# Reinforcement learning

Gergő Orbán

# Recap: Learning frameworks

# Recap: Learning frameworks

- Unsupervised

# Recap: Learning frameworks

- Unsupervised

- Supervised learning: y=f(x) — essentially a mapping from input to output
  -> task specific
  -> requires labelled data points
  -> essentially optimization
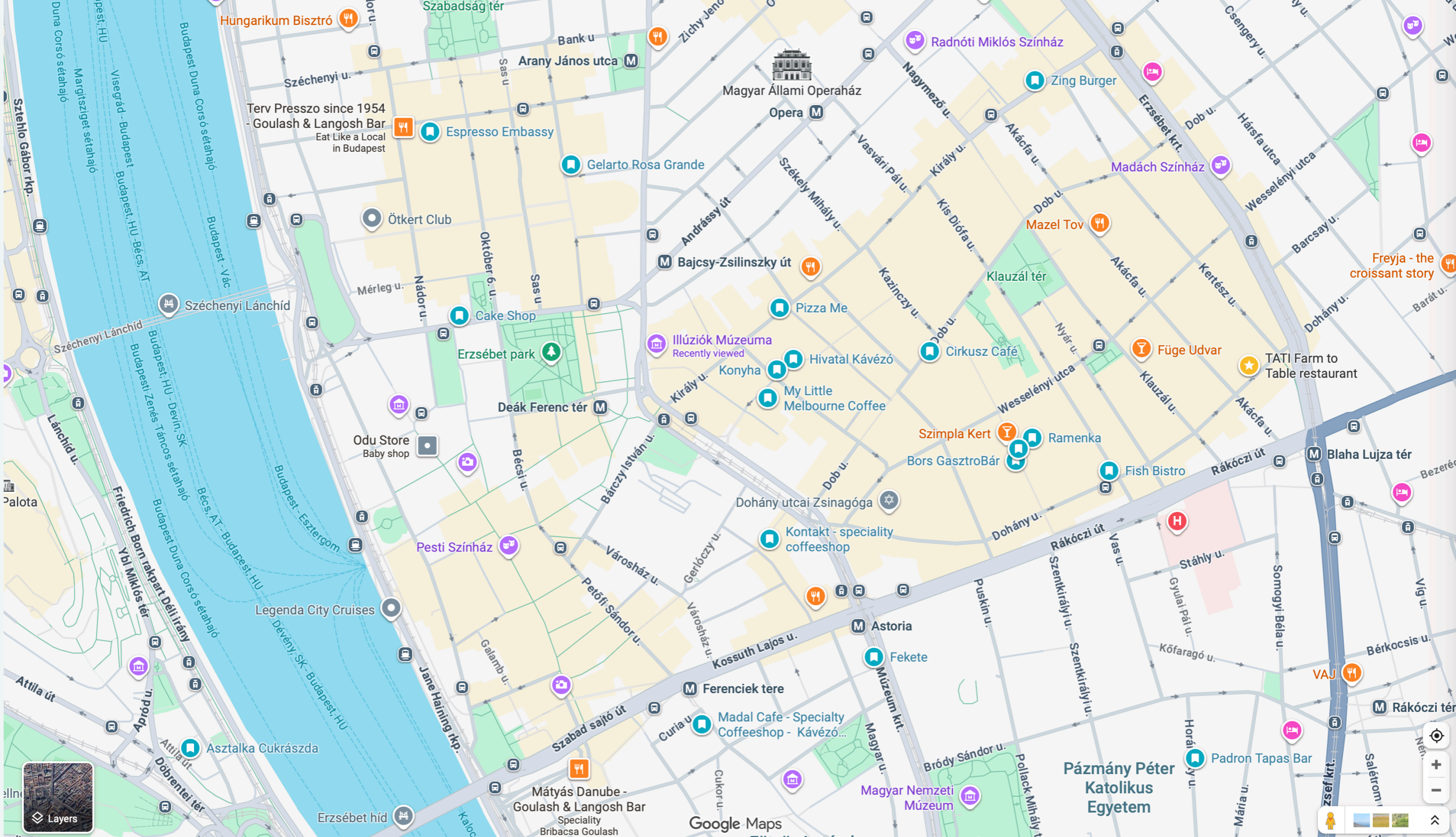  -> back-propagation could be easily obtained

# Recap: Learning frameworks

- Unsupervised

- Supervised learning: y=f(x) — essentially a mapping from input to output
  -> task specific
  -> requires labelled data points
  -> essentially optimization
  -> back-propagation could be easily obtained

- Reinforcement learning
  -> phrasing learning as collection of reward
  -> sparser learning signal
  -> max E [Reward(input, action)]

This is a map image and contains no document-level body text to transcribe as prose. The content consists of map labels only.

Map of central Budapest (Google Maps), including labels:

- Szabadság tér
- Hungarikum Bisztró
- Bank u.
- Zichy Jenő u.
- Arany János utca
- Radnóti Miklós Színház
- Magyar Állami Operaház
- Opera
- Zing Burger
- Széchenyi u.
- Sas u.
- Nagymező u.
- Erzsébet krt.
- Terv Presszo since 1954 - Goulash & Langosh Bar
- Eat Like a Local in Budapest
- Espresso Embassy
- Madách Színház
- Gelarto Rosa Grande
- Akácfa u.
- Dob u.
- Hársfa utca
- Andrássy út
- Mazel Tov
- Freyja - the croissant story
- Ötkert Club
- Székely Mihály u.
- Király u.
- Kis Diófa u.
- Klauzál tér
- Kazinczy u.
- Barcsay u.
- Mérleg u.
- Bajcsy-Zsilinszky út
- Nyár u.
- Kertész u.
- Wesselényi utca
- Dohány u.
- Pizza Me
- Széchenyi Lánchíd
- Erzsébet park
- Cake Shop
- Illúziók Múzeuma — Recently viewed
- Hivatal Kávézó
- Cirkusz Café
- Füge Udvar
- TATI Farm to Table restaurant
- Konyha
- My Little Melbourne Coffee
- Klauzál u.
- Deák Ferenc tér
- Szimpla Kert
- Ramenka
- Bors GasztroBár
- Odu Store — Baby shop
- Bécsi u.
- Bárczy István u.
- Dob u.
- Fish Bistro
- Rákóczi út
- Blaha Lujza tér
- Bezerédi
- Dohány utcai Zsinagóga
- Pesti Színház
- Városház u.
- Gerlóczy u.
- Kontakt - speciality coffeeshop
- Stáhly u.
- Petőfi Sándor u.
- Legenda City Cruises
- Kossuth Lajos u.
- Astoria
- Szentkirályi u.
- Gyulai Pál u.
- Somogyi Béla u.
- Ferenciek tere
- Fekete
- Múzeum krt.
- Curia
- Magyar u.
- Kőfaragó u.
- Bérkocsis u.
- VAJ
- Szabad sajtó út
- Madal Cafe - Specialty Coffeeshop - Kávézó...
- Bródy Sándor u.
- Rákóczi tér
- Asztalka Cukrászda
- Cukor u.
- Magyar Nemzeti Múzeum
- Pollack Mihály tér
- Padron Tapas Bar
- Mátyás Danube - Goulash & Langosh Bar — Speciality Bribacsa Goulash
- Pázmány Péter Katolikus Egyetem
- Erzsébet híd
- Google Maps
- Layers

**Dual challenge**

reward is not immediate:

most of the actions are not rewarding by themselves, as rewards are distal

## Dual challenge

reward is not immediate:
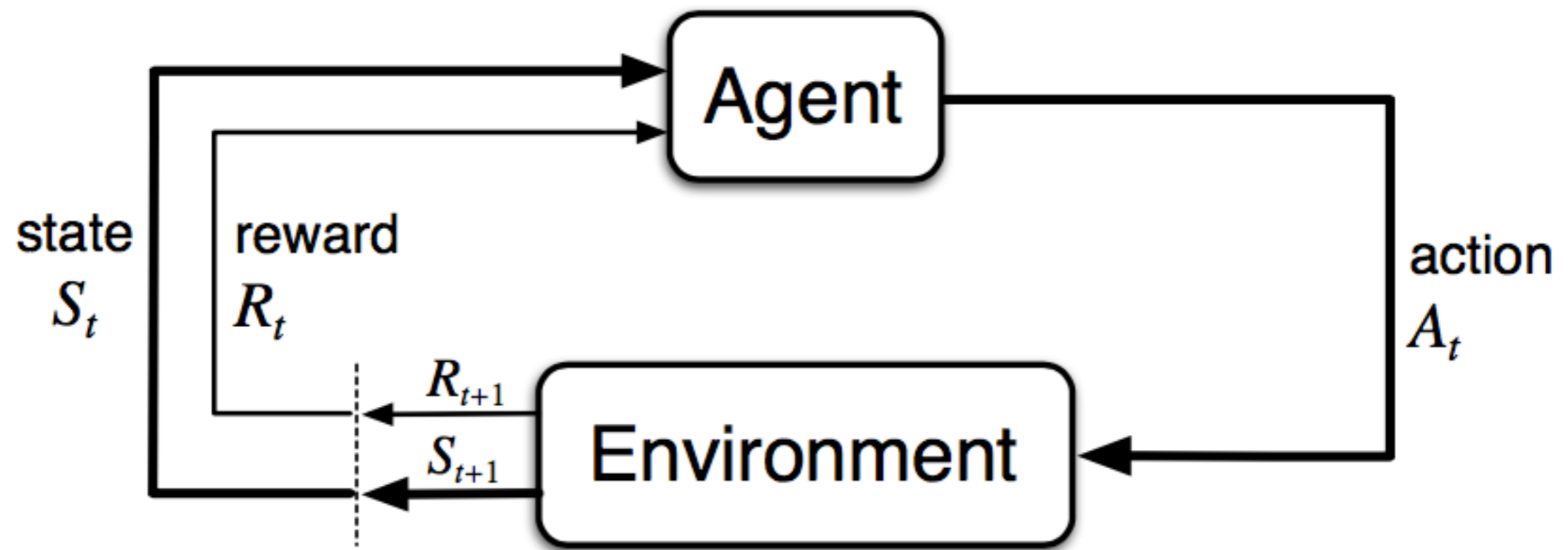most of the actions are not rewarding by themselves, as rewards are distal

options are not readily available:
rewards are not known *and* one experience does not tell exactly how rewarding a state is
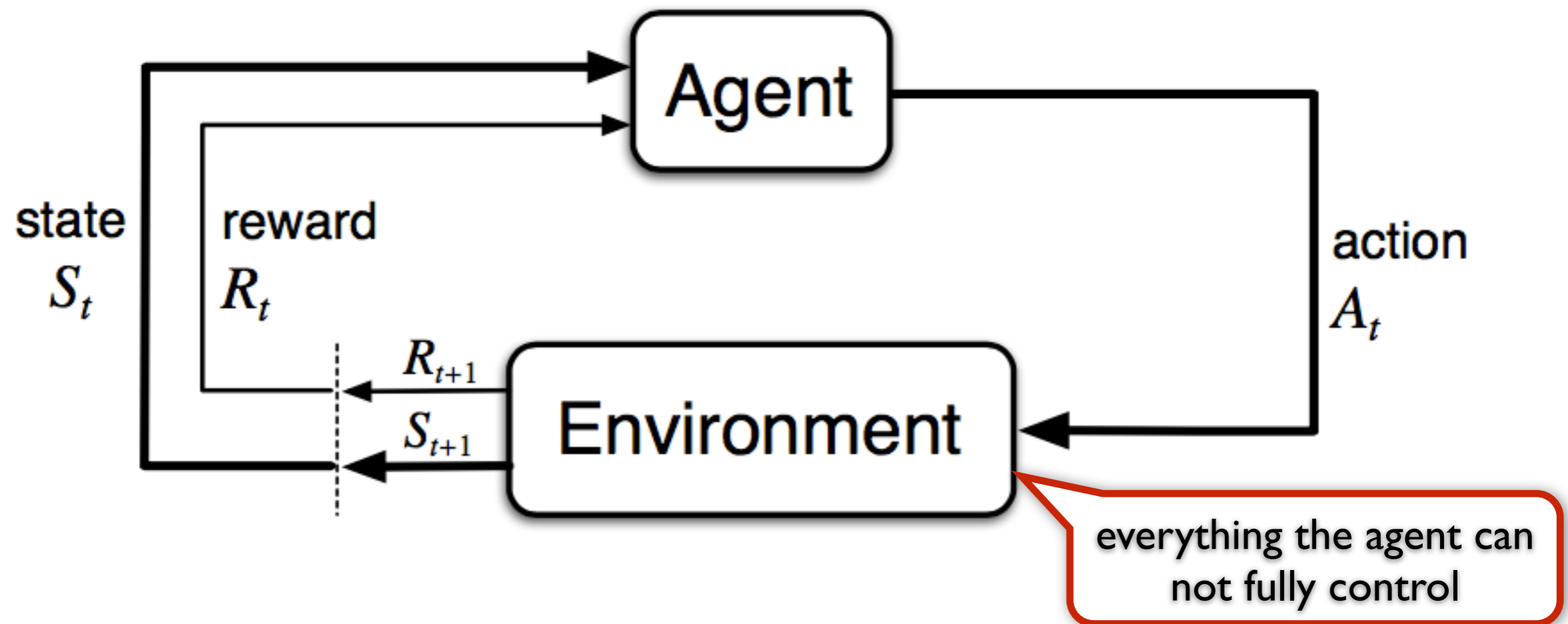
# Reinforcement learning
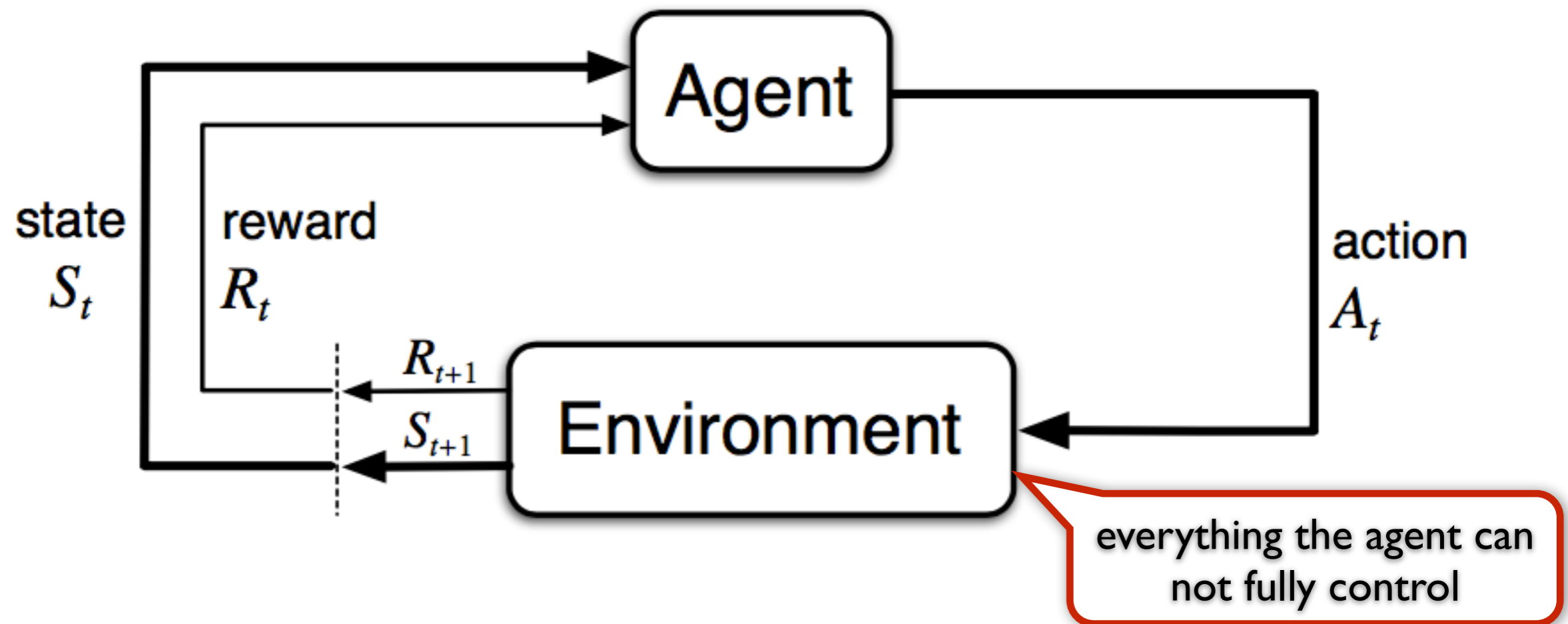
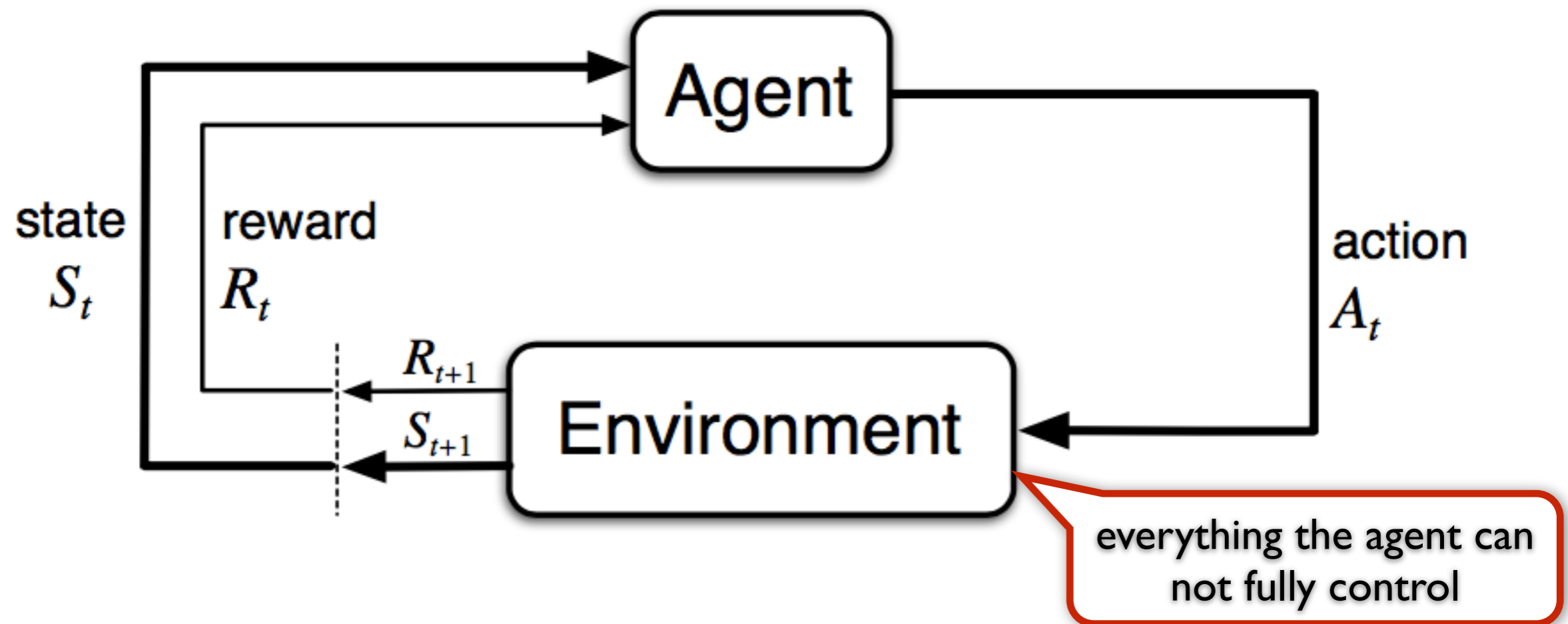# Reinforcement learning

