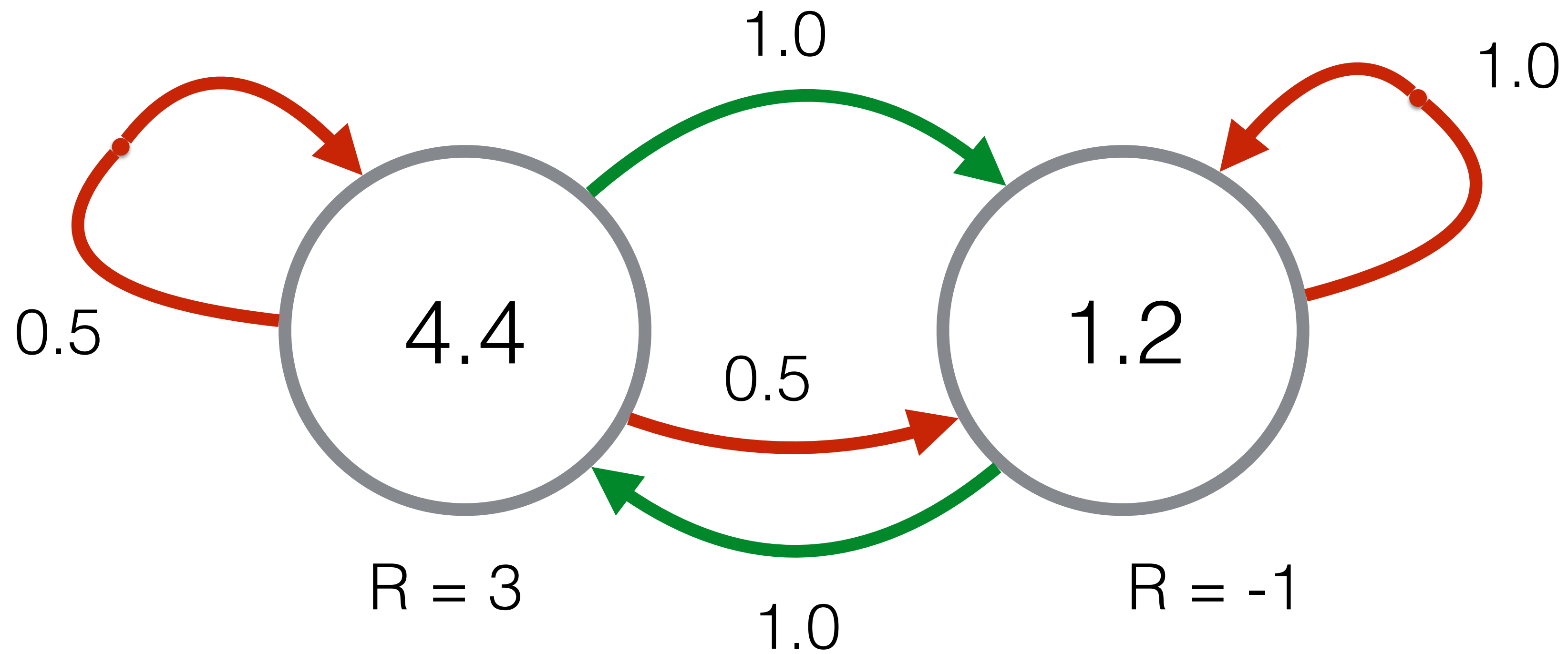
$$\gamma = 0.5$$



$$3 + 0.5 \max \{ 1.0 * 1.2, 0.5 * 4.4 + 0.5 * 1.2 \}$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

