

CS230: Lecture 9

Deep Reinforcement Learning

Kian Katanforoosh
Menti code: 80 24 08

Today's outline

I. Motivation

II. Recycling is good: an introduction to RL

III. Deep Q-Networks

IV. Application of Deep Q-Network: Breakout (Atari)

V. Tips to train Deep Q-Network

VI. Advanced topics

I. Motivation



AlphaGo



Human Level Control through
Deep Reinforcement Learning

[Silver et al. (2017): Mastering the game of Go without human knowledge]

[Mnih et al. (2015): Human Level Control through Deep Reinforcement Learning]

Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri

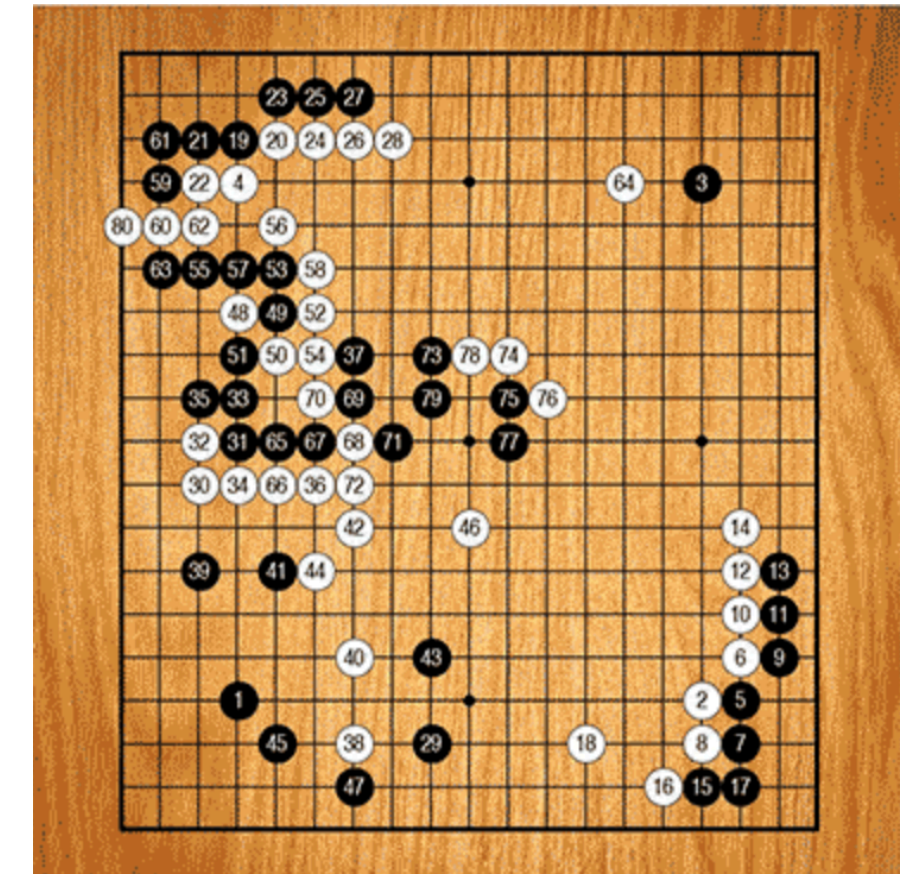
I. Motivation

Why RL?

- Delayed labels
- Making sequences of decisions

What is RL?

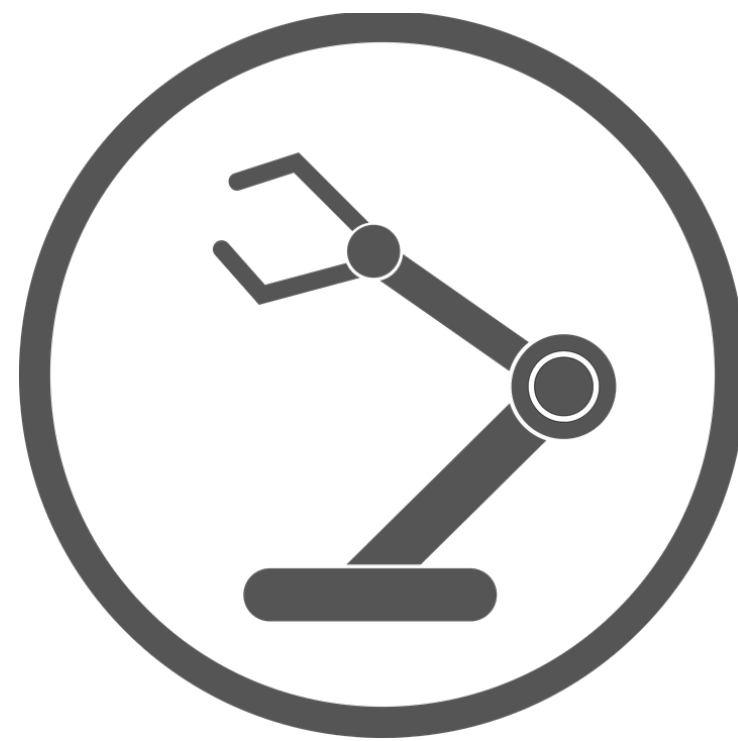
- Automatically learn to make good sequences of decision



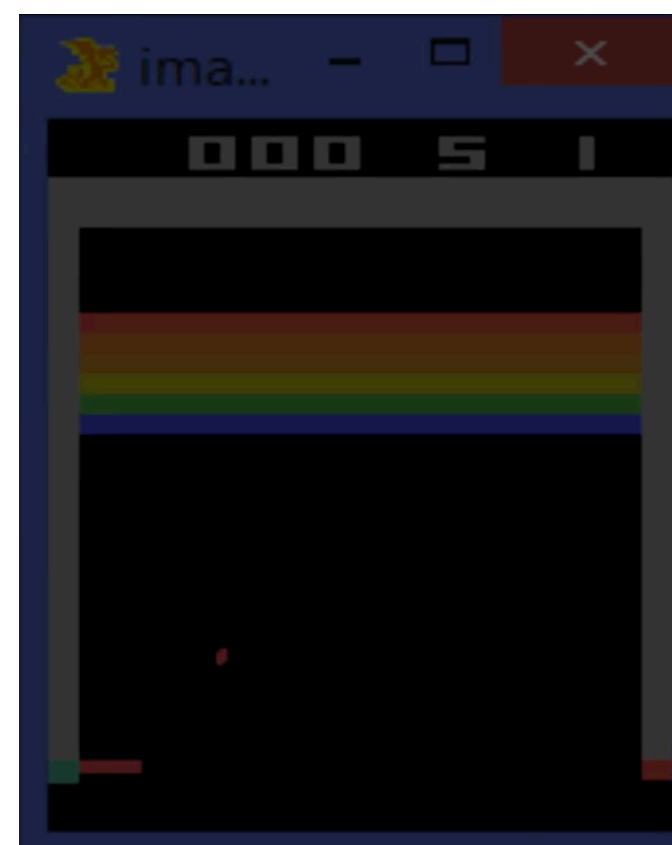
Source: <https://deepmind.com/blog/alphago-zero-learning-scratch/>

Examples of RL applications

Robotics



Games



Advertisement



Today's outline

I. Motivation

II. Recycling is good: an introduction to RL

III. Deep Q-Networks

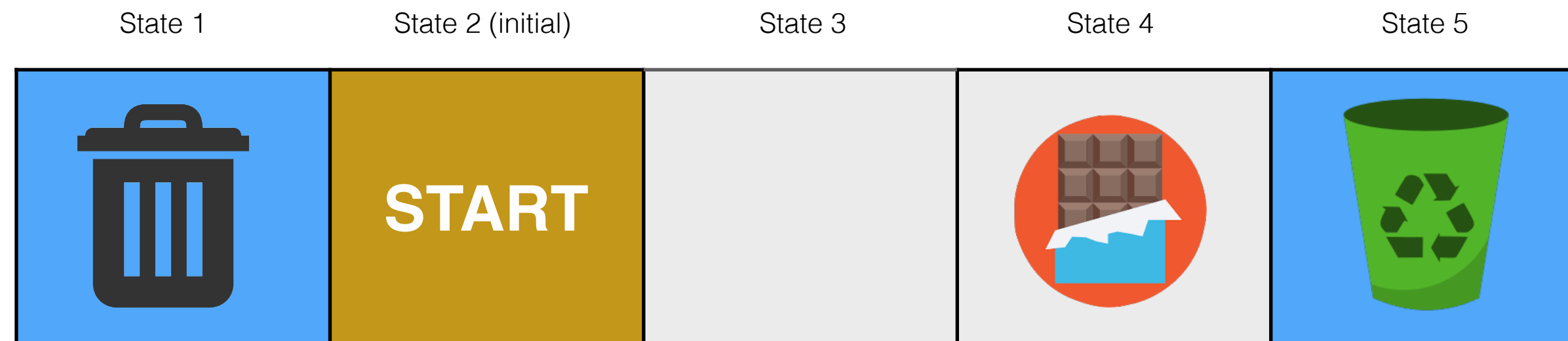
IV. Application of Deep Q-Network: Breakout (Atari)

V. Tips to train Deep Q-Network

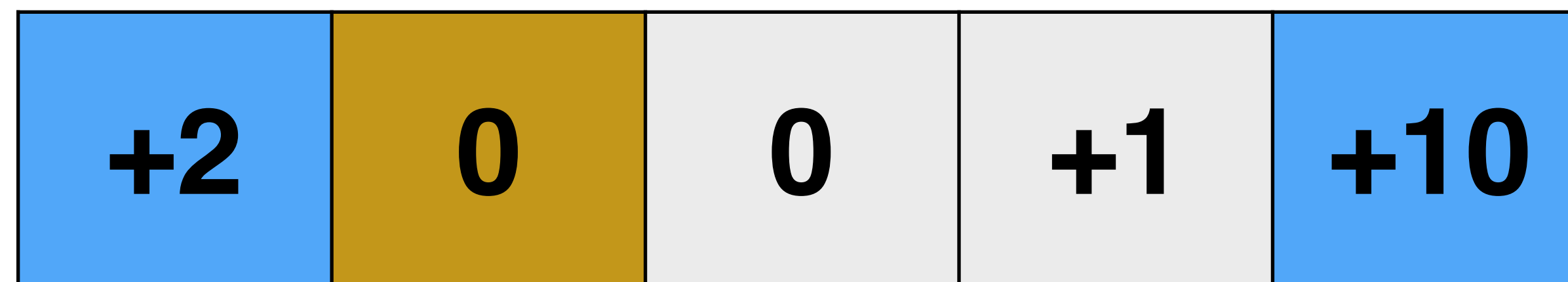
VI. Advanced topics

II. Recycling is good: an introduction to RL

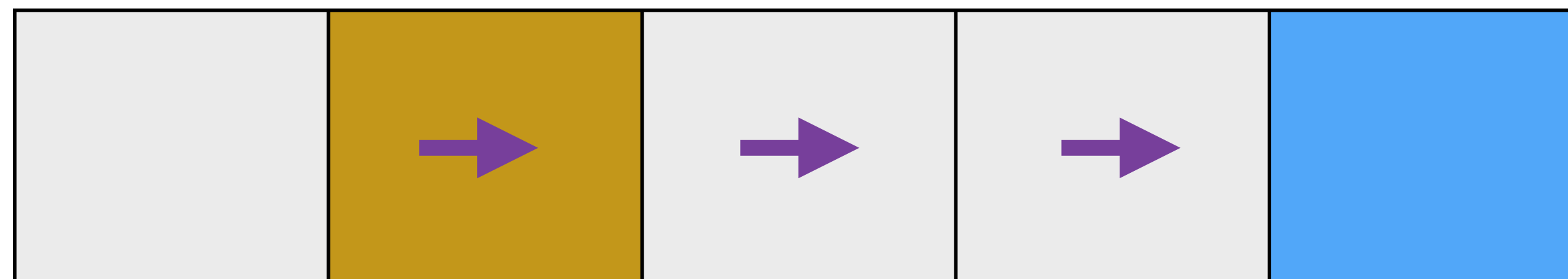
Problem statement



Define reward “ r ” in every state



Best strategy to follow if $\gamma = 1$



Goal: maximize the return (rewards)

Number of states: 5

Types of states: 

Agent's Possible actions: 

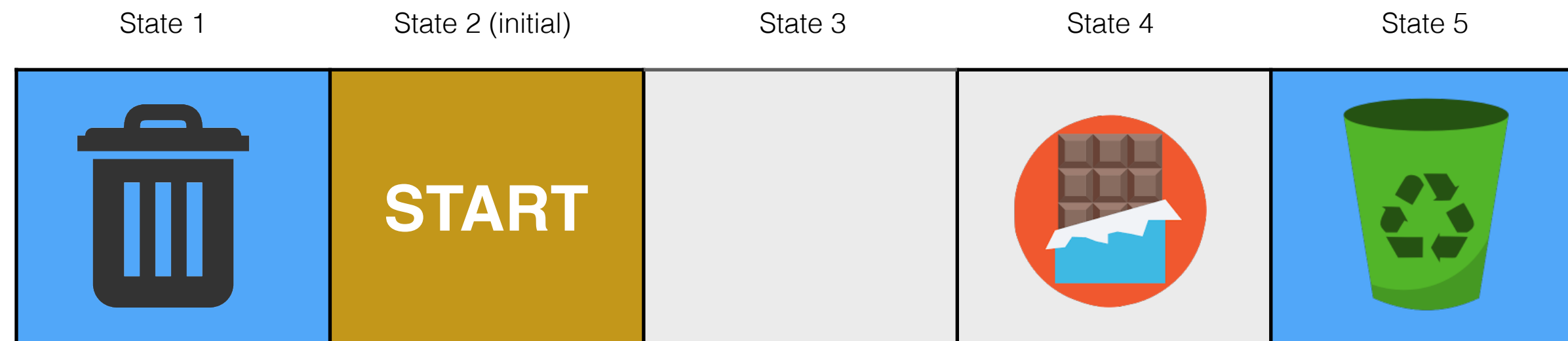
Additional rule: garbage collector coming in 3min, it takes 1min to move between states

How to define the long-term return?

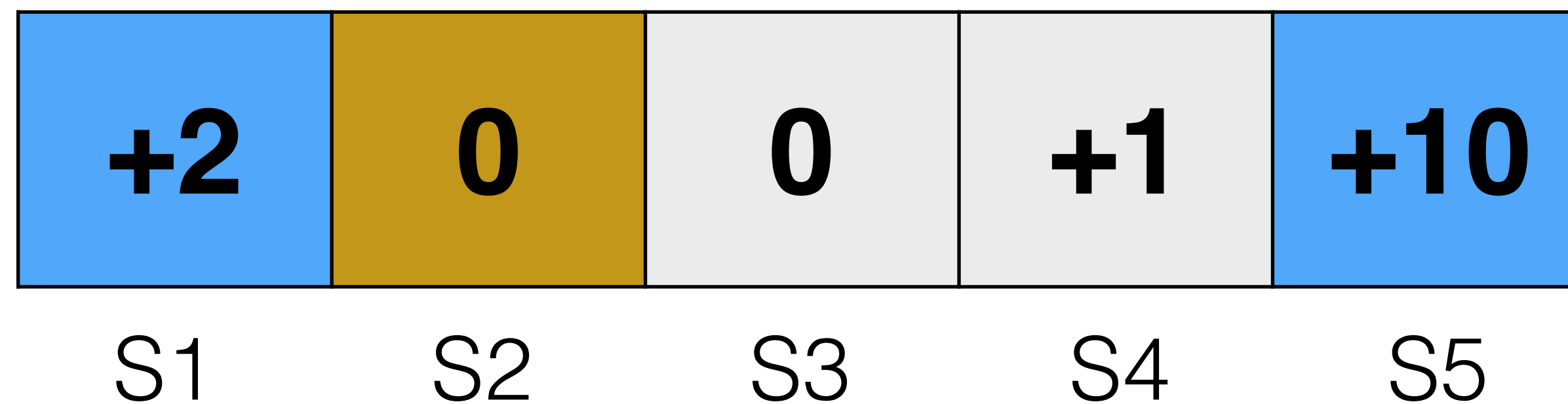
Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

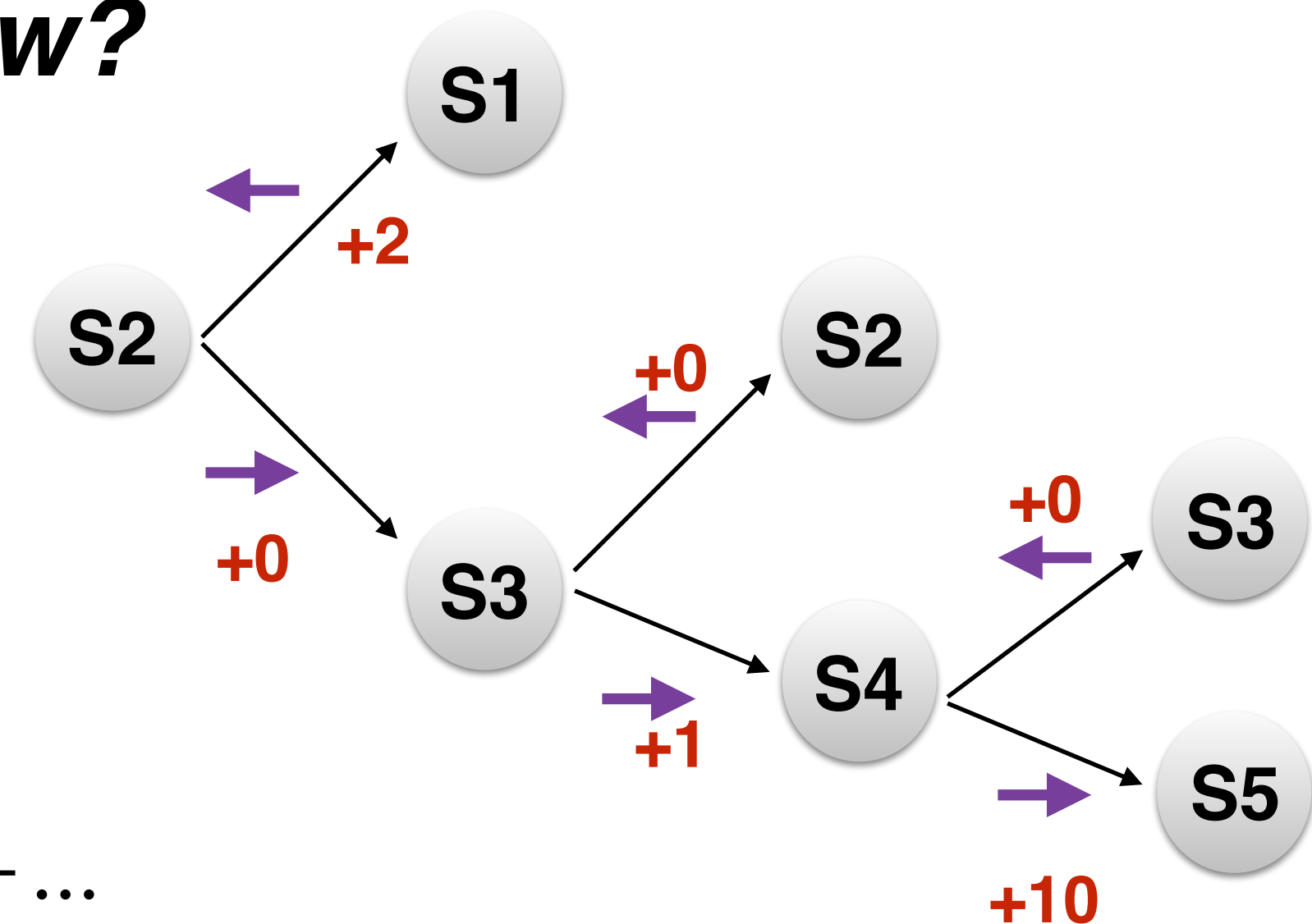
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

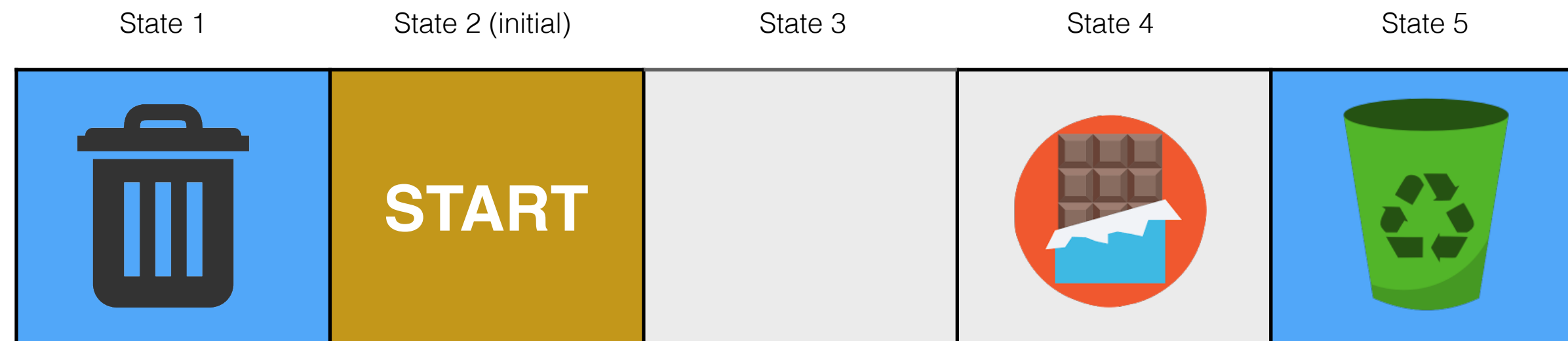
$$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} \downarrow \text{\#states} \\ \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{pmatrix} \end{matrix} \end{matrix}$$

How?

