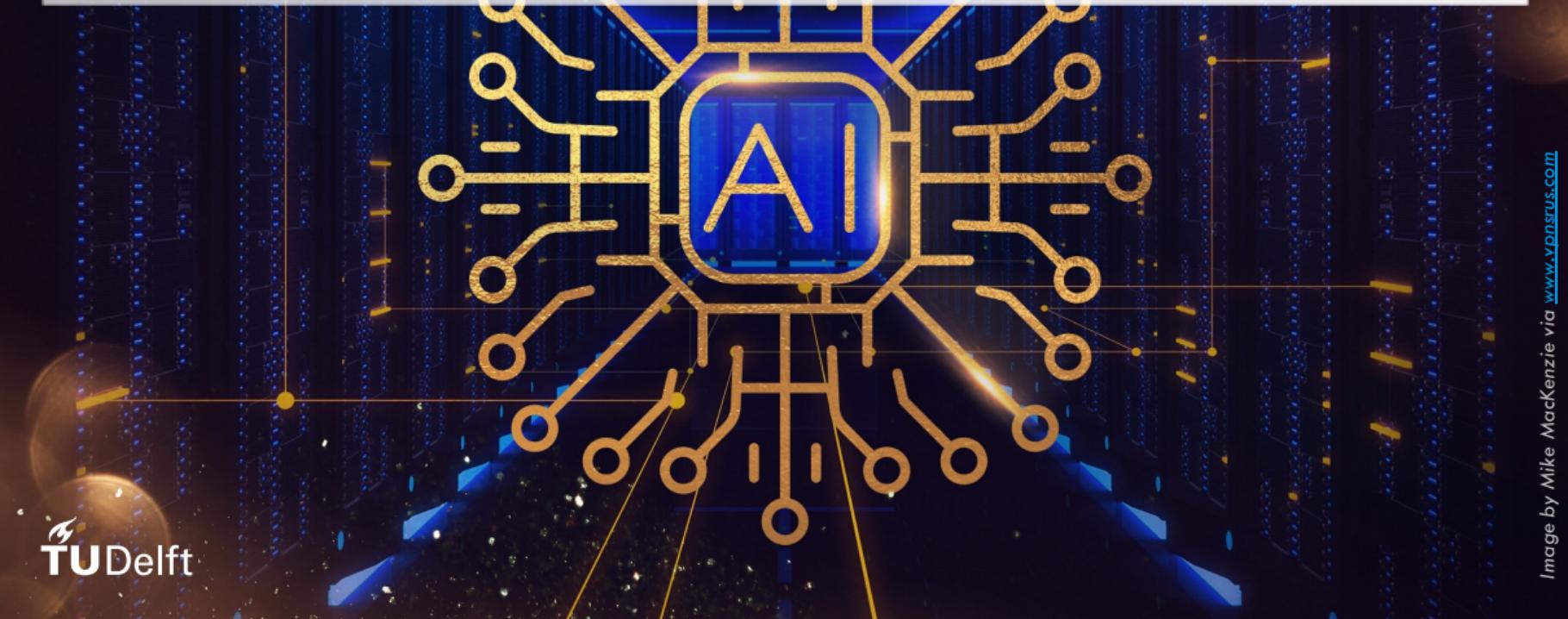


# Fundamentals of AI Programme

## Reinforcement Learning

Dr. Wendelin Böhmer <[j.w.bohmer@tudelft.nl](mailto:j.w.bohmer@tudelft.nl)>



# Content of this lecture

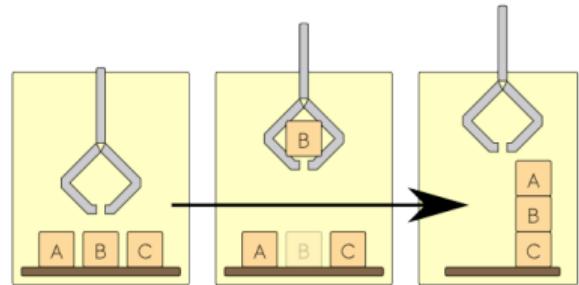
- 1 Learning to Interact
- 2 Tabular Q-learning
- 3 Deep Reinforcement Learning
- 4 State-of-the-Art and Limits

1

# Reinforcement Learning: Learning to Interact

# AI methods for interaction

- Search & optimization
  - task: reach goal with minimal costs
  - given: dynamics model of environment
  - prediction: action sequence that reaches goal



# AI methods for interaction

- Search & optimization
  - task: reach goal with minimal costs
  - given: dynamics model of environment
  - prediction: action sequence that reaches goal
- But what if we do **not know** the environment?
  - go boldly where no one has gone before
  - uncertain dynamics or observations
  - unknown dynamics, e.g. humans in the loop
- Sometimes we need to **learn** from interaction!

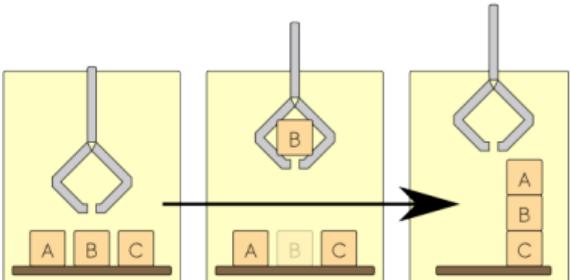
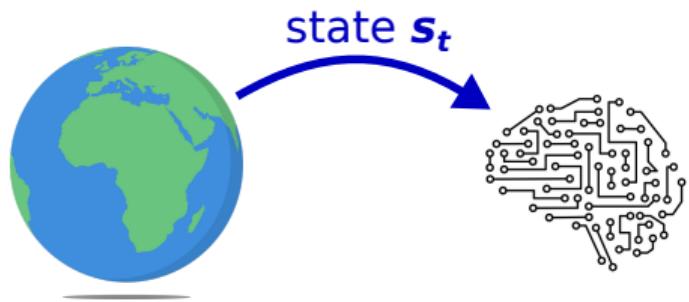


image source: Steven Shorrock

# Markov decision processes

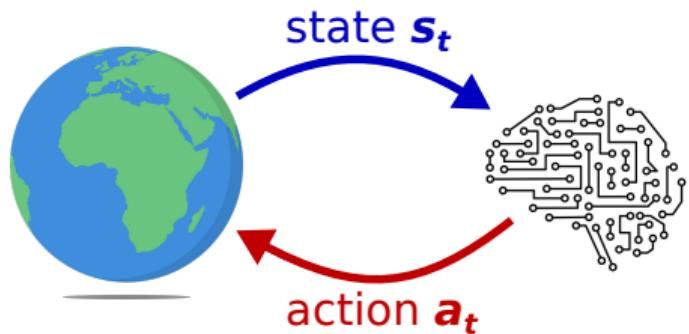
- MDP  $M = \langle \mathcal{S}, \mathcal{A}, \rho, P, r \rangle$  consists of:
  - state space  $\mathcal{S}$  of environment configurations
  - distribution  $\rho(s_0)$  of initial states  $s_0 \in \mathcal{S}$



see Chapter 3 in Sutton and Barto (2018) for a more formal introduction

# Markov decision processes

- MDP  $M = \langle \mathcal{S}, \mathcal{A}, \rho, P, r \rangle$  consists of:
  - state space  $\mathcal{S}$  of environment configurations
  - distribution  $\rho(s_0)$  of initial states  $s_0 \in \mathcal{S}$
  - action space  $\mathcal{A}$  of possible agent choices



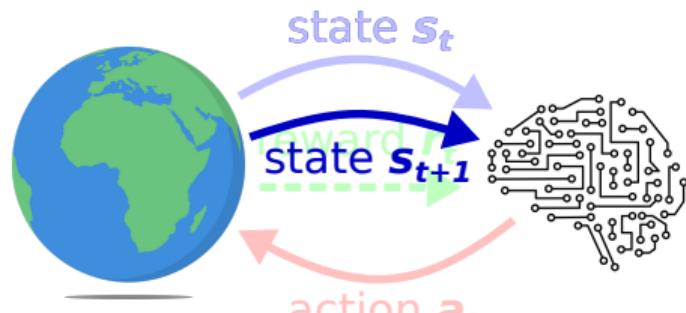
# Markov decision processes

- MDP  $M = \langle \mathcal{S}, \mathcal{A}, \rho, P, r \rangle$  consists of:
  - state space  $\mathcal{S}$  of environment configurations
  - distribution  $\rho(s_0)$  of initial states  $s_0 \in \mathcal{S}$
  - action space  $\mathcal{A}$  of possible agent choices
  - reward function  $r(s_t, a_t)$  for choosing good or bad actions  $a_t \in \mathcal{A}$  in  $s_t \in \mathcal{S}$



# Markov decision processes

- MDP  $M = \langle \mathcal{S}, \mathcal{A}, \rho, P, r \rangle$  consists of:
  - state space  $\mathcal{S}$  of environment configurations
  - distribution  $\rho(s_0)$  of initial states  $s_0 \in \mathcal{S}$
  - action space  $\mathcal{A}$  of possible agent choices
  - reward function  $r(s_t, a_t)$  for choosing good or bad actions  $a_t \in \mathcal{A}$  in  $s_t \in \mathcal{S}$
  - distribution  $P(s_{t+1}|s_t, a_t)$  of transitions from  $s_t \in \mathcal{S}$  to  $s_{t+1} \in \mathcal{S}$  after choosing  $a_t \in \mathcal{A}$