$$\gamma = 0.5$$

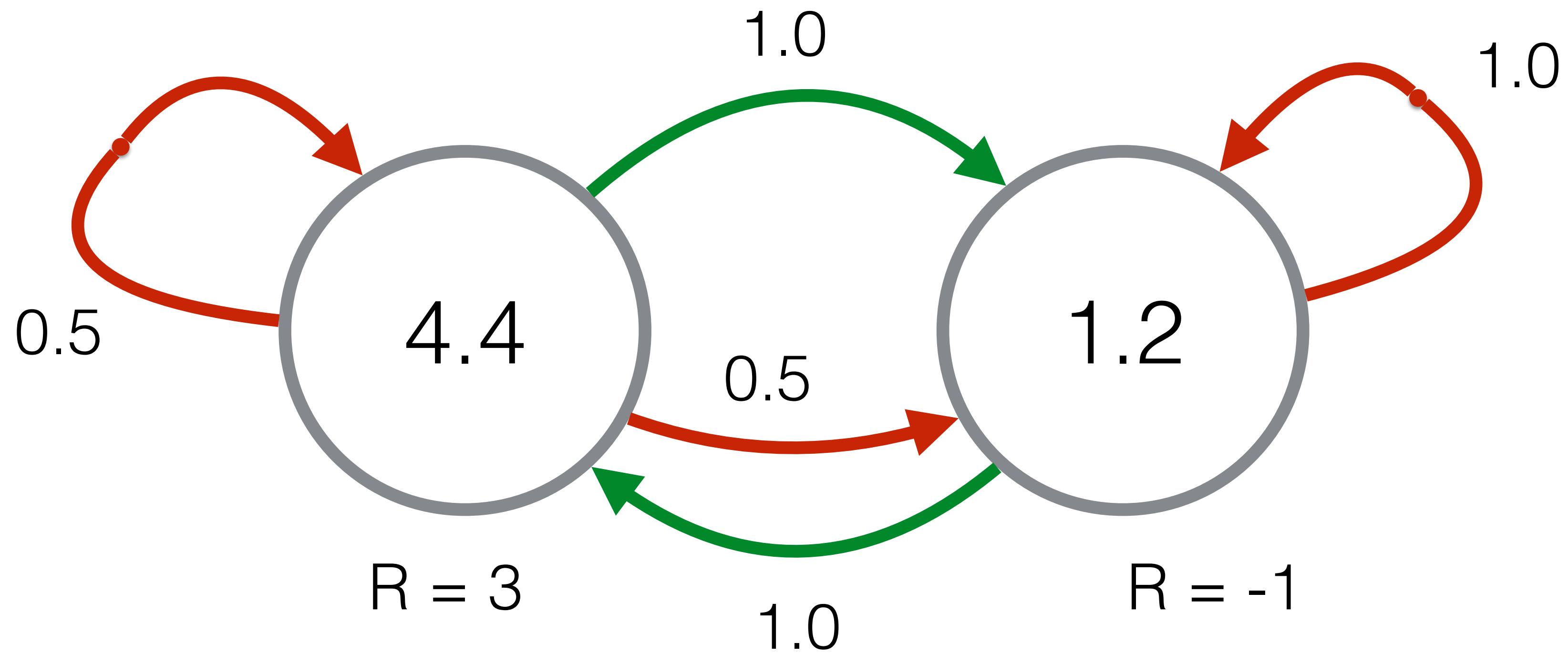


$$3 + 0.5 \max \{ 1.0 * 1.2, \quad 0.5 * 4.4 + 0.5 * 1.2 \}$$

$$3 + 0.5 \max \{ 1.2, \quad 2.2 + 0.6 \}$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$$