**State**

# Reinforcement learning

## State

- state of the environment

# Reinforcement learning

## State

- state of the environment
- partially observed



everything the agent can not fully control

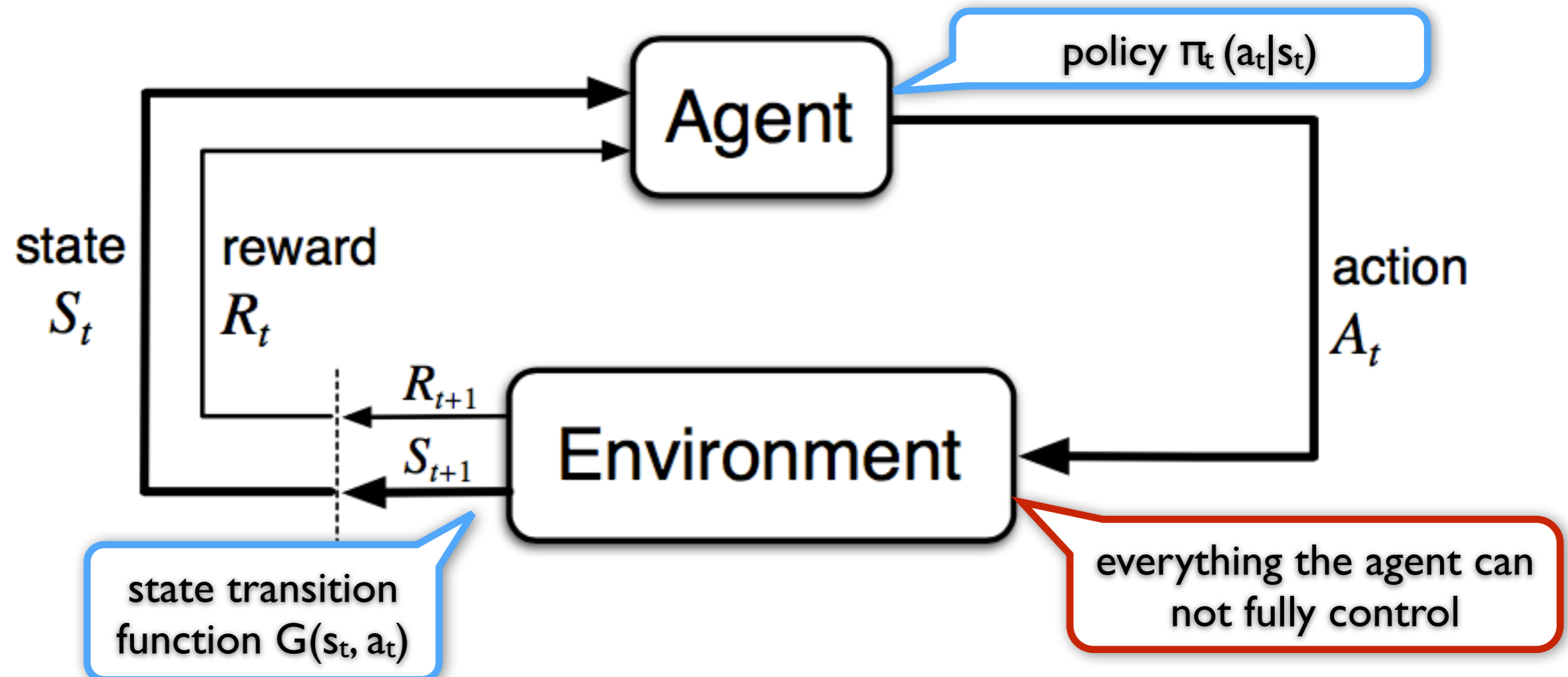## State

- state of the environment

- partially observed

- Markov property

# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action



state $S_t$

reward $R_t$

$R_{t+1}$

$S_{t+1}$

Agent

Environment

action $A_t$

state transition function G($s_t$, $a_t$)

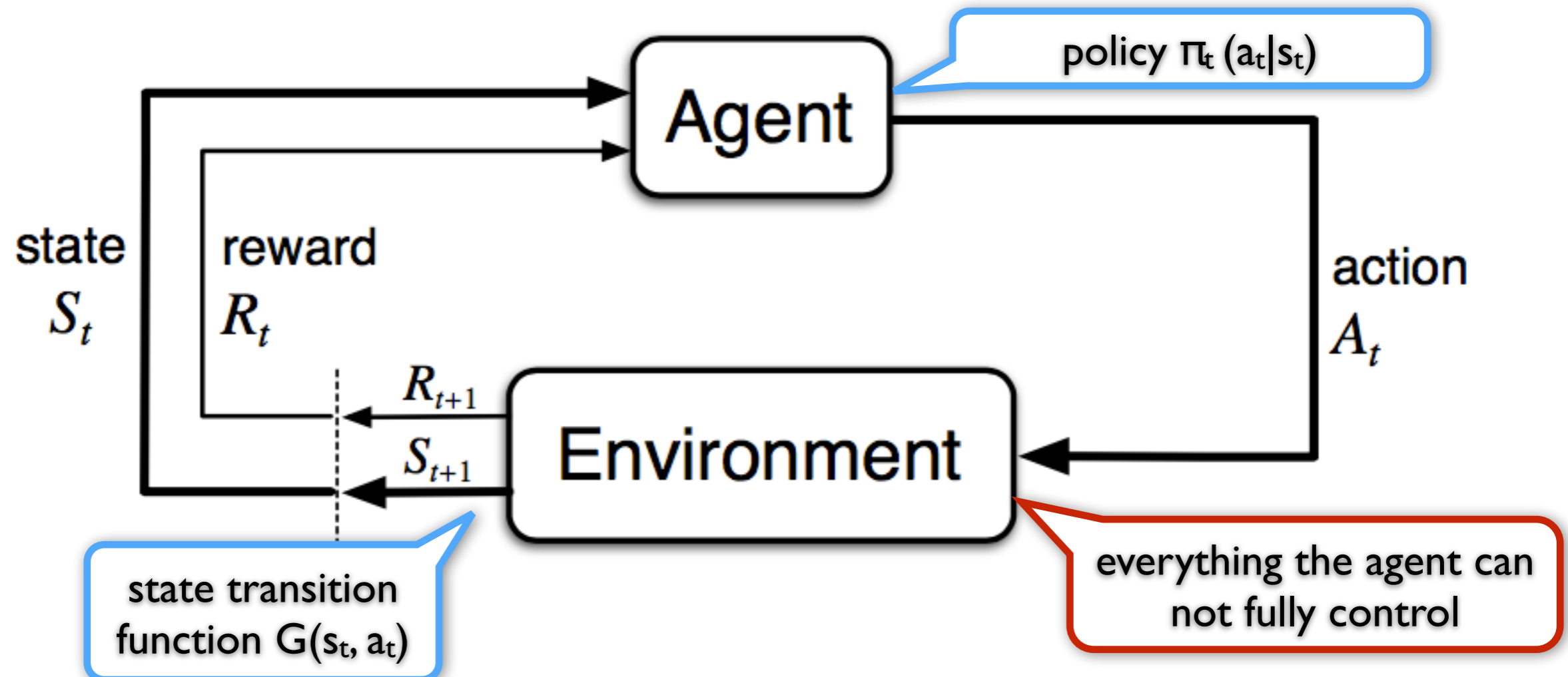everything the agent can not fully control

# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions



state $S_t$

reward $R_t$

action $A_t$

$R_{t+1}$

$S_{t+1}$

state transition function $G(s_t, a_t)$

everything the agent can not fully control

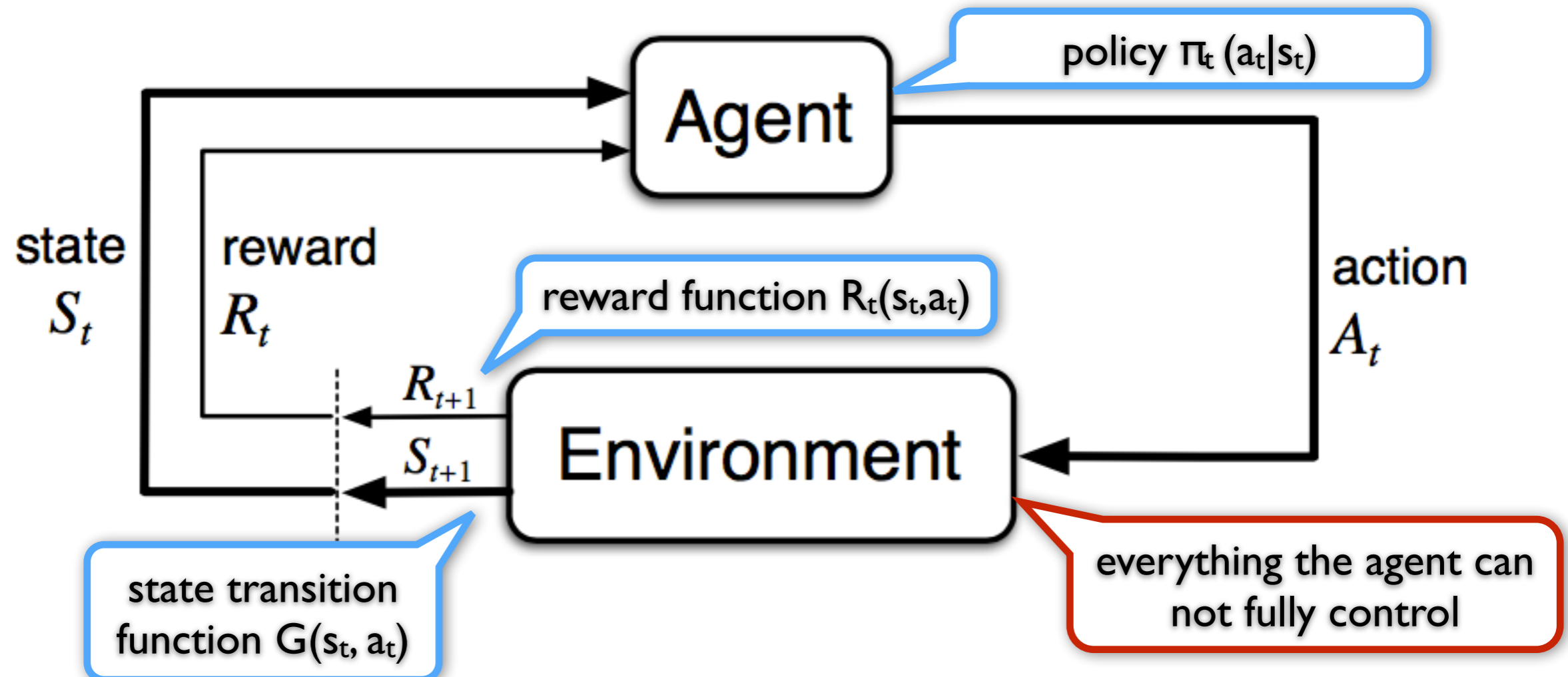# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
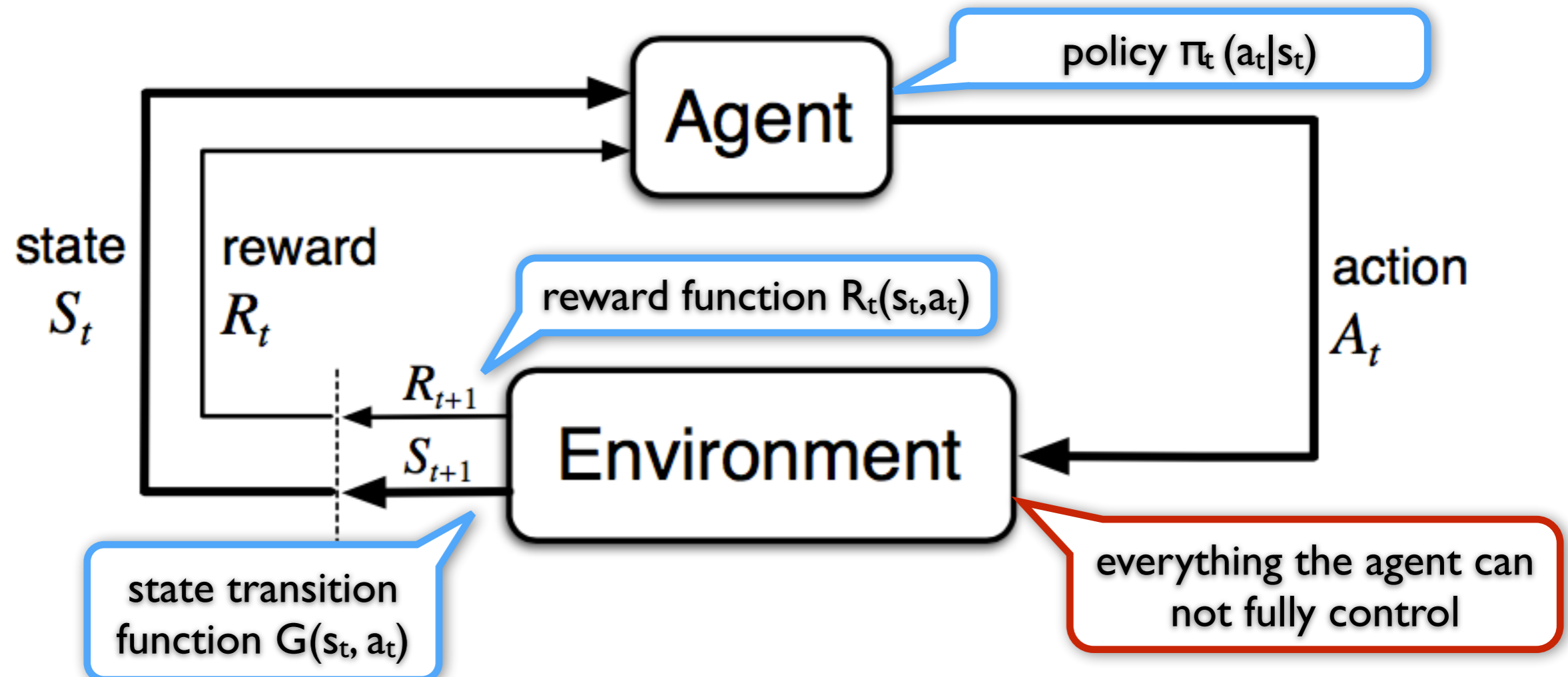
# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
- actions may change the state and/or lead to reward