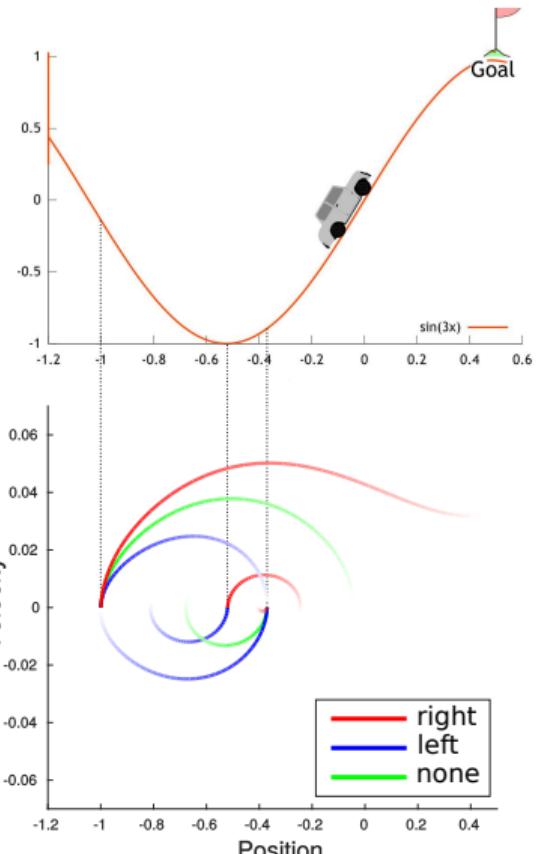
# Markov decision processes

- MDP  $M = \langle \mathcal{S}, \mathcal{A}, \rho, P, r \rangle$  consists of:
  - state space  $\mathcal{S}$  of environment configurations
  - distribution  $\rho(s_0)$  of initial states  $s_0 \in \mathcal{S}$
  - action space  $\mathcal{A}$  of possible agent choices
  - reward function  $r(s_t, a_t)$  for choosing good or bad actions  $a_t \in \mathcal{A}$  in  $s_t \in \mathcal{S}$
  - distribution  $P(s_{t+1}|s_t, a_t)$  of transitions from  $s_t \in \mathcal{S}$  to  $s_{t+1} \in \mathcal{S}$  after choosing  $a_t \in \mathcal{A}$
- Reinforcement learning aims for **actions** that maximize **all future reward**
  - future reward depends where we end up in the future → complex problem
  - no prior knowledge of **reward** or **transition** model



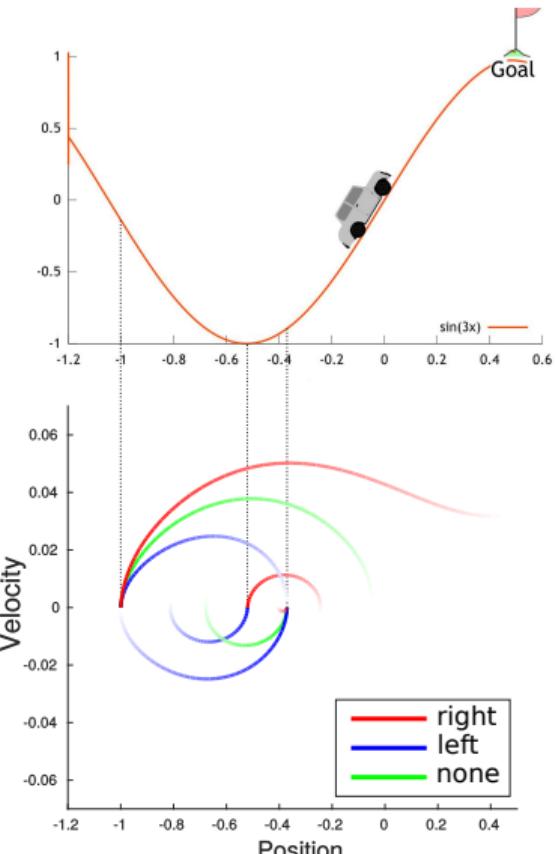
# Example: MountainCar

- Underactuated car in a valley between mountains
  - state is position and velocity:  $s_t = [x_t, v_t] \in \mathcal{S} \subset \mathbb{R}^2$
  - actions accelerate car  $a_t \in \mathcal{A} = \{\text{left}, \text{none}, \text{right}\}$
  - physics dynamics  $P$ : gravitation, but no friction



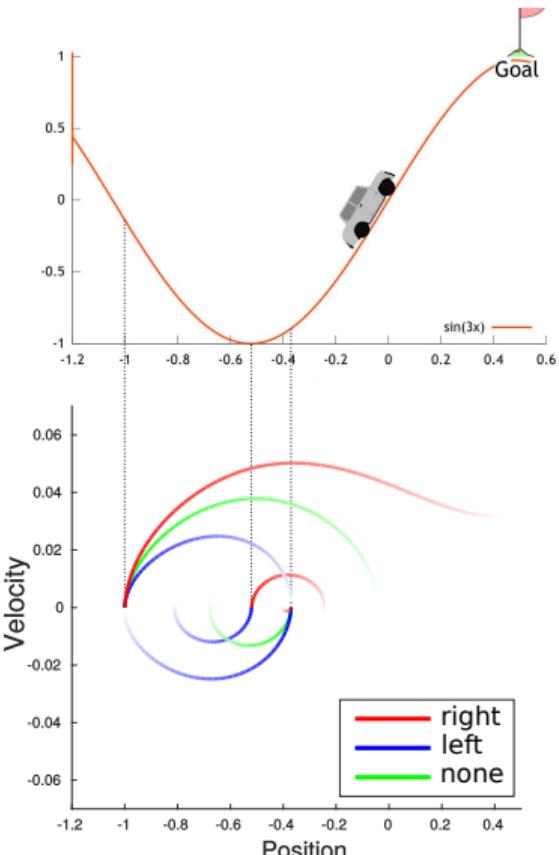
# Example: MountainCar

- Underactuated car in a valley between mountains
  - state is position and velocity:  $s_t = [x_t, v_t] \in \mathcal{S} \subset \mathbb{R}^2$
  - actions accelerate car  $a_t \in \mathcal{A} = \{\text{left}, \text{none}, \text{right}\}$
  - physics dynamics  $P$ : gravitation, but no friction
  - reward of 1 for reaching the flag, otherwise 0



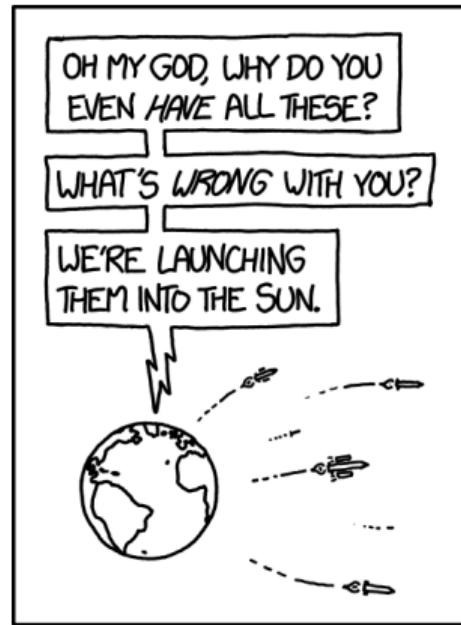
# Example: MountainCar

- Underactuated car in a valley between mountains
  - state is position and velocity:  $s_t = [x_t, v_t] \in \mathcal{S} \subset \mathbb{R}^2$
  - actions accelerate car  $a_t \in \mathcal{A} = \{\text{left}, \text{none}, \text{right}\}$
  - physics dynamics  $P$ : gravitation, but no friction
  - reward of 1 for reaching the flag, otherwise 0
  - underactuated car needs to pick up speed first



# Learning to interact

- In this video we have:
  - discussed under which situations we need to learn interaction
  - learned how a stochastic environment is formalized as a MDP
  - familiarized ourselves with the MountainCar environment



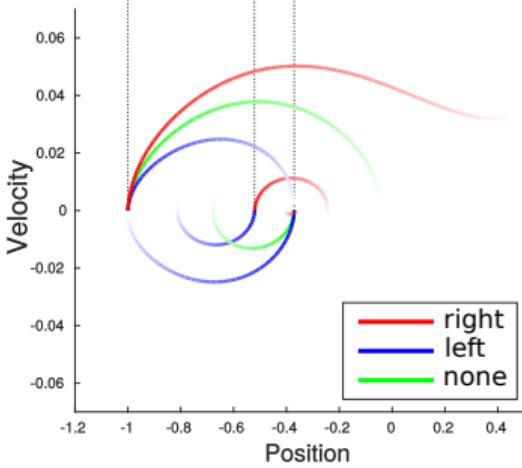
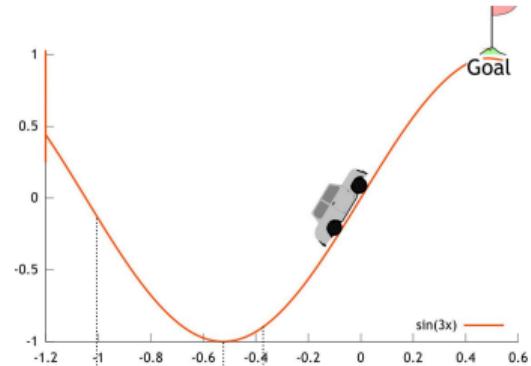
THE MOMENT THE COMPUTERS  
CONTROLLING OUR NUCLEAR  
ARSENALS BECAME SENTIENT

2

# Reinforcement Learning: Tabular Q-learning

# Agent policies

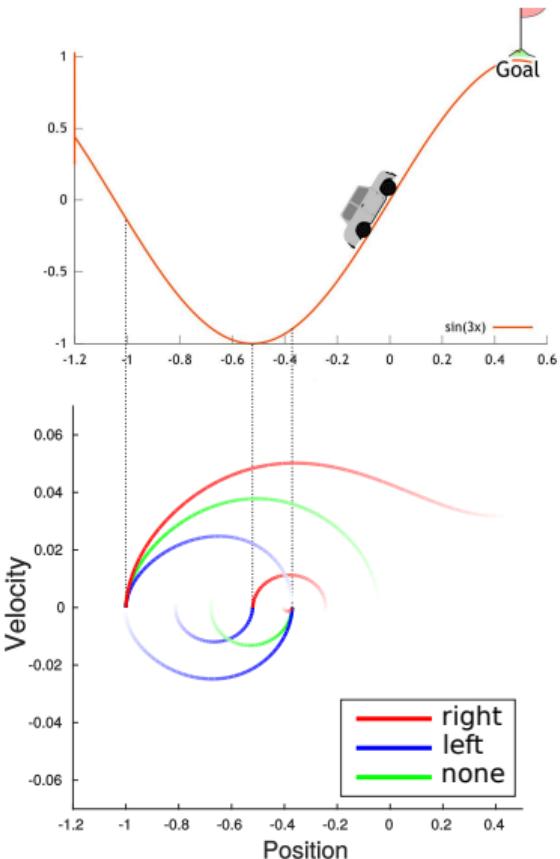
- Policy  $\pi$  defines agent's behavior
  - $\pi(a|s)$  probability the agent chooses  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$