$$\gamma = 0.5$$

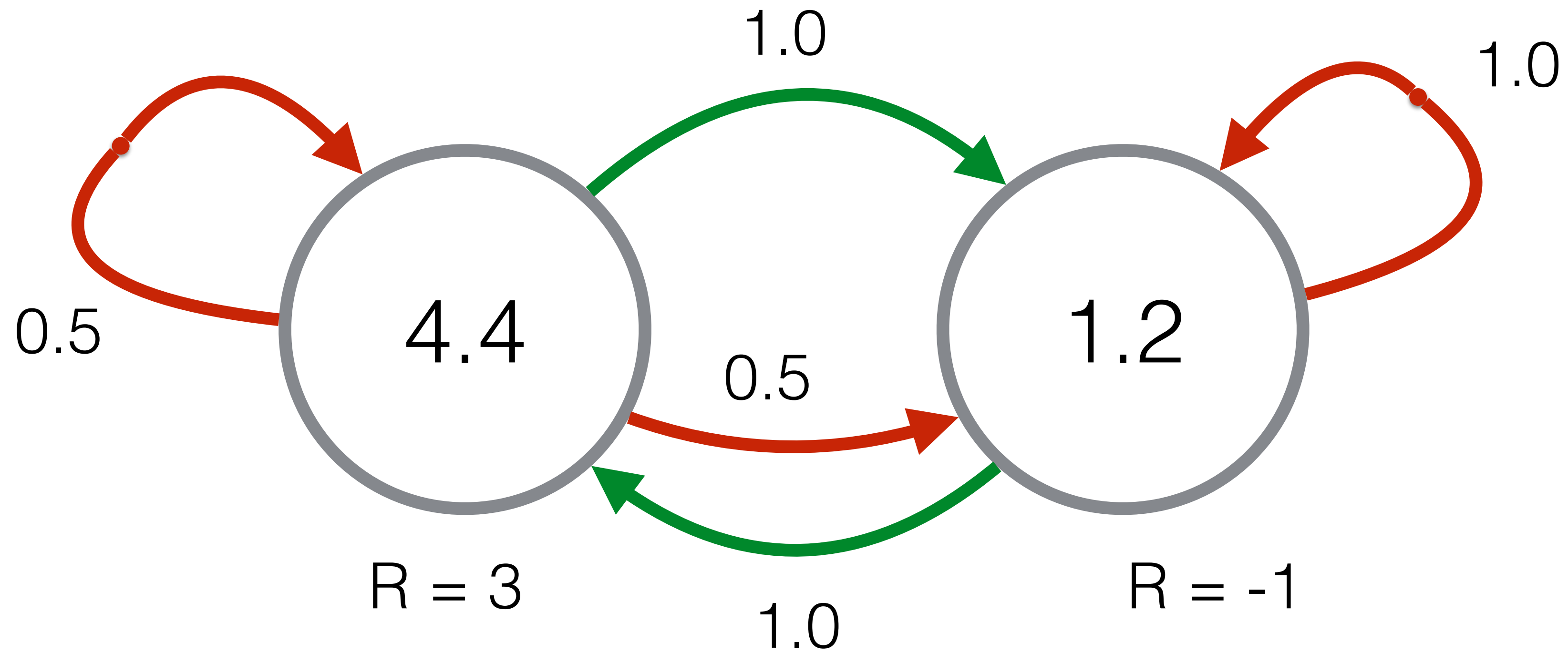


$$3 + 0.5 \max \{ 1.0 * 1.2, \quad 0.5 * 4.4 + 0.5 * 1.2 \}$$

$$3 + 0.5 \max \{ 1.2, \quad 2.2 + 0.6 \}$$

$$3 + 0.5 \max \{ 1.2, \quad 2.8 \}$$

$$U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s') \quad \gamma = 0.5$$



$$3 + 0.5 \max \{ 1.0 * 1.2, \quad 0.5 * 4.4 + 0.5 * 1.2 \}$$

$$3 + 0.5 \max \{ 1.2, \quad 2.2 + 0.6 \}$$

$$3 + 0.5 \max \{ 1.2, \quad 2.8 \}$$

$$3 + 0.5 * 2.8 = 4.4$$

# Summary



# Summary

- Markov decision processes

# Summary

- Markov decision processes
  - actions have probabilistic state transitions

# Summary

- Markov decision processes
  - actions have probabilistic state transitions
- Discounted reward function

# Summary

- Markov decision processes
  - actions have probabilistic state transitions
- Discounted reward function
- Optimal policy maximizes expected reward

# Summary

- Markov decision processes
  - actions have probabilistic state transitions
- Discounted reward function
- Optimal policy maximizes expected reward
- Value iteration

# Summary

- Markov decision processes
  - actions have probabilistic state transitions
- Discounted reward function
- Optimal policy maximizes expected reward
- Value iteration
- Chapter 17 to end of 17.2