policy $\pi_t(a_t|s_t)$

Agent

state $S_t$

reward $R_t$

action $A_t$

$R_{t+1}$

$S_{t+1}$

Environment

state transition function $G(s_t, a_t)$

everything the agent can not fully control
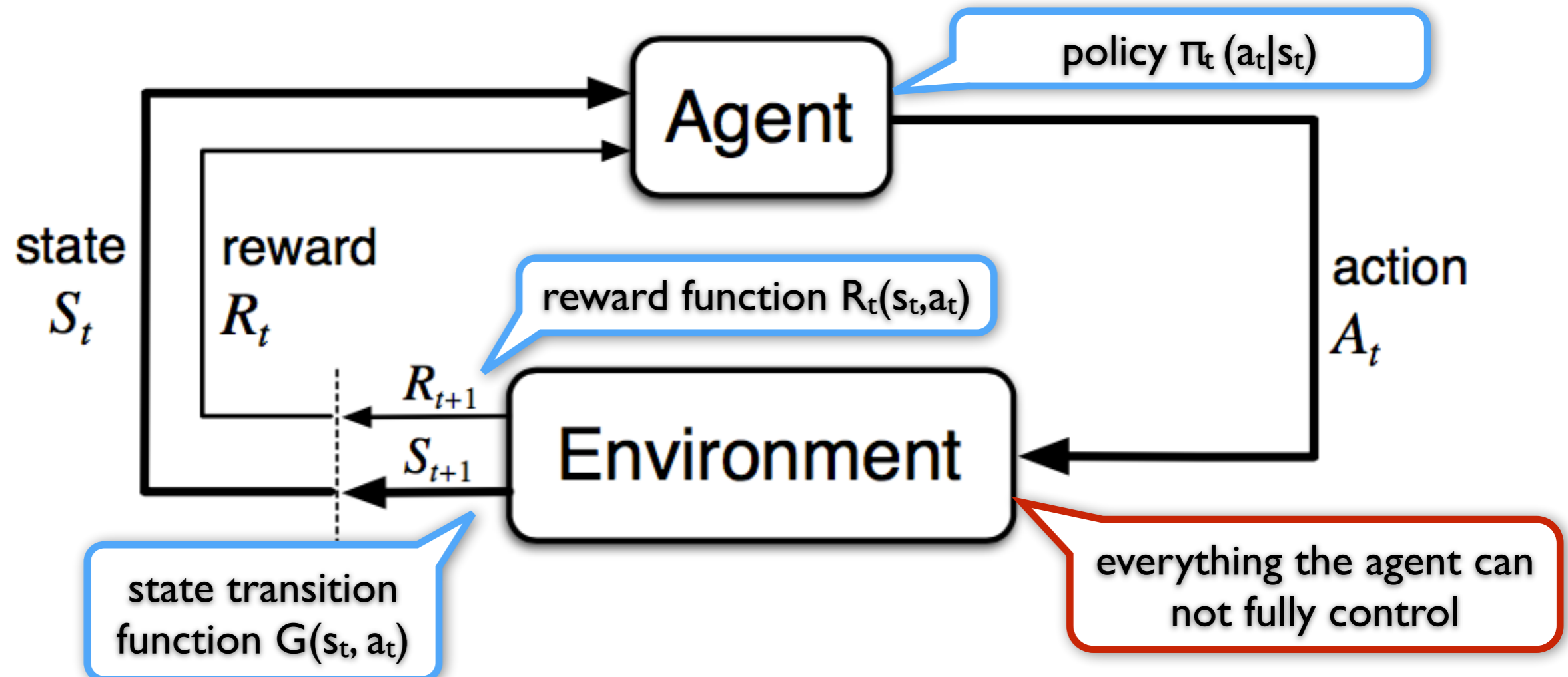
# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
- actions may change the state and/or lead to reward

## Reward



policy $\pi_t(a_t|s_t)$

Agent

state $S_t$

reward $R_t$

action $A_t$

$R_{t+1}$

$S_{t+1}$

Environment

state transition function $G(s_t, a_t)$

everything the agent can not fully control

# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
- actions may change the state and/or lead to reward

## Reward

- $R_t(s,a)$ is the reward given at time t in state s and action a

policy $\pi_t(a_t|s_t)$

reward function $R_t(s_t,a_t)$

state $S_t$

reward $R_t$

action $A_t$

$R_{t+1}$

$S_{t+1}$

state transition function $G(s_t, a_t)$

everything the agent can not fully control

Agent

Environment

# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
- actions may change the state and/or lead to reward

## Reward

- $R_t(s,a)$ is the reward given at time t in state s and action a
- negative or positive



policy $\pi_t(a_t|s_t)$

reward function $R_t(s_t,a_t)$

state transition function $G(s_t, a_t)$

everything the agent can not fully control

# Reinforcement learning

## State

- state of the environment
- partially observed
- Markov property

## Action

- agents interact with the environment with actions
- policy maps states to actions
- actions may change the state and/or lead to reward

## Reward

- $R_t(s,a)$ is the reward given at time t in state s and action a
- negative or positive
- goal: maximise total reward



policy $\pi_t(a_t|s_t)$

reward function $R_t(s_t,a_t)$

state transition function $G(s_t, a_t)$

everything the agent can not fully control

state $S_t$

reward $R_t$

Agent

$R_{t+1}$

$S_{t+1}$

Environment

action $A_t$

# Reinforcement learning

cumulative reward: $\quad \mathscr{R} = \sum_k r_k$

# Reinforcement learning

cumulative reward: $\quad \mathcal{R} = \sum_k r_k$

discounted reward: $\quad \mathcal{R} = \sum_k \alpha^k r_k$

# Reinforcement learning

cumulative reward: $\quad \mathscr{R} = \sum_k r_k$

discounted reward: $\quad \mathscr{R} = \sum_k \alpha^k r_k$

formal goal: $\quad V = \mathbb{E}[\mathscr{R}]$

# Reinforcement learning

cumulative reward: $\mathscr{R} = \sum_k r_k$

discounted reward: $\mathscr{R} = \sum_k \alpha^k r_k$

formal goal: $V = \mathbb{E}[\mathscr{R}]$

more precisely: $V(s) = \mathbb{E}[\mathscr{R}_t \,|\, s]$

# Reinforcement learning

cumulative reward: $\mathscr{R} = \sum_k r_k$

discounted reward: $\mathscr{R} = \sum_k \alpha^k r_k$

formal goal: $V = \mathbb{E}[\mathscr{R}]$

more precisely: $V(s) = \mathbb{E}[\mathscr{R}_t \,|\, s]$

more more precisely: $V_\pi(s) = \mathbb{E}_\pi[\mathscr{R}_t \,|\, s]$

# Reinforcement learning

cumulative reward: $\mathscr{R} = \sum_k r_k$

discounted reward: $\mathscr{R} = \sum_k \alpha^k r_k$

formal goal: $V = \mathbb{E}[\mathscr{R}]$

more precisely: $V(s) = \mathbb{E}[\mathscr{R}_t \,|\, s]$

more more precisely: $V_\pi(s) = \mathbb{E}_\pi[\mathscr{R}_t \,|\, s]$

**policy:**
$\pi(a_t \,|\, s_t)$

# Reinforcement learning

cumulative reward: $\mathcal{R} = \sum_k r_k$

discounted reward: $\mathcal{R} = \sum_k \alpha^k r_k$

formal goal: $V = \mathbb{E}[\mathcal{R}]$

more precisely: $V(s) = \mathbb{E}[\mathcal{R}_t \,|\, s]$

more more precisely: $V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \,|\, s]$

**policy:**
$\pi(a_t \,|\, s_t)$

thus, make explicit that not all states are equivalent: $r_{t+1} = r(s_t, a_t, s_{t+1})$

# Reinforcement learning

cumulative reward: $\mathscr{R} = \sum\limits_k r_k$

discounted reward: $\mathscr{R} = \sum\limits_k \alpha^k r_k$

**policy:**
$\pi(a_t | s_t)$

formal goal: $V = \mathbb{E}[\mathscr{R}]$

more precisely: $V(s) = \mathbb{E}[\mathscr{R}_t | s]$

more more precisely: $V_\pi(s) = \mathbb{E}_\pi[\mathscr{R}_t | s]$

thus, make explicit that not all states are equivalent: $r_{t+1} = r(s_t, a_t, s_{t+1})$

along with how the world works $- p(s_{t+1} | s_t, a_t) -$ we have all the ingredients

# Simple model environment: gridworld

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \,|\, S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\Big]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \;\Big|\; S_t = s\Big]$$

$$= \mathbb{E}_\pi\Big[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \;\Big|\; S_t = s\Big]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\right]$$

$$= \mathbb{E}_\pi\left[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s\right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\left[r(s,a,s') + \gamma\mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s'\right]\right]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \Big| S_t = s\Big]$$

$$= \mathbb{E}_\pi\Big[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \Big| S_t = s\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma\mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \Big| S_t = s'\Big]\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma V_\pi(s')\Big]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \Big| S_t = s\Big]$$

$$= \mathbb{E}_\pi\Big[r_{t+1} + \gamma\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \Big| S_t = s\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma\mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \Big| S_t = s'\Big]\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma V_\pi(s')\Big]$$

**Slight generalization: Bellmann equation for Q:**

$$Q_\pi(s,a) =$$
$$= \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \sum_a \pi(a|s')\,\gamma\,Q_\pi(s',a)\Big]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Challenges in RL

# Challenges in RL

**Why is it hard?**

# Challenges in RL

## Why is it hard?

- rewards are distal - credit assignment problem

# Challenges in RL

**Why is it hard?**

• rewards are distal - credit assignment problem

• exploration - exploitation dilemma

# Challenges in RL

**Why is it hard?**

- rewards are distal - credit assignment problem

- exploration - exploitation dilemma

- state-space is huge - searching takes a long time!

# Challenges in RL

**Why is it hard?**

- rewards are distal - credit assignment problem

- exploration - exploitation dilemma

- state-space is huge - searching takes a long time!

- state dynamics is unknown

# Challenges in RL

**Why is it hard?**

• rewards are distal **-** credit assignment problem

• exploration **-** exploitation dilemma

• state-space is huge **-** searching takes a long time!

• state dynamics is unknown

• state dynamics can be stochastic **-** noisy environment or noisy action

# Challenges in RL

**Why is it hard?**

- rewards are distal - credit assignment problem

- exploration - exploitation dilemma

- state-space is huge - searching takes a long time!

- state dynamics is unknown

- state dynamics can be stochastic - noisy environment or noisy action

- rewards can be stochastic

# Challenges in RL

## Why is it hard?

- rewards are distal - credit assignment problem

- exploration - exploitation dilemma

- state-space is huge - searching takes a long time!

- state dynamics is unknown

- state dynamics can be stochastic - noisy environment or noisy action

- rewards can be stochastic

- states are only partially observed

# Challenges in RL

## Why is it hard?

- rewards are distal **-** credit assignment problem

- exploration **-** exploitation dilemma

- state-space is huge **-** searching takes a long time!

- state dynamics is unknown

- state dynamics can be stochastic **-** noisy environment or noisy action

- rewards can be stochastic

- states are only partially observed

- rules change with time

# Challenge 1: uncertain outcomes

We consider the simple setting, where every action might be rewarding, only the expected reward for the state is changing

# Challenge 1: uncertain outcomes

We consider the simple setting, where every action might be rewarding, only the expected reward for the state is changing

# Challenge 1: uncertain outcomes

We consider the simple setting, where every action might be rewarding, only the expected reward for the state is changing

Every state has a $Q$ value

Keep up the current state as long as the value of it exceeds the value of others.
**OR**
Keep the current state as long as the value of it exceeds the average value of states

# Challenge 1: uncertain outcomes

We consider the simple setting, where every action might be rewarding, only the expected reward for the state is changing

Every state has a $Q$ value

Keep up the current state as long as the value of it exceeds the value of others.
**OR**
Keep the current state as long as the value of it exceeds the average value of states

# Challenge 1: uncertain outcomes

We consider the simple setting, where every action might be rewarding, only the expected reward for the state is changing