# Agent policies

- Policy  $\pi$  defines agent's behavior
  - $\pi(a|s)$  probability the agent chooses  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$
  - allows expectations over all possible futures

$$\mathbb{E}_\pi [\dots] := \mathbb{E} \left[ \dots \middle| \begin{array}{l} s_0 \sim \rho(\cdot), \quad r_t = r(s_t, a_t) \\ a_t \sim \pi(\cdot | s_t), \quad s_{t+1} \sim P(\cdot | s_t, a_t) \end{array}, \forall t \geq 0 \right]$$

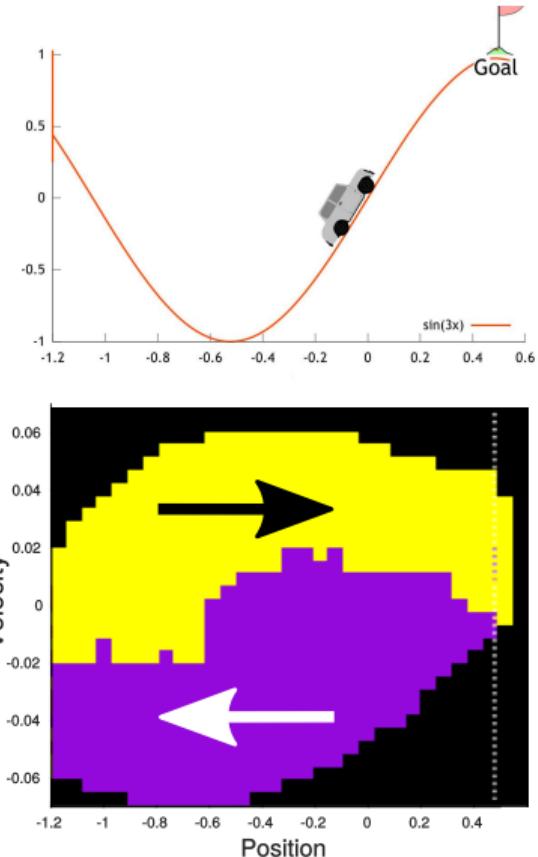


# Agent policies

- Policy  $\pi$  defines agent's behavior
  - $\pi(a|s)$  probability the agent chooses  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$
  - allows expectations over all possible futures

$$\mathbb{E}_\pi [\dots] := \mathbb{E} \left[ \dots \middle| s_0 \sim \rho(\cdot), r_t = r(s_t, a_t), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t), \forall t \geq 0 \right]$$

- Discretize continuous states  $s \in \mathcal{S}$ 
  - policy table/matrix  $\Pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  with  $\Pi_{sa} := \pi(a|s)$

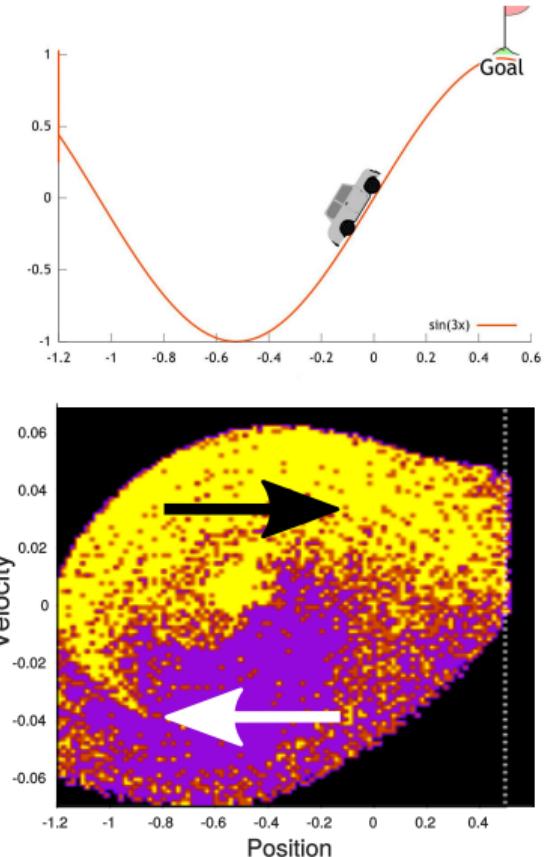


# Agent policies

- Policy  $\pi$  defines agent's behavior
  - $\pi(a|s)$  probability the agent chooses  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$
  - allows expectations over all possible futures

$$\mathbb{E}_\pi [\dots] := \mathbb{E} \left[ \dots \middle| s_0 \sim \rho(\cdot), r_t = r(s_t, a_t), a_t \sim \pi(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t), \forall t \geq 0 \right]$$

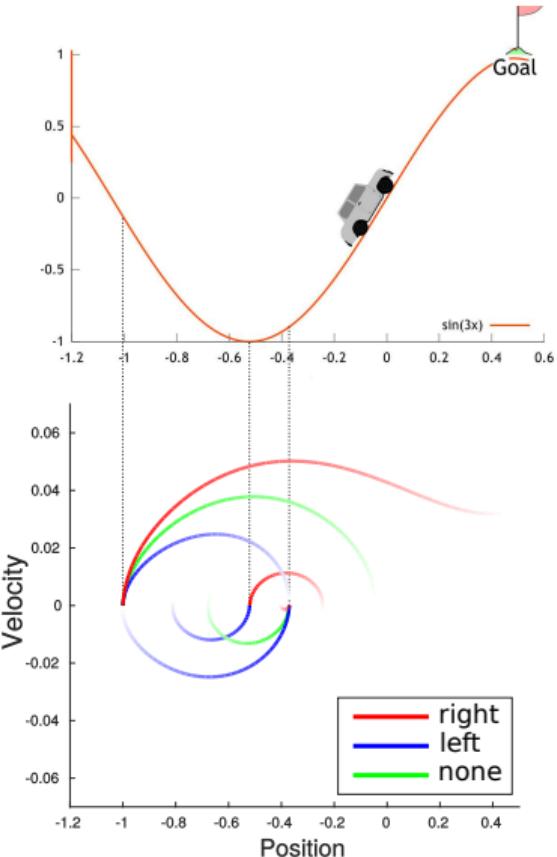
- Discretize continuous states  $s \in \mathcal{S}$ 
  - policy table/matrix  $\Pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  with  $\Pi_{sa} := \pi(a|s)$
- resolution determines control precision



# Policy evaluation

- Policy should maximize the **value**  $V^\pi$ 
  - expected sum of discounted future rewards

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t \middle| s_0=s \right] , \gamma \in [0, 1)$$

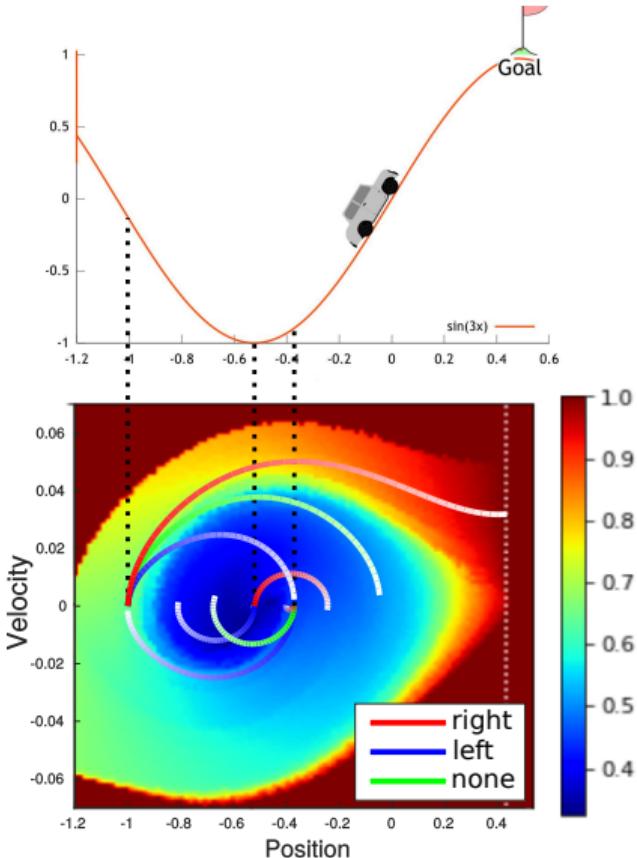


the infinite sum is a geometric series with  $\sum_{t=0}^{\infty} \gamma^t = \frac{1}{1-\gamma}$

# Policy evaluation

- Policy should maximize the **value**  $V^\pi$ 
  - expected sum of discounted future rewards
  - same discretization  $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$  as policy

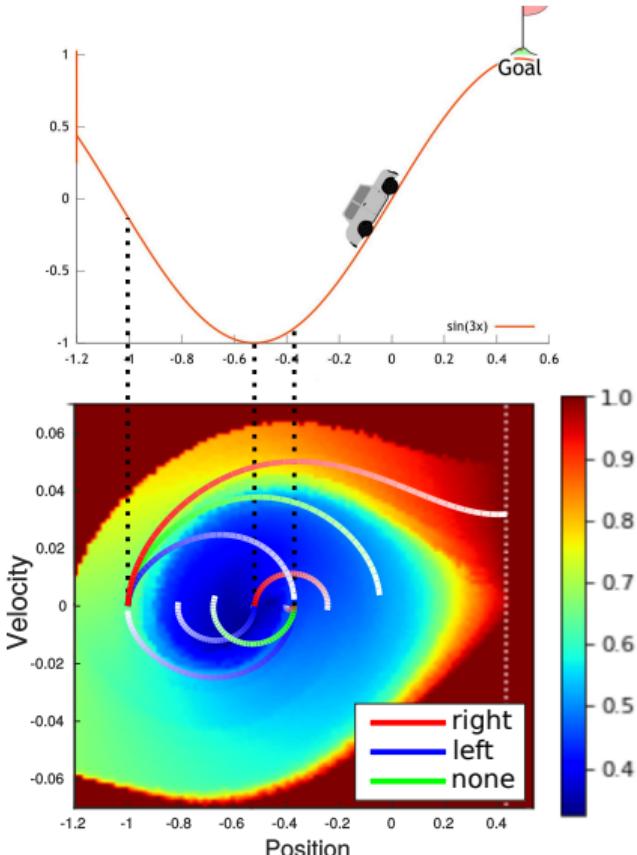
$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t \middle| s_0=s \right] , \gamma \in [0, 1)$$