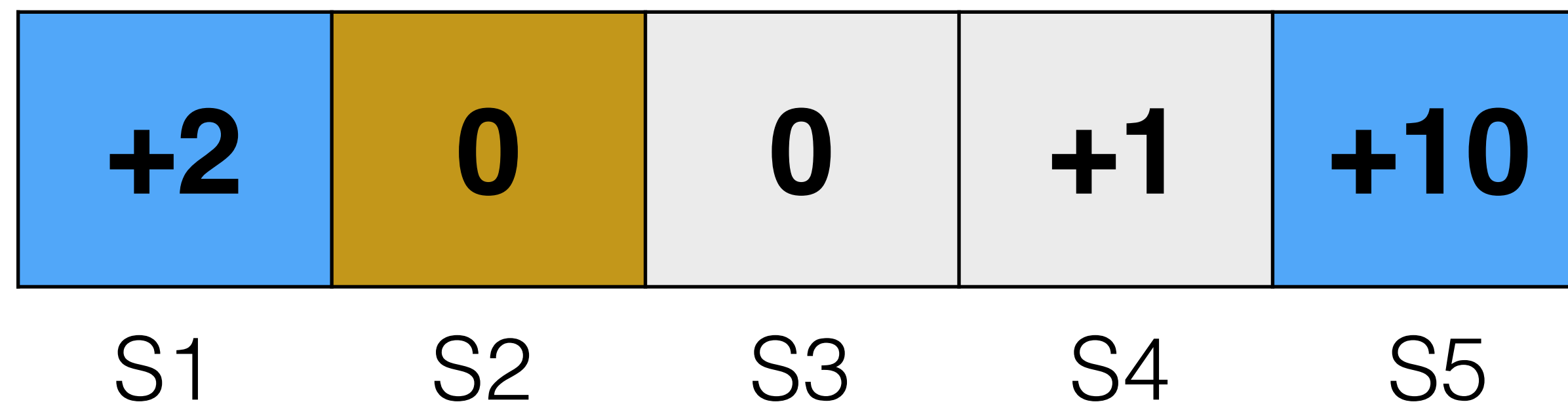


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

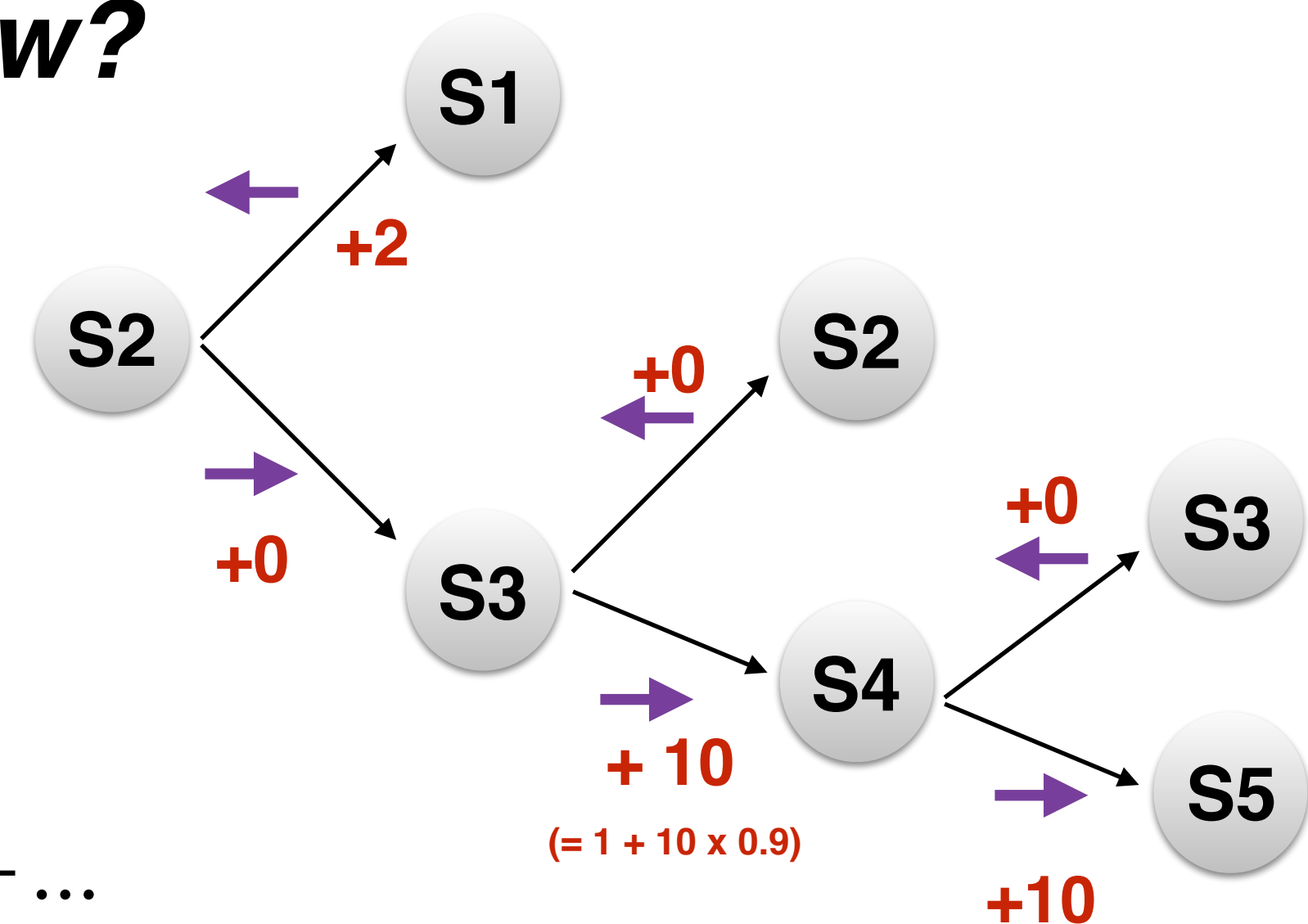
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

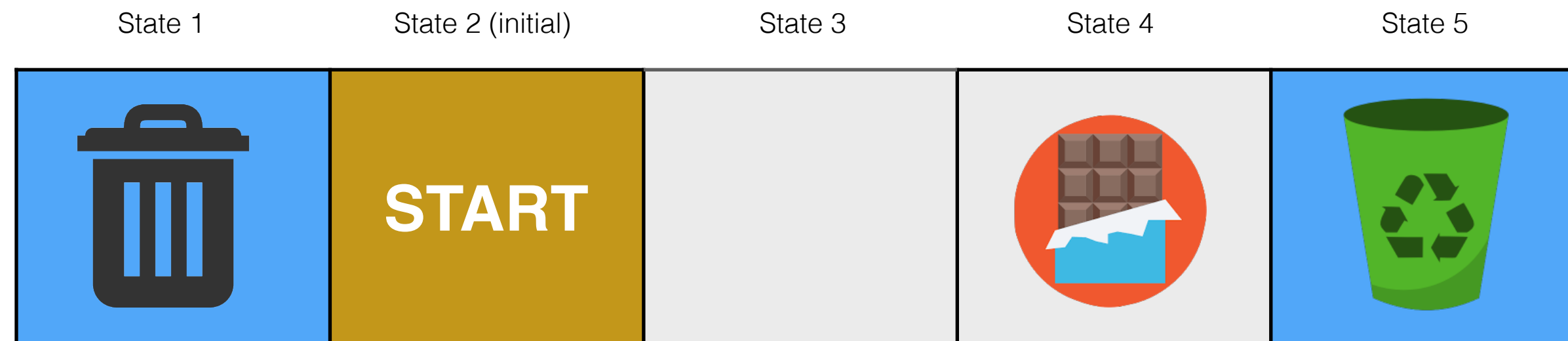
$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} \uparrow \text{\#states} \\ \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{pmatrix} \end{matrix} \end{matrix}$

How?

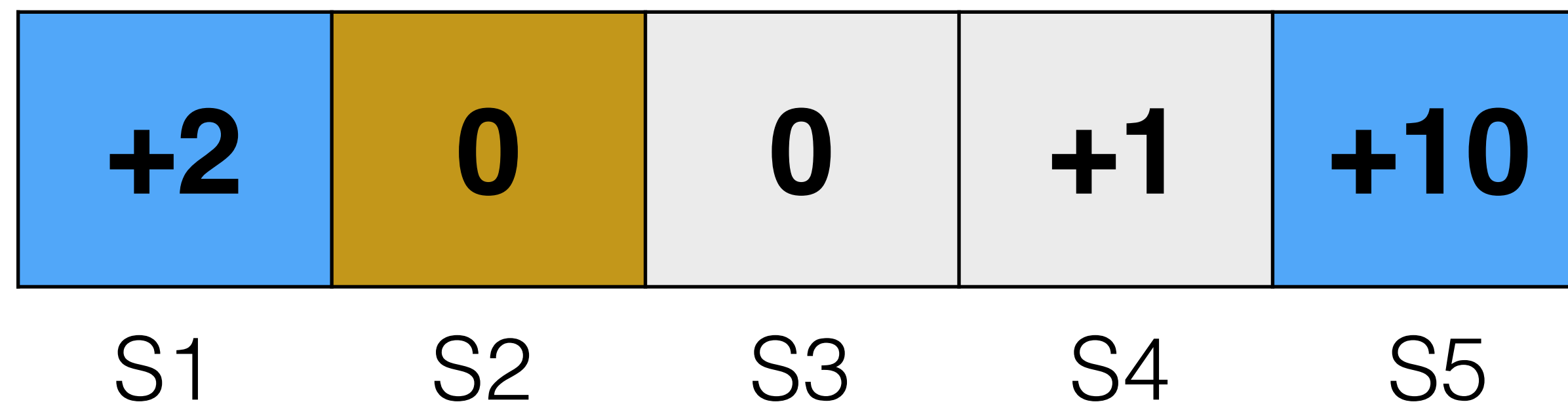


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

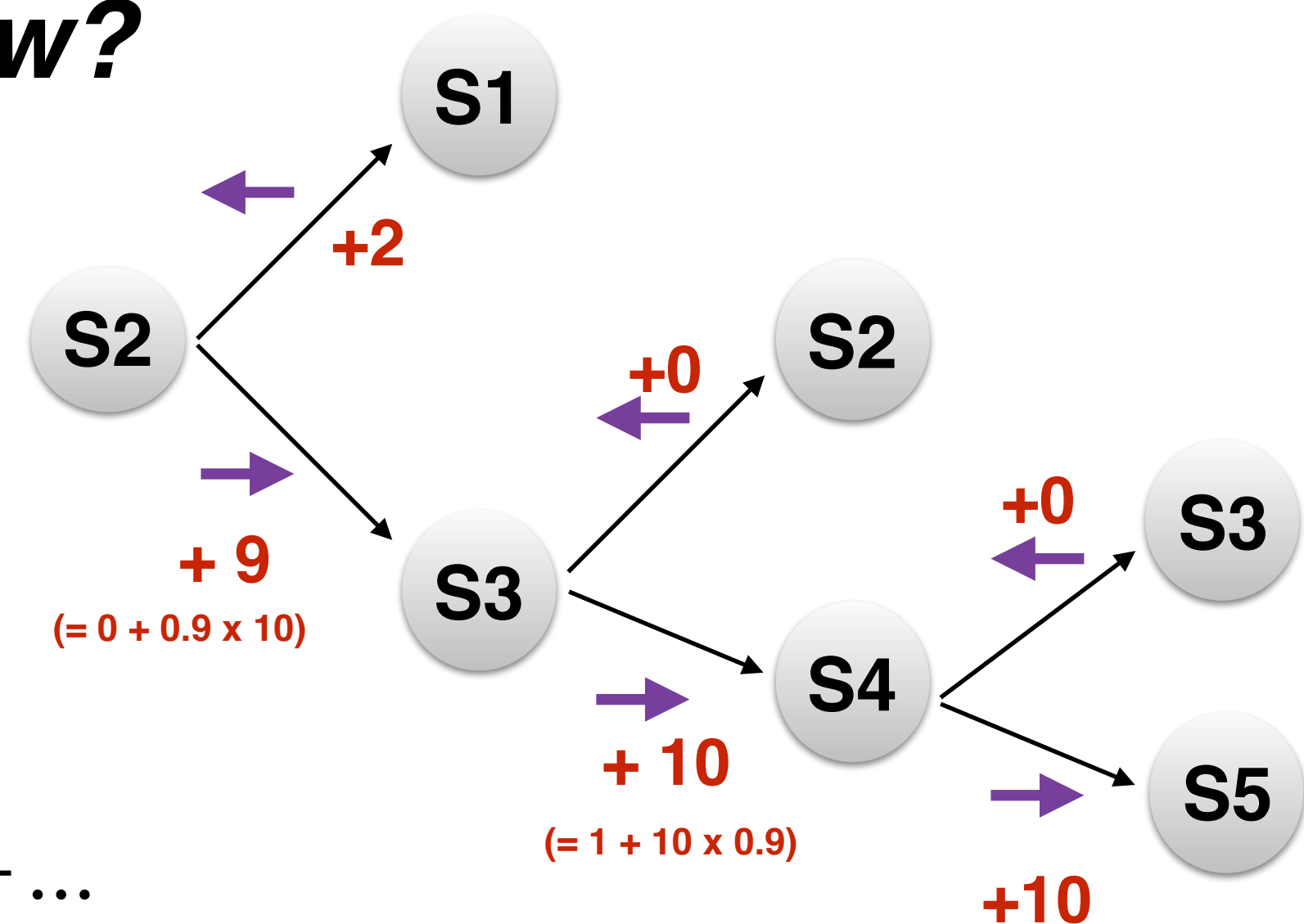
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

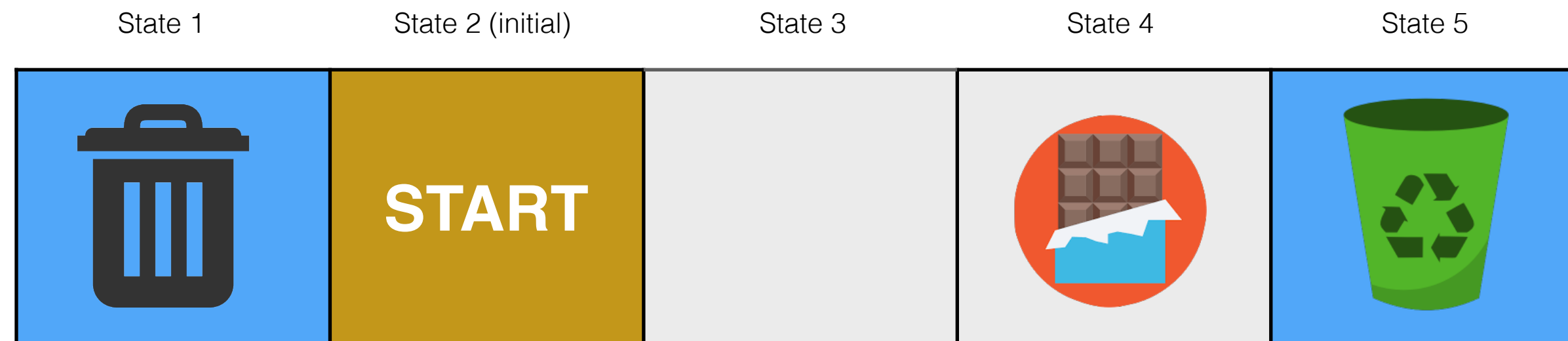
$$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} \downarrow \text{\#states} \\ \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{pmatrix} \end{matrix} \end{matrix}$$

How?

