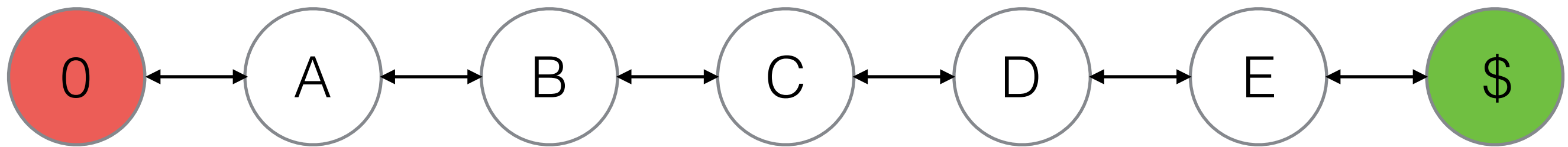
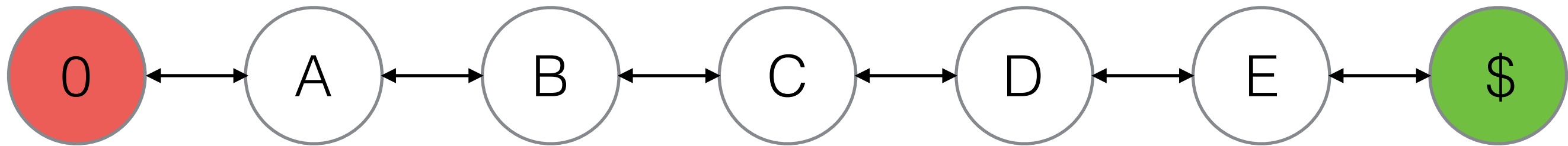


Random Walk Problem



Random Walk Problem



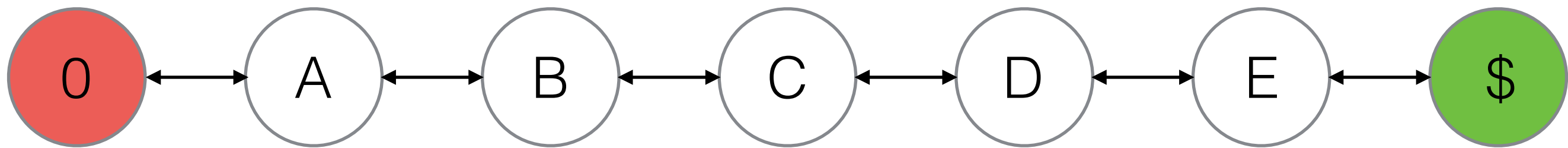
$$S := \{A, B, C, D, E\}$$

$$A := \{Left, Right\}$$

$$\phi := \{0, \$\}$$

$$R = \{-1, 0, 1\}$$

Random Walk Problem



$$S := \{A, B, C, D, E\}$$

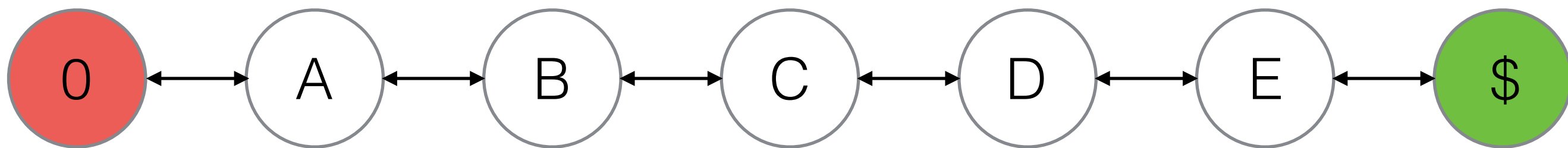
$$A := \{Left, Right\}$$

$$\phi := \{0, \$\}$$

$$R = \{-1, 0, 1\}$$

What about the transition probabilities?