# Policy evaluation

- Policy should maximize the **value**  $V^\pi$ 
  - expected sum of discounted future rewards
  - same discretization  $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$  as policy

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t \middle| s_0=s \right] \quad , \gamma \in [0, 1) \\ &= \mathbb{E}_\pi \left[ r_0 + \underbrace{\mathbb{E}_\pi [\gamma r_1 + \gamma^2 r_2 + \dots | s_1]}_{\gamma V^\pi(s_1)} \middle| s_0=s \right] \end{aligned}$$



this equation goes back to Richard Bellman (1957), see Section 3.5 in Sutton and Barto (2018)

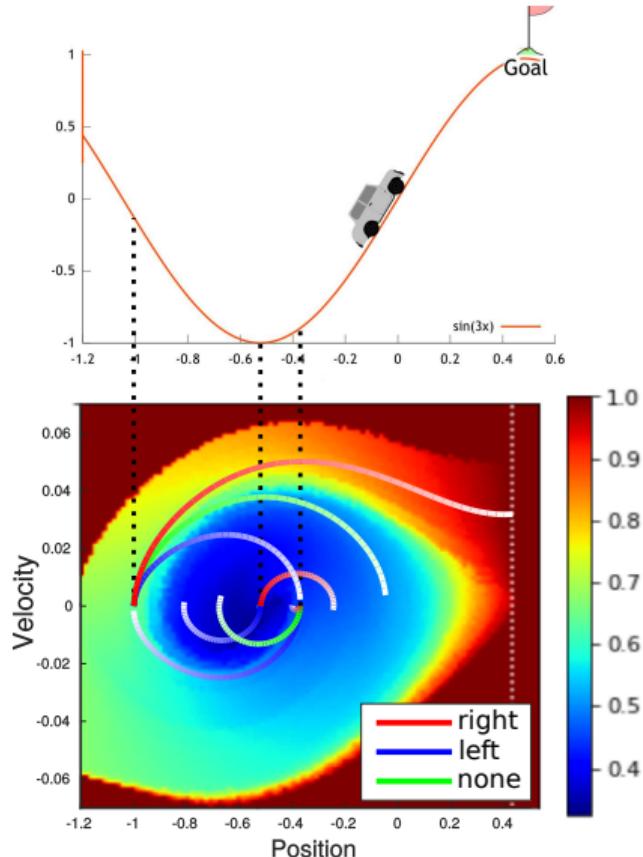
# Policy evaluation

- Policy should maximize the **value**  $V^\pi$ 
  - expected sum of discounted future rewards
  - same discretization  $\mathbf{V} \in \mathbb{R}^{|S|}$  as policy

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t r_t \middle| s_0=s \right] \quad , \gamma \in [0, 1) \\ &= \mathbb{E}_\pi \left[ r_0 + \underbrace{\mathbb{E}_\pi [\gamma r_1 + \gamma^2 r_2 + \dots | s_1]}_{\gamma V^\pi(s_1)} \middle| s_0=s \right] \\ &= \mathbb{E}_\pi \left[ r_0 + \gamma V^\pi(s_1) \middle| s_0=s \right] \end{aligned}$$

- Recursive **Bellman equation** of the state value

this equation goes back to Richard Bellman (1957), see Section 3.5 in Sutton and Barto (2018)



# The optimal Bellman equation

$$\max_{\pi} V^{\pi}(s) = \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s \right]$$

# The optimal Bellman equation

$$\max_{\pi} V^{\pi}(s) = \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right] = \max_a \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right]}_{Q^*(s,a)}$$

# The optimal Bellman equation

$$\begin{aligned} \max_{\pi} V^{\pi}(s) &= \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right] = \max_{\mathbf{a}} \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right]}_{Q^*(s,a)} \\ &\stackrel{(*)}{=} \max_{\mathbf{a}} \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{\pi} V^{\pi}(s_1) \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right] \end{aligned}$$

(\*)  $\max_{\pi} \mathbb{E}[\dots] = \mathbb{E}[\max_{\pi} \dots]$  only holds because a policy  $\pi^*$  exists that is optimal in all states:  $V^{\pi^*}(s) \geq V^{\pi}(s), \forall \pi, \forall s \in \mathcal{S}$ .

# The optimal Bellman equation

$$\begin{aligned} \max_{\pi} V^{\pi}(s) &= \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right] = \max_{\mathbf{a}} \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right]}_{Q^*(s,a)} \\ &\stackrel{(*)}{=} \max_{\mathbf{a}} \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{\pi} V^{\pi}(s_1) \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right] \\ &= \max_{\mathbf{a}} \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{\mathbf{a}'} \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s'_t, a'_t) \middle| \begin{array}{l} s'_0=s_1 \\ a'_0=a' \end{array} \right]}_{Q^*(s_1,a')} \middle| \begin{array}{l} s_0=s \\ a_0=a \end{array} \right] \end{aligned}$$

(\*)  $\max_{\pi} \mathbb{E}[\dots] = \mathbb{E}[\max_{\pi} \dots]$  only holds because a policy  $\pi^*$  exists that is optimal in all states:  $V^{\pi^*}(s) \geq V^{\pi}(s)$ ,  $\forall \pi, \forall s \in \mathcal{S}$ .

# The optimal Bellman equation

$$\begin{aligned} \max_{\pi} V^{\pi}(s) &= \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right] = \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right]}_{Q^*(s,a)} \\ &\stackrel{(*)}{=} \max_a \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{\pi} V^{\pi}(s_1) \middle| a_0=a \right] \\ &= \max_a \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{a'} \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s'_t, a'_t) \middle| a'_0=a' \right]}_{Q^*(s_1, a')} \middle| a_0=a \right] \end{aligned}$$

- Optimal state-action **Q-values**:

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E} \left[ \max_{a'} Q^*(s', a') \middle| s' \sim P(\cdot | s, a) \right]$$

# The optimal Bellman equation

$$\begin{aligned} \max_{\pi} V^{\pi}(s) &= \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right] = \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0=s \right]}_{Q^*(s,a)} \\ &\stackrel{(*)}{=} \max_a \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{\pi} V^{\pi}(s_1) \middle| a_0=a \right] \\ &= \max_a \mathbb{E} \left[ r(s_0, a_0) + \gamma \max_{a'} \underbrace{\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s'_t, a'_t) \middle| a'_0=a' \right]}_{Q^*(s_1, a')} \middle| a_0=a \right] \end{aligned}$$

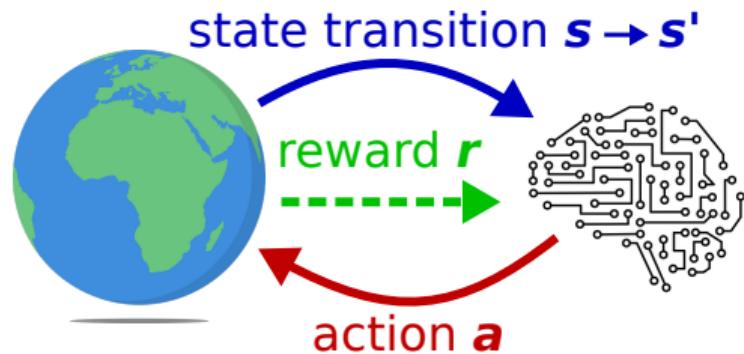
- Optimal state-action **Q-values** and **policy**:

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E} \left[ \max_{a'} Q^*(s', a') \middle| s' \sim P(\cdot | s, a) \right]$$

$$\pi^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a'} Q^*(s, a') \\ 0, & \text{otherwise} \end{cases}$$

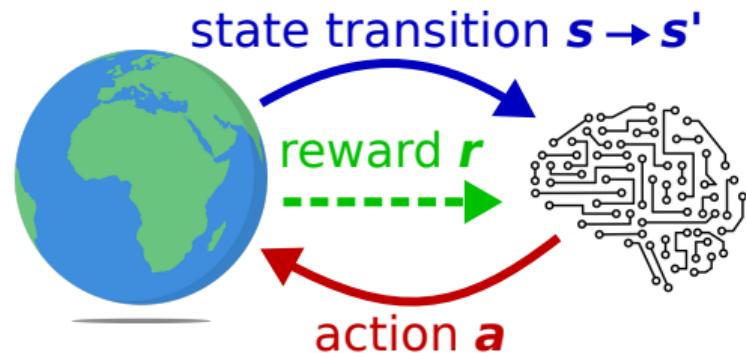
# Tabular Q-learning

- Estimate optimal Q-value  $Q^*$ 
  - while interacting with environment
  - MDP is unknown, can only be sampled



# Tabular Q-learning

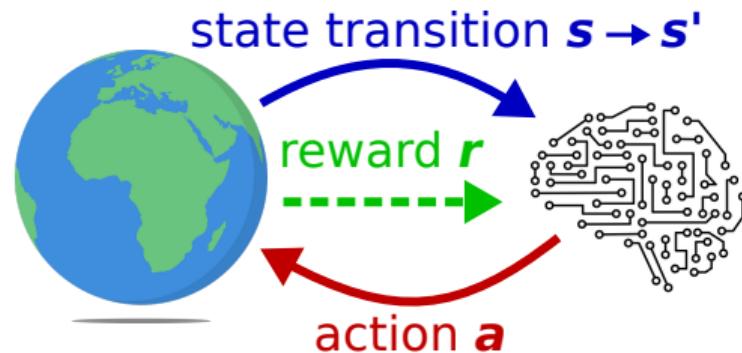
- Estimate optimal Q-value  $Q^*$ 
  - while interacting with environment
  - MDP is unknown, can only be sampled
- Table  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  of finite  $\mathcal{S} \times \mathcal{A}$ 
  - same discretization as policy



# Tabular Q-learning

- Estimate optimal Q-value  $Q^*$ 
  - while interacting with environment
  - MDP is unknown, can only be sampled
- Table  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  of finite  $\mathcal{S} \times \mathcal{A}$ 
  - same discretization as policy
- Update table  $\mathbf{Q}$  after observing transition  $s, a \rightarrow r, s'$ 
  - temporal averaging with learning rate  $\alpha \in (0, 1]$