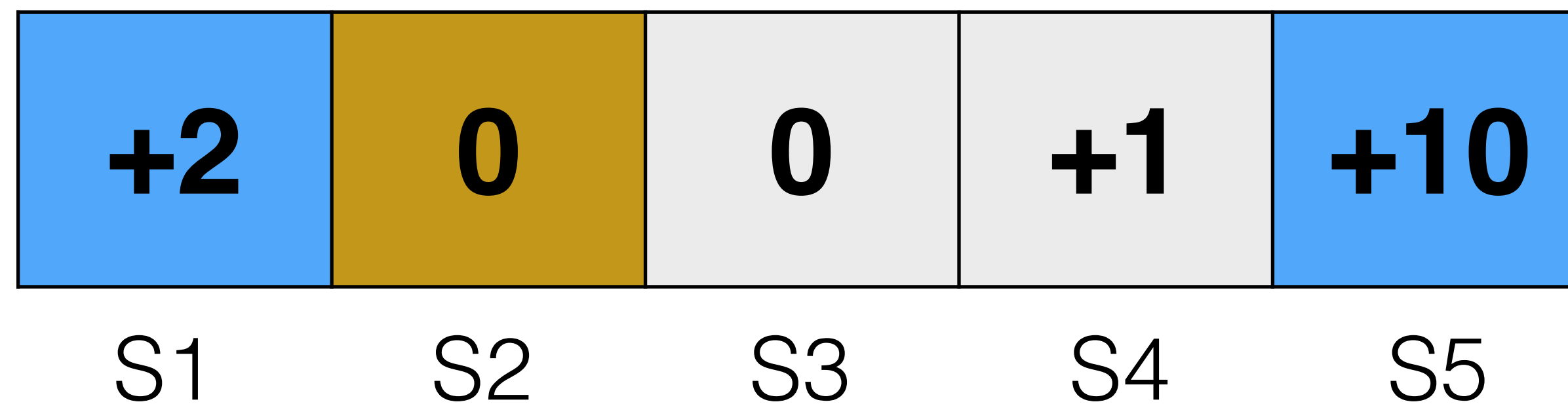


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

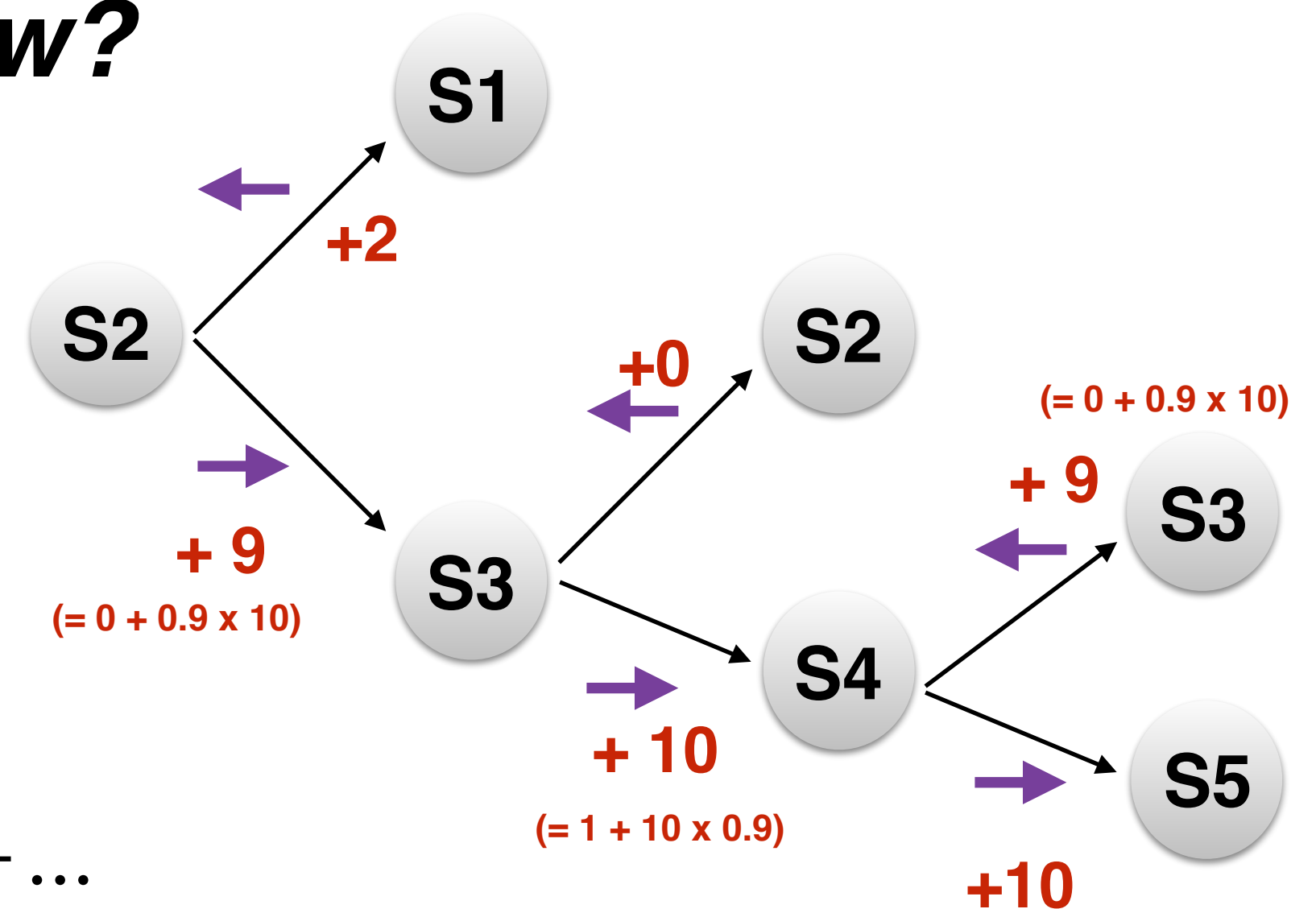
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

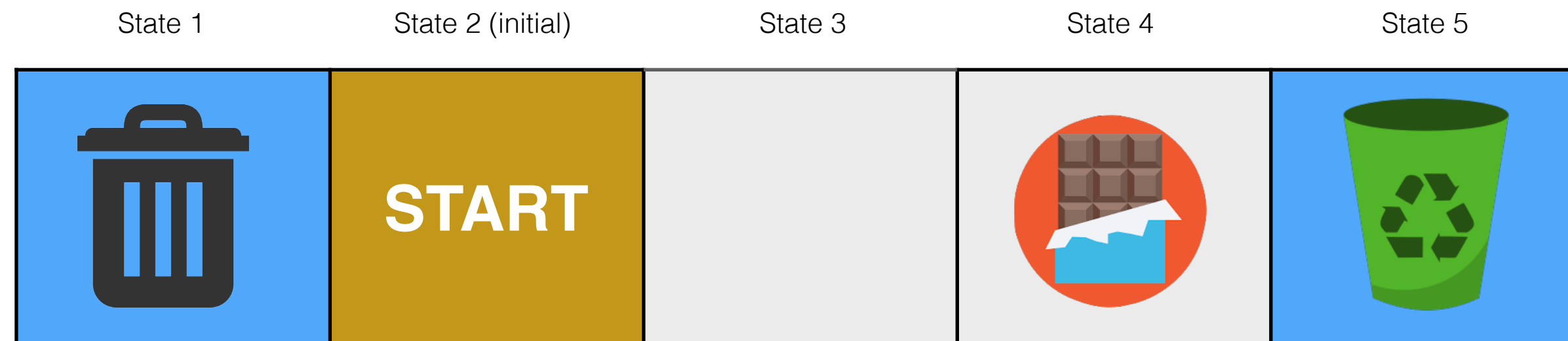
$$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{matrix} & \begin{matrix} \updownarrow \\ \text{\#states} \end{matrix} \end{matrix}$$

How?

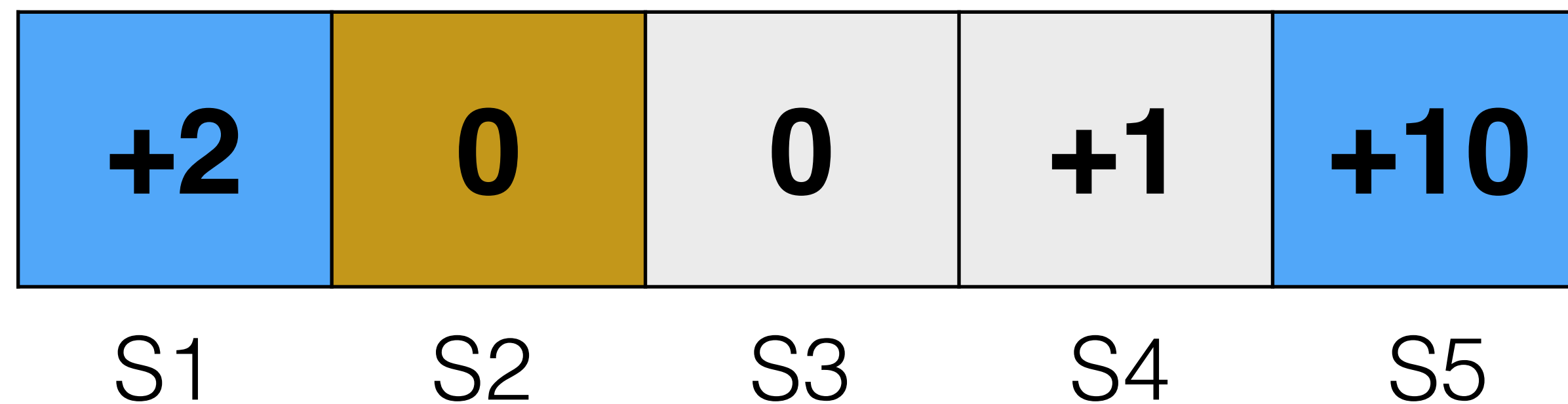


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

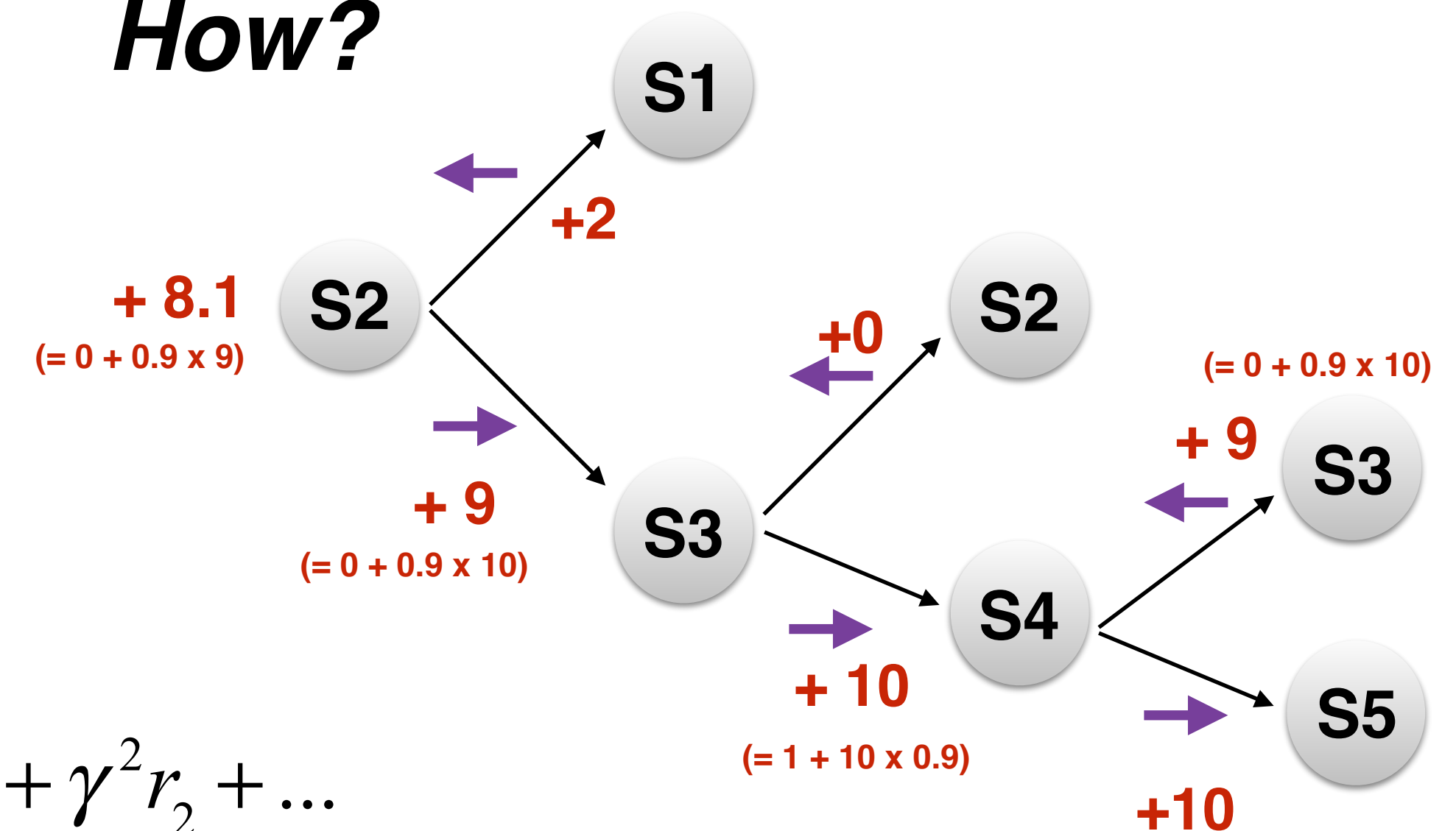
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

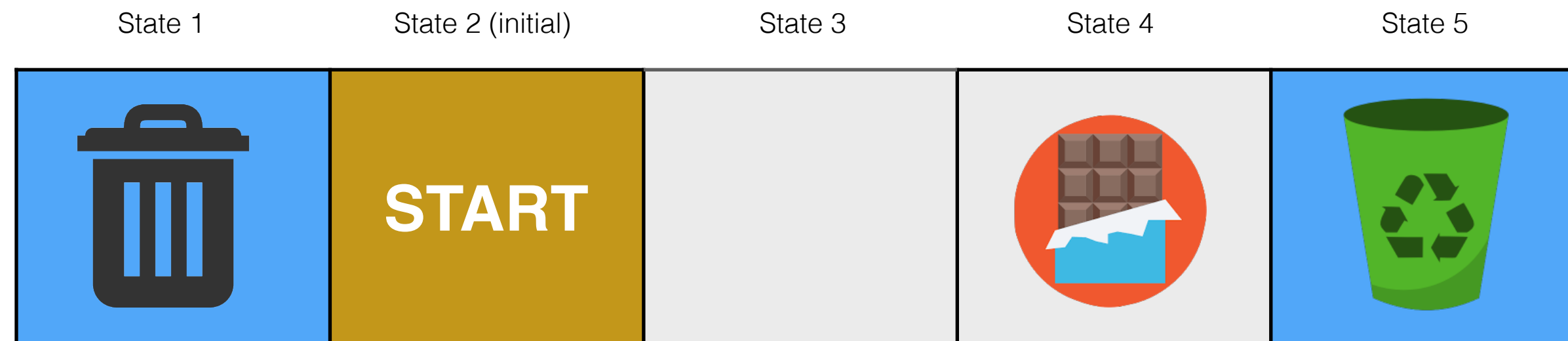
$$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{matrix} & \begin{matrix} \updownarrow \\ \text{\#states} \end{matrix} \end{matrix}$$

How?

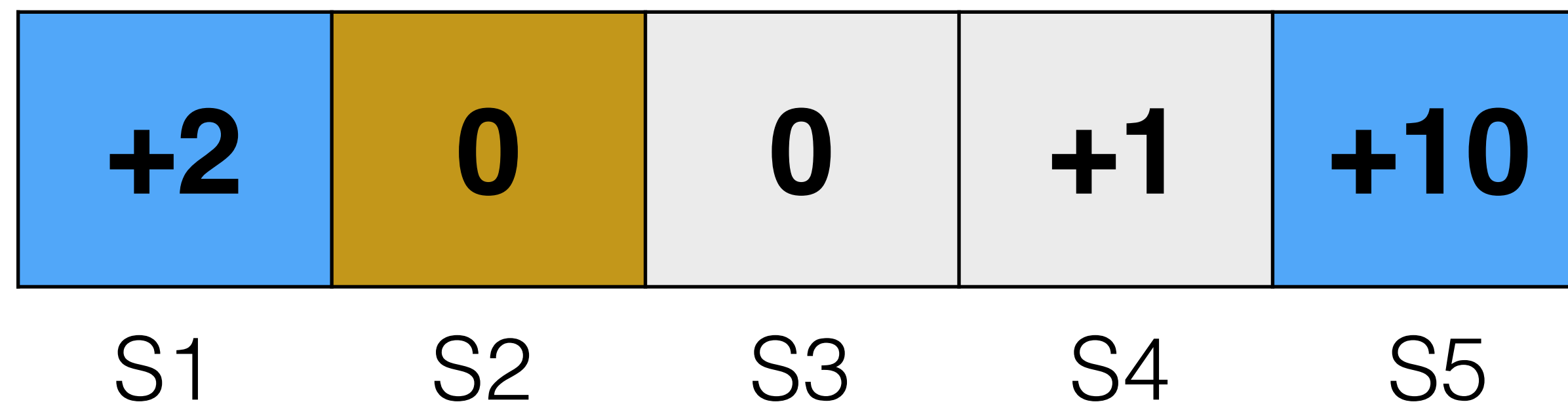


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

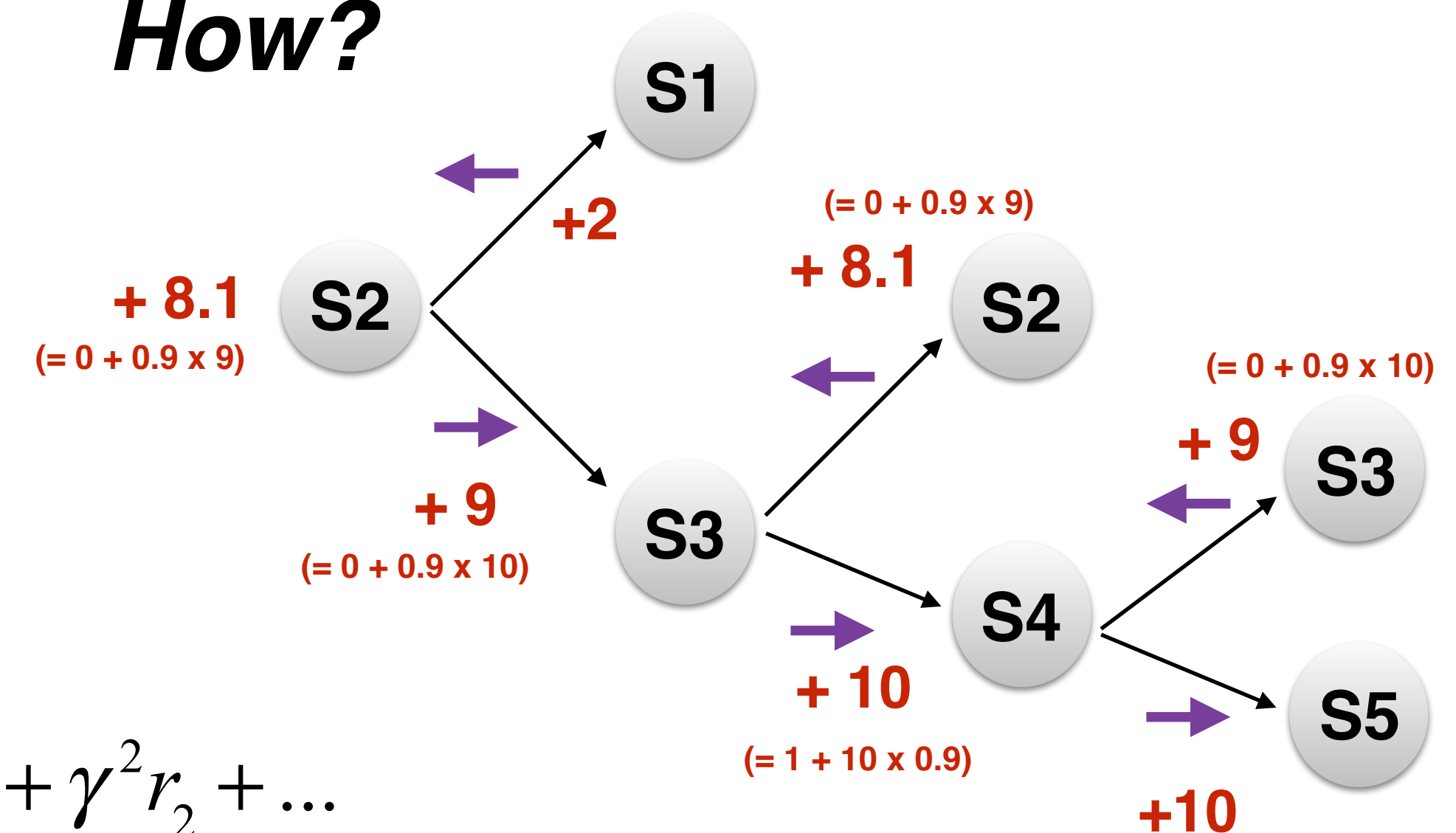
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

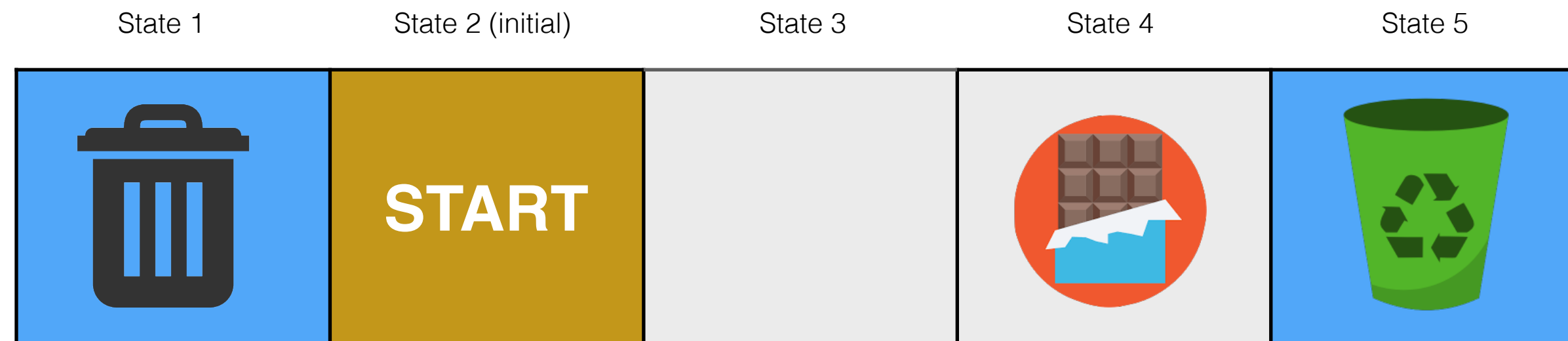
$$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} \downarrow \text{\#states} \\ \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \\ Q_{31} & Q_{32} \\ Q_{41} & Q_{42} \\ Q_{51} & Q_{52} \end{pmatrix} \end{matrix} \end{matrix}$$

How?