formal goal: $Q(s) = \mathbb{E}[\mathscr{R}]$

discounted reward: $\mathscr{R} = \sum_k \alpha^k r_k(s)$

Every state has a $Q$ value

Keep up the current state as long as the value of it exceeds the value of others.
**OR**
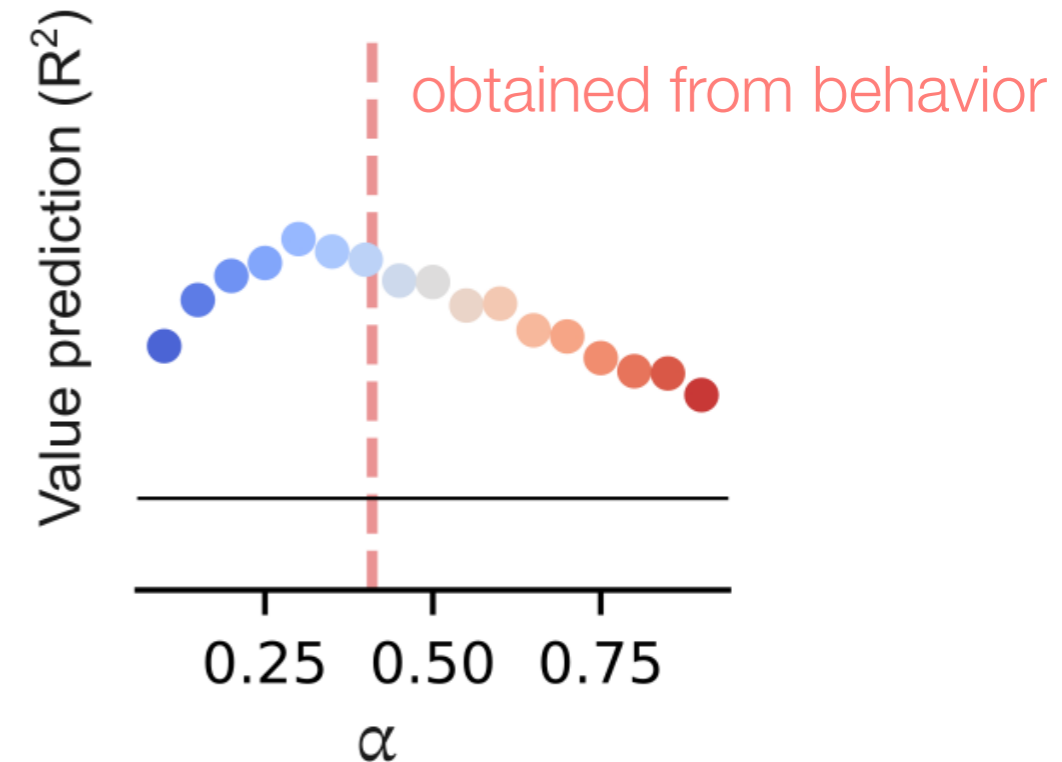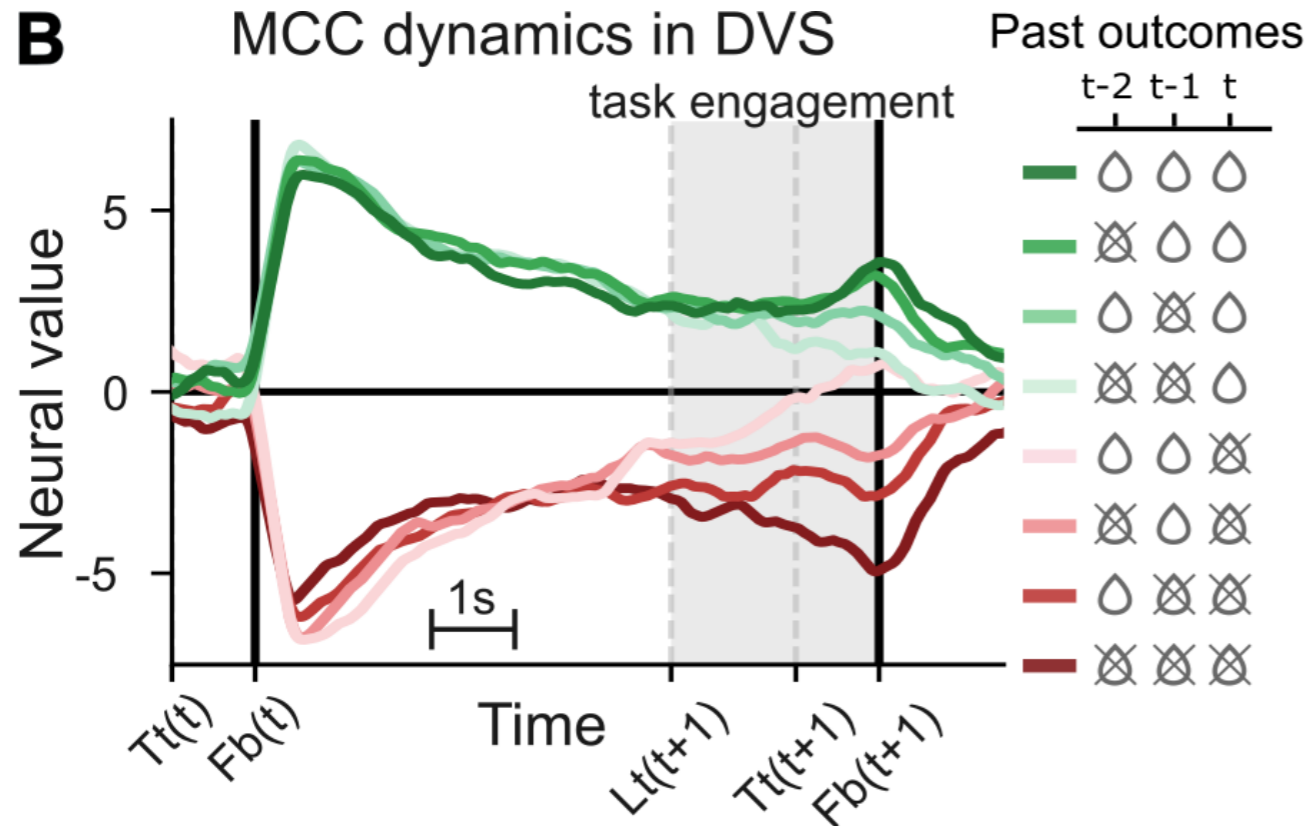Keep the current state as long as the value of it exceeds the average value of states
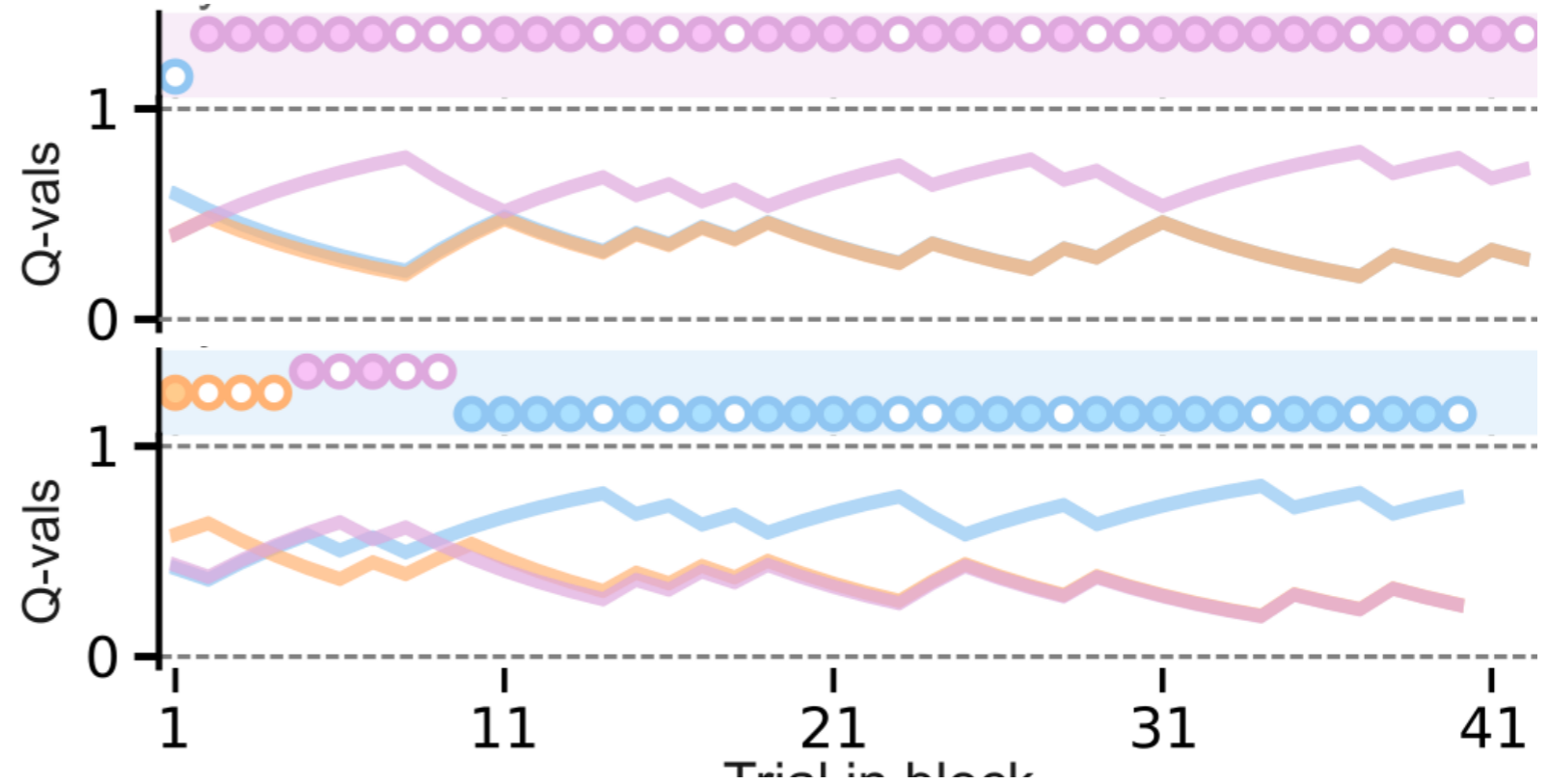
# Neural correlates of RL

# Neural correlates of RL

# Neural correlates of RL
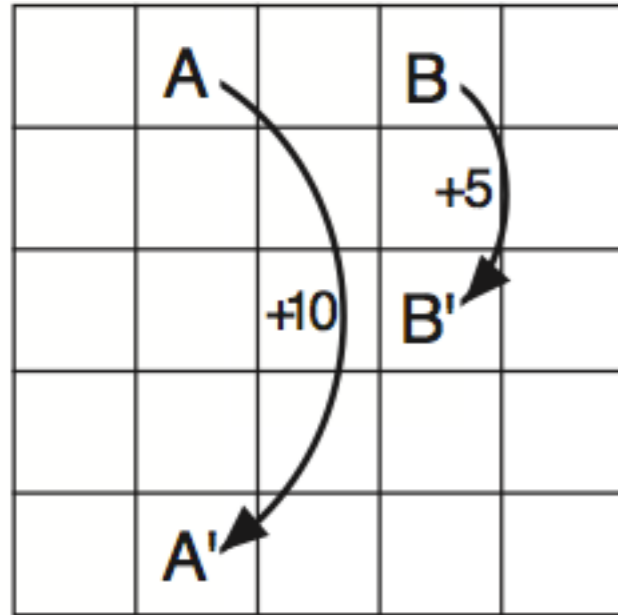


trial start
st

lever touch
Lt

*hold
1000ms*

lever validation
Lv

target touch
Tt

*hold
450ms*

target validation
Tv

*delay
500ms*

feedback
Fb

**B** MCC dynamics in DVS

task engagement

Past outcomes

t-2  t-1  t

Neural value

Time

Tt(t)  Fb(t)  Lt(t+1)  Tt(t+1)  Fb(t+1)

Q-vals

Trial in block

# Neural correlates of RL



trial start
st

lever touch
Lt

*hold
1000ms*

lever validation
Lv

target touch
Tt

*hold
450ms*

target validation
Tv

*delay
500ms*

feedback
Fb

**B** MCC dynamics in DVS

task engagement

Neural value

1s

Tt(t)  Fb(t)  Time  Lt(t+1)  Tt(t+1)  Fb(t+1)

Past outcomes

t-2  t-1  t

Q-vals

Q-vals

1    11    21    31    41

Trial in block

Value prediction ($R^2$)

obtained from behavior

0.25    0.50    0.75

$\alpha$

# Simple model environment: gridworld

# Simple model environment: gridworld
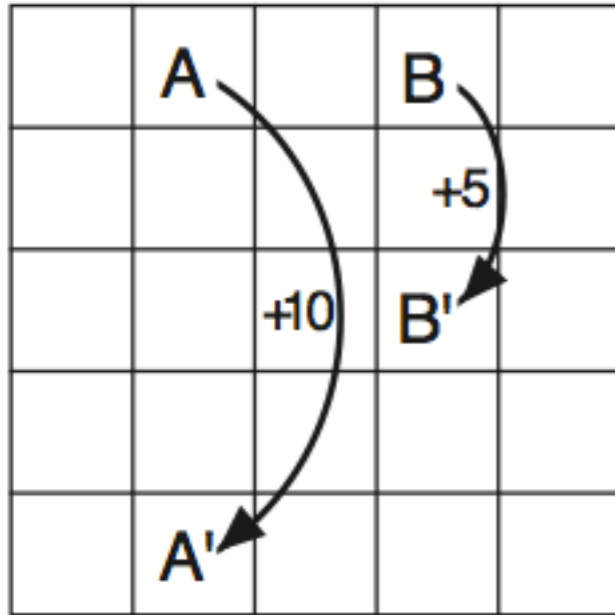
# Simple model environment: gridworld

**state transitions**     $P(s_{t+1}|s_t, a_t)$
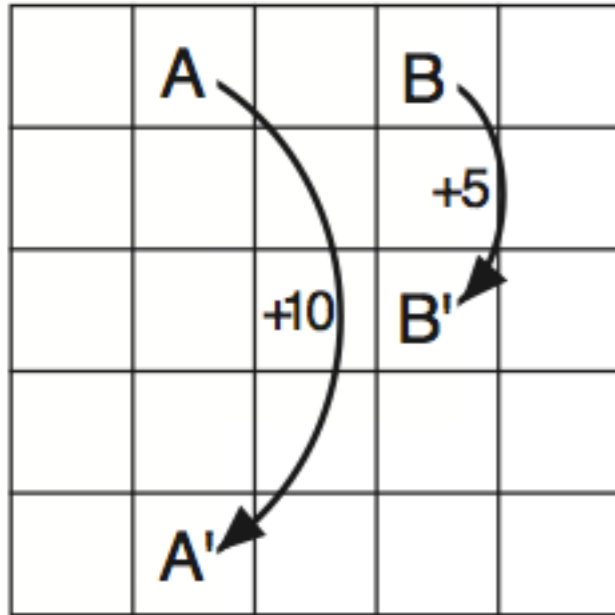
# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

# Simple model environment: gridworld
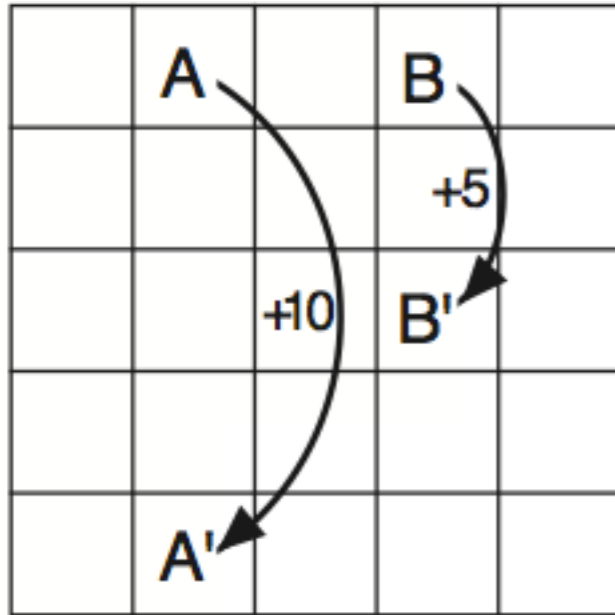


**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$
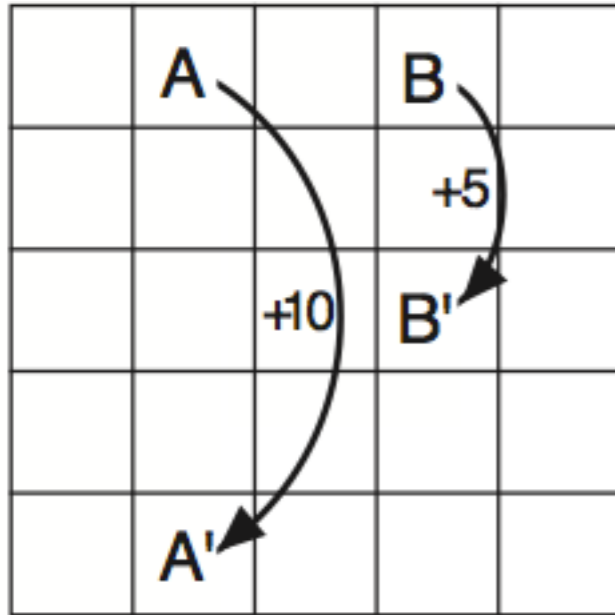
**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$
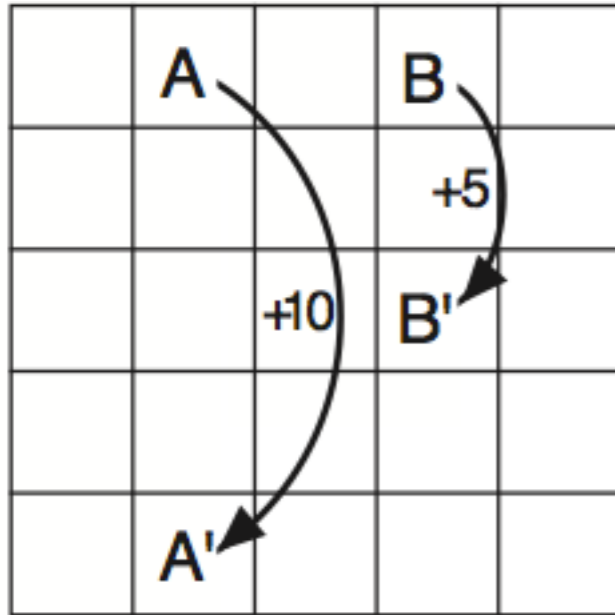
$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$
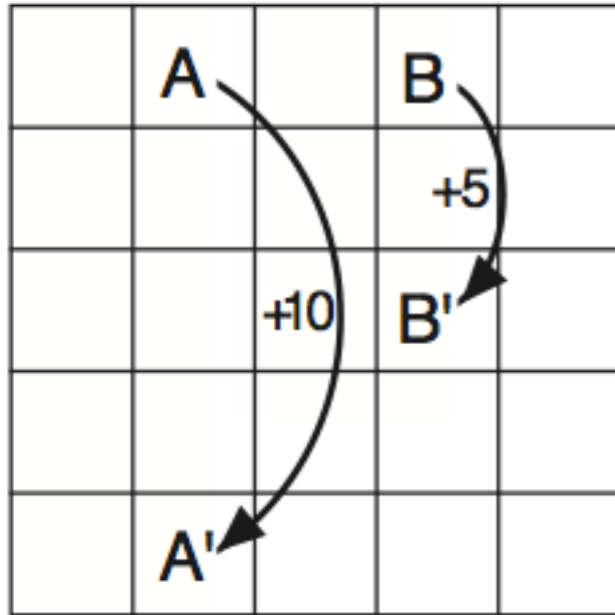
$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\Big]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \ldots$