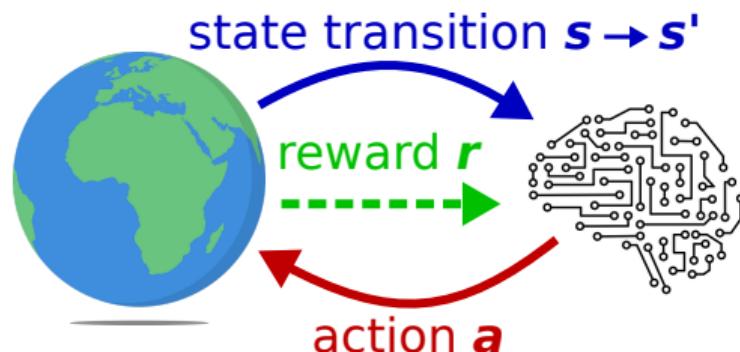
$$Q_{sa} \leftarrow (1 - \alpha) \underbrace{Q_{sa}}_{\text{old value}} + \alpha \underbrace{\left( r + \gamma \max_{a'} Q_{s'a'} \right)}_{\text{Bellman update}}$$



# Tabular Q-learning

- Estimate optimal Q-value  $Q^*$ 
  - while interacting with environment
  - MDP is unknown, can only be sampled
- Table  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  of finite  $\mathcal{S} \times \mathcal{A}$ 
  - same discretization as policy
- Update table  $\mathbf{Q}$  after observing transition  $s, a \rightarrow r, s'$ 
  - temporal averaging with learning rate  $\alpha \in (0, 1]$

$$Q_{sa} \leftarrow (1 - \alpha) \underbrace{Q_{sa}}_{\text{old value}} + \alpha \underbrace{\left( r + \gamma \max_{a'} Q_{s'a'} \right)}_{\text{Bellman update}} = Q_{sa} + \alpha \underbrace{\left( r + \gamma \max_{a'} Q_{s'a'} - Q_{sa} \right)}_{\text{temporal difference (TD) error}}$$



# Tabular Q-learning

- Estimate optimal Q-value  $Q^*$ 
  - while interacting with environment
  - MDP is unknown, can only be sampled
- Table  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$  of finite  $\mathcal{S} \times \mathcal{A}$ 
  - same discretization as policy
- Update table  $\mathbf{Q}$  after observing transition  $s, a \rightarrow r, s'$ 
  - temporal averaging with learning rate  $\alpha \in (0, 1]$

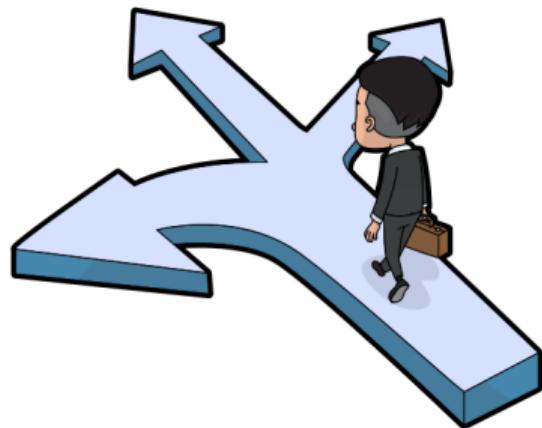
$$Q_{sa} \leftarrow (1 - \alpha) \underbrace{Q_{sa}}_{\text{old value}} + \alpha \underbrace{\left( r + \gamma \max_{a'} Q_{s'a'} \right)}_{\text{Bellman update}} = Q_{sa} + \underbrace{\alpha \left( r + \gamma \max_{a'} Q_{s'a'} - Q_{sa} \right)}_{\text{temporal difference (TD) error}}$$

tabular Q-learning has been introduced by Watkins and Dayan (1992)

```
1 s' = environment.reset()
2 while not done:
3     s, a = s', argmax(Q[s'])
4     s', r, done = environment.step(a)
5     td = r + γ * max(Q[s']) - Q[s, a]
6     Q[s, a] += α * td
```

# Exploration and exploitation

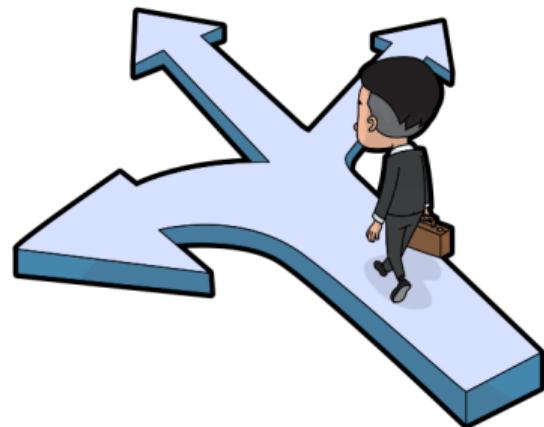
- Tabular Q-learning converges to optimal  $Q^*$ 
  - if all state-action pairs are visited often enough
  - how can we guarantee that we **explore** enough?



# Exploration and exploitation

- Tabular Q-learning converges to optimal  $Q^*$ 
  - if all state-action pairs are visited often enough
  - how can we guarantee that we **explore** enough?
- 1  $\epsilon$ -greedy: explore randomly with probability  $\epsilon$

$$\pi_\epsilon(a|s) := \underbrace{(1 - \epsilon) \pi^*(a|s)}_{\text{greedy}} + \underbrace{\epsilon \frac{1}{|\mathcal{A}|}}_{\text{random}}$$



# Exploration and exploitation

- Tabular Q-learning converges to optimal  $Q^*$

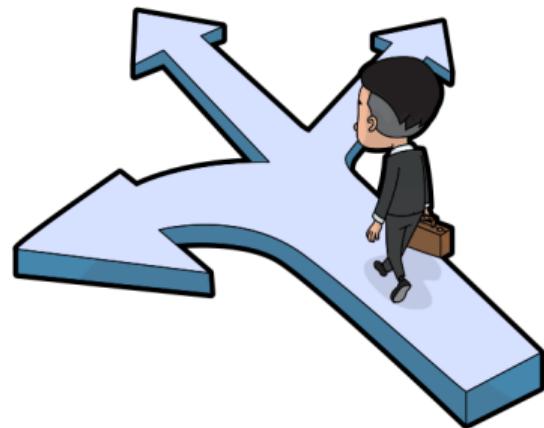
- if all state-action pairs are visited often enough
- how can we guarantee that we **explore** enough?

- ➊  $\epsilon$ -greedy: explore randomly with probability  $\epsilon$

$$\pi_\epsilon(a|s) := \underbrace{(1 - \epsilon) \pi^*(a|s)}_{\text{greedy}} + \underbrace{\epsilon \frac{1}{|\mathcal{A}|}}_{\text{random}}$$

- ➋ optimistic initialization: overestimate values

- initialize  $Q_{sa} > \max_{s', a'} Q^*(s', a')$
- updates reduce table entries towards  $Q^*(s, a)$
- under-explored actions are more attractive



# Tabular Q-learning

- In this video we have:
  - defined the agent's behavior as policy  $\pi$
  - defined how to evaluate a policy  $V^\pi$
  - derived the value  $Q^*$  of the optimal policy
  - introduced the tabular Q-learning algorithm
  - discussed how to guarantee sufficient exploration



3

# Reinforcement Learning: Deep Reinforcement Learning

# Scaling up reinforcement learning

- Tabular reinforcement learning is restricted
  - small state-action spaces
  - curse of dimensionality

# Scaling up reinforcement learning

- Tabular reinforcement learning is restricted
  - small state-action spaces
  - curse of dimensionality
- How can we learn to e.g. play ATARI games
  - no model available for planning
  - image is high dimensional input
  - no feasible discretization



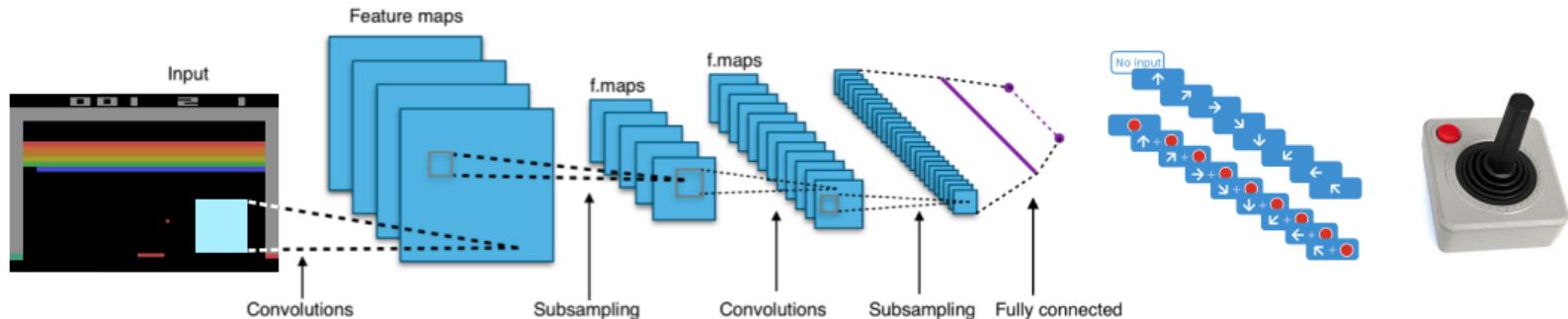
# Scaling up reinforcement learning

- Tabular reinforcement learning is restricted
  - small state-action spaces
  - curse of dimensionality
- How can we learn to e.g. play ATARI games
  - no model available for planning
  - image is high dimensional input
  - no feasible discretization
- Can we replace the Q-table with a neural network?



