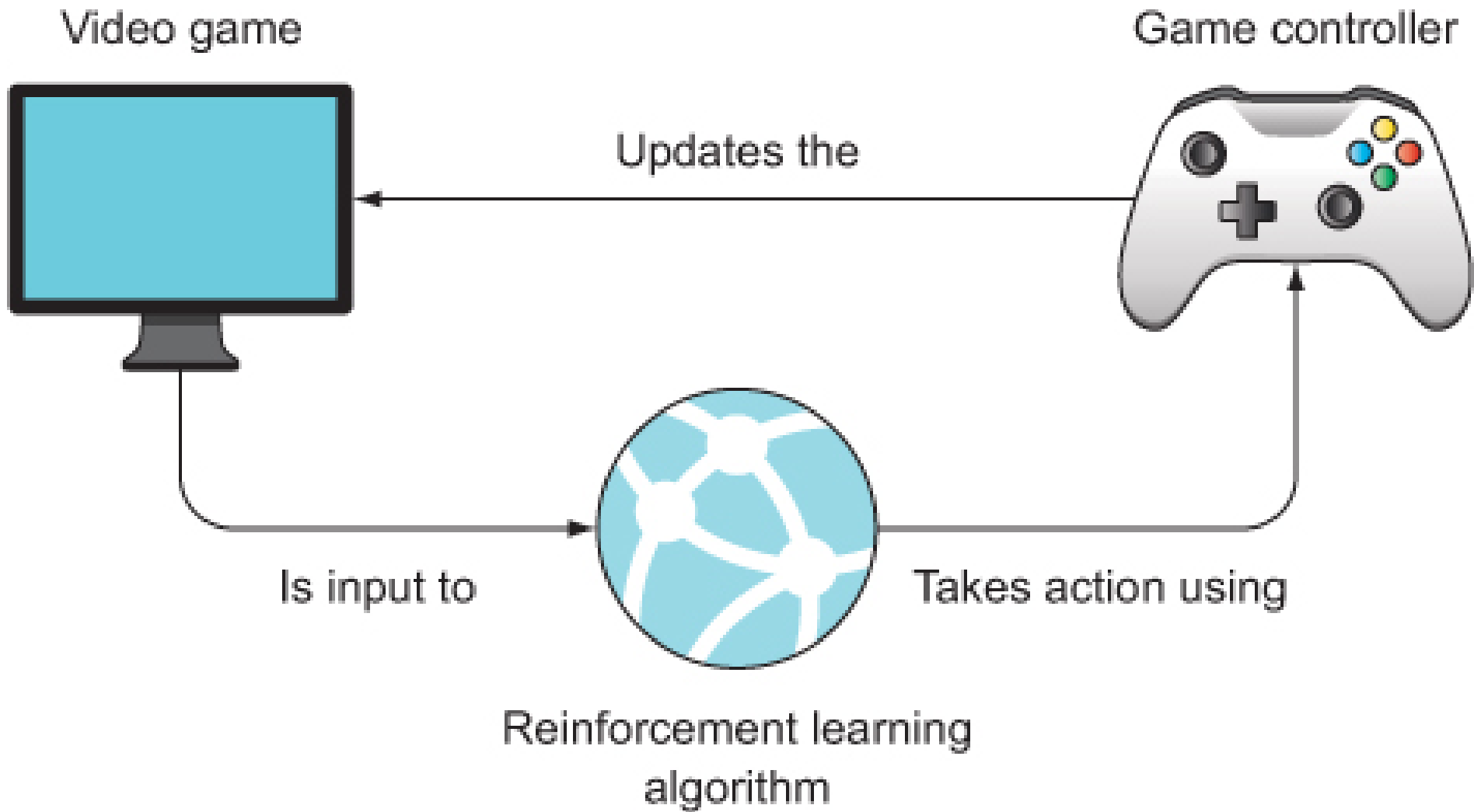# Reinforcement Learning: A Gentle Introduction
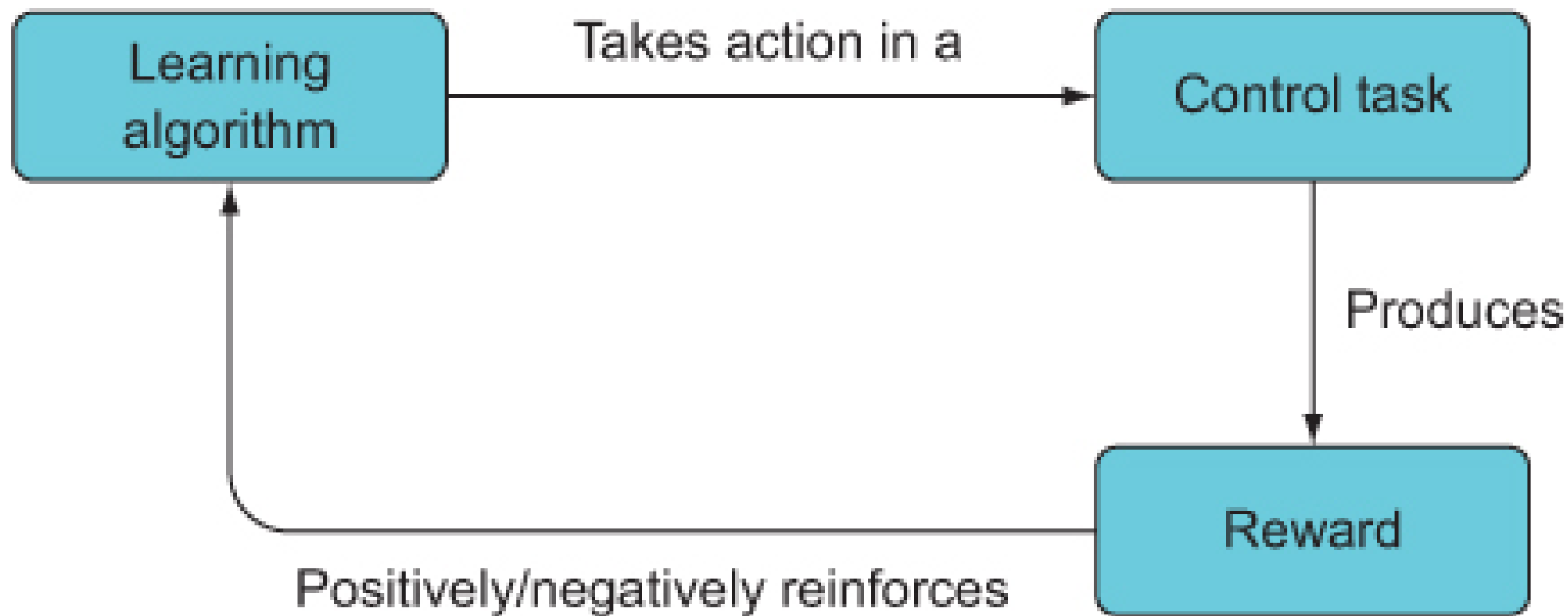