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$\pi = \epsilon - Greedy$$

S_t	C					
R_t	0					
A_t	r					

Initialize $Q(s, a)$ arbitrarily, $\forall s \in \mathcal{S}, a \in \mathcal{A}$

Initialize π to be ϵ -greedy with respect to Q , or to a fixed given policy

Parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$, a positive integer n

All store and access operations (for S_t , A_t , and R_t) can take their index mod n

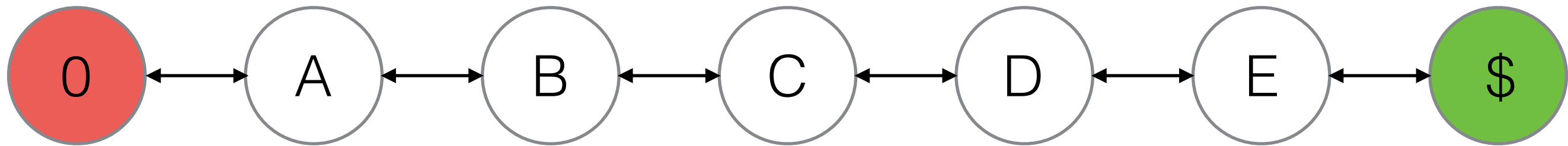
Repeat (for each episode):

Initialize and store $S_0 \neq \text{terminal}$

Select and store an action $A_0 \sim \pi(\cdot | S_0)$

$T \leftarrow \infty$

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 0$$

S_t	C	D				
R_t	0	0				
A_t	r	l				

For $t = 0, 1, 2, \dots$:

 If $t < T$, then:

 Take action A_t

 Observe and store the next reward as R_{t+1} and the next state as S_{t+1}

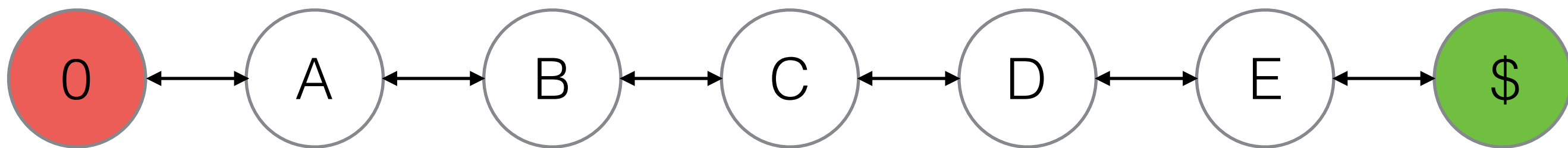
 If S_{t+1} is terminal, then:

$T \leftarrow t + 1$

 else:

 Select and store an action $A_{t+1} \sim \pi(\cdot | S_{t+1})$

Random Walk Problem



$$T = \infty$$

$$n = 2$$

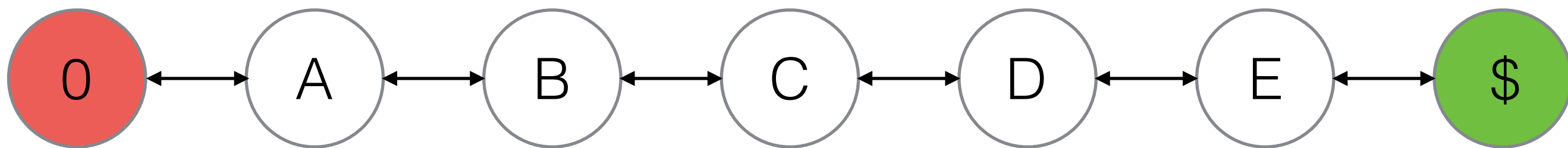
$$t = 0$$

$$\text{tau} = -1$$

St	C	D				
Rt	0	0				
At	r	l				

```
|  $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
| If  $\tau \geq 0$ :
|    $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|   If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|    $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
```