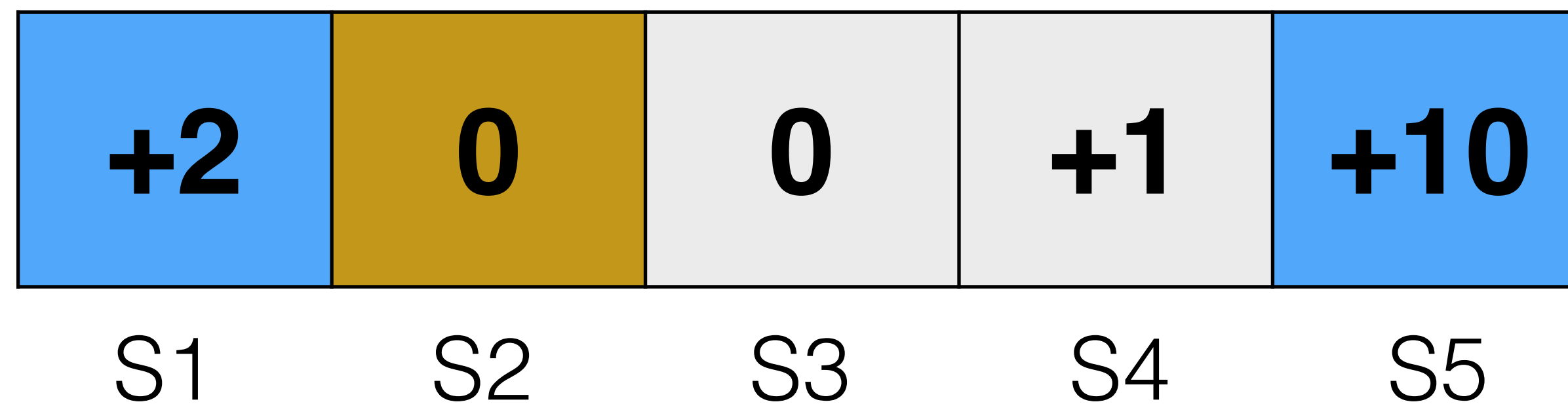


II. Recycling is good: an introduction to RL

Problem statement



Define reward “**r**” in every state



Assuming $\gamma = 0.9$

Discounted return $R = \sum_{t=0} \gamma^t r_t = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

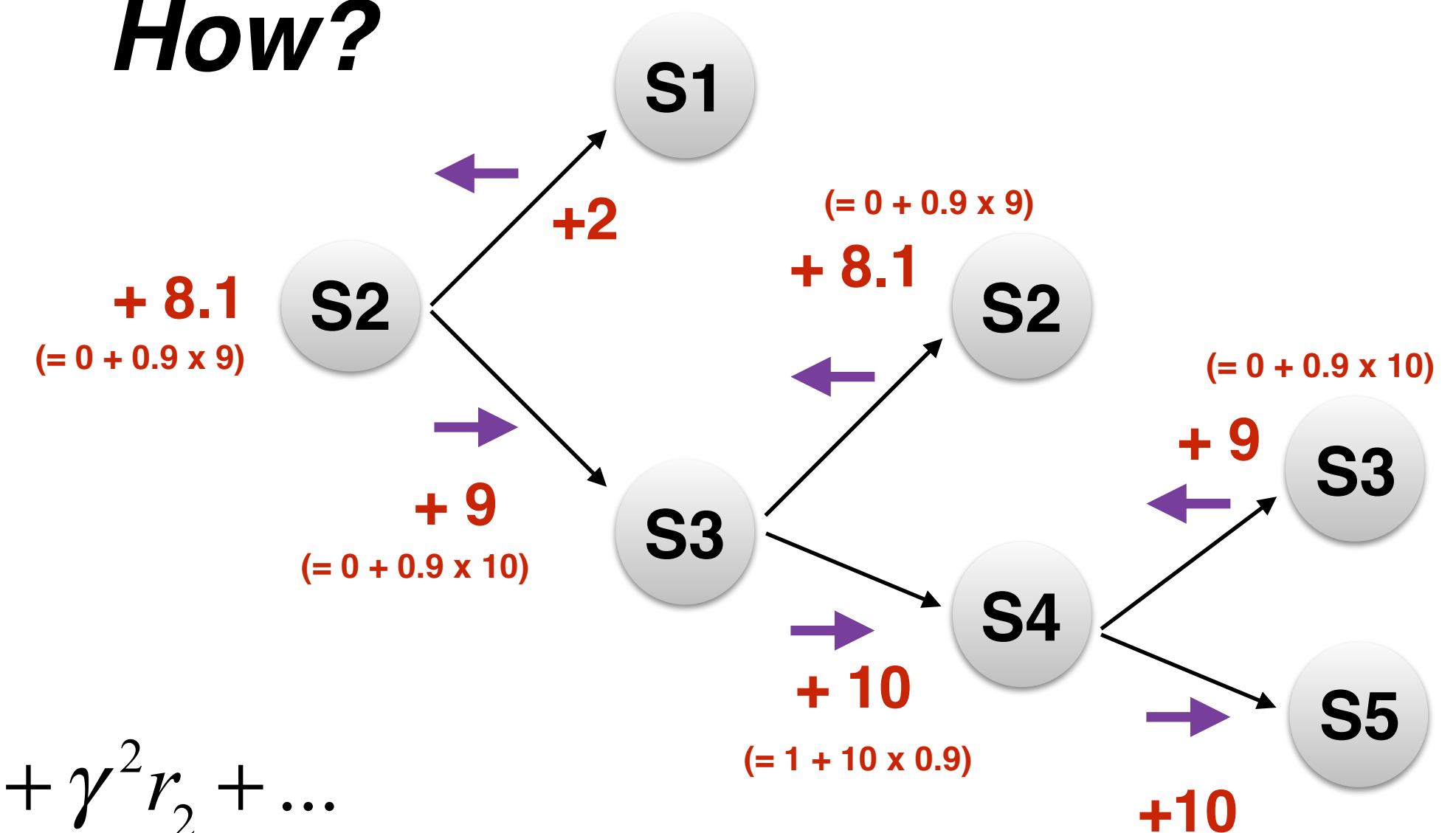
What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

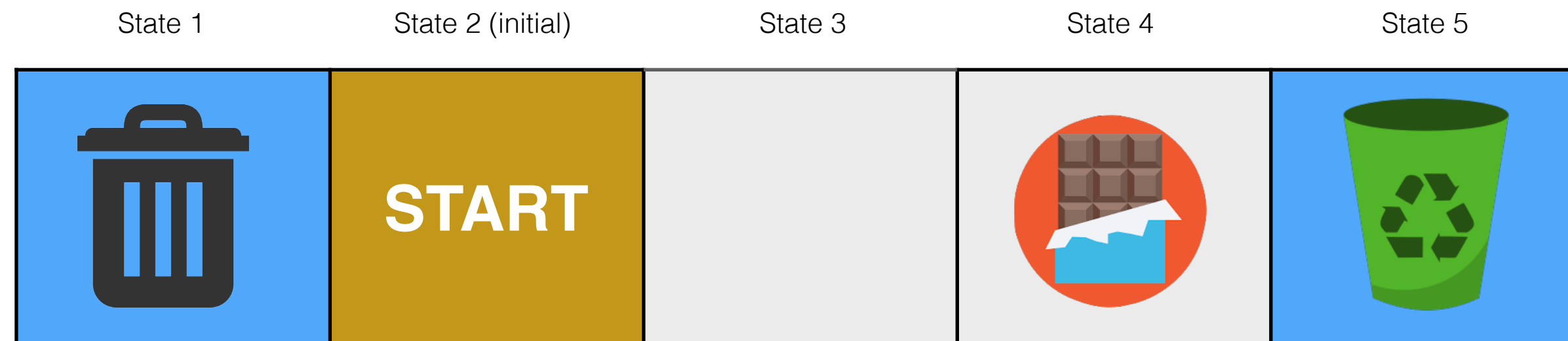
$Q = \begin{matrix} & \xrightarrow{\text{\#actions}} \\ \begin{matrix} \downarrow \text{\#states} \\ \begin{pmatrix} 0 & 0 \\ 2 & 9 \\ 8.1 & 10 \\ 9 & 10 \\ 0 & 0 \end{pmatrix} \end{matrix} \end{matrix}$

How?

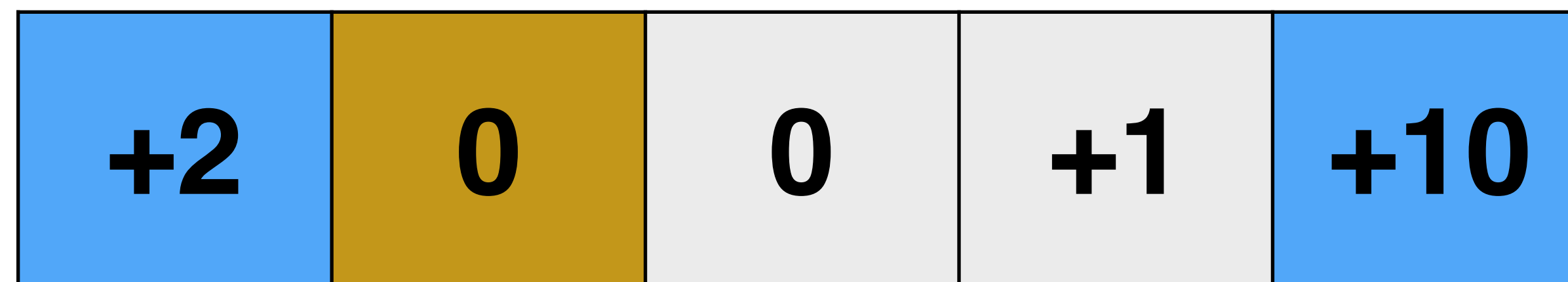


II. Recycling is good: an introduction to RL

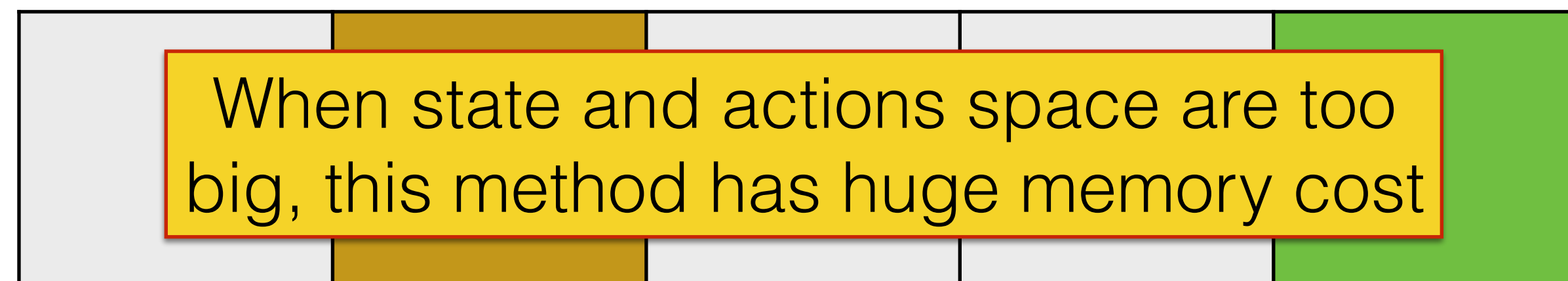
Problem statement



Define reward “*r*” in every state



Best strategy to follow if $\gamma = 0.9$



Function telling us our best strategy

What do we want to learn?

how good is it to take
action 1 in state 2

Q-table

$Q =$

#actions		#states
0	0	
2	9	
8.1	10	
9	10	
0	0	

*Bellman equation
(optimality equation)*

$$Q^*(s, a) = r + \gamma \max_{a'} (Q^*(s', a'))$$

Policy $\pi(s) = \arg \max_a (Q^*(s, a))$

What we've learned so far:

- *Vocabulary: environment, agent, state, action, reward, total return, discount factor.*
- *Q-table: matrix of entries representing “how good is it to take action a in state s ”*
- *Policy: function telling us what's the best strategy to adopt*
- *Bellman equation satisfied by the optimal Q-table*

Today's outline

I. Motivation

II. Recycling is good: an introduction to RL

III. Deep Q-Networks

IV. Application of Deep Q-Network: Breakout (Atari)

V. Tips to train Deep Q-Network

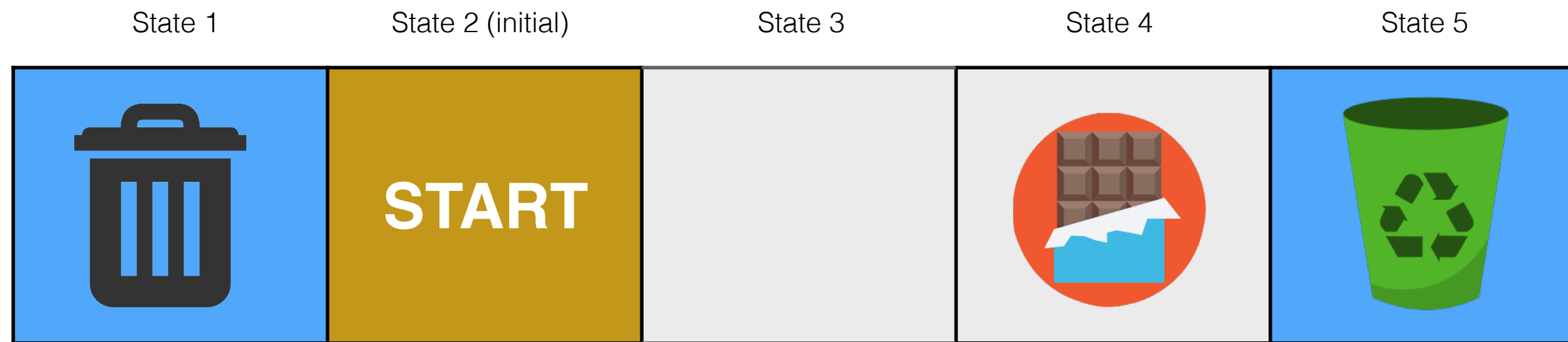
VI. Advanced topics

III. Deep Q-Networks

Main idea: find a Q-function to replace the Q-table

Problem statement

Neural Network

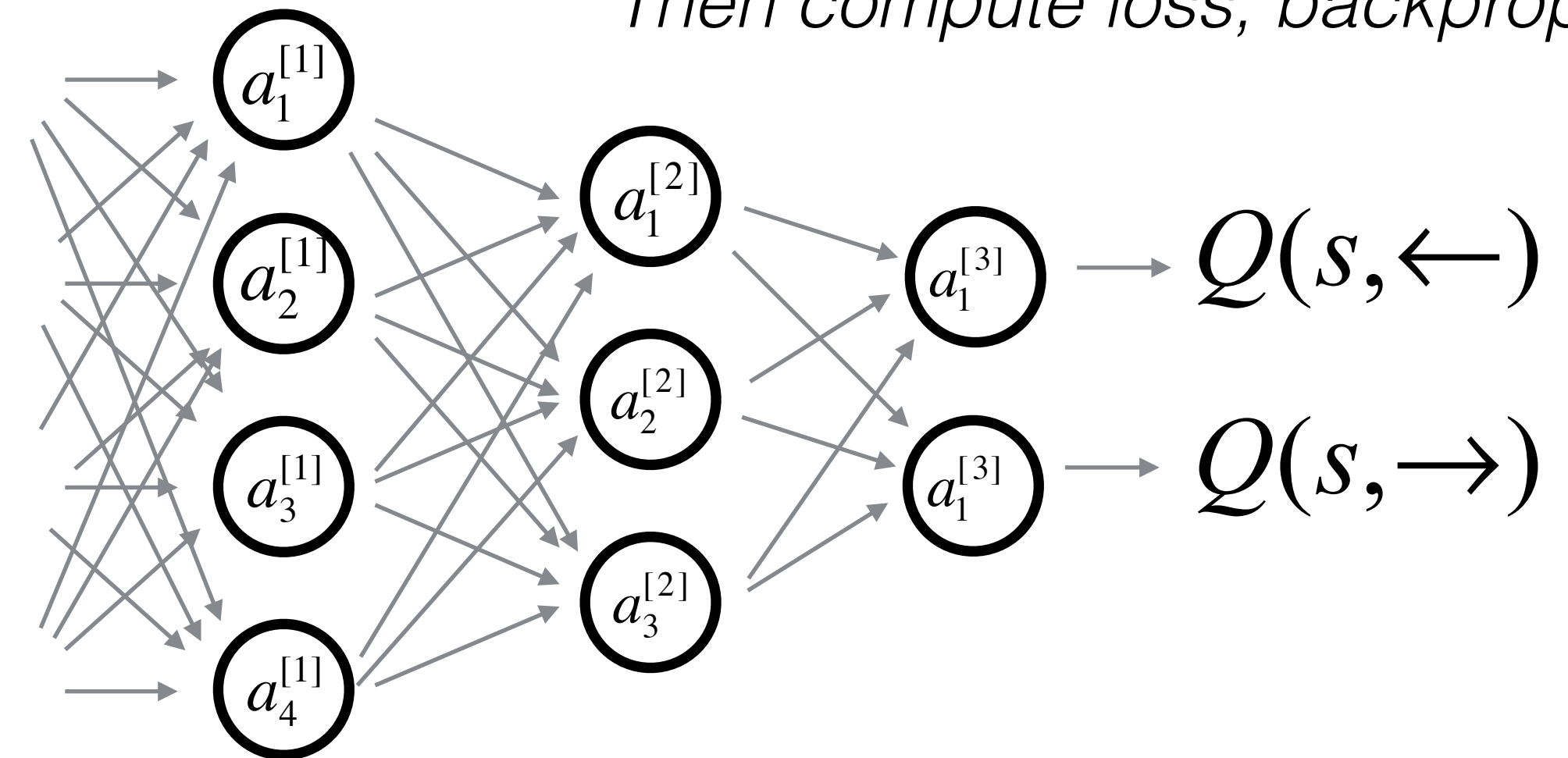


Q-table

#actions		#states
$Q = \begin{pmatrix} 0 & 0 \\ 2 & 9 \\ 8.1 & 10 \\ 9 & 10 \\ 0 & 0 \end{pmatrix}$	0	
	2	
	8.1	
	9	
	0	
	0	