PUBH 7475/8475
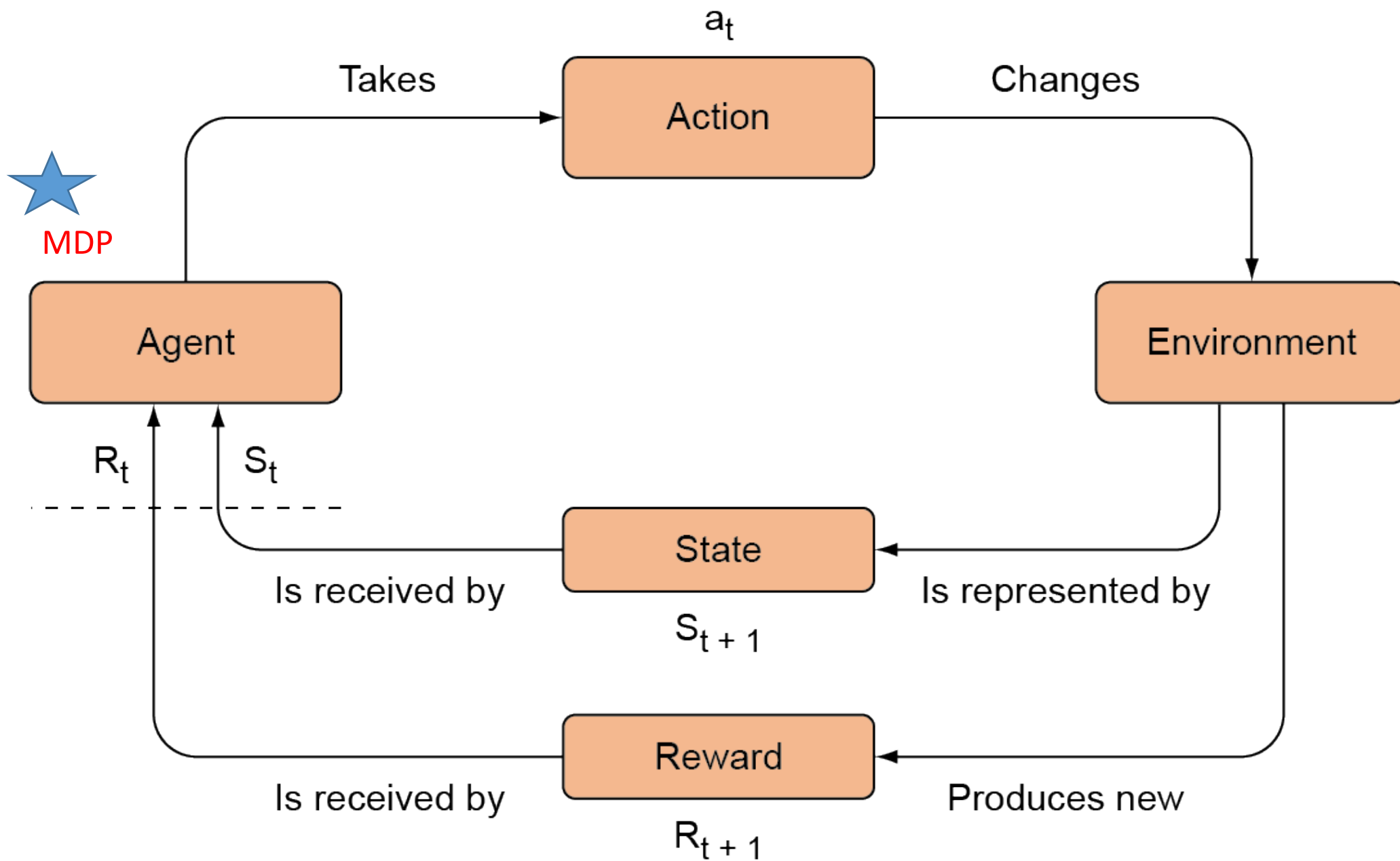
Based on "Deep Reinforcement Learning in Action" by A. Zai & B. Brown