Random Walk Problem



$$T = \infty$$

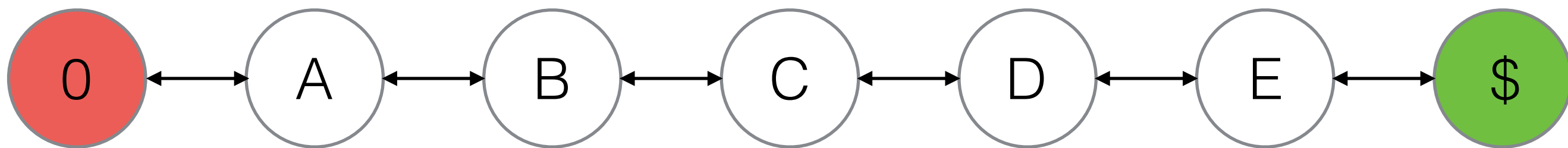
$$n = 2$$

$$t = 1$$

S_t	C	D	C			
R_t	0	0	0			
A_t	r	l	r			

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 1$$

$$\text{tau} = 0$$

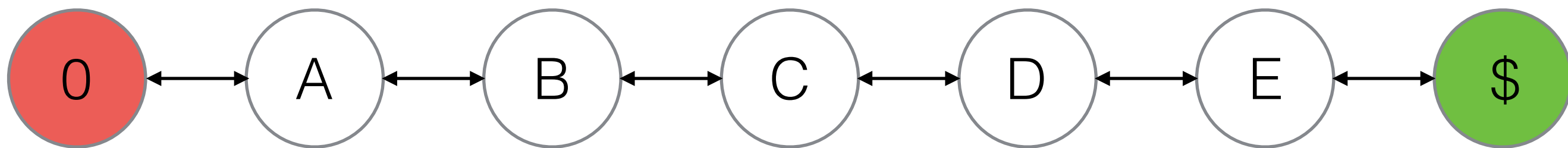
St	C	D	C			
Rt	0	0	0			
At	r	l	r			

$$G = \gamma^0 R_1 + \gamma^1 R_2$$

```

|   $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|  If  $\tau \geq 0$ :
|       $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|      If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|       $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```


Random Walk Problem



$$T = \infty$$

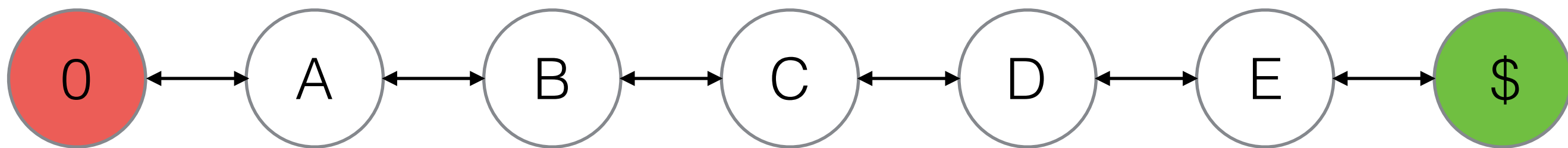
$$n = 2$$

$$t = 2$$

S_t	C	D	C	D		
R_t	0	0	0	0		
A_t	r	l	r	r		

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 2$$

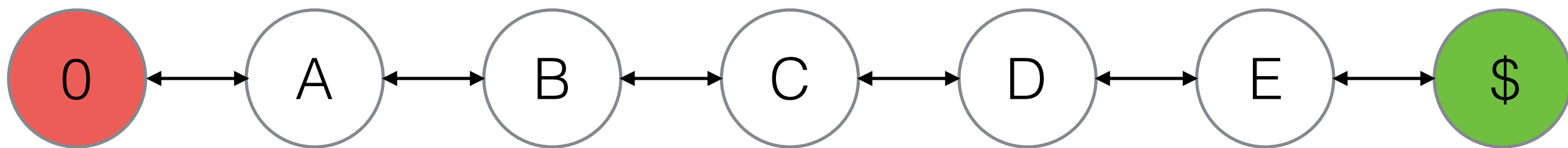
$$\text{tau} = 1$$

St	C	D	C	D		
Rt	0	0	0	0		
At	r	l	r	r		

```

|    $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|   If  $\tau \geq 0$ :
|        $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|       If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|        $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```