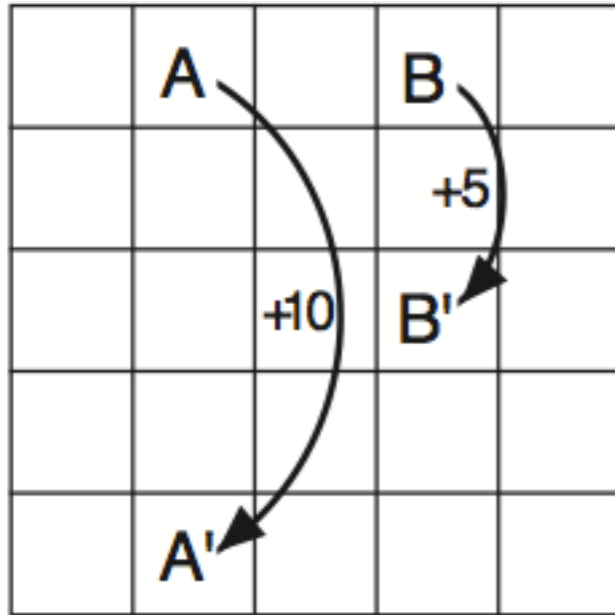$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\right]$$

$$= \mathbb{E}_\pi\left[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s\right]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

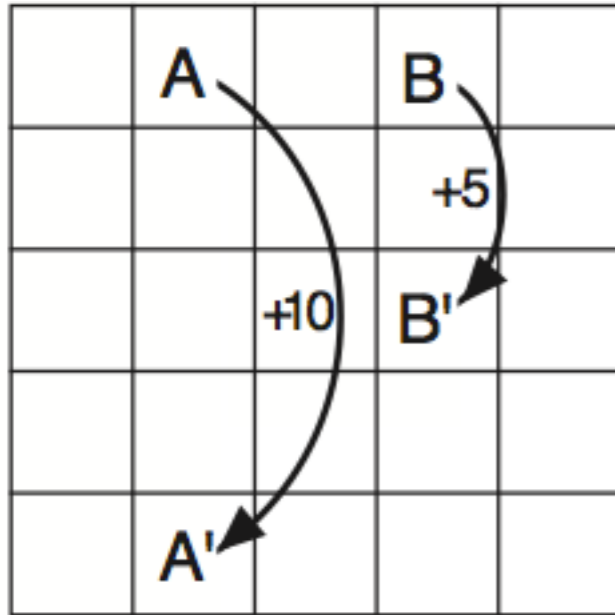**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\Big]$$

$$= \mathbb{E}_\pi\Big[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\Big[r(s, a, s') + \gamma \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s'\Big]\Big]$$

# Simple model environment: gridworld



**state transitions** $\quad P(s_{t+1}|s_t, a_t)$

**rewards:** $\quad r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\quad \pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
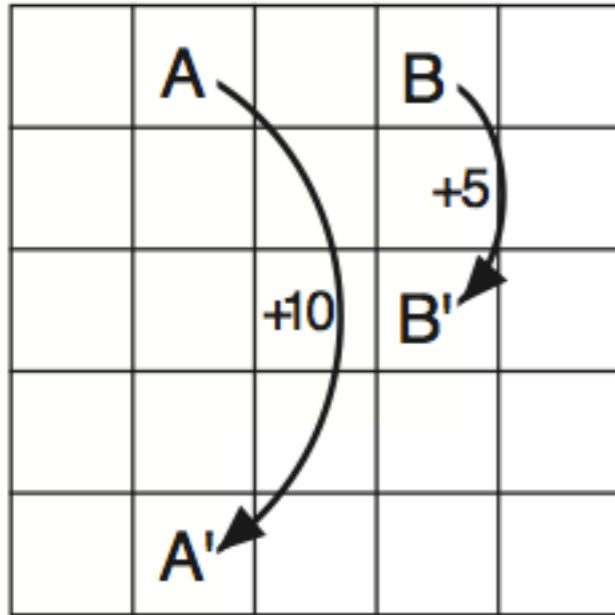**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s \Big]$$

$$= \mathbb{E}_\pi\Big[ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s \Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\Big[ r(s, a, s') + \gamma \mathbb{E}_\pi\Big[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s' \Big]\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\Big[ r(s, a, s') + \gamma V_\pi(s')\Big]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \ldots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**

**Bellmann equation:**

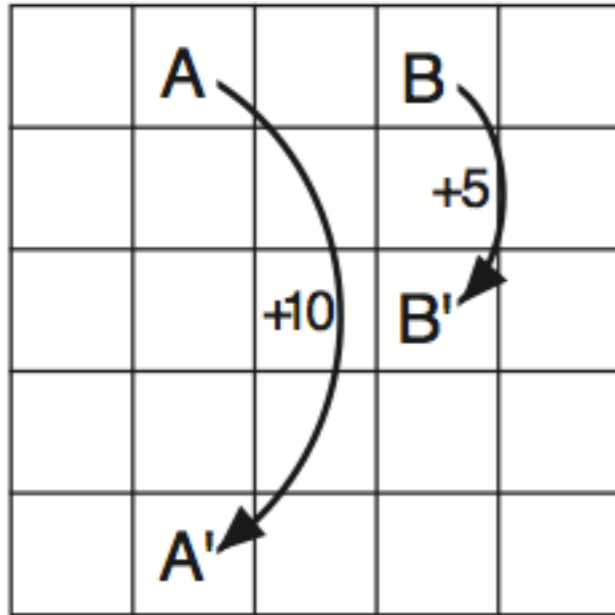$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s\right]$$

$$= \mathbb{E}_\pi\left[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s\right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\left[r(s,a,s') + \gamma\mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s'\right]\right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\left[r(s,a,s') + \gamma V_\pi(s')\right]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

state transitions $P(s_{t+1}|s_t, a_t)$

rewards: $r(s, a, s')$

**Bellmann equation for Q:**

$$Q_\pi(s, a) = \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \sum_a \pi(a|s') \, \gamma \, Q_\pi(s', a) \right]$$

policy: $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \,|\, S_t = s]$$

$$= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s \right]$$

$$= \mathbb{E}_\pi \left[ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \gamma \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s' \right] \right]$$
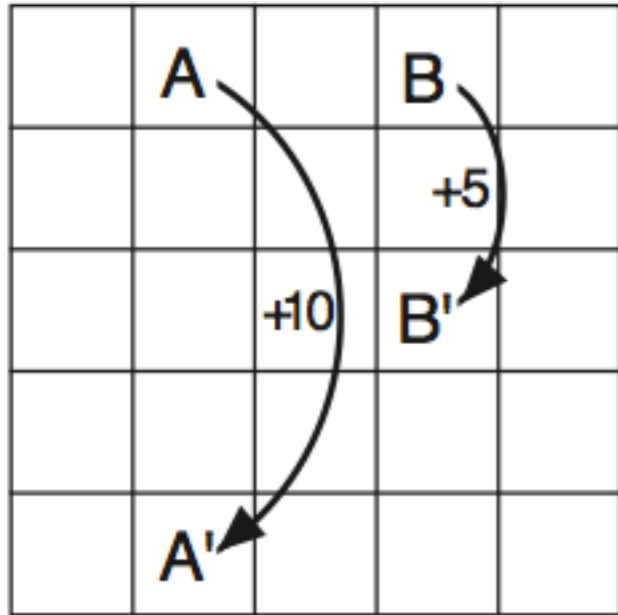
$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \gamma V_\pi(s') \right]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

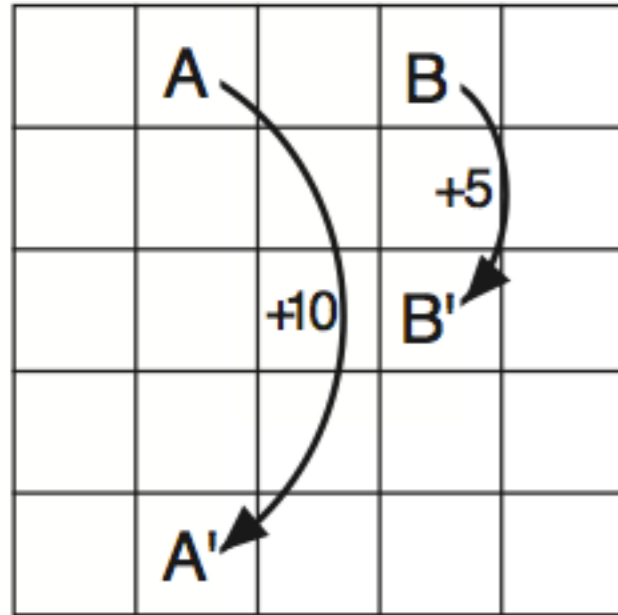# Simple model environment: gridworld

# Simple model environment: gridworld

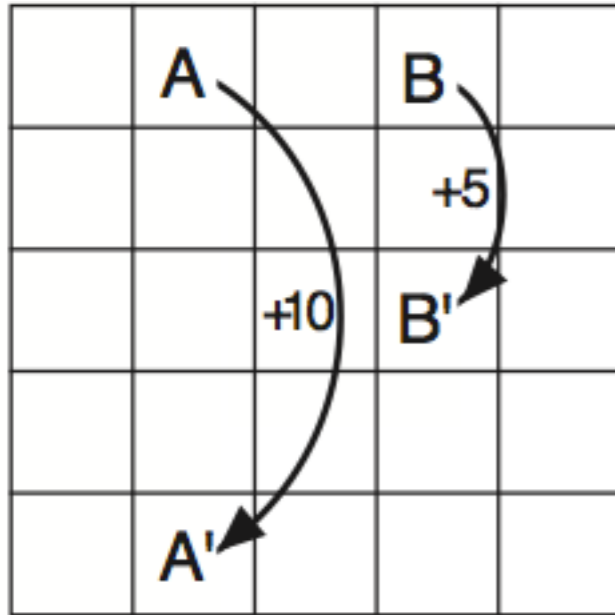# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**

# Simple model environment: gridworld



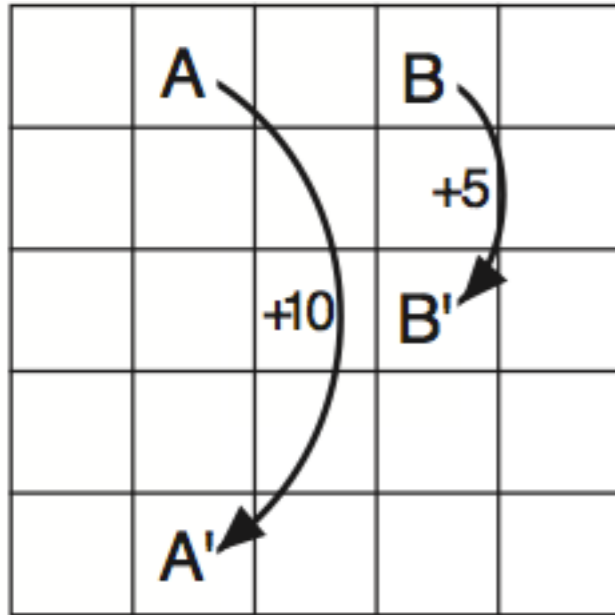**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

# Simple model environment: gridworld
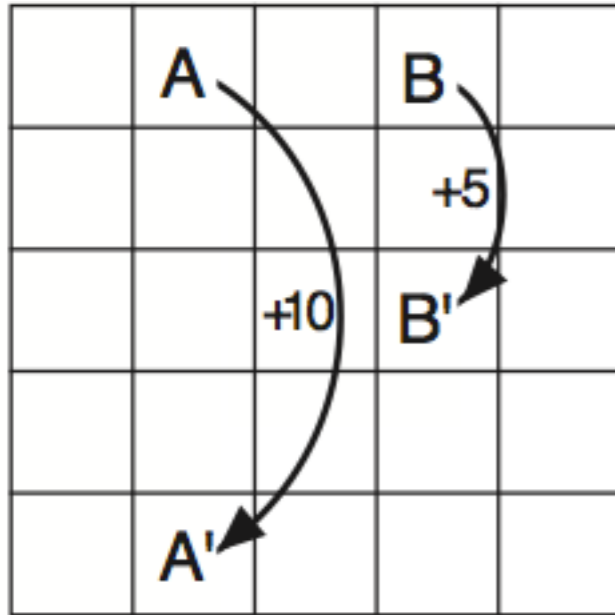


**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

# Simple model environment: gridworld



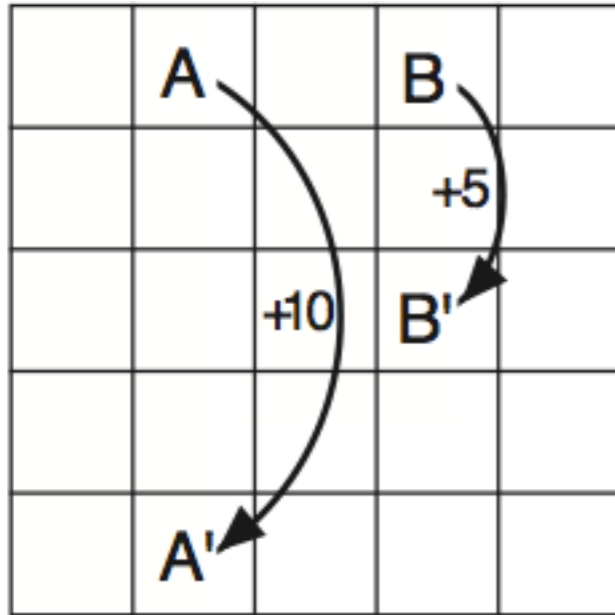**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s\Big]$$

# Simple model environment: gridworld



**state transitions** $\quad P(s_{t+1}|s_t, a_t)$

**rewards:** $\quad r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$
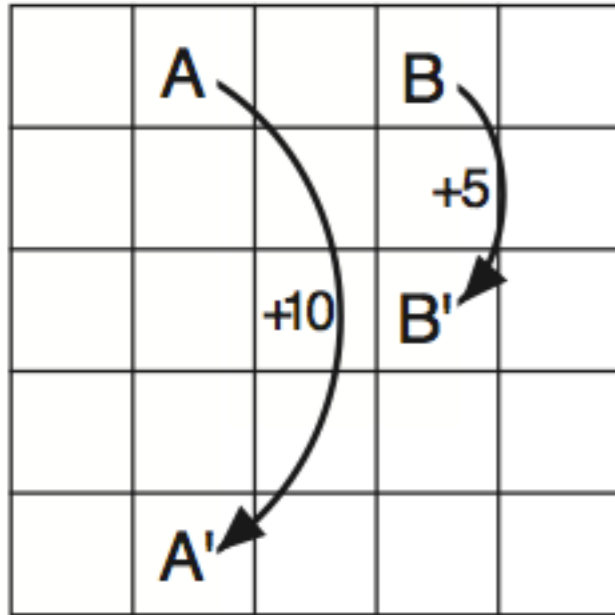
$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