https://www.manning.com/books/deep-reinforcement-learning-in-action

# Outline

- What is RL?
- Q-Learning: DeepMind's DQN
- Policy methods
- Actor-critic methods

- General comments:
  +/-: highly mathematical
  +: most similar to human learning; impressive applications (e.g.
      DeepMind's AlphaGo, AlphaZero); rapid development
  -: time-consuming, not efficient; unstable, hard to train

Video game

Game controller

Updates the

Is input to

Takes action using

Reinforcement learning
algorithm

# The action-value (Q) function

Math

English

$Q_\pi: (s|a) \rightarrow E(R|a,s,\pi),$

$Q_\pi$ is a function that maps a pair, $(s, a)$, of a state, $s$, and an action, $a$, to the expected reward of taking action $a$ in state $s$, given that we're using the policy (or "strategy") $\pi$.

# The policy function

Math

English

$\pi; s \rightarrow Pr(A|s),$ where $s \in S$

A policy, $\pi$, is a mapping from states to the (probabilistically) best actions for those states.
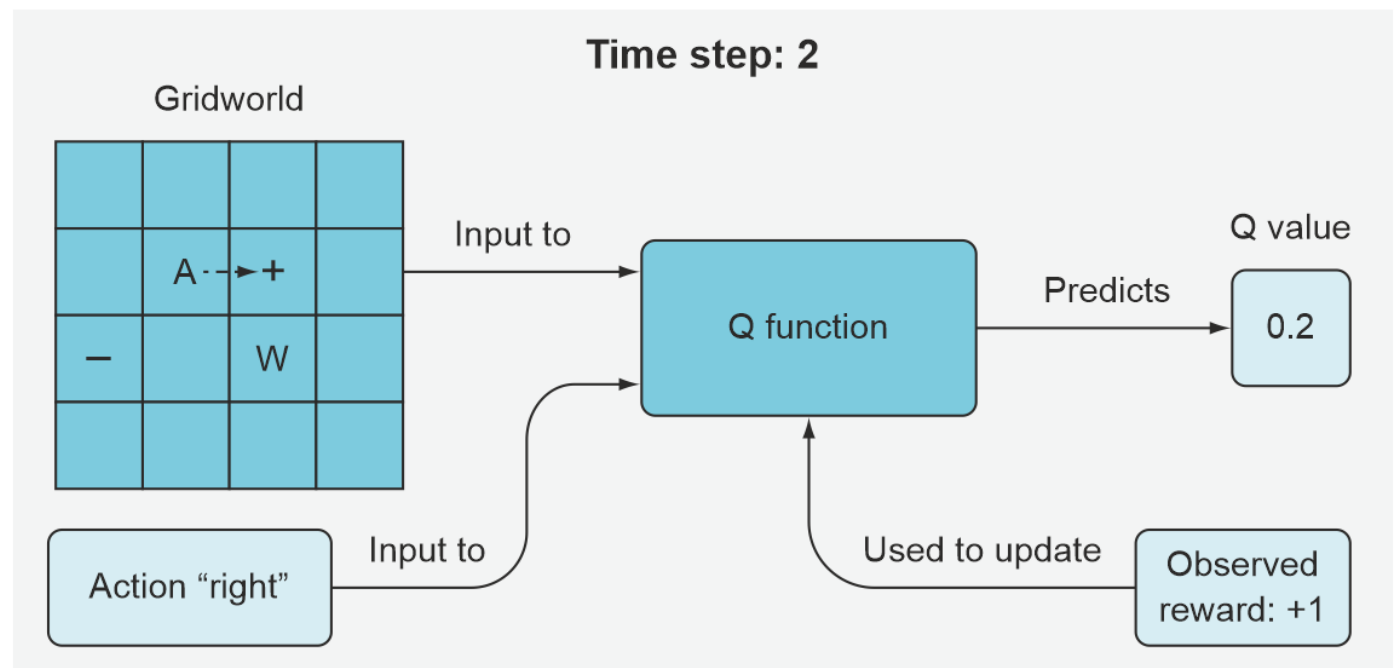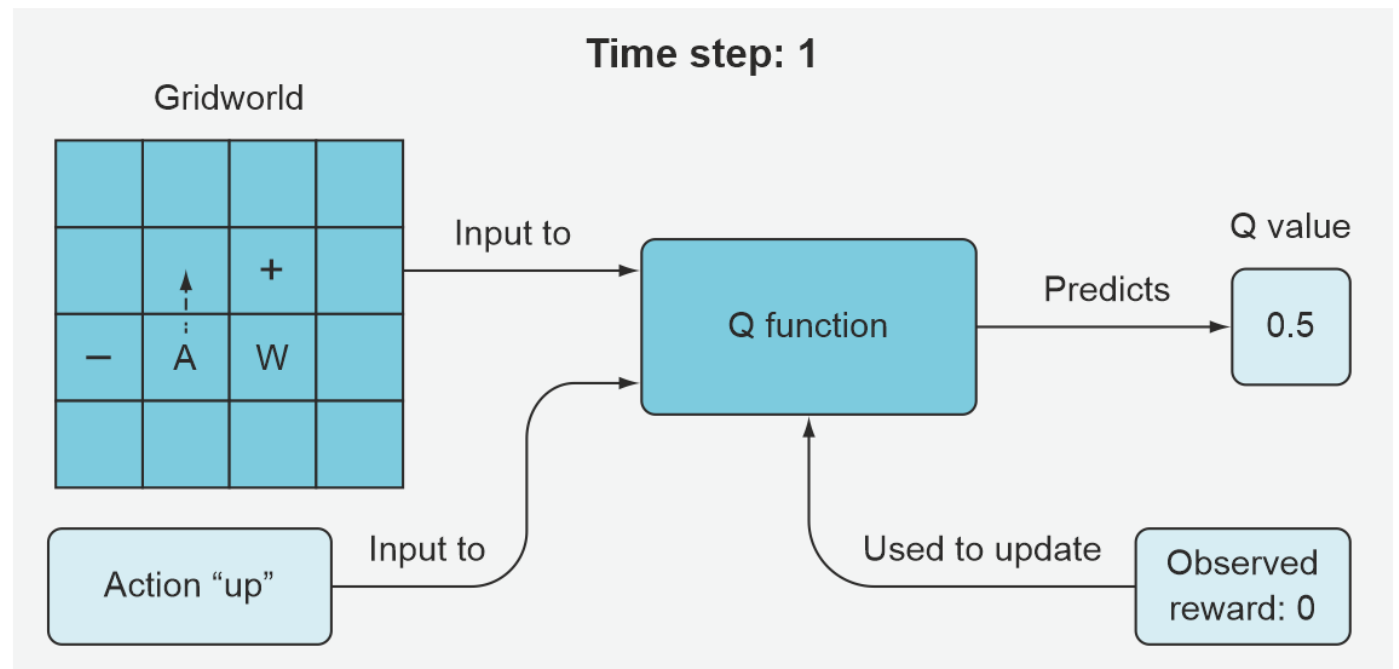
## The state-value function