Random Walk Problem



$$T = \infty$$

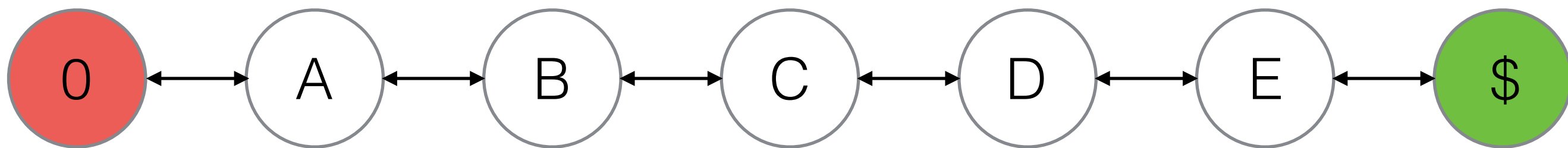
$$n = 2$$

$$t = 3$$

S_t	C	D	C	D	E	
R_t	0	0	0	0	0	
A_t	r	l	r	r	r	

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = \infty$$

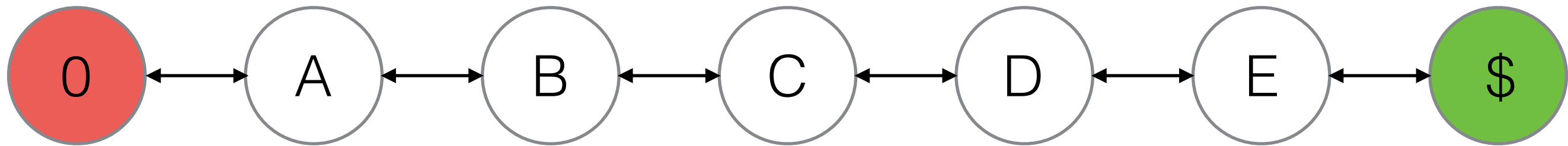
$$n = 2$$

$$t = 4$$

S_t	C	D	C	D	E	\$
R_t	0	0	0	0	0	1
A_t	r	l	r	r	r	

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = 5$$

$$n = 2$$

$$t = 4$$

S_t	C	D	C	D	E	\$
R_t	0	0	0	0	0	1
A_t	r	l	r	r	r	

For $t = 0, 1, 2, \dots$:

| If $t < T$, then:

| Take action A_t

| Observe and store the next reward as R_{t+1} and the next state as S_{t+1}

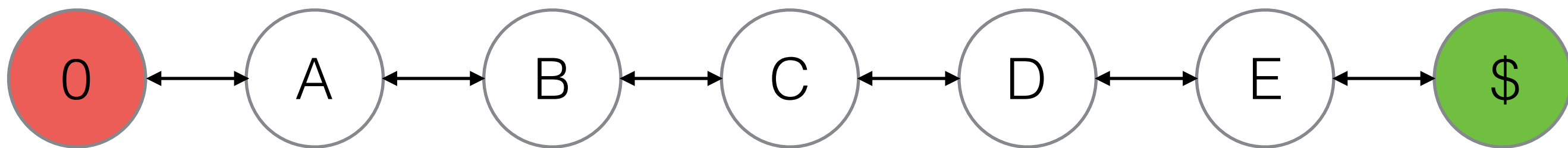
| If S_{t+1} is terminal, then:

| $T \leftarrow t + 1$

| else:

| Select and store an action $A_{t+1} \sim \pi(\cdot | S_{t+1})$

Random Walk Problem



$$T = 5$$

$$n = 2$$

$$t = 4$$

$$\text{tau} = 3$$

St	C	D	C	D	E	\$
Rt	0	0	0	0	0	1
At	r	l	r	r	r	

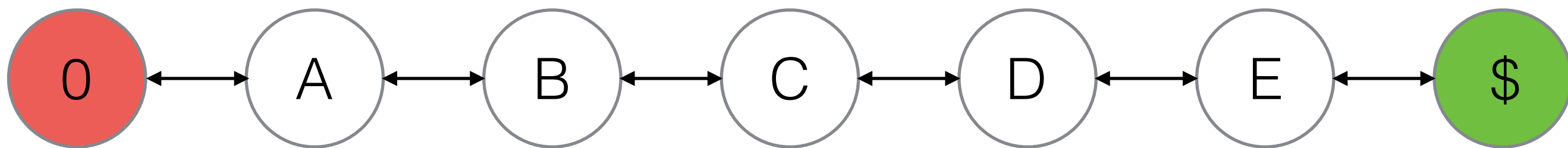
$$G = \gamma^0 R_4 + \gamma^1 R_5 = 0 + 1$$