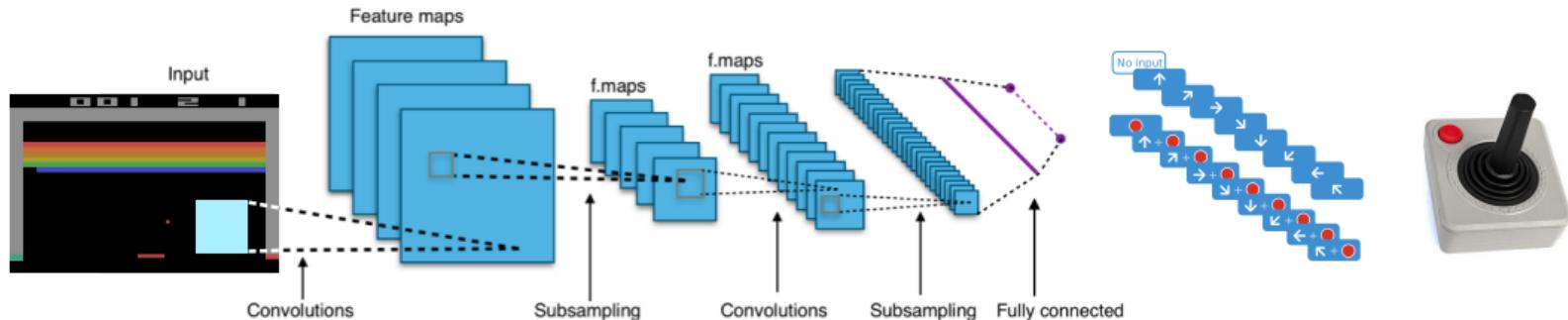
# Deep Q-networks (DQN)

- Convolutional neural network  $q_{\theta}(s) \in \mathbb{R}^{|\mathcal{A}|}$  with one output per action



# Deep Q-networks (DQN)

- Convolutional neural network  $q_\theta(s) \in \mathbb{R}^{|\mathcal{A}|}$  with one output per action



- For transitions  $s, a \rightarrow r, s'$  we update parameters  $\theta$ 
  - gradient descent on mean-squared TD error

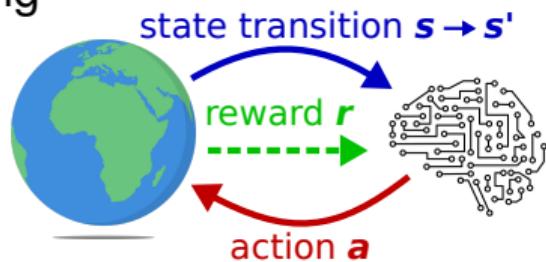
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \underbrace{\left( r + \gamma \max_{a'} (q_{\theta}(s'))_{a'} - (q_{\theta}(s))_a \right)^2}_{\text{TD-error}}$$



image modified from Apex34 and Bilby (wikipedia.com)

# RL is not supervised learning!

- Q-learning violates assumptions of supervised learning
  - TD-error: regression targets are *non-stationary*
  - exploration: the state-distribution is *non-stationary*
  - sampling: transitions are not visited *i.i.d.*



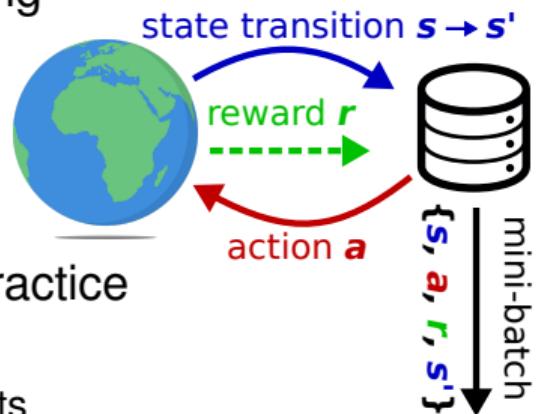
$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left( \underbrace{r + \gamma \max_{a'} (q_{\theta}(s'))_{a'}}_{\text{TD-error}} - (q_{\theta}(s))_a \right)^2$$

regression target

TD-error

# RL is not supervised learning!

- Q-learning violates assumptions of supervised learning
  - TD-error: regression targets are *non-stationary*
  - exploration: the state-distribution is *non-stationary*
  - sampling: transitions are not visited *i.i.d.*
- DQN must be stabilized with many tricks to work in practice
  - draw mini-batches i.i.d. from *experience replay buffer*  $\mathcal{D}$
  - use slowly changing *target network*  $q_{\theta'}$  for regression targets



$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathbb{E} \left[ \underbrace{\left( r + \gamma \max_{a'} \left( \overbrace{q_{\theta'}(s')}^{\text{target network}} \right)_{a'} - (q_{\theta}(s))_a \right)^2}_{\text{TD error}} \mid \underbrace{s, a, r, s' \sim \mathcal{D}}_{\text{from experience replay buffer}} \right]$$

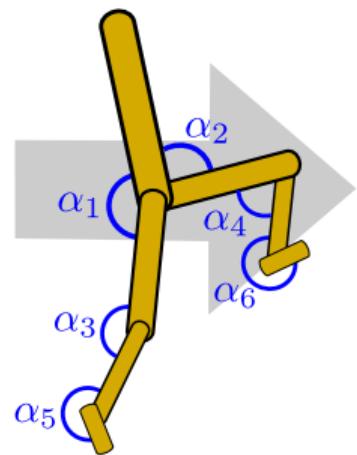


see Mnih et al. (2013, 2015) for details on stabilization techniques; see Badia et al. (2020) for a modern version; image modified from Designmodo

# Continuous actions

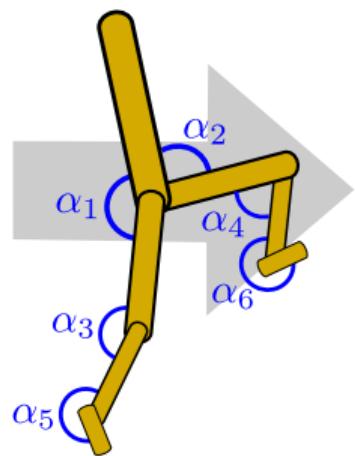
- Applications like robotics have continuous actions
  - states are limb angles, momentum, pressure sensors, etc.  
e.g.  $s = [\alpha_1, \dot{\alpha}_1, \dots, \alpha_6, \dot{\alpha}_6]$
  - actions are motor torques, e.g.  $a = [\ddot{\alpha}_1, \dots, \ddot{\alpha}_6] \in \mathbb{R}^{|a|}$

$$\max_{\mathbf{a}} q_{\theta}(\mathbf{s}, \mathbf{a}) = ???$$



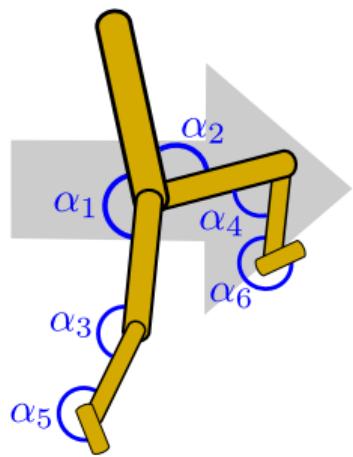
# Continuous actions

- Applications like robotics have continuous actions
  - states are limb angles, momentum, pressure sensors, etc.  
e.g.  $s = [\alpha_1, \dot{\alpha}_1, \dots, \alpha_6, \dot{\alpha}_6]$
  - actions are motor torques, e.g.  $a = [\ddot{\alpha}_1, \dots, \ddot{\alpha}_6] \in \mathbb{R}^{|a|}$
- Neural network outputs distribution over actions
  - e.g. a Gaussian distribution  $\pi_\theta(a|s) := \mathcal{N}(a|\mu_\theta(s), \sigma_\theta^2(s))$
  - $\mu_\theta(s) \in \mathbb{R}^{|a|}$  and  $\sigma_\theta(s) \in \mathbb{R}^{|a|}$  are neural networks



# Continuous actions

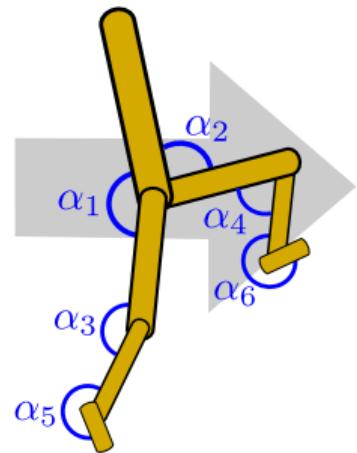
- Applications like robotics have continuous actions
  - states are limb angles, momentum, pressure sensors, etc.  
e.g.  $s = [\alpha_1, \dot{\alpha}_1, \dots, \alpha_6, \dot{\alpha}_6]$
  - actions are motor torques, e.g.  $a = [\ddot{\alpha}_1, \dots, \ddot{\alpha}_6] \in \mathbb{R}^{|a|}$
- Neural network outputs distribution over actions
  - e.g. a Gaussian distribution  $\pi_\theta(a|s) := \mathcal{N}(a|\mu_\theta(s), \sigma_\theta^2(s))$
  - $\mu_\theta(s) \in \mathbb{R}^{|a|}$  and  $\sigma_\theta(s) \in \mathbb{R}^{|a|}$  are neural networks
- How can we learn the parameters  $\theta$  of the policy?



# Policy gradients

- Directly maximize expected return  $R_t := \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ 
  - starting at the first state
  - gradient reinforce actions that led to good returns

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R_0] = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \right]$$



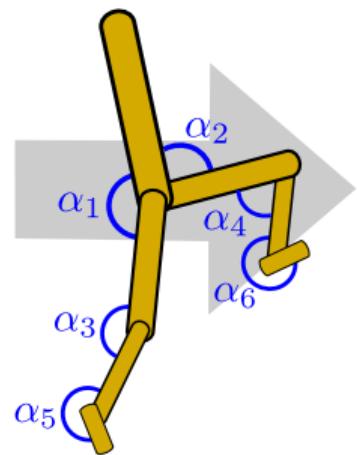
see the *policy gradient theorem* (originally by Sutton et al., 1999) in Section 13.2 of Sutton and Barto (2018)

# Policy gradients

- Directly maximize expected return  $R_t := \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ 
  - starting at the first state
  - gradient reinforce actions that led to good returns

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R_0] = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \right]$$

- REINFORCE algorithm ([Williams, 1992](#))
  - sample  $n$  trajectories per policy gradient
  - slow and not very stable

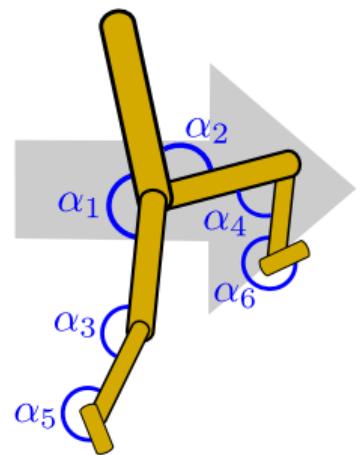


# Policy gradients

- Directly maximize expected return  $R_t := \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ 
  - starting at the first state
  - gradient reinforce actions that led to good returns