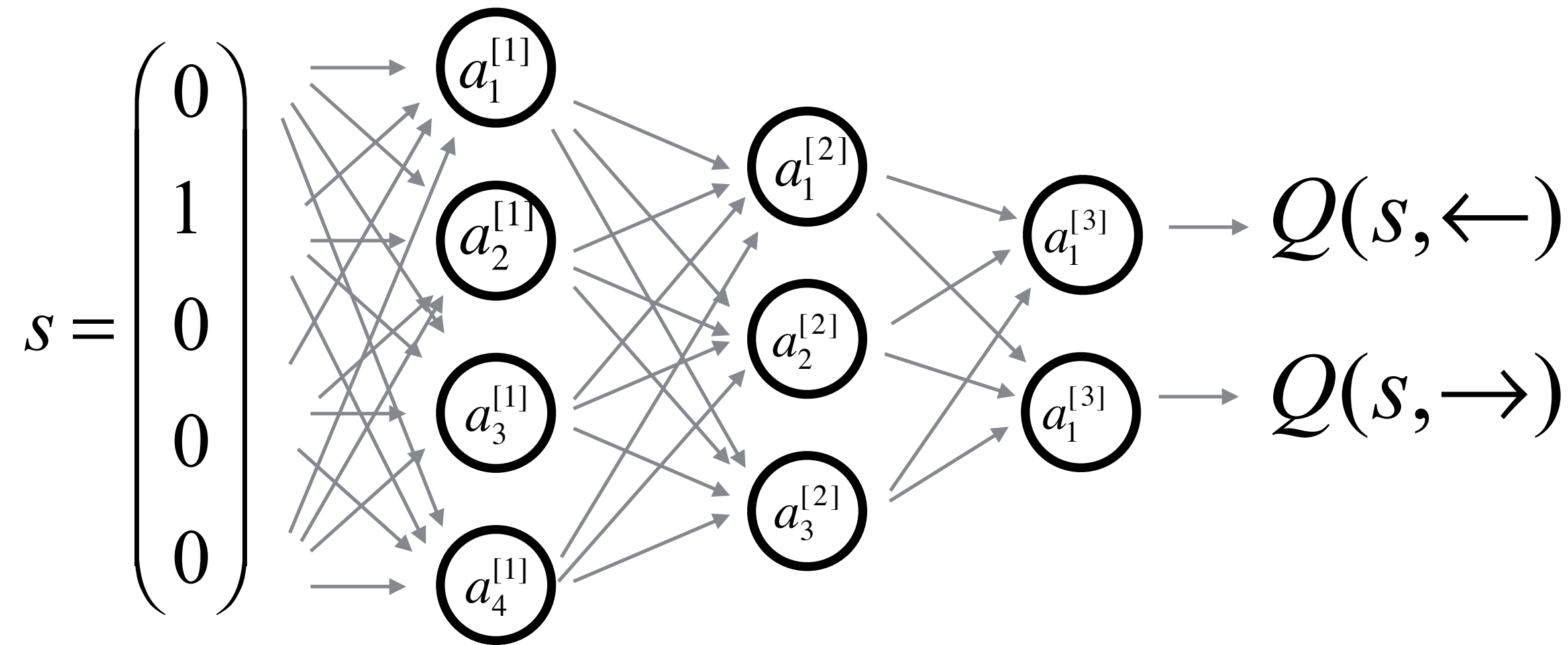
$$s = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$



How to compute the loss?

III. Deep Q-Networks

$$Q^*(s,a) = r + \gamma \max_{a'} (Q^*(s',a'))$$



Loss function

$$L = (y - Q(s, \leftarrow))^2$$

Target value

Case: $Q(s, \leftarrow) > Q(s, \rightarrow)$

$$y = r_{\leftarrow} + \gamma \max_{a'} (Q(s_{\leftarrow}^{next}, a'))$$

Hold fixed for backprop

Immediate reward for taking action \leftarrow in state s

Discounted maximum future reward when you are in state s_{\leftarrow}^{next}

Case: $Q(s, \leftarrow) < Q(s, \rightarrow)$

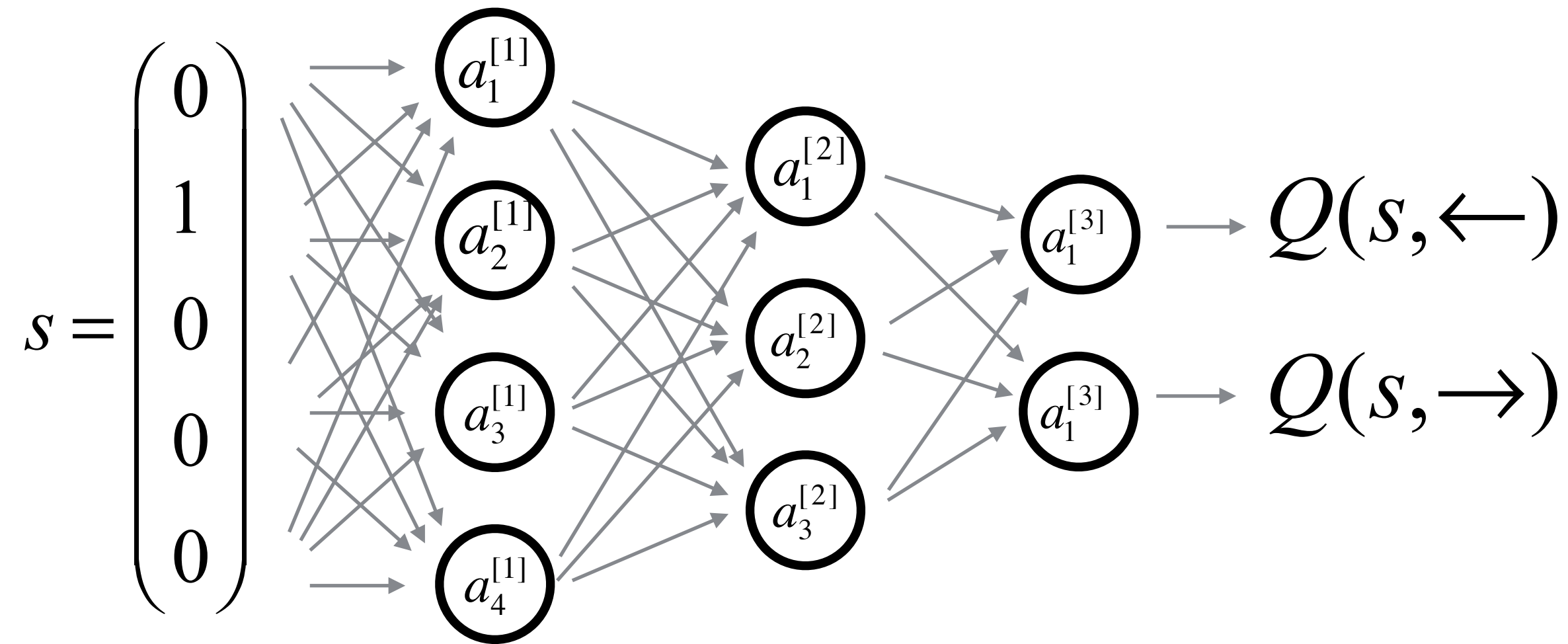
$$y = r_{\rightarrow} + \gamma \max_{a'} (Q(s_{\rightarrow}^{next}, a'))$$

Hold fixed for backprop

Immediate Reward for taking action \rightarrow in state s

Discounted maximum future reward when you are in state s_{\rightarrow}^{next}

III. Deep Q-Networks



Loss function (regression)

$$L = (y - Q(s, \rightarrow))^2$$

Target value

Case: $Q(s, \leftarrow) > Q(s, \rightarrow)$

$$y = r_{\leftarrow} + \gamma \max_{a'} (Q(s_{\leftarrow}^{next}, a'))$$

Case: $Q(s, \leftarrow) < Q(s, \rightarrow)$

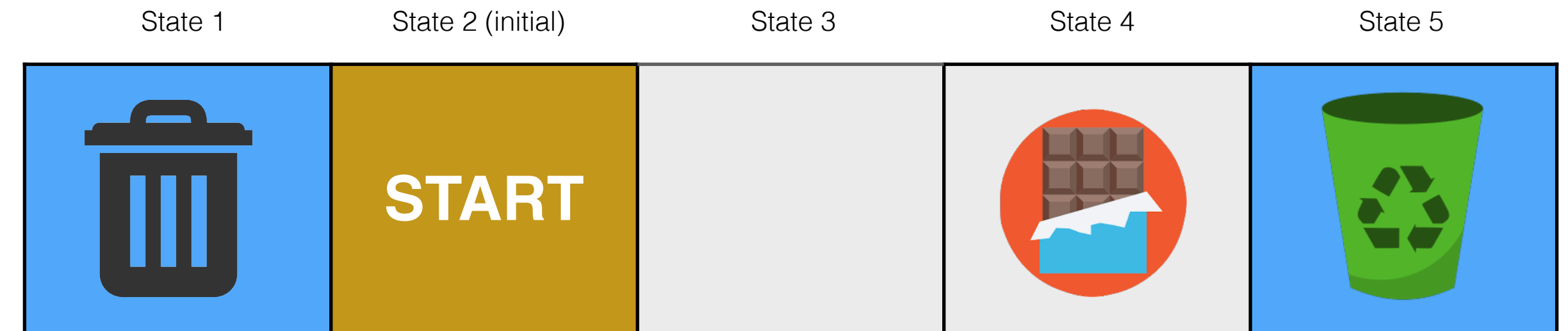
$$y = r_{\rightarrow} + \gamma \max_{a'} (Q(s_{\rightarrow}^{next}, a'))$$

Backpropagation Compute $\frac{\partial L}{\partial W}$ and update W using stochastic gradient descent

Recap'

DQN Implementation:

- Initialize your Q-network parameters
- Loop over episodes:
 - Start from initial state s
 - Loop over time-steps:
 - Forward propagate s in the Q-network
 - Execute action a (that has the maximum $Q(s,a)$ output of Q-network)
 - Observe rewards r and next state s'
 - Compute targets y by forward propagating state s' in the Q-network, then compute loss.
 - Update parameters with gradient descent



Today's outline

I. Motivation

II. Recycling is good: an introduction to RL

III. Deep Q-Networks

IV. Application of Deep Q-Network: Breakout (Atari)

V. Tips to train Deep Q-Network

VI. Advanced topics

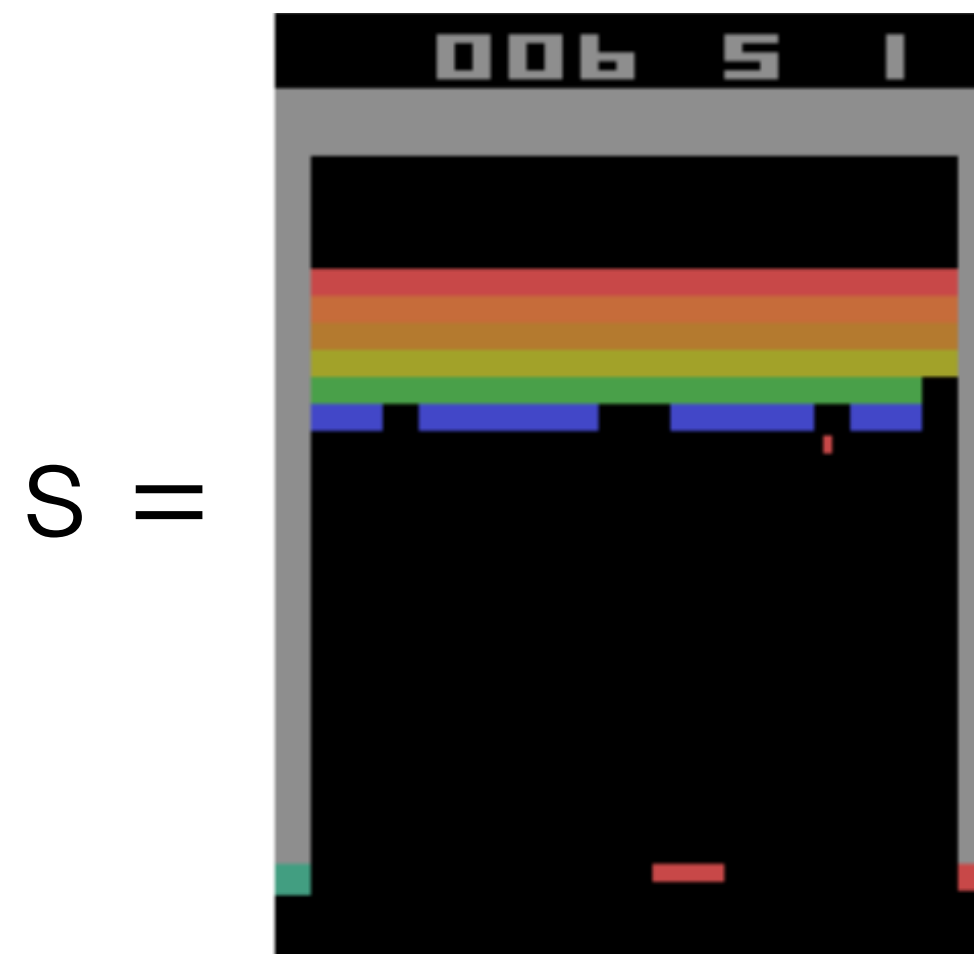
IV. Deep Q-Networks application: Breakout (Atari)

Goal: play breakout, i.e. destroy all the bricks.

Demo



input of Q-network



Output of Q-network

Q-values

$$\begin{pmatrix} Q(s, \leftarrow) \\ Q(s, \rightarrow) \\ Q(s, -) \end{pmatrix}$$

Would that work?

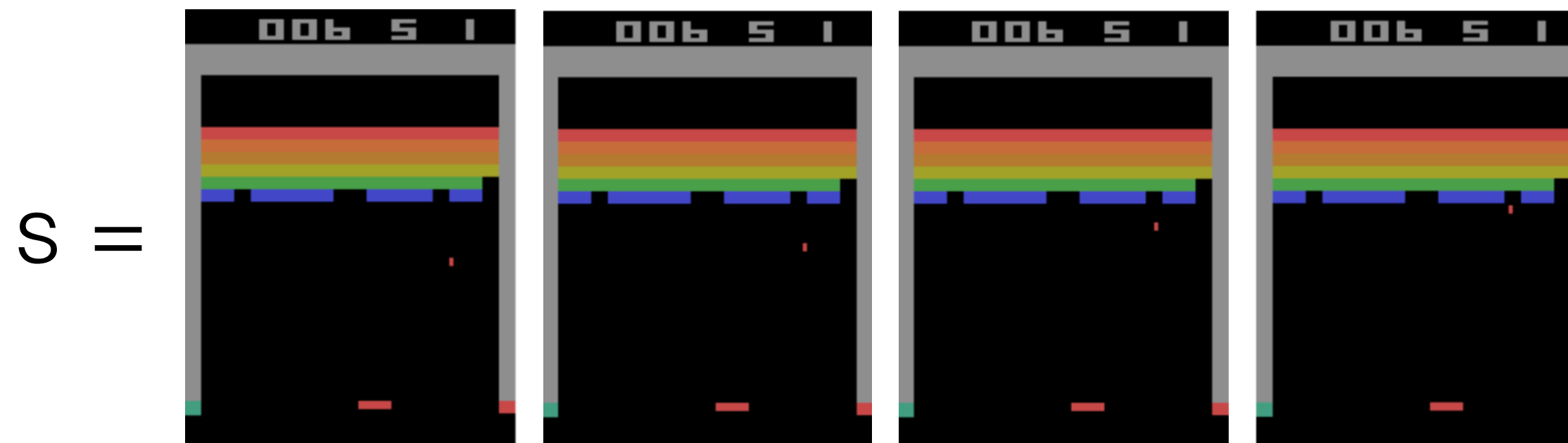
IV. Deep Q-Networks application: Breakout (Atari)

Goal: play breakout, i.e. destroy all the bricks.

Demo



input of Q-network



Output of Q-network

Q-values

$$\begin{pmatrix} Q(s, \leftarrow) \\ Q(s, \rightarrow) \\ Q(s, -) \end{pmatrix}$$

Preprocessing

$\phi(s)$

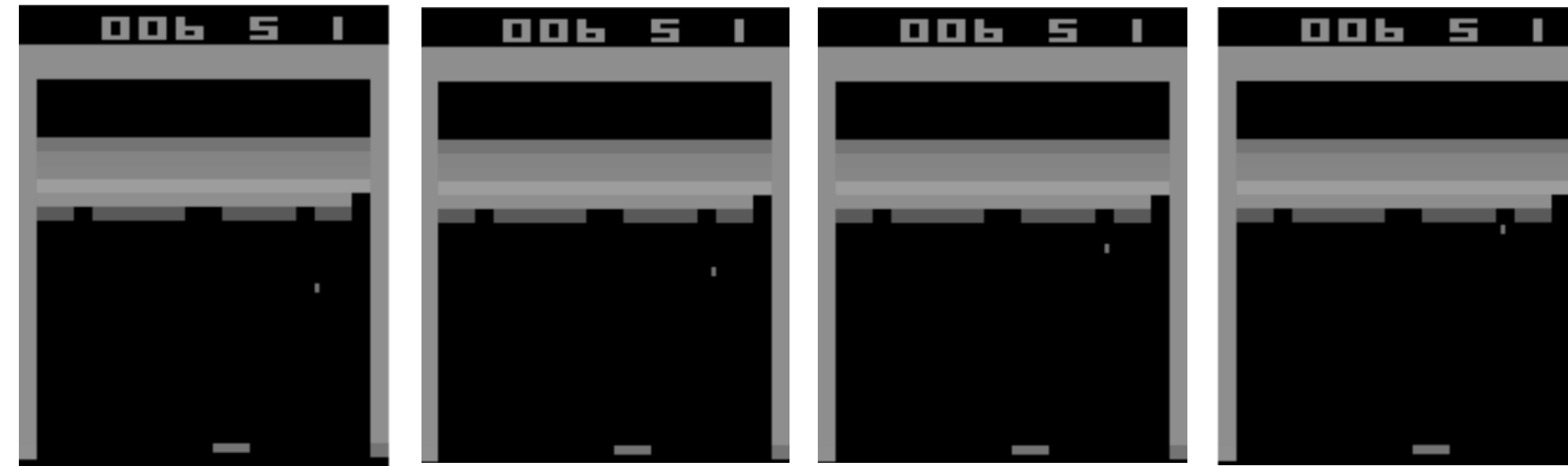
- Convert to grayscale
- Reduce dimensions (h,w)
- History (4 frames)

What is done in preprocessing?