Math

English

$V_\pi: s \rightarrow E(R|s,\pi),$

A value function, $V_\pi$, is a function that maps a state, $s$, to the expected rewards, given that we start in state $s$ and follow some policy, $\pi$.

# Q learning: learn the Q function



**Time step: 1**

Gridworld

Input to → Q function → Predicts → Q value: 0.5

Action "up" — Input to → Q function

Q function — Used to update → Observed reward: 0

**Time step: 2**

Gridworld

Input to → Q function → Predicts → Q value: 0.2

Action "right" — Input to → Q function

Q function — Used to update → Observed reward: +1

Updated
Q value

Current
Q value

Observed
reward

Max Q value
for all actions

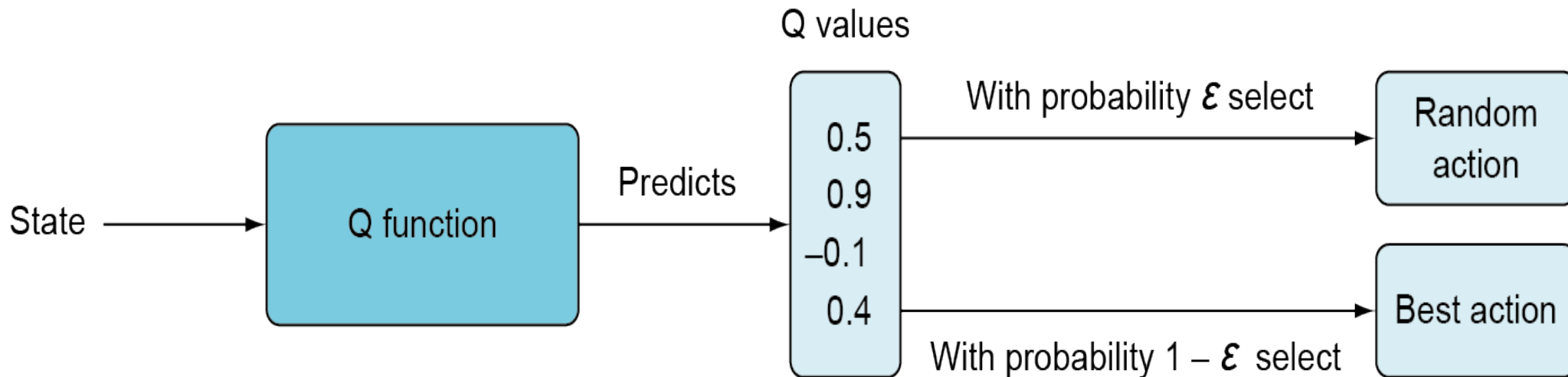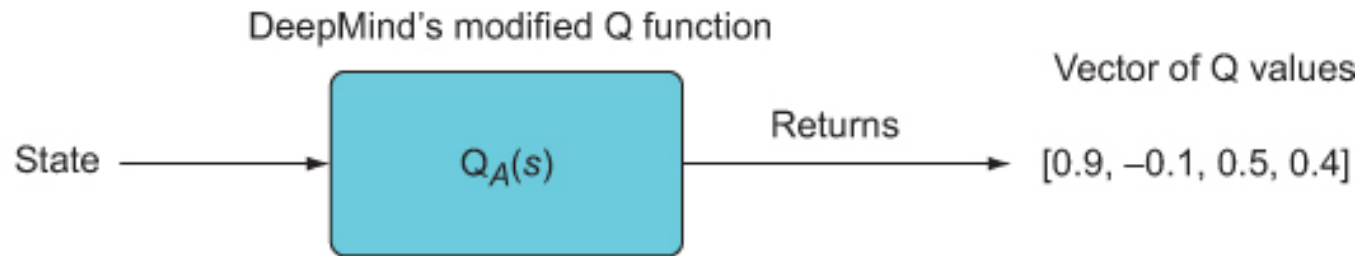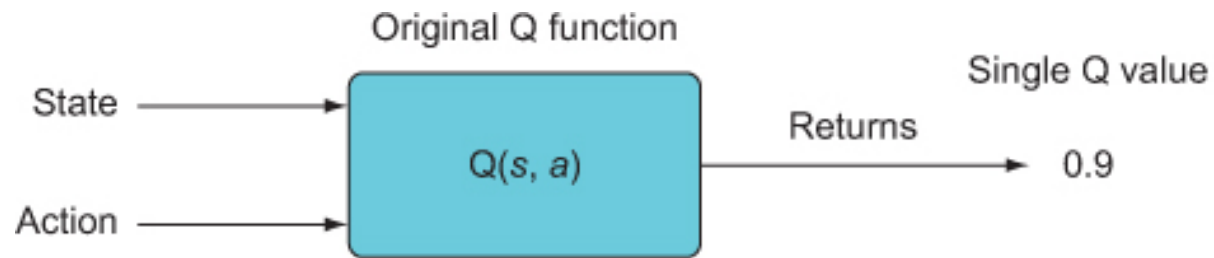$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max Q(S_{t+1}, a) - Q(S_t, A_t)]$$