**policy:** $\quad \pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \,|\, S_t = s]$$

$$= \mathbb{E}_\pi\Big[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s \Big]$$

$$= \mathbb{E}_\pi\Big[ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s \Big]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \ldots$
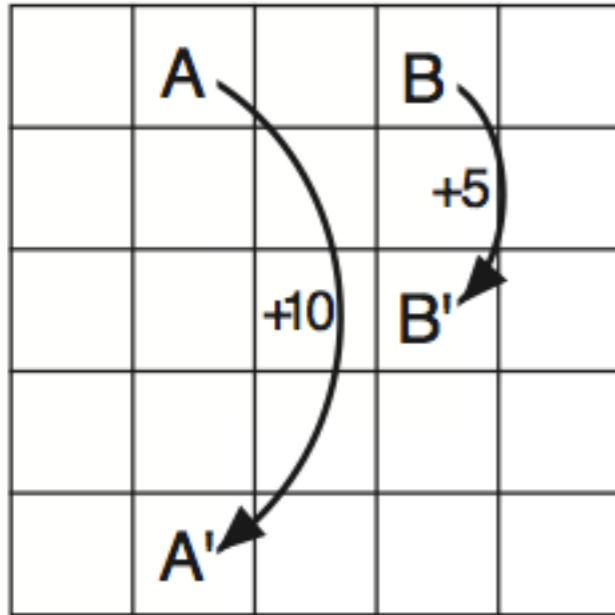
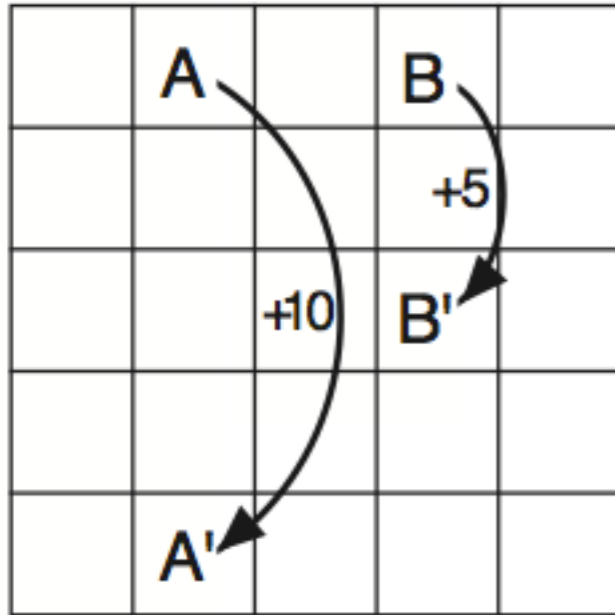$$= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s\right]$$

$$= \mathbb{E}_\pi\left[r_{t+1} + \gamma\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s\right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\left[r(s, a, s') + \gamma\mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s'\right]\right]$$

**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

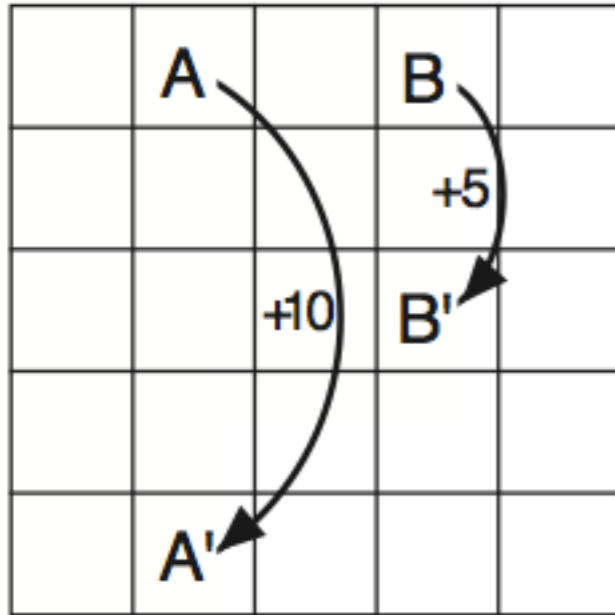**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid S_t = s\Big]$$

$$= \mathbb{E}_\pi\Big[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma\mathbb{E}_\pi\Big[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid S_t = s'\Big]\Big]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a,s)\Big[r(s,a,s') + \gamma V_\pi(s')\Big]$$

# Simple model environment: gridworld



**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s_t, a_t, s_{t+1})$

discounting: $\mathcal{R}_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+2} + \dots$

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

**policy:** $\pi(a_t|s_t)$

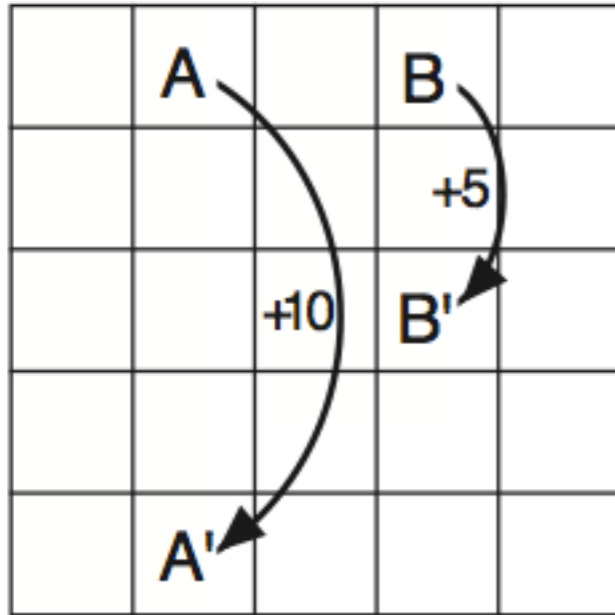**what is the value associated with a given state under a policy?**
**Bellmann equation:**

$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \mid S_t = s]$$

$$= \mathbb{E}_\pi\left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s \right]$$

$$= \mathbb{E}_\pi\left[ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\left[ r(s, a, s') + \gamma \mathbb{E}_\pi\left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s' \right] \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s)\left[ r(s, a, s') + \gamma V_\pi(s') \right]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

**state transitions** $P(s_{t+1}|s_t, a_t)$

**rewards:** $r(s, a, s')$

**Bellmann equation for Q:**

$$Q_\pi(s, a) = \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \sum_a \pi(a|s') \gamma Q_\pi(s', a) \right]$$

**policy:** $\pi(a_t|s_t)$

**what is the value associated with a given state under a policy?**

**Bellmann equation:**

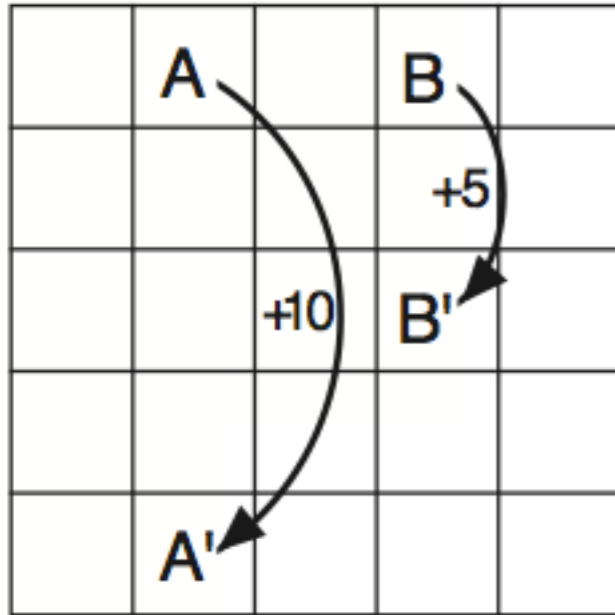$$V_\pi(s) = \mathbb{E}_\pi[\mathcal{R}_t \,|\, S_t = s]$$

$$= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \,\Big|\, S_t = s \right]$$

$$= \mathbb{E}_\pi \left[ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \gamma \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \,\Big|\, S_t = s' \right] \right]$$

$$= \sum_a \pi(a|s) \sum_{s'} P(s'|a, s) \left[ r(s, a, s') + \gamma V_\pi(s') \right]$$

- consistency relationship between states
- depends on policy
- optimal policy: highest value
- learning: find the optimal policy

# Simple model environment: gridworld

Computing the value function, V(s)

# Simple model environment: gridworld

## Computing the value function, V(s)

random policy

# Simple model environment: gridworld

## Computing the value function, V(s)

random policy



| 3.3 | 8.8 | 4.4 | 5.3 | 1.5 |
|-----|-----|-----|-----|-----|
| 1.5 | 3.0 | 2.3 | 1.9 | 0.5 |
| 0.1 | 0.7 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.4 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

# Simple model environment: gridworld

Computing the value function, V(s)

random policy



| 3.3 | 8.8 | 4.4 | 5.3 | 1.5 |
|-----|-----|-----|-----|-----|
| 1.5 | 3.0 | 2.3 | 1.9 | 0.5 |
| 0.1 | 0.7 | 0.7 | 0.4 | -0.4 |
| -1.0 | -0.4 | -0.4 | -0.6 | -1.2 |
| -1.9 | -1.3 | -1.2 | -1.4 | -2.0 |

optimal policy

| 22.0 | 24.4 | 22.0 | 19.4 | 17.5 |
|------|------|------|------|------|
| 19.8 | 22.0 | 19.8 | 17.8 | 16.0 |
| 17.8 | 19.8 | 17.8 | 16.0 | 14.4 |
| 16.0 | 17.8 | 16.0 | 14.4 | 13.0 |
| 14.4 | 16.0 | 14.4 | 13.0 | 11.7 |

# Alternative solutions to Bellmann equation

# Alternative solutions to Bellmann equation

**dynamic programming**

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment

# Alternative solutions to Bellmann equation

**dynamic programming**

- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: $\mathbf{V}_\pi(s) \mid \boldsymbol{\pi}(a|s)$

# Alternative solutions to Bellmann equation

**dynamic programming**

- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: $\mathbf{V}_\pi(s) \mid \boldsymbol{\pi}(a|s)$
- policy improvement: $\boldsymbol{\pi}(a|s) \mid \mathbf{V}_\pi(s)$

# Alternative solutions to Bellmann equation

**dynamic programming**

- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: $V_\pi(s) \mid \pi(a|s)$
- policy improvement: $\pi(a|s) \mid V_\pi(s)$

**Monte Carlo techniques**

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: $\mathbf{V}_\pi(s) \mid \boldsymbol{\pi}(a|s)$
- policy improvement: $\boldsymbol{\pi}(a|s) \mid \mathbf{V}_\pi(s)$

**Monte Carlo techniques**
- wait until the reward arrives

# Alternative solutions to Bellmann equation

**dynamic programming**

- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: $\mathbf{V}_\pi(s) \mid \boldsymbol{\pi}(a|s)$
- policy improvement: $\boldsymbol{\pi}(a|s) \mid \mathbf{V}_\pi(s)$

**Monte Carlo techniques**

- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \mathcal{R}_t - Q_\pi(s_t, r_t) \Big]$$

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**$_\pi$(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**$_\pi$(s)

**Monte Carlo techniques**
- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \mathcal{R}_t - Q_\pi(s_t, r_t) \right]$$

**temporal difference (TD-) learning**

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**$_\pi$(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**$_\pi$(s)

**Monte Carlo techniques**
- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \mathcal{R}_t - Q_\pi(s_t, r_t) \Big]$$

**temporal difference (TD-) learning**
- don't wait with the updates until rewards!

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**$_\pi$(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**$_\pi$(s)

**Monte Carlo techniques**
- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \mathcal{R}_t - Q_\pi(s_t, r_t) \Big]$$

**temporal difference (TD-) learning**
- don't wait with the updates until rewards!
- use intermediate value estimates to update the action-values!

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**$_\pi$(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**$_\pi$(s)

**Monte Carlo techniques**
- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \mathcal{R}_t - Q_\pi(s_t, r_t) \right]$$

**temporal difference (TD-) learning**
- don't wait with the updates until rewards!
- use intermediate value estimates to update the action-values!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \underbrace{r_{t+1} + \gamma\, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \right]$$

# Alternative solutions to Bellmann equation

**dynamic programming**
- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**$_\pi$(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**$_\pi$(s)

**model-based**

**Monte Carlo techniques**
- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \mathcal{R}_t - Q_\pi(s_t, r_t) \Big]$$

**temporal difference (TD-) learning**
- don't wait with the updates until rewards!
- use intermediate value estimates to update the action-values!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \underbrace{r_{t+1} + \gamma\, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \Big]$$

# Alternative solutions to Bellmann equation

**dynamic programming**

**model-based**

- use simulations to
- solve the Bellmann equations iteratively
- need an accurate model of the environment
- policy evaluation: **V**π(s) | **π**(a|s)
- policy improvement: **π**(a|s) | **V**π(s)

**Monte Carlo techniques**

- wait until the reward arrives
- update value functions based on average returns

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \mathcal{R}_t - Q_\pi(s_t, r_t) \Big]$$

**temporal difference (TD-) learning**

**model-free**

- don't wait with the updates until rewards!
- use intermediate value estimates to update the action-values!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \underbrace{r_{t+1} + \gamma \, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \Big]$$

# Intuition for Temporal Difference Learning

# Intuition for Temporal Difference Learning

**temporal difference learning**

# Intuition for Temporal Difference Learning

**temporal difference learning**

- don't wait with the updates until the reward!

# Intuition for Temporal Difference Learning

**temporal difference learning**

- don't wait with the updates until the reward!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \underbrace{r_{t+1} + \gamma \, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \right]$$

# Intuition for Temporal Difference Learning

**temporal difference learning**
- don't wait with the updates until the reward!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \underbrace{r_{t+1} + \gamma \, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \right]$$

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

# Intuition for Temporal Difference Learning

**temporal difference learning**

- don't wait with the updates until the reward!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \underbrace{r_{t+1} + \gamma\, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \Big]$$

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |



temporal difference

# Intuition for Temporal Difference Learning

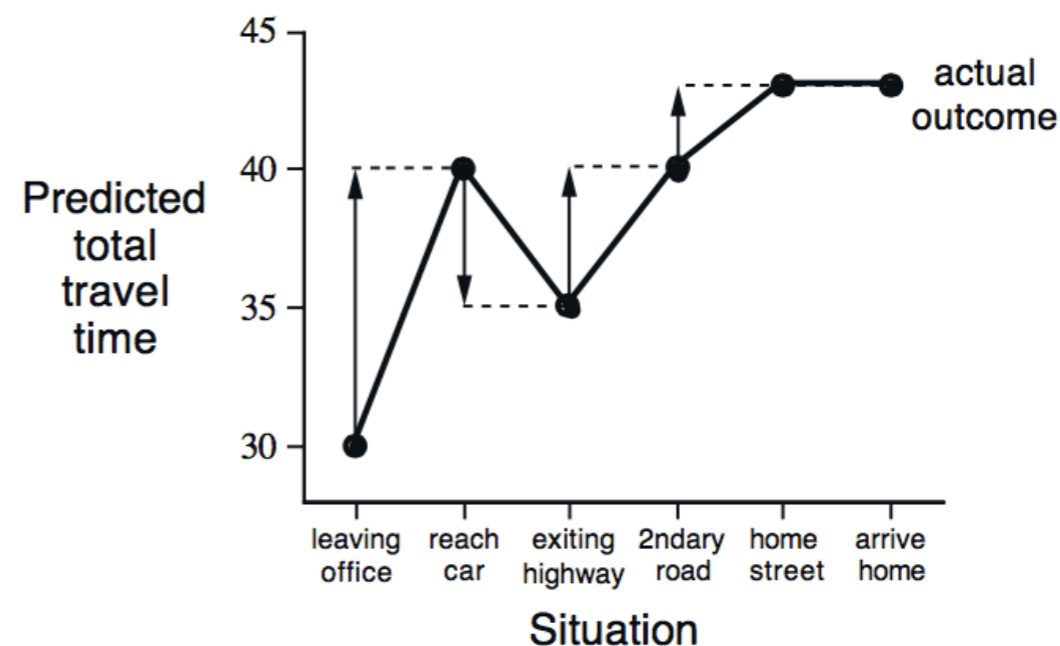**temporal difference learning**

- don't wait with the updates until the reward!