$$Q(D, r) = 0 + \alpha(1 - 0)$$

```

|    $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|   If  $\tau \geq 0$ :
|        $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|       If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|        $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```

Random Walk Problem



$$T = 5$$

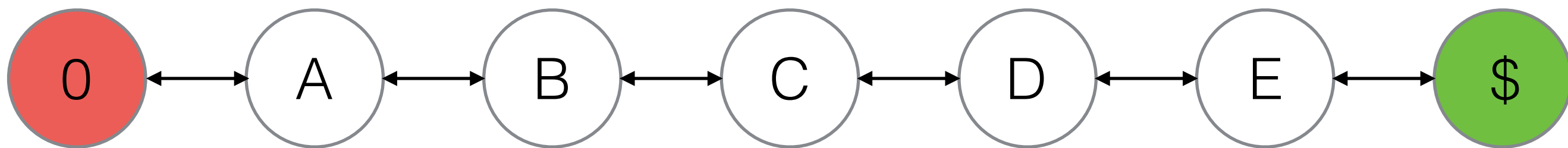
$$n = 2$$

$$t = 5$$

S_t	C	D	C	D	E	\$
R_t	0	0	0	0	0	1
A_t	r	l	r	r	r	

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = 5$$

$$n = 2$$

$$t = 5$$

$$\text{tau} = 4$$

St	C	D	C	D	E	\$
Rt	0	0	0	0	0	1
At	r	l	r	r	r	

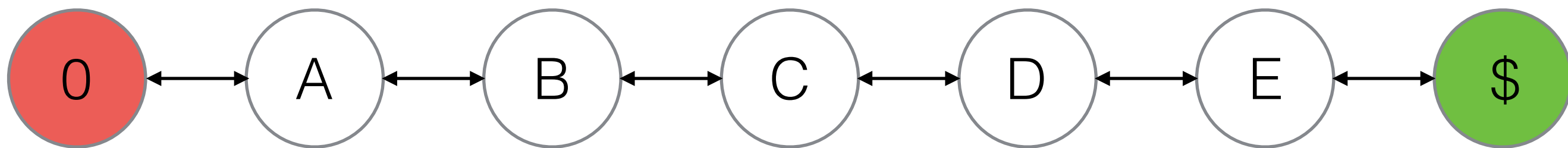
$$G = \gamma^0 R_5$$

$$Q(E, r) = 0 + \alpha(1 - 0)$$

```

|    $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|   If  $\tau \geq 0$ :
|        $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|       If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|        $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```


Random Walk Problem



$$T = 5$$

$$n = 2$$

$$t = 5$$

$$\text{tau} = 4$$

St	C	D	C	D	E	\$
Rt	0	0	0	0	0	1
At	r	l	r	r	r	

$$G = \gamma^0 R_5$$

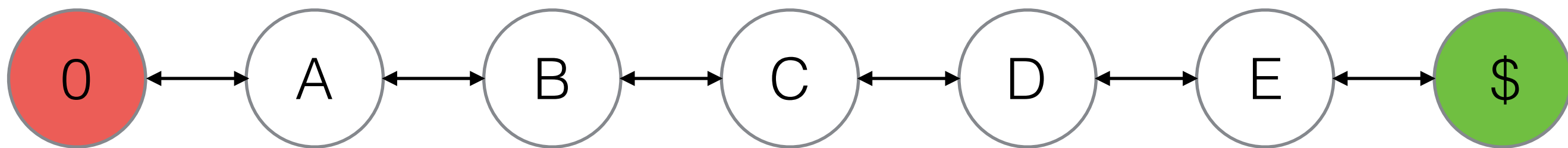
$$Q(E, r) = 0 + \alpha(1 - 0)$$

```

|    $\tau \leftarrow t - n + 1$    ( $\tau$  is the time whose estimate is being updated)
|   If  $\tau \geq 0$ :
|        $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|       If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|        $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```

Until $\tau = T - 1$

Random Walk Problem



$$T = \infty$$

$$n = 2$$

S_t	C					
R_t	0					
A_t	r					

Initialize $Q(s, a)$ arbitrarily, $\forall s \in \mathcal{S}, a \in \mathcal{A}$

Initialize π to be ε -greedy with respect to Q , or to a fixed given policy

Parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$, a positive integer n

All store and access operations (for S_t , A_t , and R_t) can take their index mod n

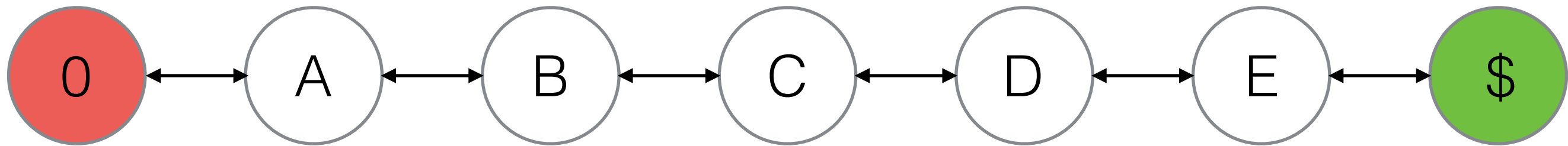
Repeat (for each episode):

Initialize and store $S_0 \neq \text{terminal}$

Select and store an action $A_0 \sim \pi(\cdot | S_0)$

$T \leftarrow \infty$

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 0$$

S_t	C	D				
R_t	0	0				
A_t	r	r				

For $t = 0, 1, 2, \dots$:

| If $t < T$, then:

| Take action A_t

| Observe and store the next reward as R_{t+1} and the next state as S_{t+1}

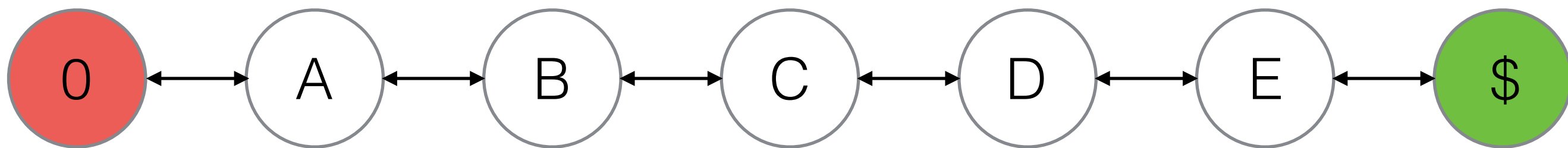
| If S_{t+1} is terminal, then:

| $T \leftarrow t + 1$

| else:

| Select and store an action $A_{t+1} \sim \pi(\cdot | S_{t+1})$

Random Walk Problem



$$T = \infty$$

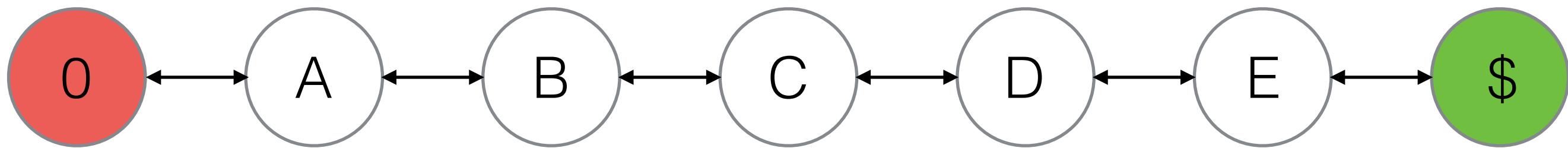
$$n = 2$$

$$t = 1$$

S_t	C	D	E			
R_t	0	0	0			
A_t	r	r	r			