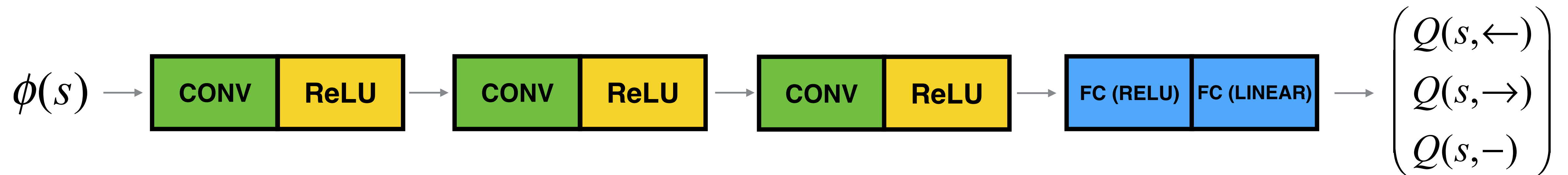
IV. Deep Q-Networks application: Breakout (Atari)

input of Q-network

$\phi(s) =$



Deep Q-network architecture?



Recap' (+ preprocessing + terminal state)

DQN Implementation:

- Initialize your Q-network parameters
- Loop over episodes:
 - Start from initial state ~~s~~ $\phi(s)$
 - Loop over time-steps: $\phi(s)$
 - Forward propagate ~~s~~ in the Q-network $\phi(s)$
 - Execute action a (that has the maximum $Q(\phi(s), a)$ output of Q-network)
 - Observe rewards r and next state s'
 - **Use s' to create $\phi(s')$**
 - Compute targets y by forward propagating state ~~s~~ in the Q-network, then compute loss. $\phi(s')$
 - Update parameters with gradient descent

Some training challenges:

- Keep track of terminal step
- Experience replay
- Epsilon greedy action choice (Exploration / Exploitation tradeoff)

Recap' (+ preprocessing + terminal state)

DQN Implementation:

- Initialize your Q-network parameters
- Loop over episodes:
 - Start from initial state ~~s~~ $\phi(s)$
 - Loop over time-steps: $\phi(s)$
 - Forward propagate ~~s~~ in the Q-network $\phi(s)$
 - Execute action a (that has the maximum $Q(\phi(s), a)$ output of Q-network)
 - Observe rewards r and next state s'
 - **Use s' to create $\phi(s')$**
 - Compute targets y by forward propagating state ~~s~~ in the Q-network, then compute loss. $\phi(s')$
 - Update parameters with gradient descent

Some training challenges:

- Keep track of terminal step
- Experience replay
- Epsilon greedy action choice (Exploration / Exploitation tradeoff)

Recap' (+ preprocessing + terminal state)

DQN Implementation:

- Initialize your Q-network parameters
- Loop over episodes:
 - Start from initial state ~~s~~ $\phi(s)$
 - **Create a boolean to detect terminal states: $terminal = False$**
 - Loop over time-steps:
 - Forward propagate ~~s~~ in the Q-network $\phi(s)$
 - Execute action a (that has the maximum $Q(\text{ ~~$s$~~ , a)}$ output of Q-network)
 - Observe rewards r and next state s'
 - Use s' to create $\phi(s')$
 - **Check if s' is a terminal state.** Compute targets y by forward propagating state ~~s~~ in the Q-network, then compute loss.
 - Update parameters with gradient descent

Some training challenges:

- Keep track of terminal step
- Experience replay
- Epsilon greedy action choice (Exploration / Exploitation tradeoff)

$$\begin{cases} \text{if } terminal = False & : y = r + \gamma \max_{a'} (Q(s', a')) \\ \text{if } terminal = True & : y = r \quad (break) \end{cases}$$

IV - DQN training challenges

Experience replay

1 experience (leads to one iteration of gradient descent)

Current method is to start from initial state s and follow:

E1 $\phi(s) \rightarrow a \rightarrow r \rightarrow \phi(s')$
E2 $\phi(s') \rightarrow a' \rightarrow r' \rightarrow \phi(s'')$
E3 $\phi(s'') \rightarrow a'' \rightarrow r'' \rightarrow \phi(s''')$
...

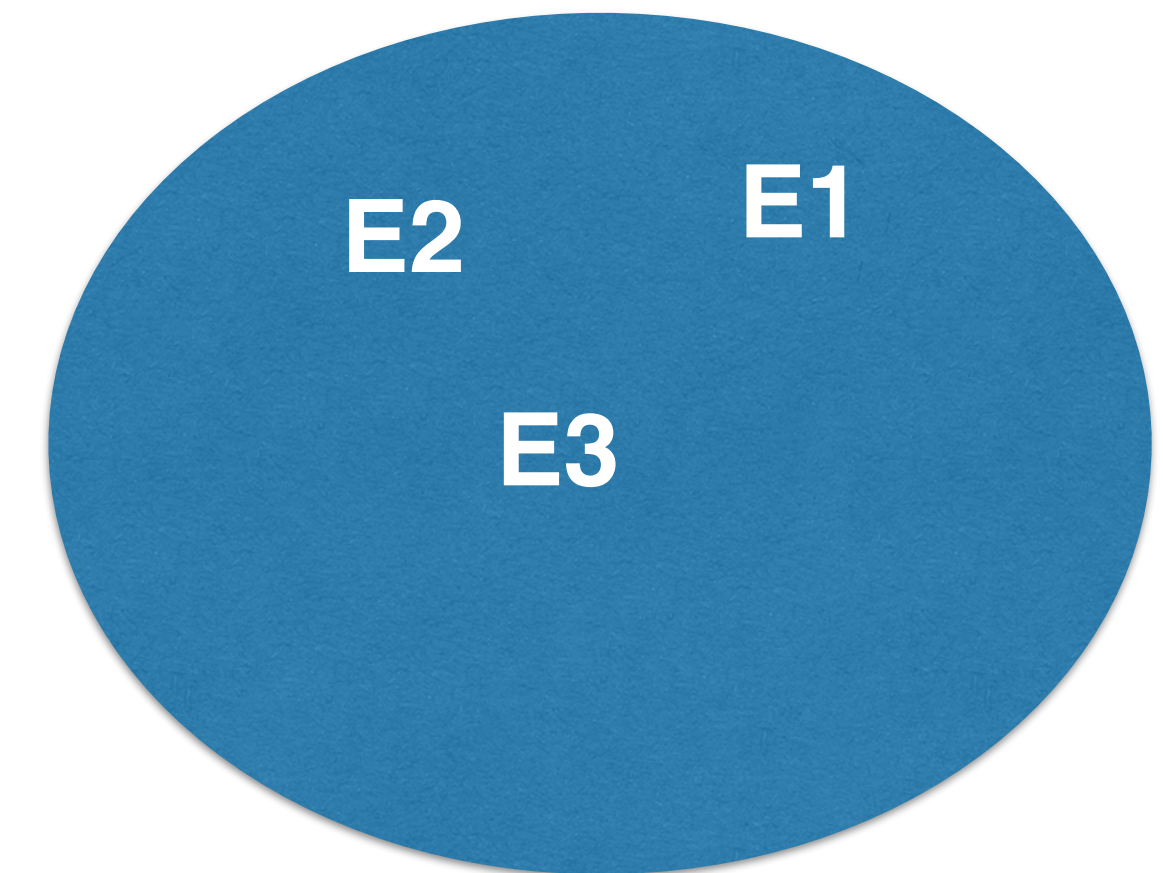
Experience Replay

E1

E2

E3

...



Replay memory (D)

Training: $E1 \rightarrow E2 \rightarrow E3$

Training: $E1 \rightarrow \text{sample}(E1, E2) \rightarrow \text{sample}(E1, E2, E3) \rightarrow \text{sample}(E1, E2, E3, E4) \rightarrow \dots$

Can be used with mini batch gradient descent

Advantages of experience replay?

Recap' (+ experience replay)

DQN Implementation:

- Initialize your Q-network parameters
- **Initialize replay memory D**
- Loop over episodes:
 - Start from initial state $\phi(s)$
 - Create a boolean to detect terminal states: $terminal = False$
 - Loop over time-steps:
 - Forward propagate $\phi(s)$ in the Q-network
 - Execute action a (that has the maximum $Q(\phi(s), a)$ output of Q-network)
 - Observe rewards r and next state s'
 - Use s' to create $\phi(s')$
 - **Add experience $(\phi(s), a, r, \phi(s'))$ to replay memory (D)**
- **Sample random mini-batch of transitions from D**
 - Check if s' is a terminal state. Compute targets y by forward propagating state $\phi(s')$ in the Q-network, then compute loss.
 - Update parameters with gradient descent

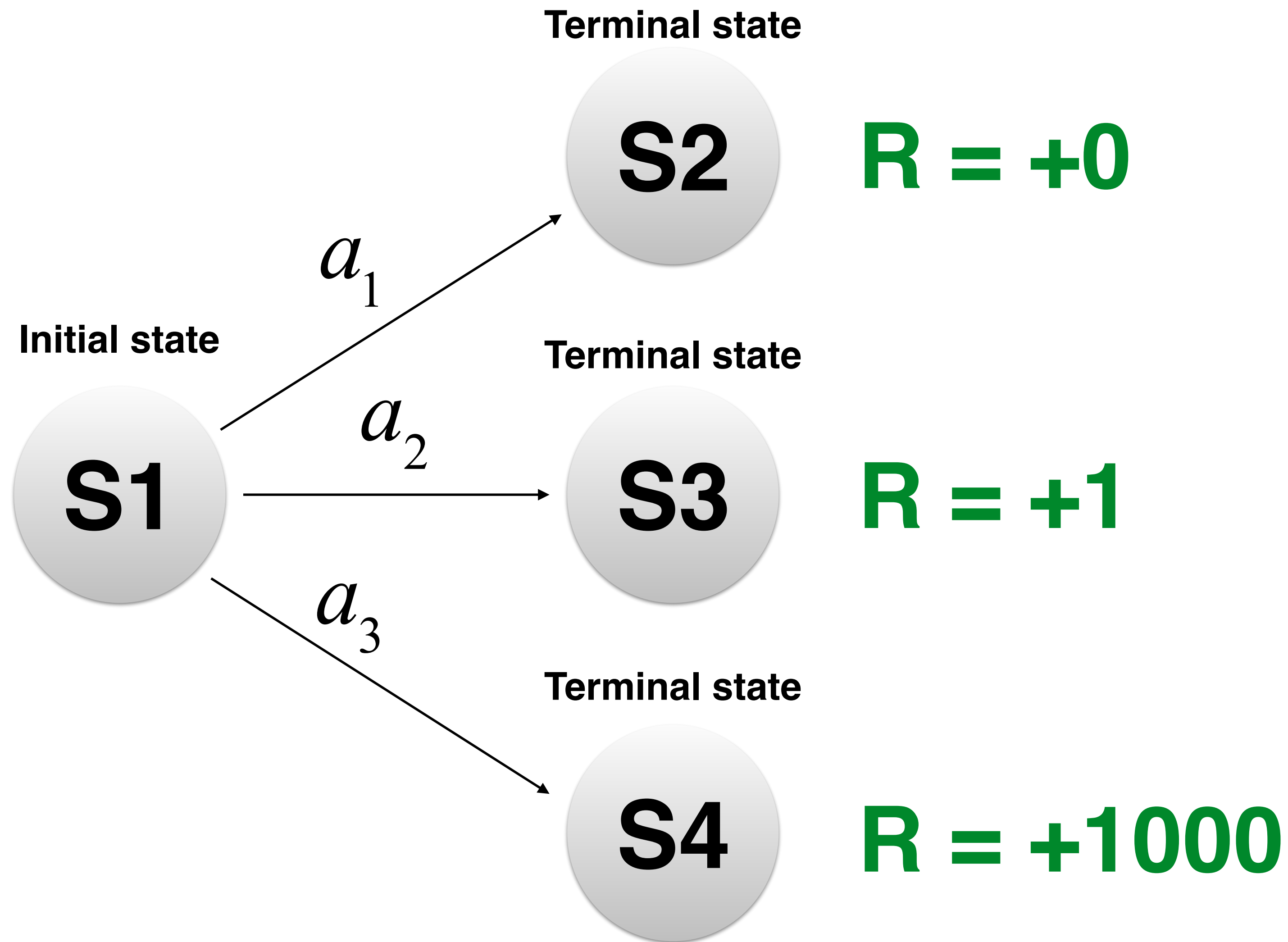
Some training challenges:

- Keep track of terminal step
- Experience replay
- Epsilon greedy action choice (Exploration / Exploitation tradeoff)

The transition resulting from this is added to D , and will not always be used in this iteration's update!

Update using sampled transitions

Exploration vs. Exploitation



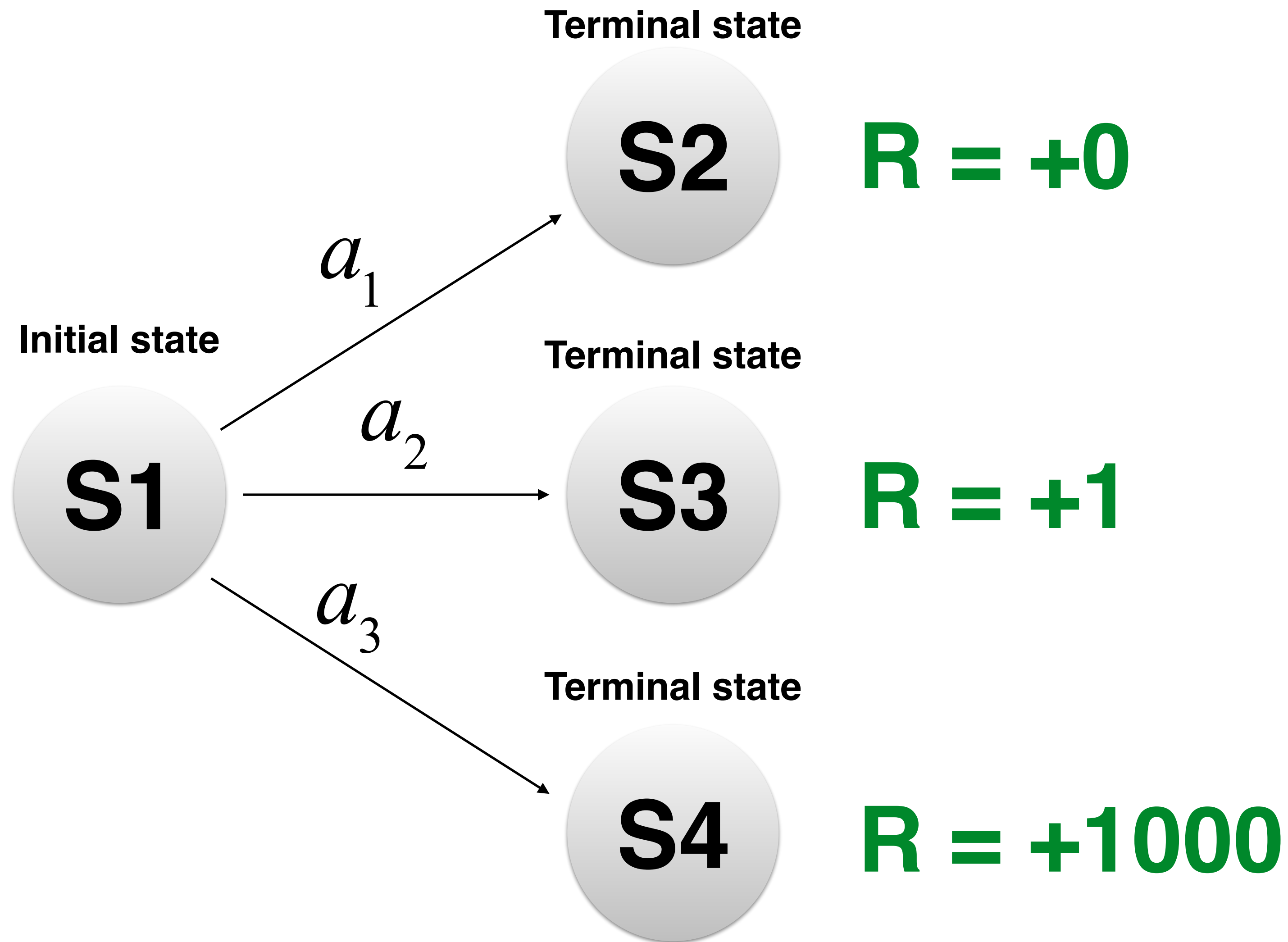
Just after initializing the Q-network, we get:

$$Q(S1, a_1) = 0.5$$

$$Q(S1, a_2) = 0.4$$

$$Q(S1, a_3) = 0.3$$

Exploration vs. Exploitation



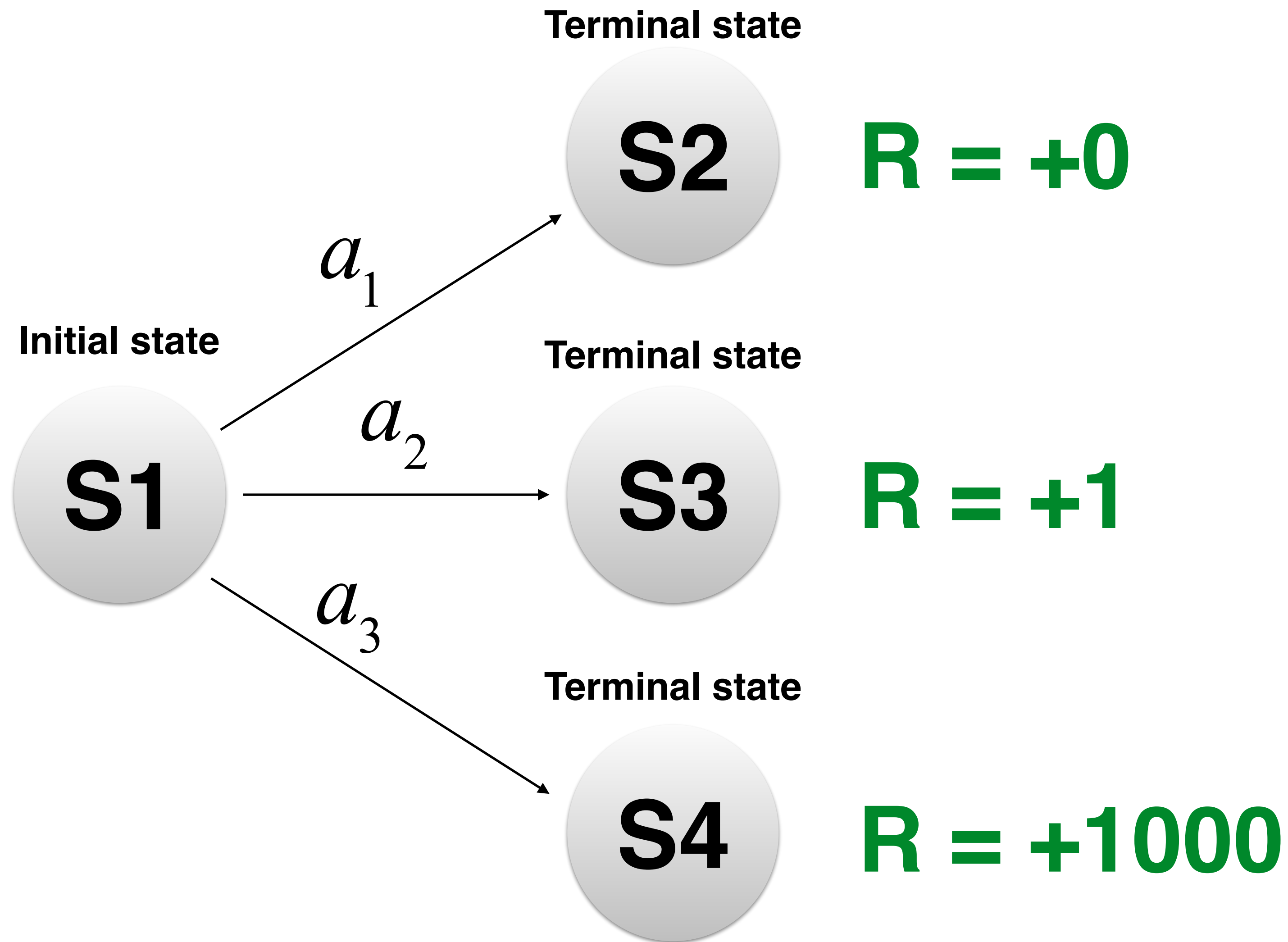
Just after initializing the Q-network, we get:

$$Q(S1, a_1) = \cancel{0.5} \quad 0$$

$$Q(S1, a_2) = 0.4$$

$$Q(S1, a_3) = 0.3$$

Exploration vs. Exploitation



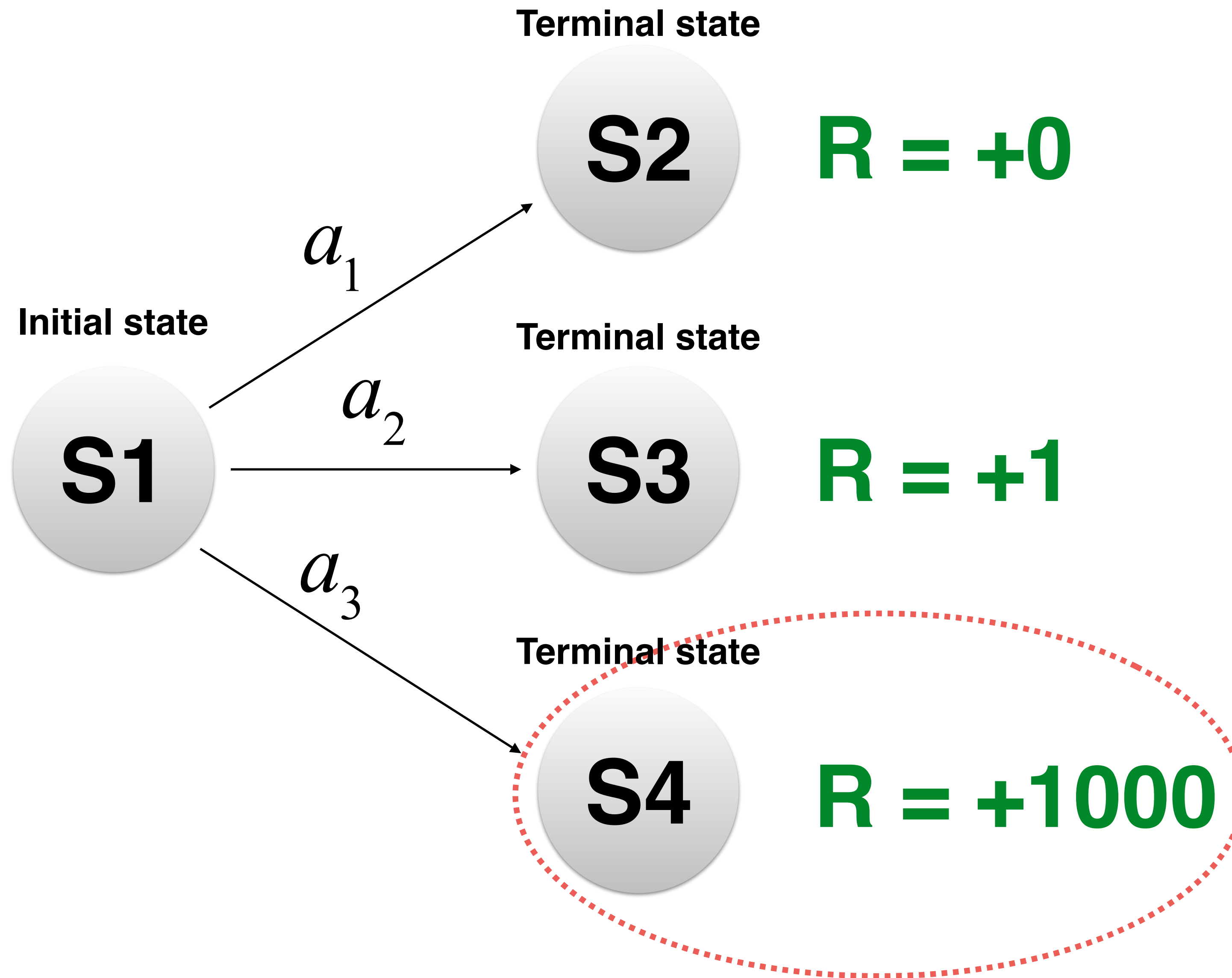
Just after initializing the Q-network, we get:

$$Q(S1, a_1) = \cancel{0.5} \quad 0$$

$$Q(S1, a_2) = \cancel{0.4} \quad 1$$

$$Q(S1, a_3) = 0.3$$

Exploration vs. Exploitation



Just after initializing the Q-network, we get:

$$Q(S1, a_1) = \cancel{0.5} \quad 0$$

$$Q(S1, a_2) = \cancel{0.4} \quad 1$$

$$Q(S1, a_3) = 0.3$$

Will never be visited, because $Q(S1, a_3) < Q(S1, a_2)$

Recap' (+ epsilon greedy action)

DQN Implementation:

- Initialize your Q-network parameters
- Initialize replay memory D
- Loop over episodes:
 - Start from initial state $\phi(s)$
 - Create a boolean to detect terminal states: $terminal = False$
 - Loop over time-steps:
 - **With probability epsilon, take random action a .**
 - **Otherwise:**
 - Forward propagate $\phi(s)$ in the Q-network
 - Execute action a (that has the maximum $Q(\phi(s), a)$ output of Q-network).
 - Observe rewards r and next state s'
 - Use s' to create $\phi(s')$
 - Add experience $(\phi(s), a, r, \phi(s'))$ to replay memory (D)
 - Sample random mini-batch of transitions from D
 - Check if s' is a terminal state. Compute targets y by forward propagating state $\phi(s')$ in the Q-network, then compute loss.
 - Update parameters with gradient descent

Overall recap'

DQN Implementation:

- Initialize your Q-network parameters
- **Initialize replay memory D**
- Loop over episodes:
 - Start from initial state $\phi(s)$
 - **Create a boolean to detect terminal states: $terminal = False$**
 - Loop over time-steps:
 - **With probability ϵ , take random action a .**
 - **Otherwise:**
 - Forward propagate $\phi(s)$ in the Q-network
 - Execute action a (that has the maximum $Q(\phi(s), a)$ output of Q-network).
 - Observe rewards r and next state s'
 - **Use s' to create $\phi(s')$**
 - **Add experience $(\phi(s), a, r, \phi(s'))$ to replay memory (D)**
 - **Sample random mini-batch of transitions from D**
 - **Check if s' is a terminal state.** Compute targets y by forward propagating state $\phi(s')$ in the Q-network, then compute loss.
 - Update parameters with gradient descent

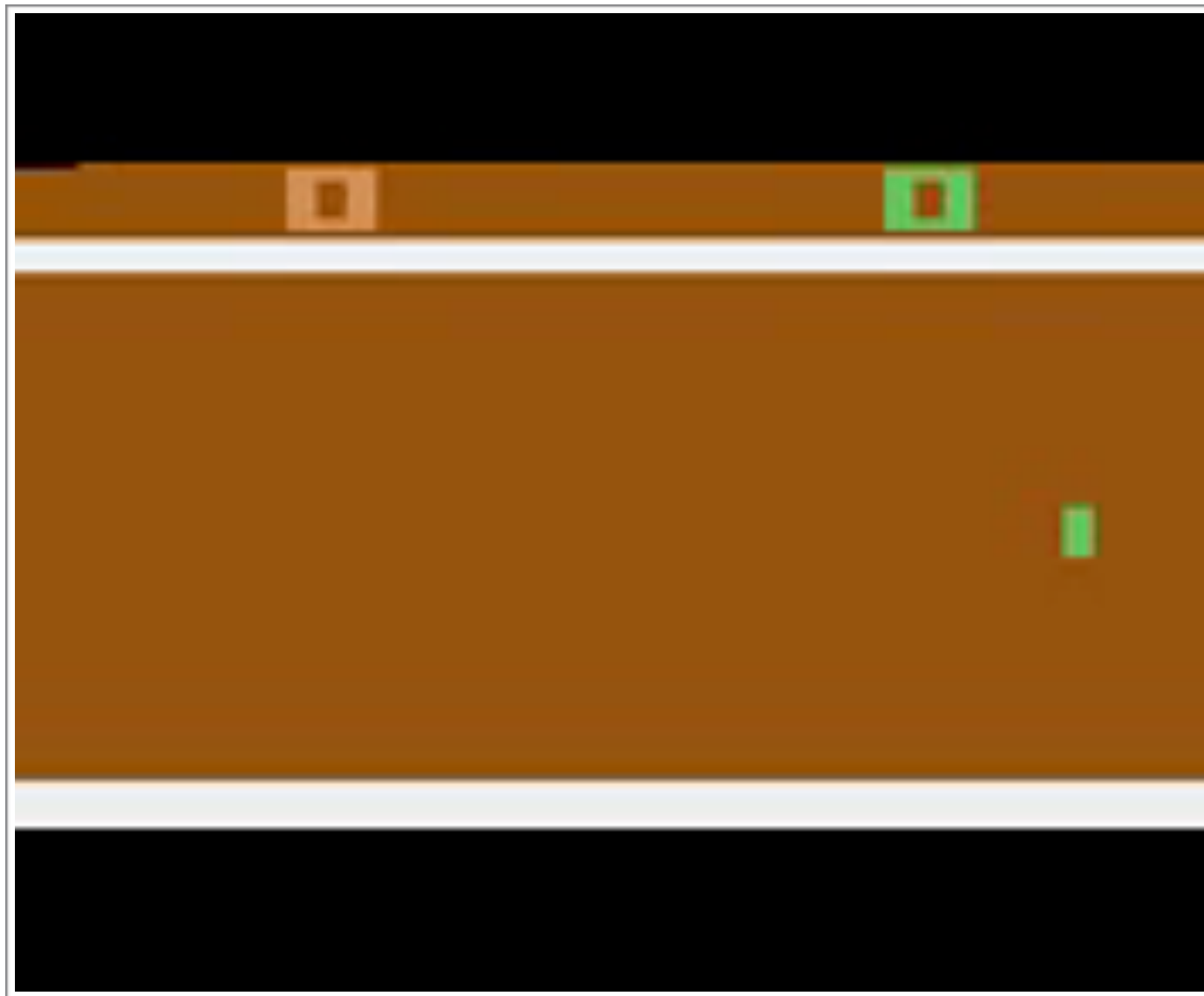
- **Preprocessing**
- **Detect terminal state**
- **Experience replay**
- **Epsilon greedy action**

Results



Other Atari games

Pong



SeaQuest



Space Invaders



[<https://www.youtube.com/watch?v=NirMkC5uvWU>]

[https://www.youtube.com/watch?v=p88R2_3yWPA]

[<https://www.youtube.com/watch?v=W2CAghUiofY&t=2s>]

Today's outline

I. Motivation

II. Recycling is good: an introduction to RL

III. Deep Q-Networks

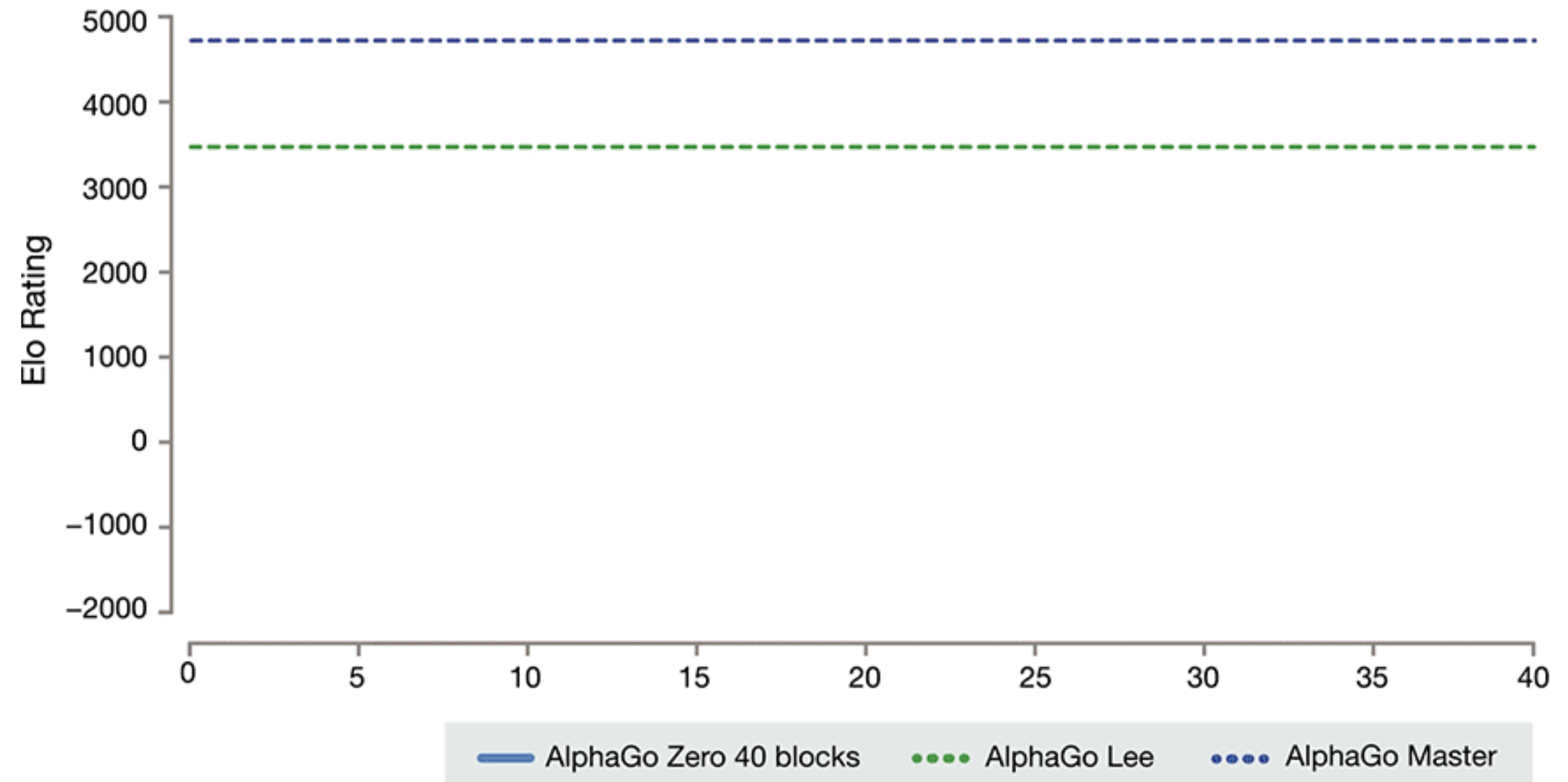
IV. Application of Deep Q-Network: Breakout (Atari)

V. Tips to train Deep Q-Network

VI. Advanced topics

VI - Advanced topics

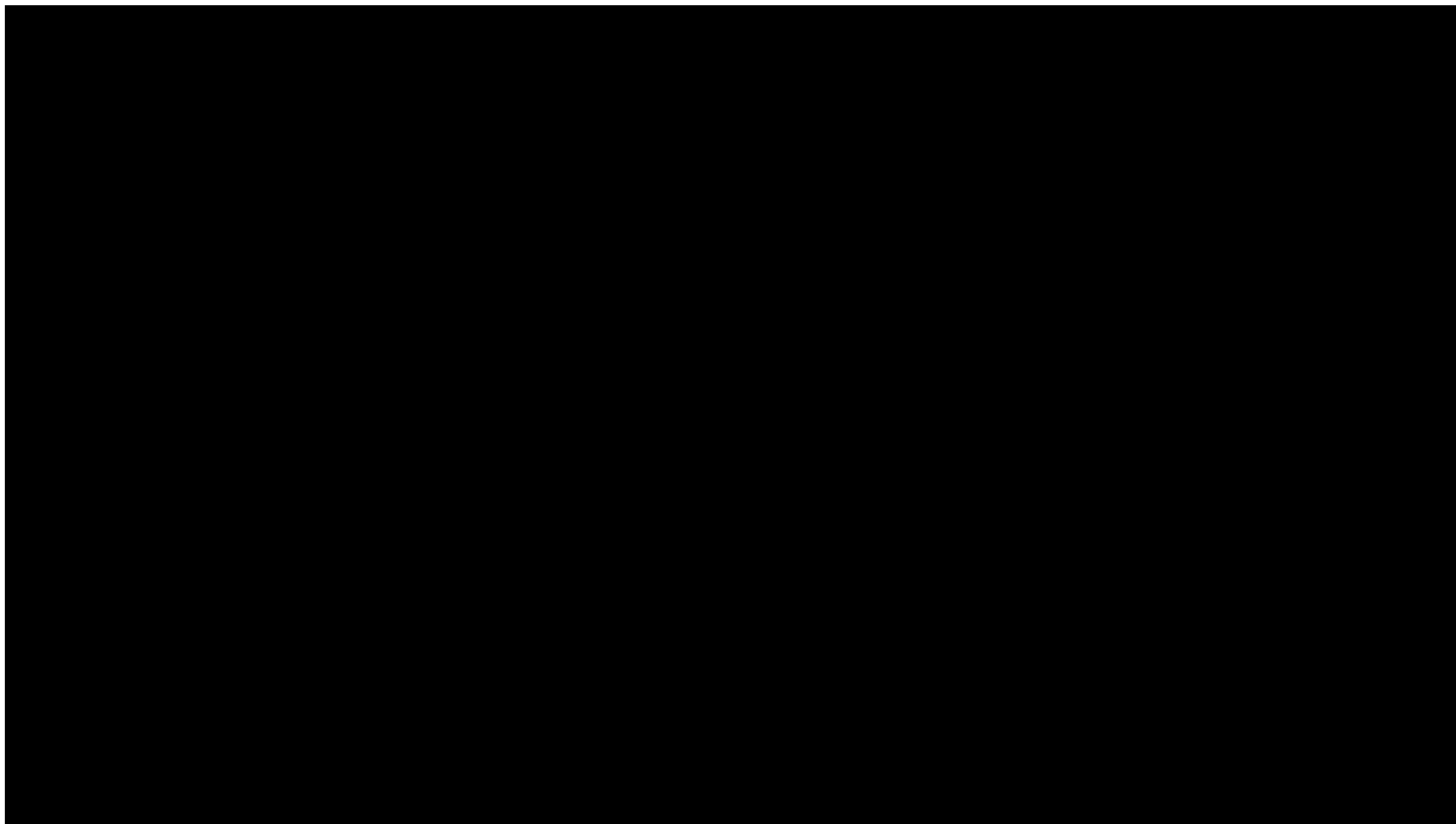
Alpha Go



[[DeepMind Blog](#)]