$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \left[ \underbrace{r_{t+1} + \gamma\, Q_\pi(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q_\pi(s_t, a_t) \right]$$

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office, friday at 6 | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exiting highway | 20 | 15 | 35 |
| 2ndary road, behind truck | 30 | 10 | 40 |
| entering home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |



Monte Carlo



temporal difference
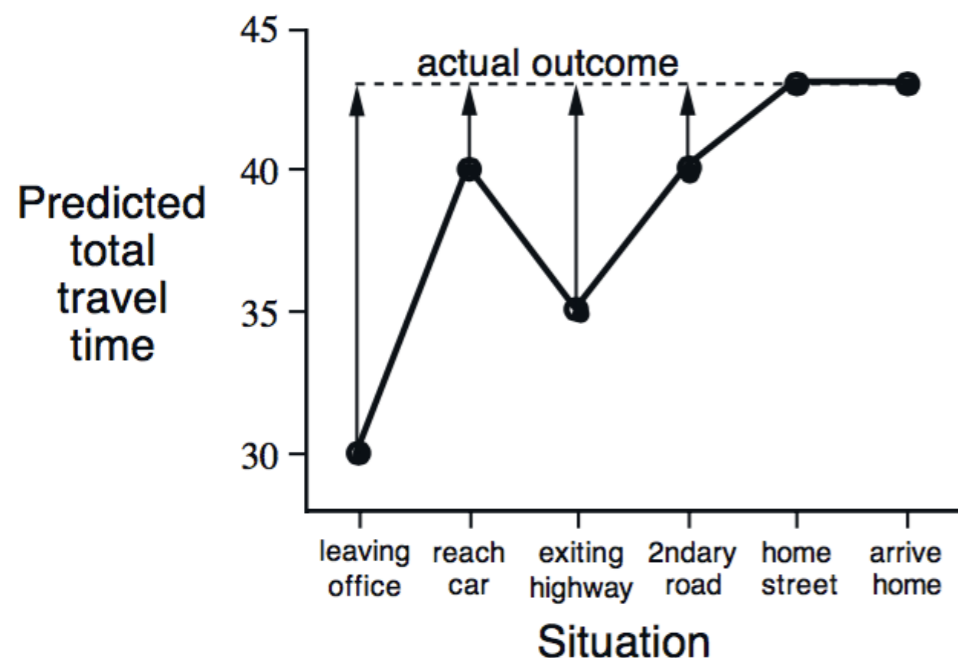
# RL in practice

**temporal difference learning**
- don't wait with the updates until the end of the trial!

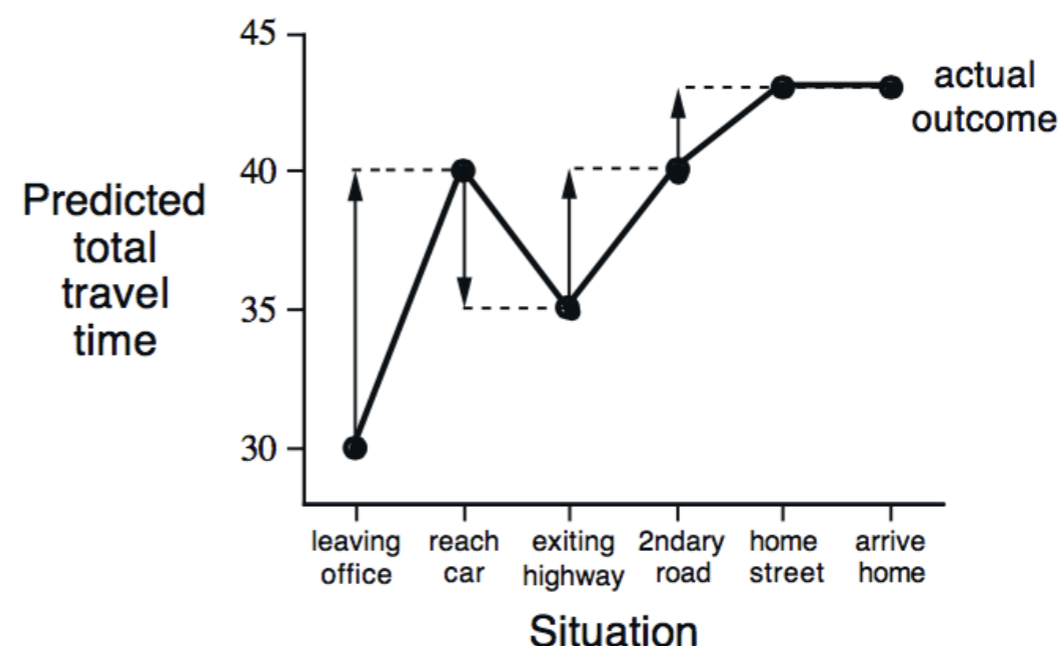$$Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha \Big[ \underbrace{r_{t+1} + \gamma\, Q(s_{t+1}, a_{t+1})}_{\text{estimate}} - Q(s_t, a_t) \Big]$$

- Q-learning:
  a powerful algorithm that has been applied to many different real-word problems
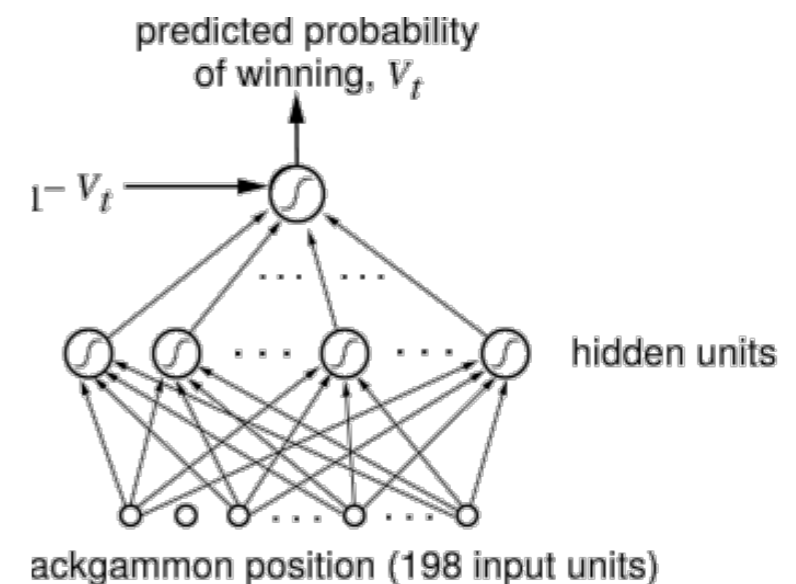
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \Big[ \underbrace{r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)}_{\text{prediction error}} \Big]$$

neuronal implementation:
- learn the state space - representational learning
- tabular vs. function approximation
- learning is based on prediction error
- is reward prediction error calculated by the brain?

# Decision making with a neural network and TD learning

- Gerard Tesauro TD-backgammon

  - Multi-layer neural network

  - Input: possible states achieved by potential moves

  - Output: the probability of winning from an actual state

- Based on these, a policy can be established

- Result: performance is compatible with the best human players

- Training the algorithm takes about 5s

## AlphaZero (Silver et al., 2018)

# Atari games (Mnih et al., 2015)



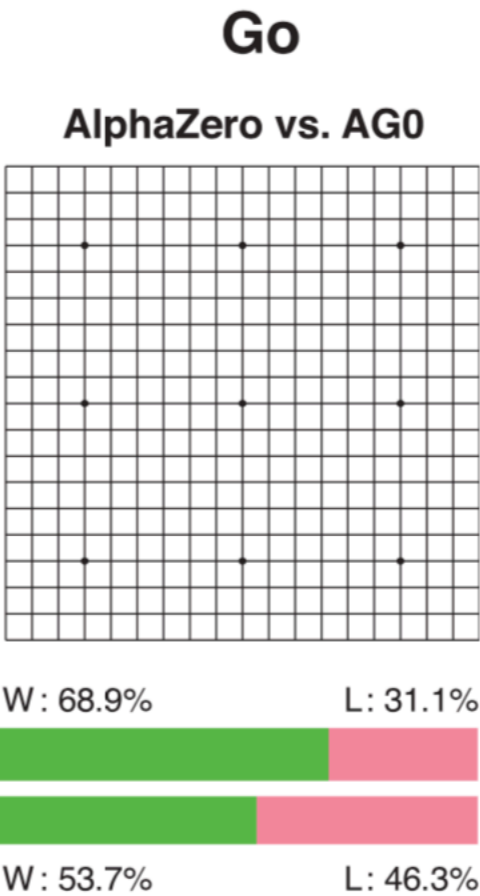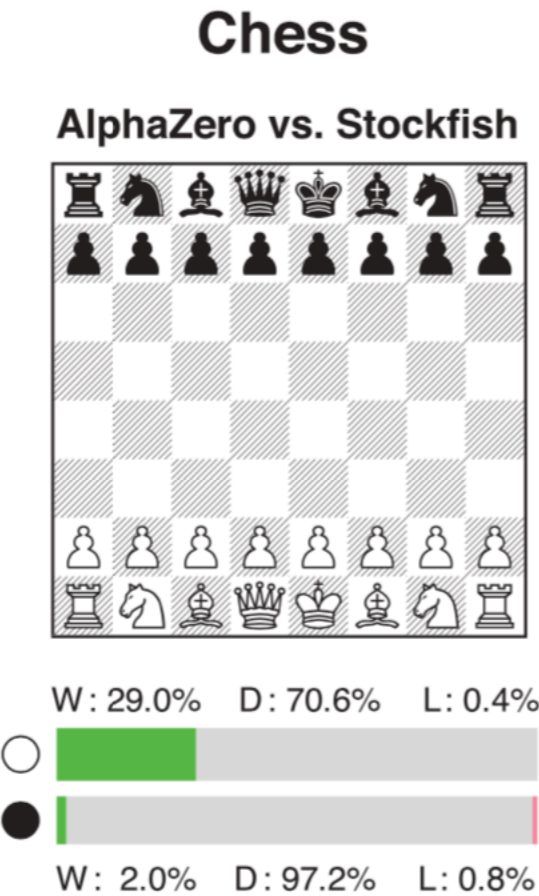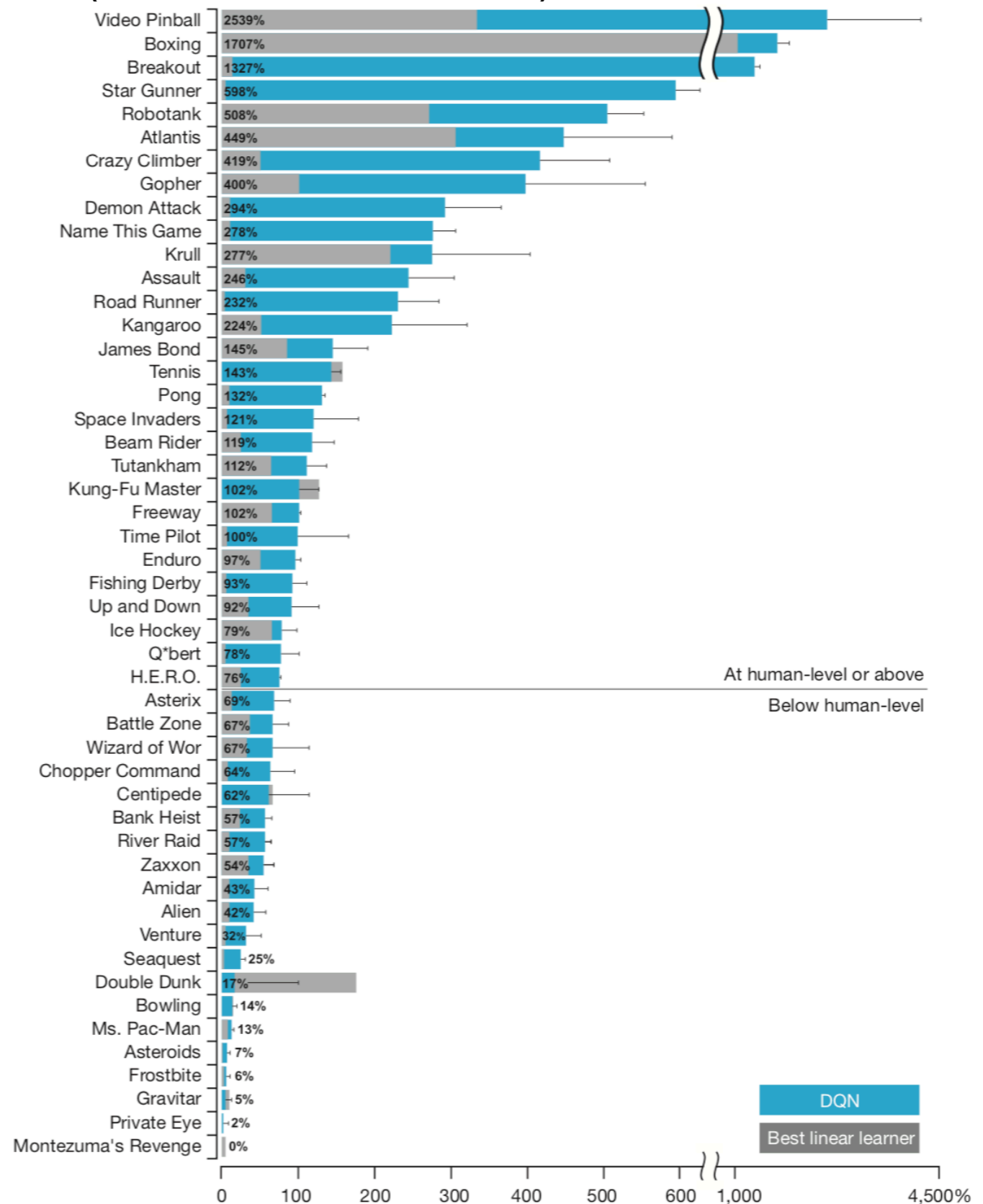| Game | Value |
|---|---|
| Video Pinball | 2539% |
| Boxing | 1707% |
| Breakout | 1327% |
| Star Gunner | 598% |
| Robotank | 508% |
| Atlantis | 449% |
| Crazy Climber | 419% |
| Gopher | 400% |
| Demon Attack | 294% |
| Name This Game | 278% |
| Krull | 277% |
| Assault | 246% |
| Road Runner | 232% |
| Kangaroo | 224% |
| James Bond | 145% |
| Tennis | 143% |
| Pong | 132% |
| Space Invaders | 121% |
| Beam Rider | 119% |
| Tutankham | 112% |
| Kung-Fu Master | 102% |
| Freeway | 102% |
| Time Pilot | 100% |
| Enduro | 97% |
| Fishing Derby | 93% |
| Up and Down | 92% |
| Ice Hockey | 79% |
| Q*bert | 78% |
| H.E.R.O. | 76% |
| Asterix | 69% |
| Battle Zone | 67% |
| Wizard of Wor | 67% |
| Chopper Command | 64% |
| Centipede | 62% |
| Bank Heist | 57% |
| River Raid | 57% |
| Zaxxon | 54% |
| Amidar | 43% |
| Alien | 42% |
| Venture | 32% |
| Seaquest | 25% |
| Double Dunk | 17% |
| Bowling | 14% |
| Ms. Pac-Man | 13% |
| Asteroids | 7% |
| Frostbite | 6% |
| Gravitar | 5% |
| Private Eye | 2% |
| Montezuma's Revenge | 0% |