Step size

Discount
factor
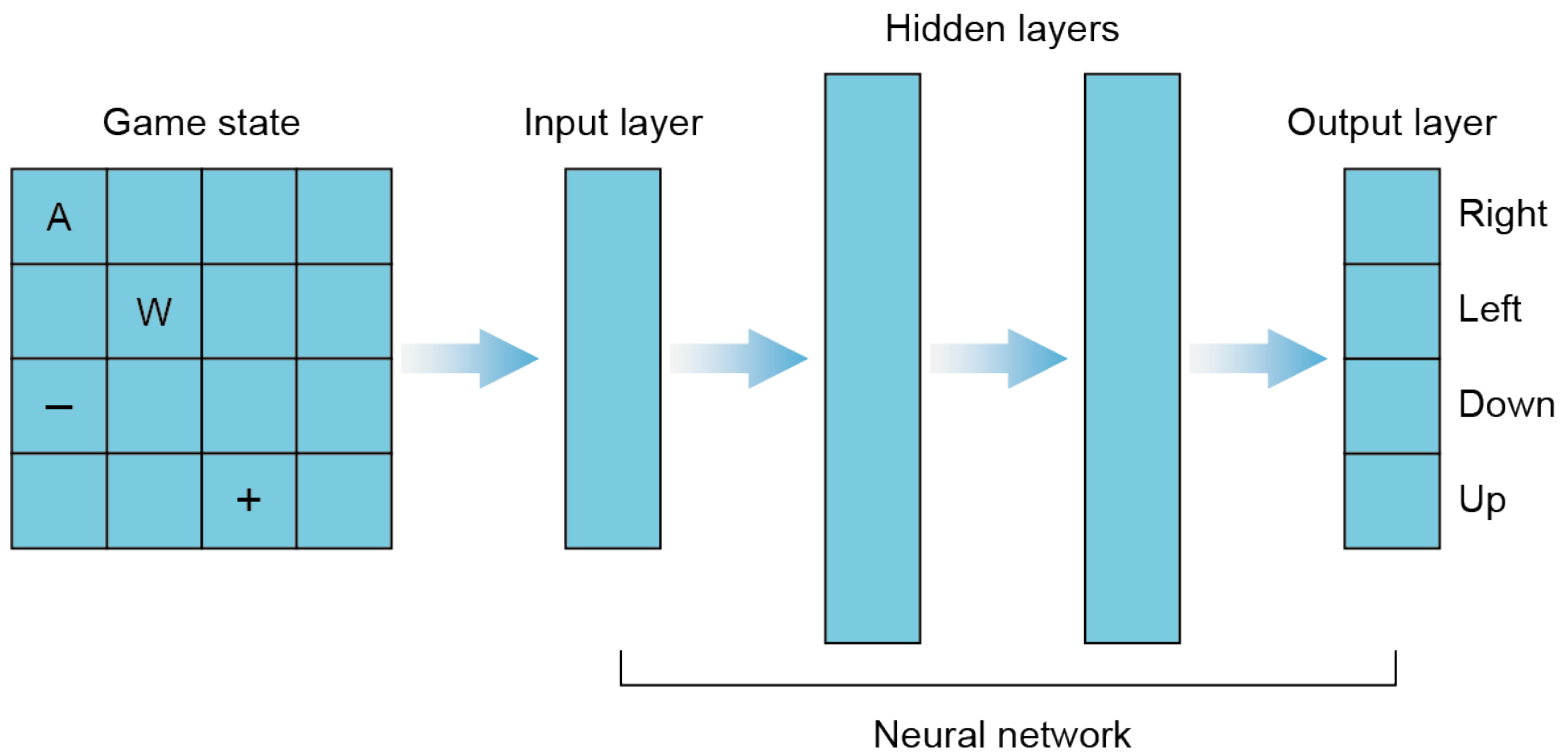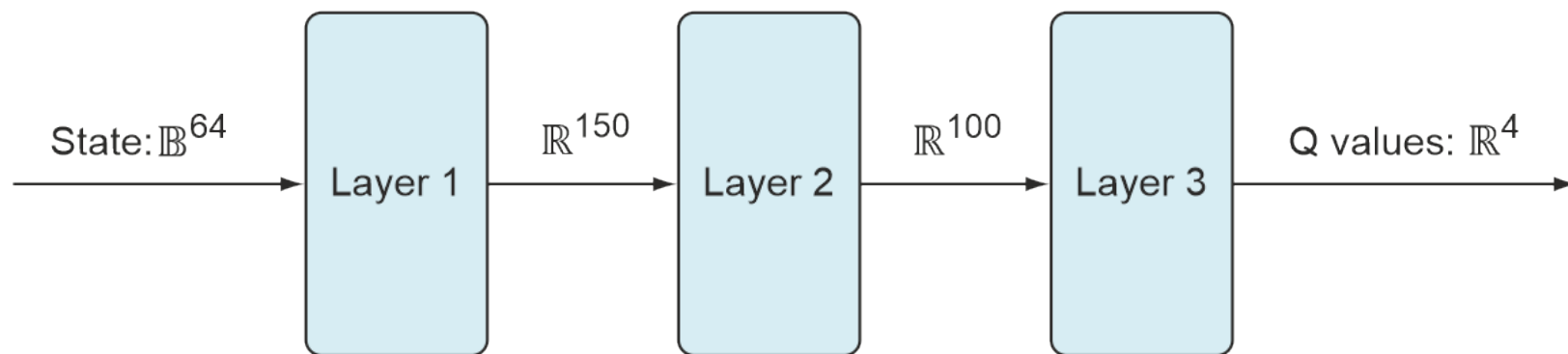
The Bellman equation:

$$Q_\pi(s_t, a_t) \leftarrow r_t + \gamma \times max[Q_\pi(s_{t+1}, a)]$$

Original Q function

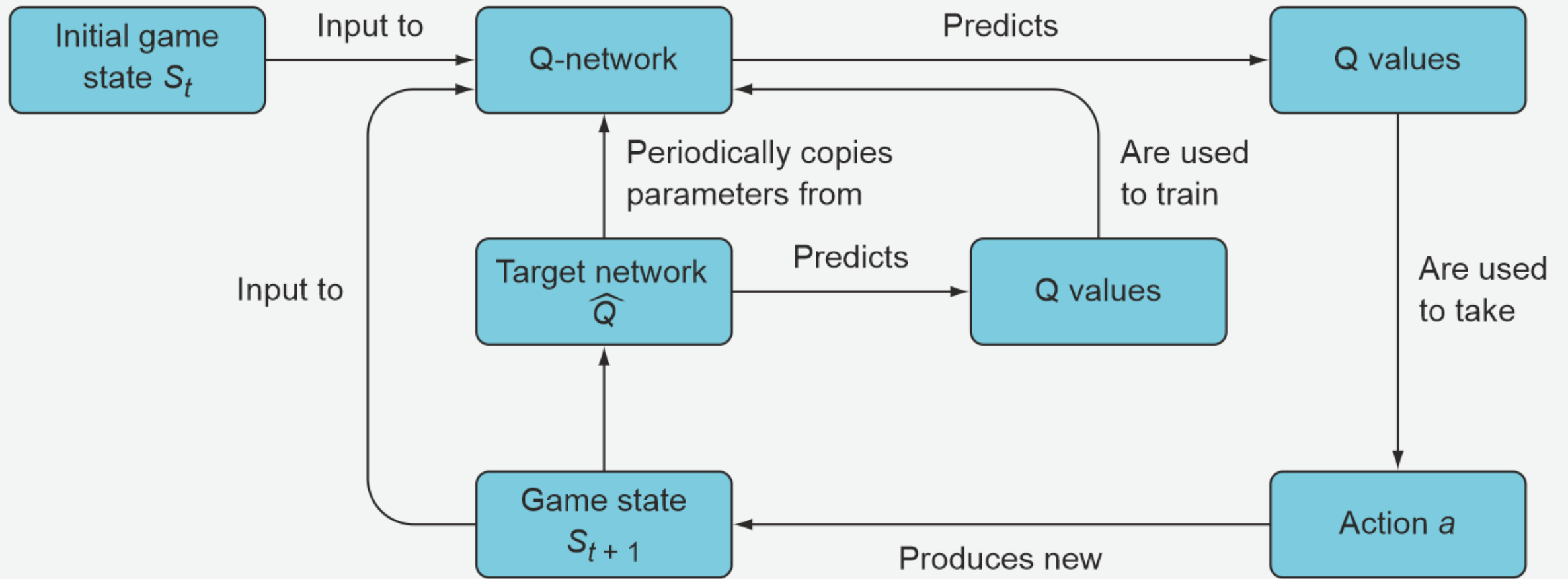State ⟶ [ Q(s, a) ] ⟶ Returns ⟶ 0.9

Single Q value

Action ⟶

DeepMind's modified Q function

State ⟶ [ $Q_A(s)$ ] ⟶ Returns ⟶ [0.9, –0.1, 0.5, 0.4]