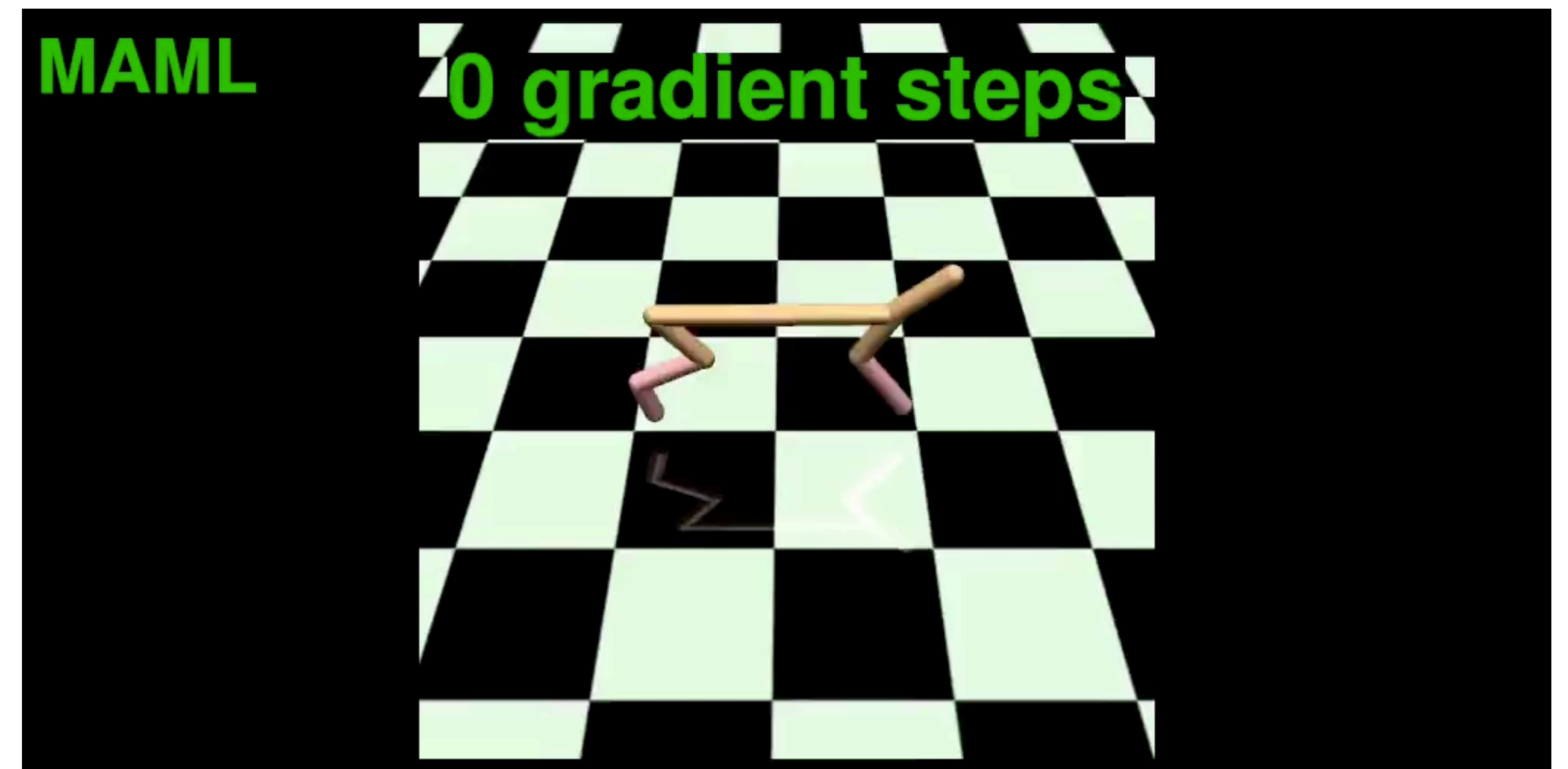
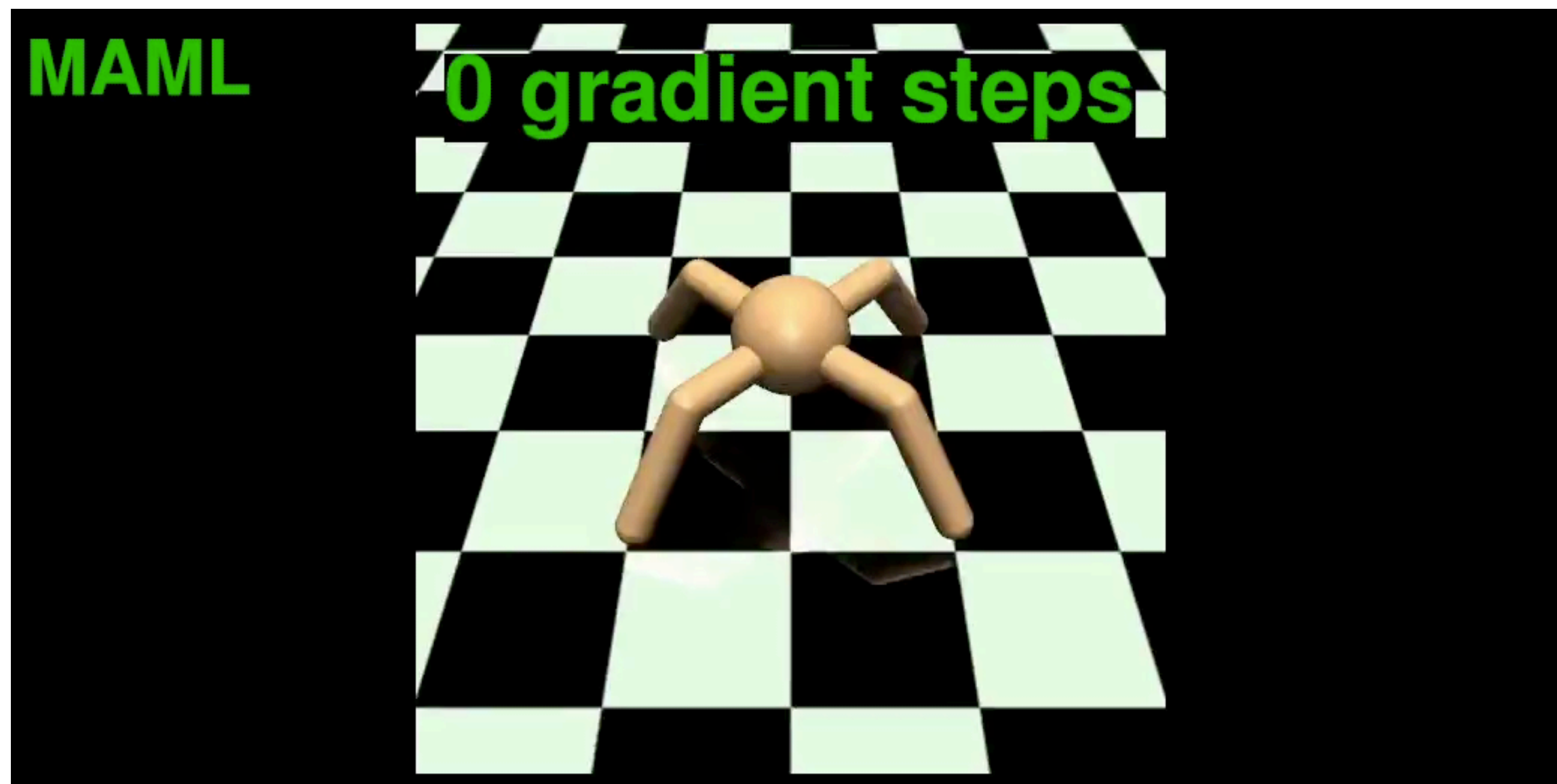
[Silver et al. (2017): Mastering the game of Go without human knowledge]

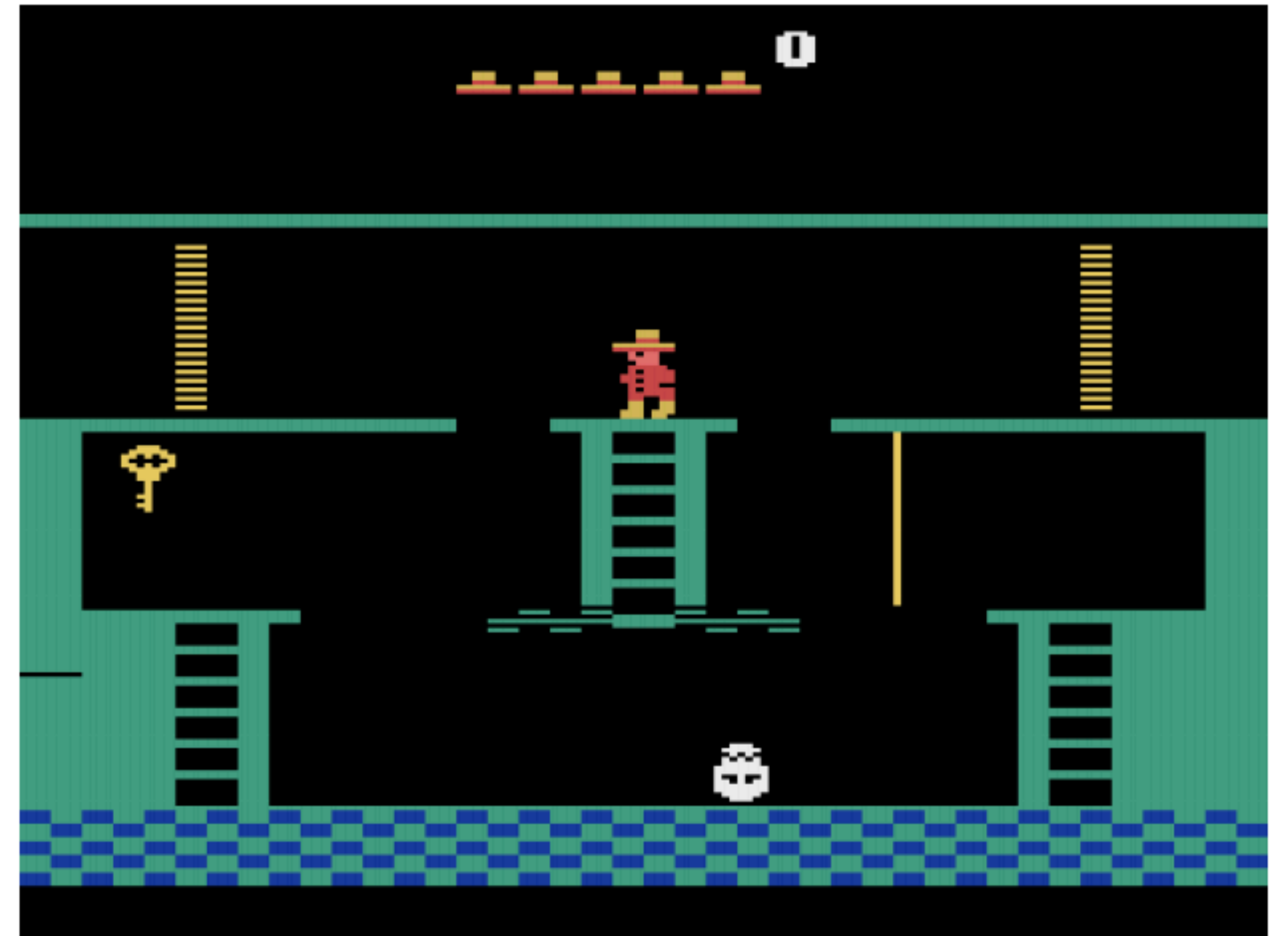
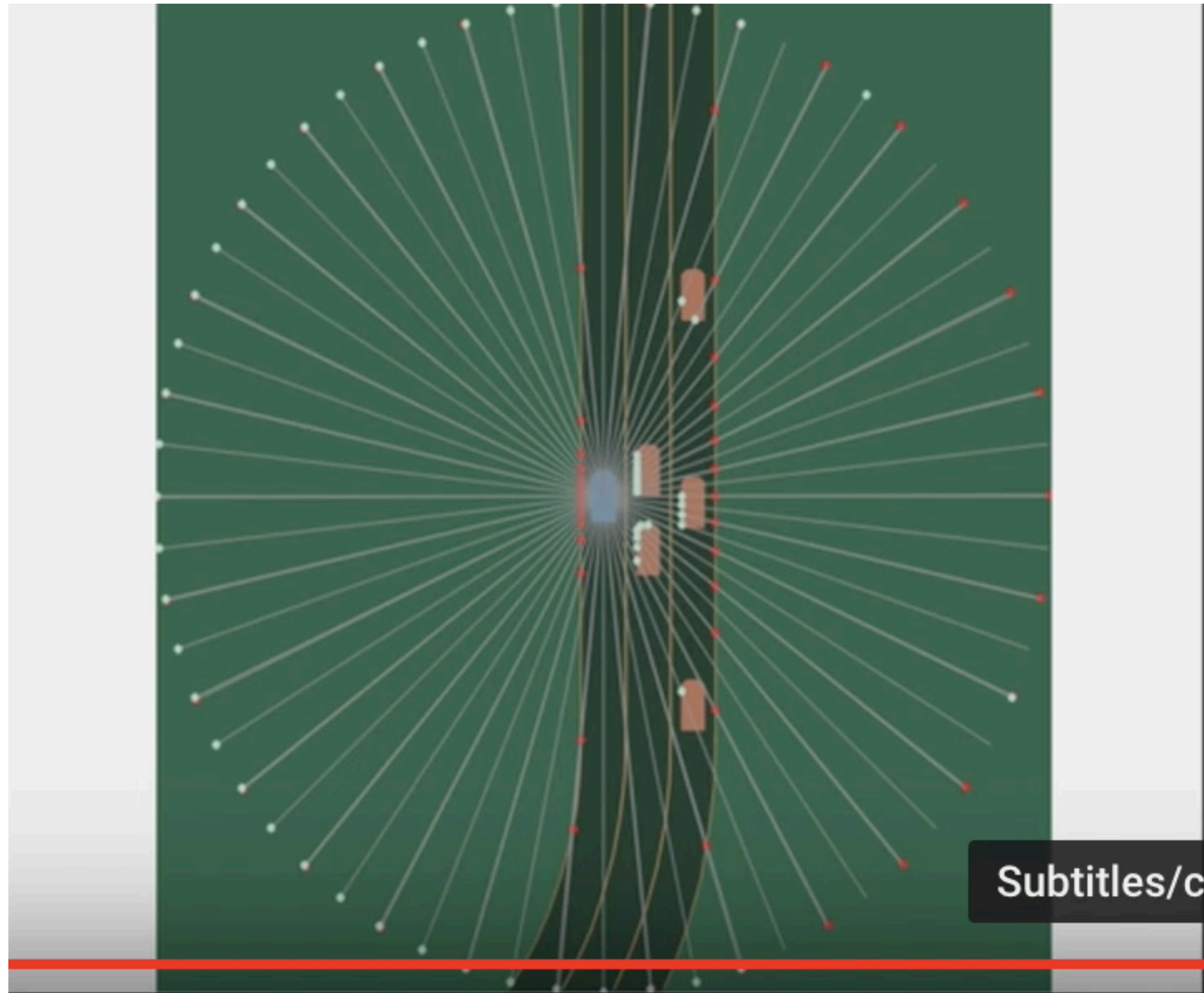
Kian Katanforoosh, Andrew Ng, Younes Bensouda Mourri



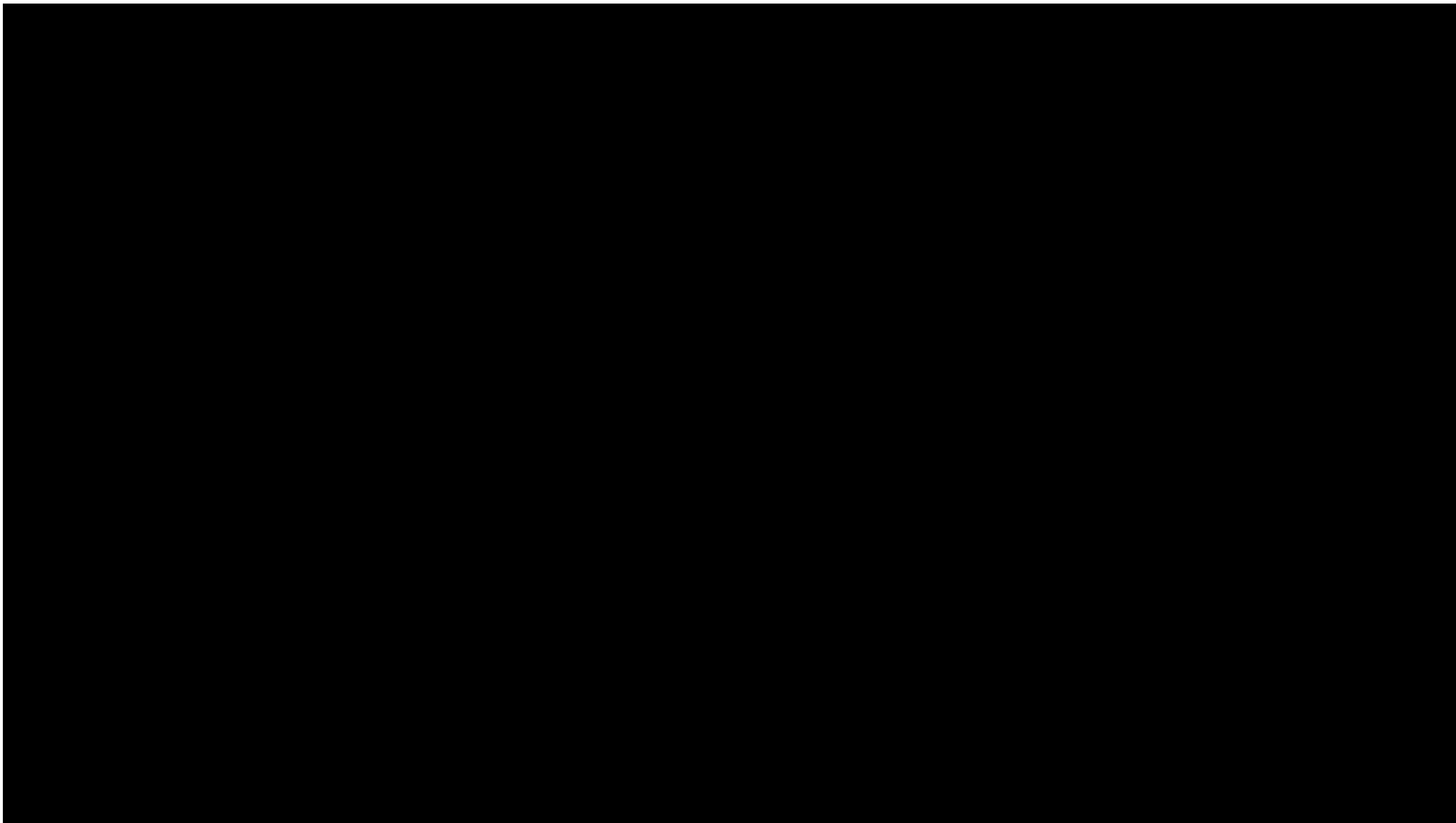
VI - Advanced topics

Meta learning





[Source: Bellemare et al. (2016): Unifying Count-Based Exploration and Intrinsic Motivation]



Announcements

For Tuesday 06/05, 9am:

This Friday:

- **TA Sections:**
 - **How to have a great final project write-up.**
 - **Advices on: How to write a great report.**
 - **Advices on: How to build a super poster.**
 - **Advices on: Final project grading criteria.**
 - **Going through examples of great projects and why they were great.**
 - **Small competitive quiz in section.**