At human-level or above
Below human-level

DQN
Best linear learner
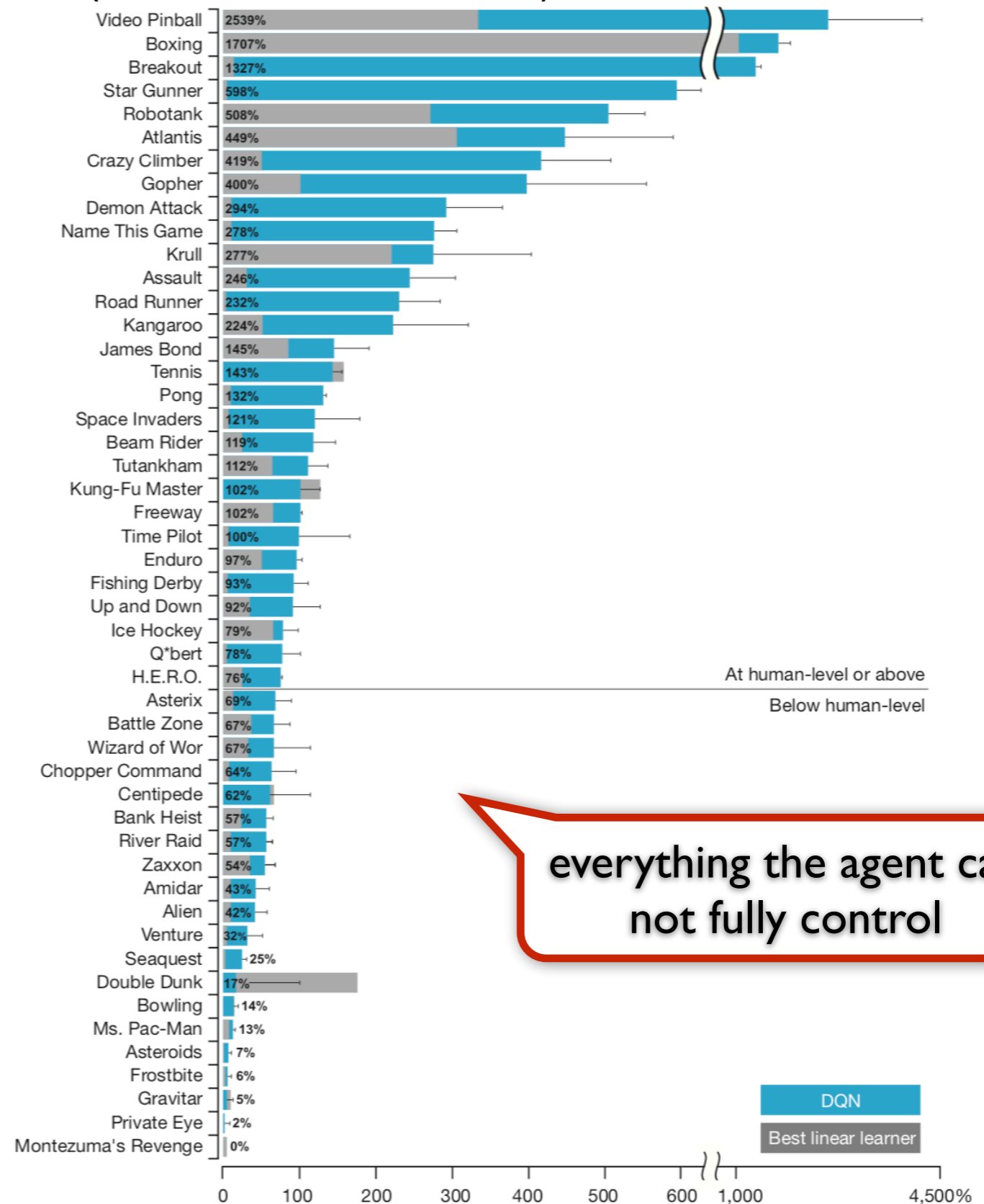
# Decision making with a neural network and TD learning Deep Q learning
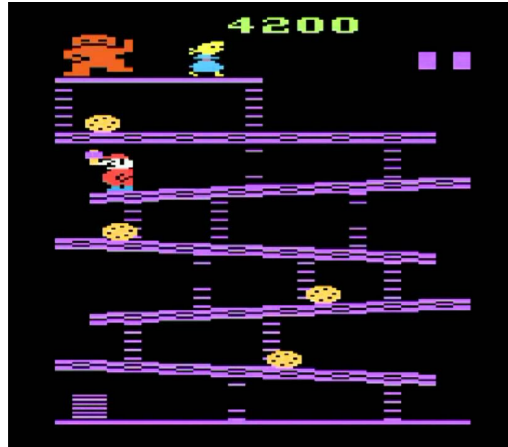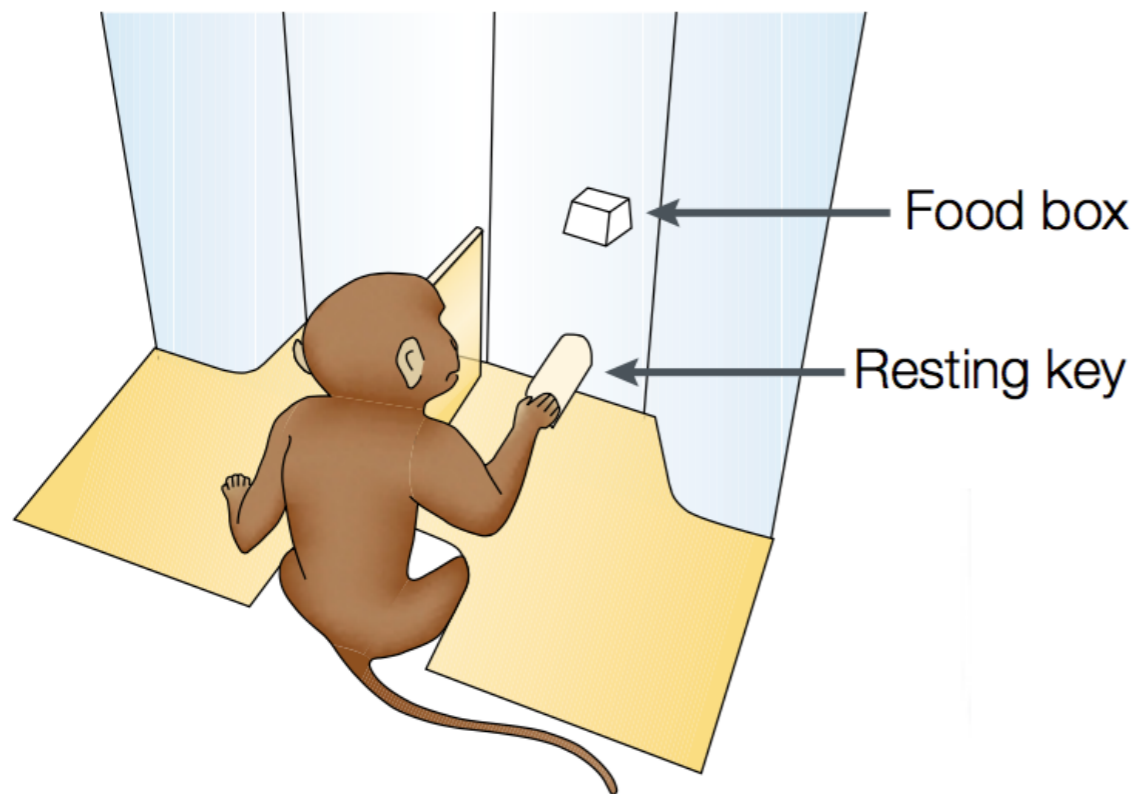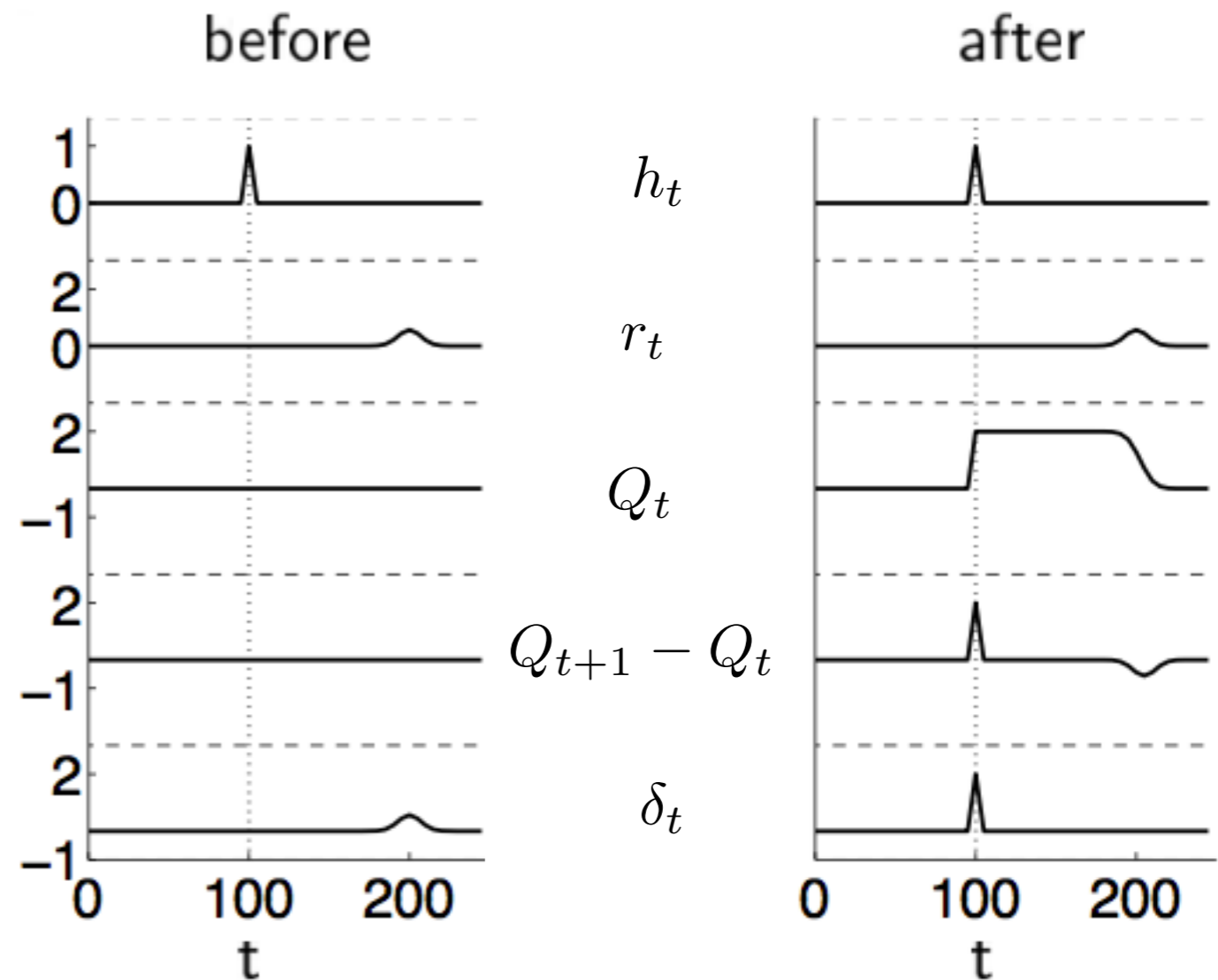
## Atari games (Mnih et al., 2015)



Video Pinball 2539%
Boxing 1707%
Breakout 1327%
Star Gunner 598%
Robotank 508%
Atlantis 449%
Crazy Climber 419%
Gopher 400%
Demon Attack 294%
Name This Game 278%
Krull 277%
Assault 246%
Road Runner 232%
Kangaroo 224%
James Bond 145%
Tennis 143%
Pong 132%
Space Invaders 121%
Beam Rider 119%
Tutankham 112%
Kung-Fu Master 102%
Freeway 102%
Time Pilot 100%
Enduro 97%
Fishing Derby 93%
Up and Down 92%
Ice Hockey 79%
Q*bert 78%
H.E.R.O. 76%
Asterix 69%
Battle Zone 67%
Wizard of Wor 67%
Chopper Command 64%
Centipede 62%
Bank Heist 57%
River Raid 57%
Zaxxon 54%
Amidar 43%
Alien 42%
Venture 32%
Seaquest 25%
Double Dunk 17%
Bowling 14%
Ms. Pac-Man 13%
Asteroids 7%
Frostbite 6%
Gravitar 5%
Private Eye 2%
Montezuma's Revenge 0%

At human-level or above
Below human-level

everything the agent can not fully control

DQN
Best linear learner

0   100   200   300   400   500   600   1,000   4,500%

# Neural representation: dopamine signal



Food box

Resting key

before

after

$h_t$

$r_t$

$Q_t$

$Q_{t+1} - Q_t$

$\delta_t$

Dopaminergic
Pathways in the Brain

Mesolimbic

Nigrostriatal

Hypothalamic
Infundibular

$$\alpha \underbrace{\left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right]}_{\delta_t}$$

# Neural representation: dopamine signal



before     after

$h_t$

$r_t$

$Q_t$

$Q_{t+1} - Q_t$

$\delta_t$

A   early   late   stimulus   reward   $t(\mathrm{s})$
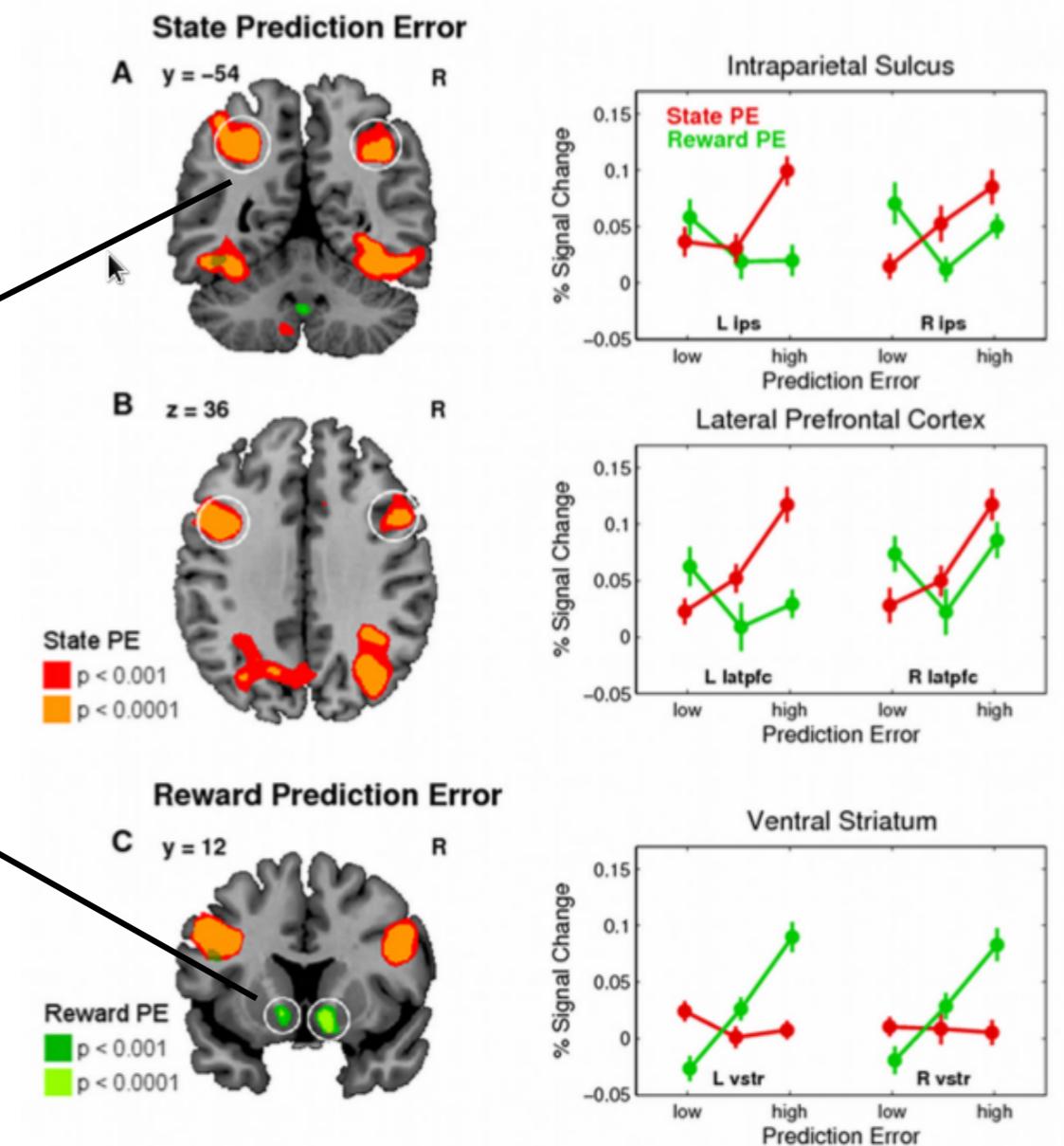
B   reward   no reward   $t(\mathrm{s})$

# Model-based RL in the brain



Deliberative, model-based RL
prefrontal and parietal cortices

Reactive, model-free RL
subcortical structures

# Conclusions

# Conclusions

- reinforcement learning:

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment

# Conclusions

- reinforcement learning:
    - learning from interaction with the environment
- **computation**: learning by interaction

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:
    - midbrain dopamine signals prediction error (PE)

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:
    - midbrain dopamine signals prediction error (PE)
    - relevant for TD learning

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:
    - midbrain dopamine signals prediction error (PE)
    - relevant for TD learning
    - PE is communicated to many brain regions

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:
    - midbrain dopamine signals prediction error (PE)
    - relevant for TD learning
    - PE is communicated to many brain regions
    - error based learning is widespread in the brain

# Conclusions

- reinforcement learning:
  - learning from interaction with the environment
- **computation**: learning by interaction
- **algorithm**:
  - TD-learning
  - model-based
- **implementation**:
  - TD-learning:
    - midbrain dopamine signals prediction error (PE)
    - relevant for TD learning
    - PE is communicated to many brain regions
    - error based learning is widespread in the brain