Vector of Q values

Q values

State ⟶ [ Q function ] ⟶ Predicts ⟶ [ 0.5 / 0.9 / –0.1 / 0.4 ]

With probability $\varepsilon$ select ⟶ [ Random action ]

With probability $1 – \varepsilon$ select ⟶ [ Best action ]

Game state

Input layer

Hidden layers

Output layer

A

W

−

+

Right

Left

Down

Up

Neural network

**Deep Q-network**

State: $\mathbb{B}^{64}$

Layer 1

$\mathbb{R}^{150}$

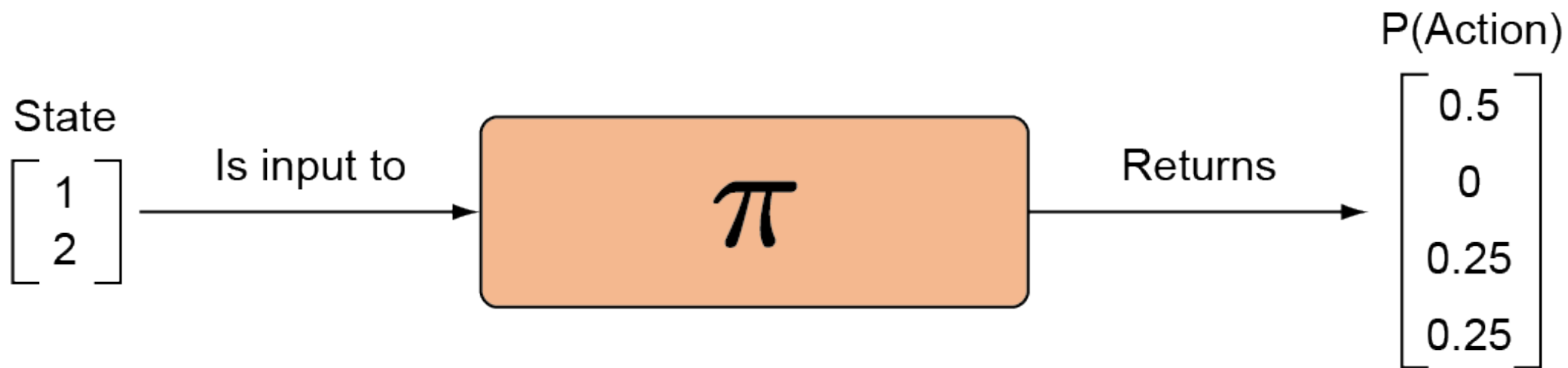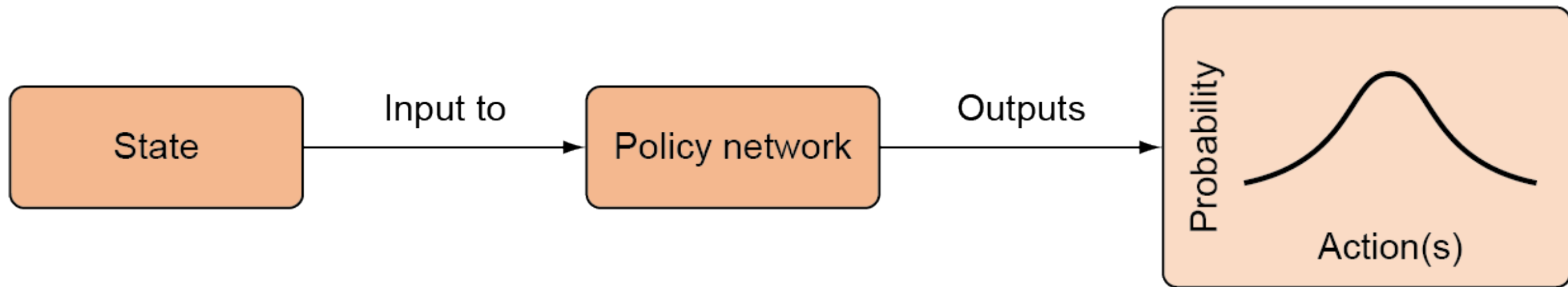Layer 2