```
For  $t = 0, 1, 2, \dots$  :  
|   If  $t < T$ , then:  
|       Take action  $A_t$   
|       Observe and store the next reward as  $R_{t+1}$  and the next state as  $S_{t+1}$   
|       If  $S_{t+1}$  is terminal, then:  
|            $T \leftarrow t + 1$   
|       else:  
|           Select and store an action  $A_{t+1} \sim \pi(\cdot | S_{t+1})$ 
```

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 1$$

$$\text{tau} = 0$$

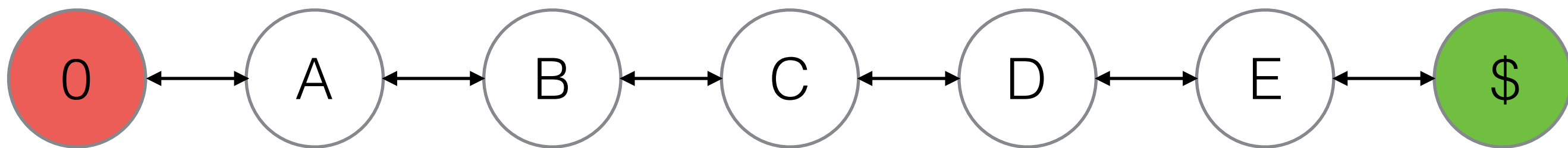
St	C	D	E			
Rt	0	0	0			
At	r	r	r			

$$G = \gamma^0 R_1 + \gamma^1 R_2$$

```

|   $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|  If  $\tau \geq 0$ :
|     $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|    If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|     $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```

Random Walk Problem



$$T = \infty$$

$$n = 2$$

$$t = 1$$

$$\text{tau} = 0$$

St	C	D	E			
Rt	0	0	0			
At	r	r	r			

$$G = \gamma^0 R_1 + \gamma^1 R_2$$

$$G = G + \gamma^2 Q(E, r)$$

```

|   $\tau \leftarrow t - n + 1$  ( $\tau$  is the time whose estimate is being updated)
|  If  $\tau \geq 0$ :
|       $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$ 
|      If  $\tau + n < T$ , then  $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$ 
|       $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha [G - Q(S_\tau, A_\tau)]$ 
  
```