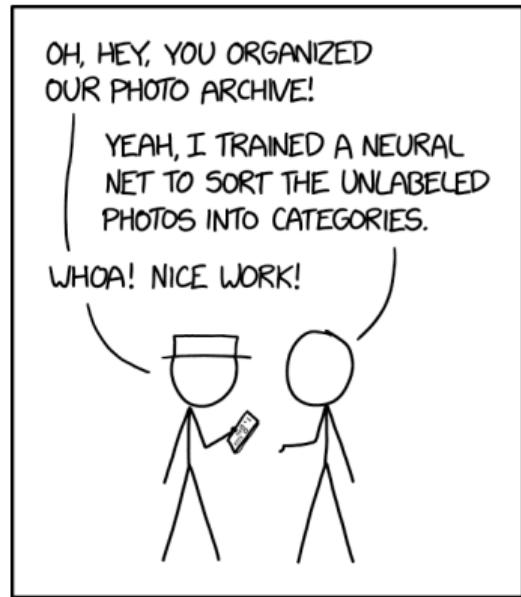
$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R_0] = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \right]$$

- REINFORCE algorithm ([Williams, 1992](#))
  - sample  $n$  trajectories per policy gradient
  - slow and not very stable
- Modern variations faster and more stable
  - PPO ([Schulman et al., 2017](#))
  - SAC ([Haarnoja et al., 2018a,b](#))



# Deep reinforcement learning

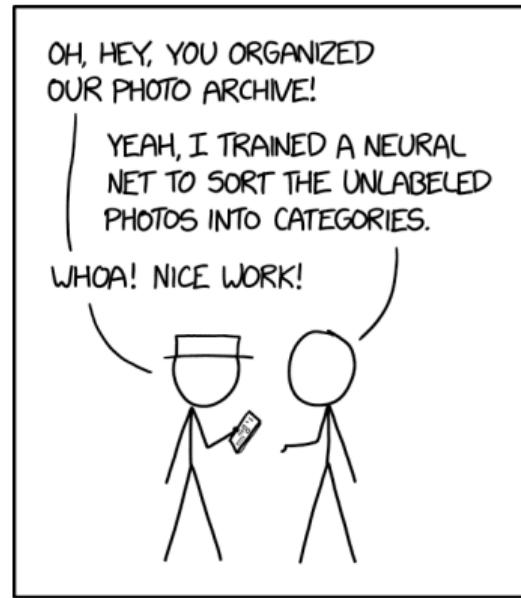
- In this video we have:
  - discussed why tabular RL does not scale
  - seen that Q-learning can be scaled up to play video games
  - learned how neural networks represent continuous policies
  - seen how continuous policies can be learned



ENGINEERING TIP:  
WHEN YOU DO A TASK BY HAND,  
YOU CAN TECHNICALLY SAY YOU  
TRAINED A NEURAL NET TO DO IT.

# Deep reinforcement learning

- In this video we have:
  - discussed why tabular RL does not scale
  - seen that Q-learning can be scaled up to play video games
  - learned how neural networks represent continuous policies
  - seen how continuous policies can be learned
- But be careful, RL is not supervised learning!
- Deep RL is an active research topic



ENGINEERING TIP:  
WHEN YOU DO A TASK BY HAND,  
YOU CAN TECHNICALLY SAY YOU  
TRAINED A NEURAL NET TO DO IT.

4

# Reinforcement Learning: State-of-the-Art and Limits

# Recent RL success stories

- Playing video games (2013/2015)
  - *challenge*: high dimensional pixel input
  - one algorithm can learn each Atari game

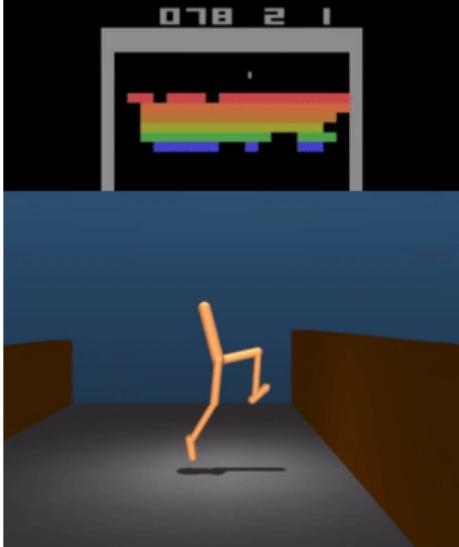
# Recent RL success stories



- Playing video games (2013/2015)
- Fast moving simulated robots (2016/2017)
  - *challenge*: continuous actions with complex dynamics
  - running robots pass random obstacle courses
  - only in simulators so far

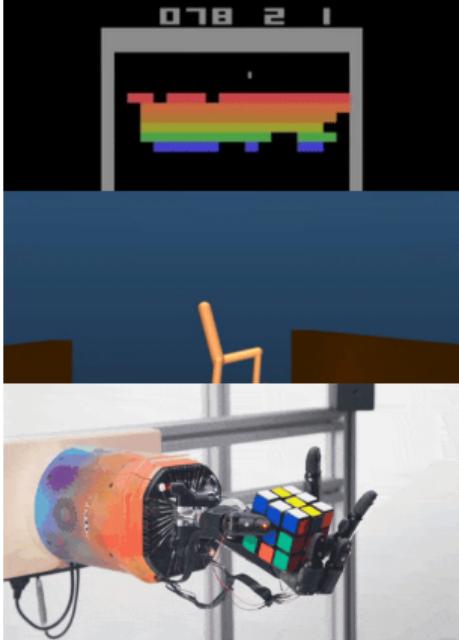
# Recent RL success stories

- Playing video games (2013/2015)
- Fast moving simulated robots (2016/2017)
- Real robotic hand solves Rubik's cubes (2019)
  - *challenge*: transfer from simulator to reality
  - works only in the lab so far



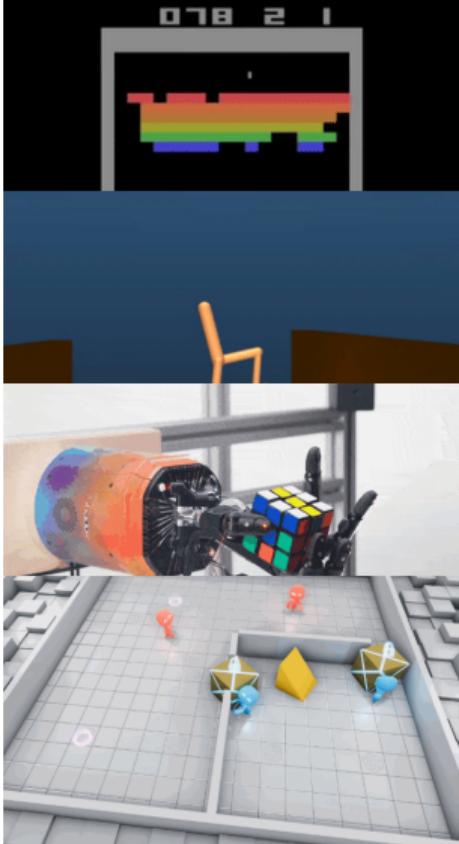
# Recent RL success stories

- Playing video games (2013/2015)
- Fast moving simulated robots (2016/2017)
- Real robotic hand solves Rubik's cubes (2019)
- Multi-agent RL in competitive games (2019/2020)
  - *challenge*: many agents interact with each other
  - surprisingly complex and seemingly intelligent behavior
  - RL beats professional players in StarCraft and MOBA's



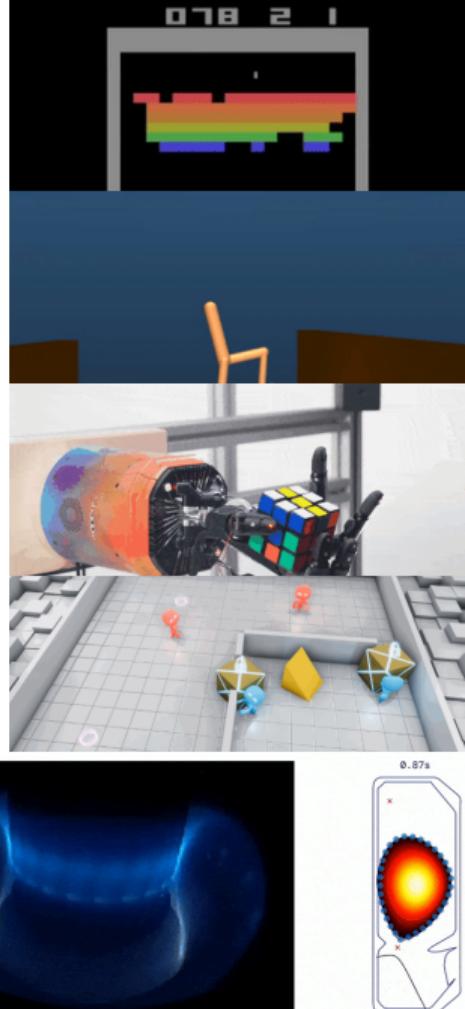
# Recent RL success stories

- Playing video games (2013/2015)
- Fast moving simulated robots (2016/2017)
- Real robotic hand solves Rubik's cubes (2019)
- Multi-agent RL in competitive games (2019/2020)
- Magnetic control of a Tokamak fusion reactor (2022)
  - *challenge*: keep the plasma bubble from bursting
  - simulations transferred to real fusion reactor



# Recent RL success stories

- Playing video games (2013/2015)
- Fast moving simulated robots (2016/2017)
- Real robotic hand solves Rubik's cubes (2019)
- Multi-agent RL in competitive games (2019/2020)
- Magnetic control of a Tokamak fusion reactor (2022)
- Most experts expect breakthrough in RL soon
  - industry applications within a decade or so
  - might impact **your** future job!



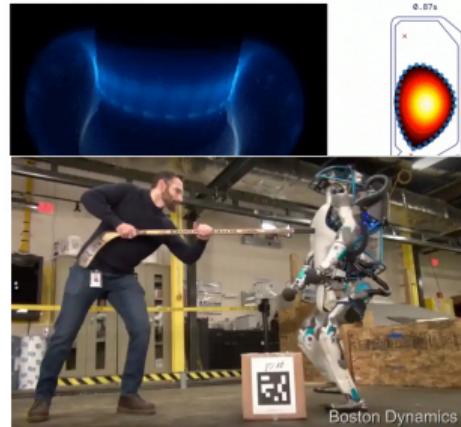
# What is possible?