$\mathbb{R}^{100}$

Layer 3

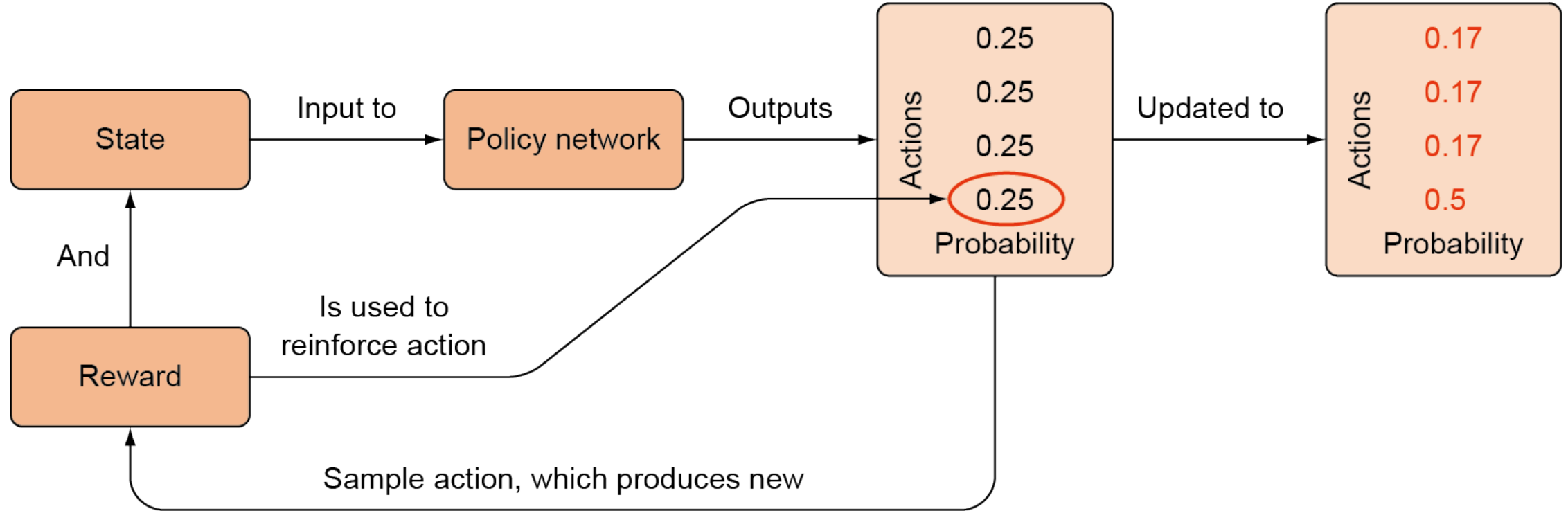Q values: $\mathbb{R}^{4}$

Q-learning with a target network

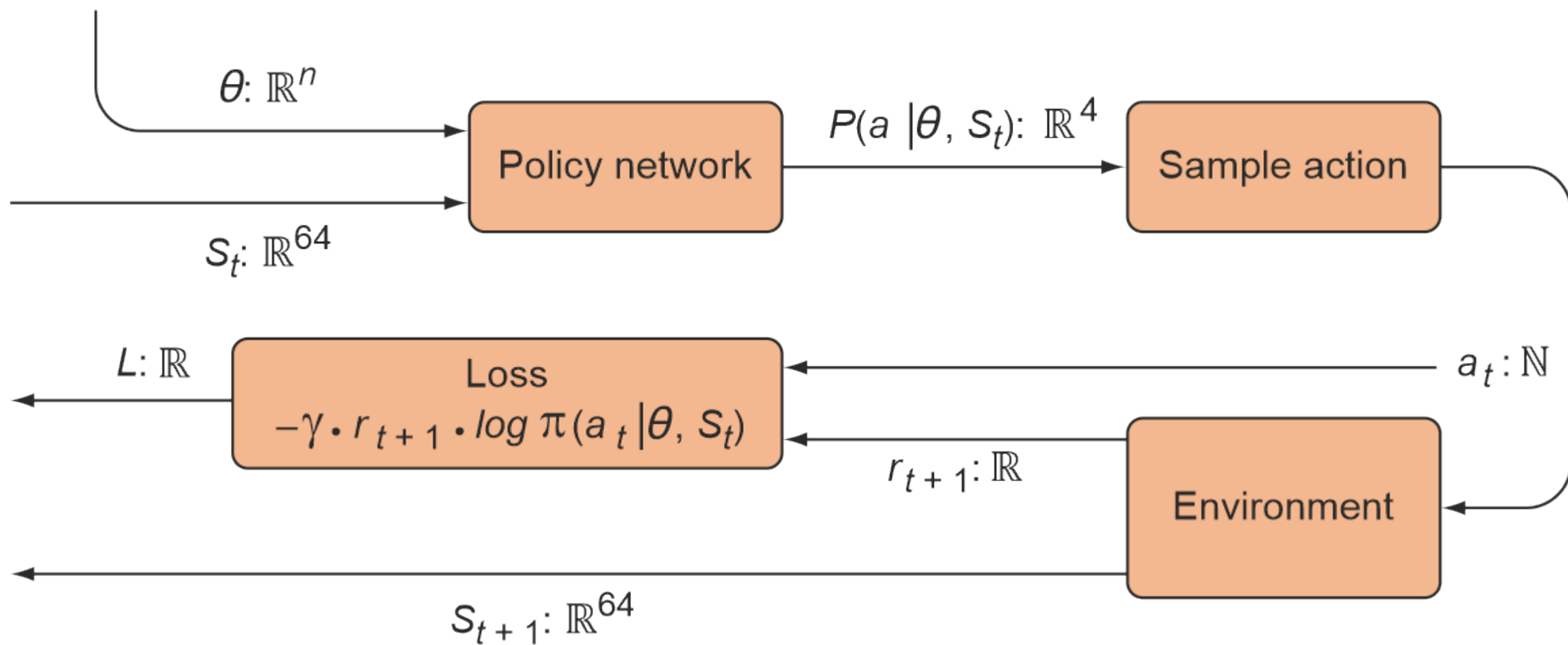Before: $Q_{new} = R_t + d*\max(Q(S_{t+1}))$, not stable
Now: $Q_{new} = R_t + d*\max(Q'(S_{t+1}))$

# Policy methods



Why? Compared to Q-learning, generally regarded advantageous, even necessary (if the action space is continuous), but more mathematical...

**State-of-the-art:**

Loss $= -log\,(\pi(a\,|\,S)) \cdot (R - V_\pi(S))$