- Now: learn to control unknown dynamics
  - environment must be stationary and completely explorable
  - many interactions with the environment (e.g. simulator/game)
  - no catastrophic/irreversible actions (e.g. crashing the car)



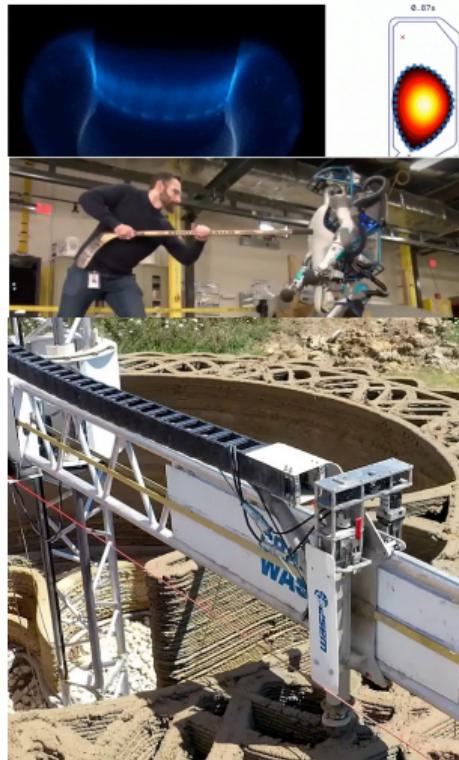
# What is possible?

- Now: learn to control unknown dynamics
  - environment must be stationary and completely explorable
  - many interactions with the environment (e.g. simulator/game)
  - no catastrophic/irreversible actions (e.g. crashing the car)
- Soon: safe transfer of learned control
  - learn in a simulation, act safely in the real world
  - learn to anticipate behavior of other learning agents
  - integrate existing knowledge/constraints into the training



# What is possible?

- Now: learn to control unknown dynamics
  - environment must be stationary and completely explorable
  - many interactions with the environment (e.g. simulator/game)
  - no catastrophic/irreversible actions (e.g. crashing the car)
- Soon: safe transfer of learned control
  - learn in a simulation, act safely in the real world
  - learn to anticipate behavior of other learning agents
  - integrate existing knowledge/constraints into the training
- Possible applications
  - infinite sampling: computer games, web-crawling, video processing
  - limited autonomy: construction, factories, power plants/grids, driver assistance



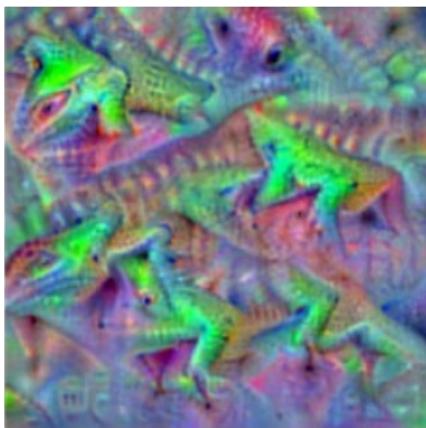
# What is impossible?

- RL agents do **not represent** the world as we do
  - only patterns: no objects, relationships, movements
  - only training data: no intuition of what is *similar*

What you see:

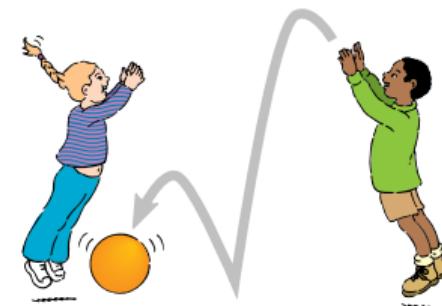
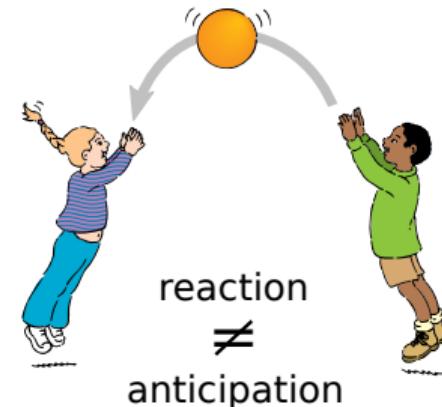


What a neural net sees:



# What is impossible?

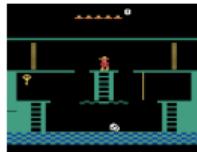
- RL agents do **not represent** the world as we do
  - only patterns: no objects, relationships, movements
  - only training data: no intuition of what is *similar*
- RL agents do **not anticipate** new situations
  - reinforce actions: all situations have to be trained for
  - no imagination: all reaction, no prediction



# What is impossible?

- RL agents do **not represent** the world as we do
  - only patterns: no objects, relationships, movements
  - only training data: no intuition of what is *similar*
- RL agents do **not anticipate** new situations
  - reinforce actions: all situations have to be trained for
  - no imagination: all reaction, no prediction
- RL agents **only learn** the task they are given
  - only looks smart: memorizes what yielded reward
  - misspecification: RL can act dumb to get reward

1. start



2. go get  
the key



3. die



4. respawn  
near door



5. open door  
get points



# What is impossible?

- RL agents do **not represent** the world as we do
  - only patterns: no objects, relationships, movements
  - only training data: no intuition of what is *similar*
- RL agents do **not anticipate** new situations
  - reinforce actions: all situations have to be trained for
  - no imagination: all reaction, no prediction
- RL agents **only learn** the task they are given
  - only looks smart: memorizes what yielded reward
  - misspecification: RL can act dumb to get reward
- RL agents are **not intelligent** like humans! Don't expect them to be!



image source: *The Simpsons* (20th Television Animation)

# When to use reinforcement learning?

- ① Is there no heuristic that already solves the task almost optimal?
- ② Is your task virtual or do you have a simulator of your tasks?
- ③ Are there no “forbidden” actions that would destroy the system?
- ④ Can the system cause **no** serious harm to people?
- ⑤ Is the variety of situations the agent can face limited?
- ⑥ Do you have a lot of compute to train the agent for all situations?
- ⑦ Are you willing to spend 6 month fiddling with hyper-parameters?

# State-of-the-art and limits

- In this video we have:
  - seen some recent success stories of RL
  - discussed what RL *can* learn
  - discussed what RL *cannot* learn
  - discussed when to use RL

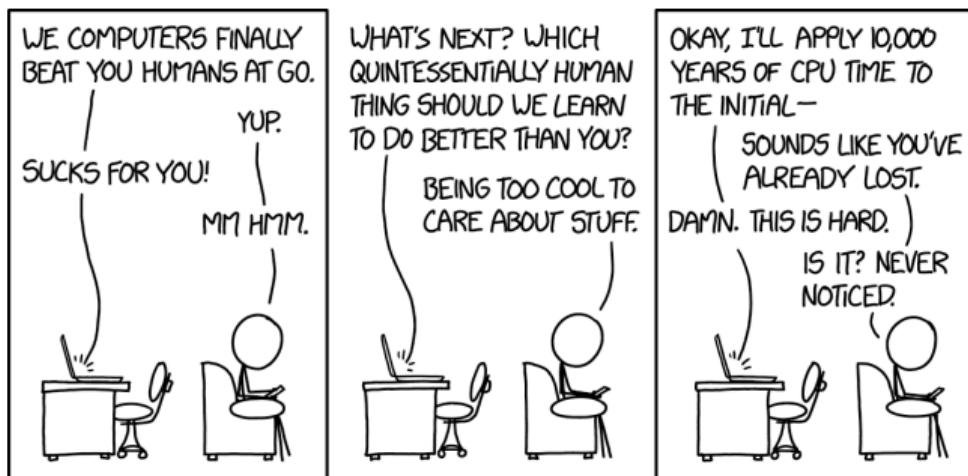


image source: [xkcd.com](http://xkcd.com)

# References |

- Adriá Puigdoménech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi, Zhaohan Daniel Guo, and Charles Blundell. Agent57: Outperforming the Atari human benchmark. In *CoRR*, abs/2003.13350, 2020. URL <https://arxiv.org/abs/2003.13350>.
- Bowen Baker, Ingmar Kanitscheider, Todor M. Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. Emergent tool use from multi-agent autocurricula. In *8th International Conference on Learning Representations (ICLR)*, 2020. URL <https://openreview.net/forum?id=SkxpxJBKwS>.
- Richard Ernest Bellman. *Dynamic programming*. Princeton University Press, 1957.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. *CoRR*, abs/1912.06680, 2019. URL <http://arxiv.org/abs/1912.06680>.
- Jonas Degrave, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego Casas, Craig Donner, Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie Kay, Antoine Merle, Jean-Marc Moret, and Martin Riedmiller. Magnetic control of tokamak plasmas through deep reinforcement learning. *Nature*, 602:414–419, 02 2022. doi: 10.1038/s41586-021-04301-9.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of Machine Learning Research (ICML)*, volume 80, pages 1861–1870, 2018a. URL <http://proceedings.mlr.press/v80/haarnoja18b.html>.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications. *CoRR*, abs/1812.05905, 2018b. URL <http://arxiv.org/abs/1812.05905>.
- Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, S. M. Ali Eslami, Martin A. Riedmiller, and David Silver. Emergence of locomotion behaviours in rich environments. *CoRR*, abs/1707.02286, 2017. URL <http://arxiv.org/abs/1707.02286>.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>. NIPS Deep Learning Workshop 2013.

## References II

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015.
- OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang. Solving rubik's cube with a robot hand. *CoRR*, abs/1910.07113, 2019. URL <http://arxiv.org/abs/1910.07113>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- Richard S. Sutton, David McAllester, and Satinder Singh and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 12, pages 1057–1063. MIT Press, 1999. URL <https://papers.nips.cc/paper/1999/hash/464d828b85b0bed98e80ade0a5c43b0f-Abstract.html>.
- Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575:350–354, 2019.
- Christopher Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992.
- Ronald J. Williams. Simple statistical gradient algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992. URL <https://link.springer.com/article/10.1007/BF00992696>.