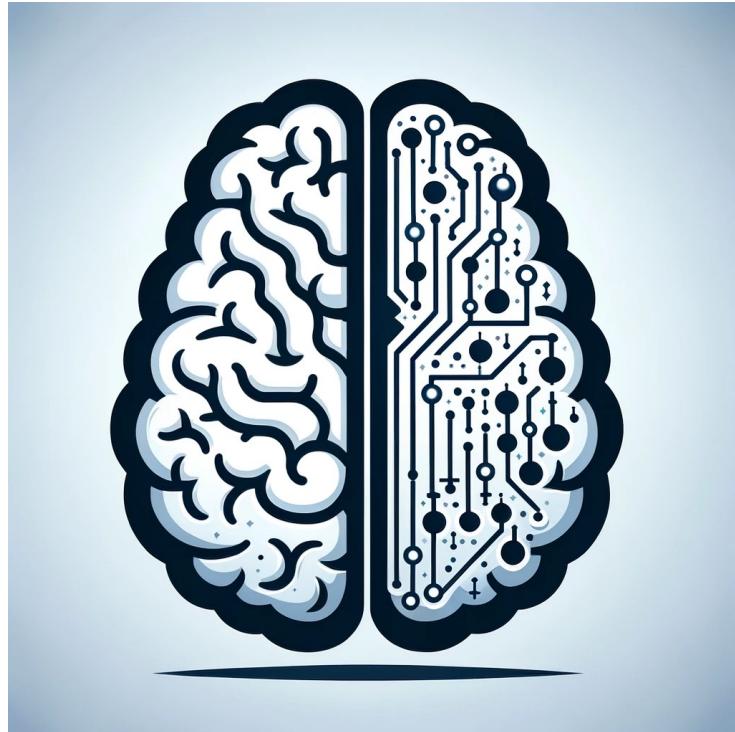


# **EN 601.473/601.673: Cognitive Artificial Intelligence (CogAI)**



**Lecture 20:  
Theory of Mind**

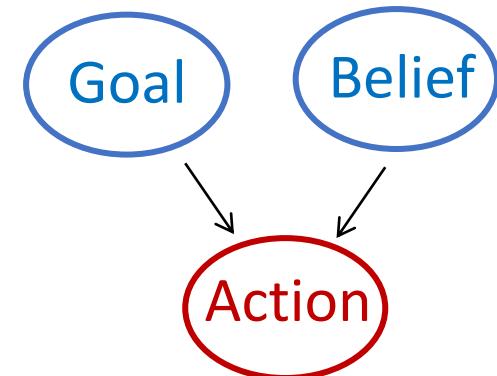
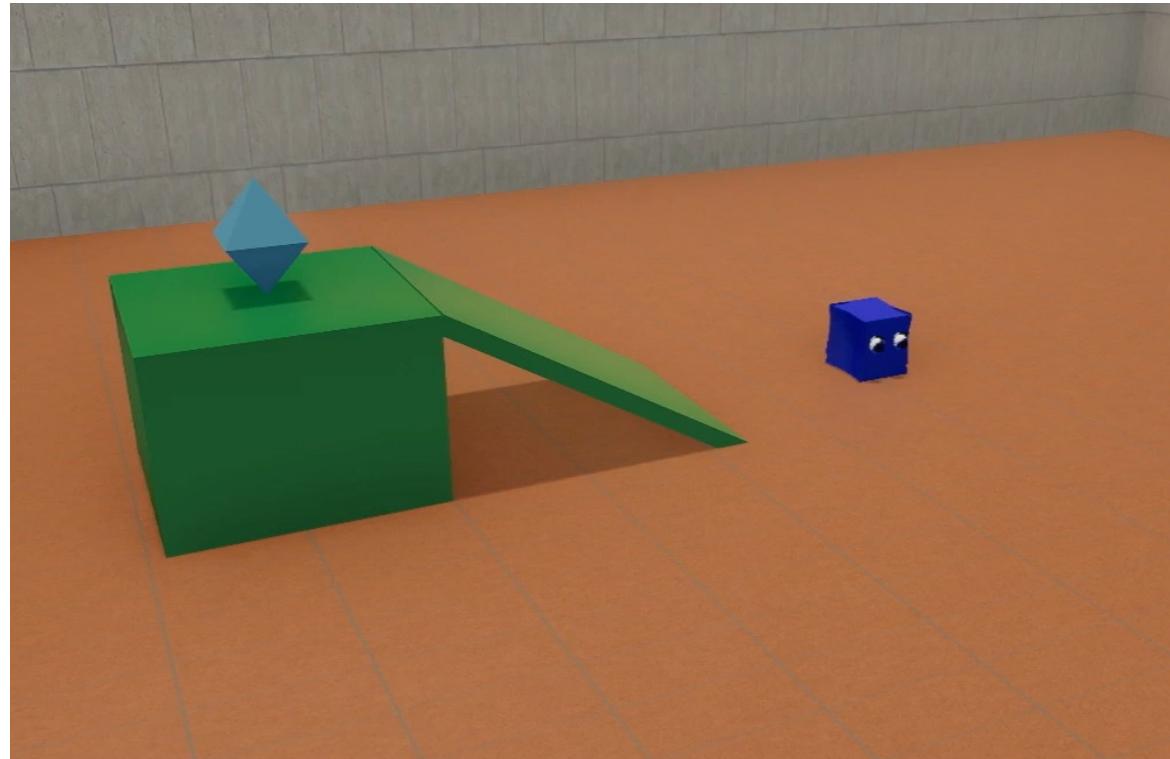
**Tianmin Shu**

# Theory of Mind

- Basic concepts and evaluation
- Single-agent decision making
- Inverse single-agent decision making

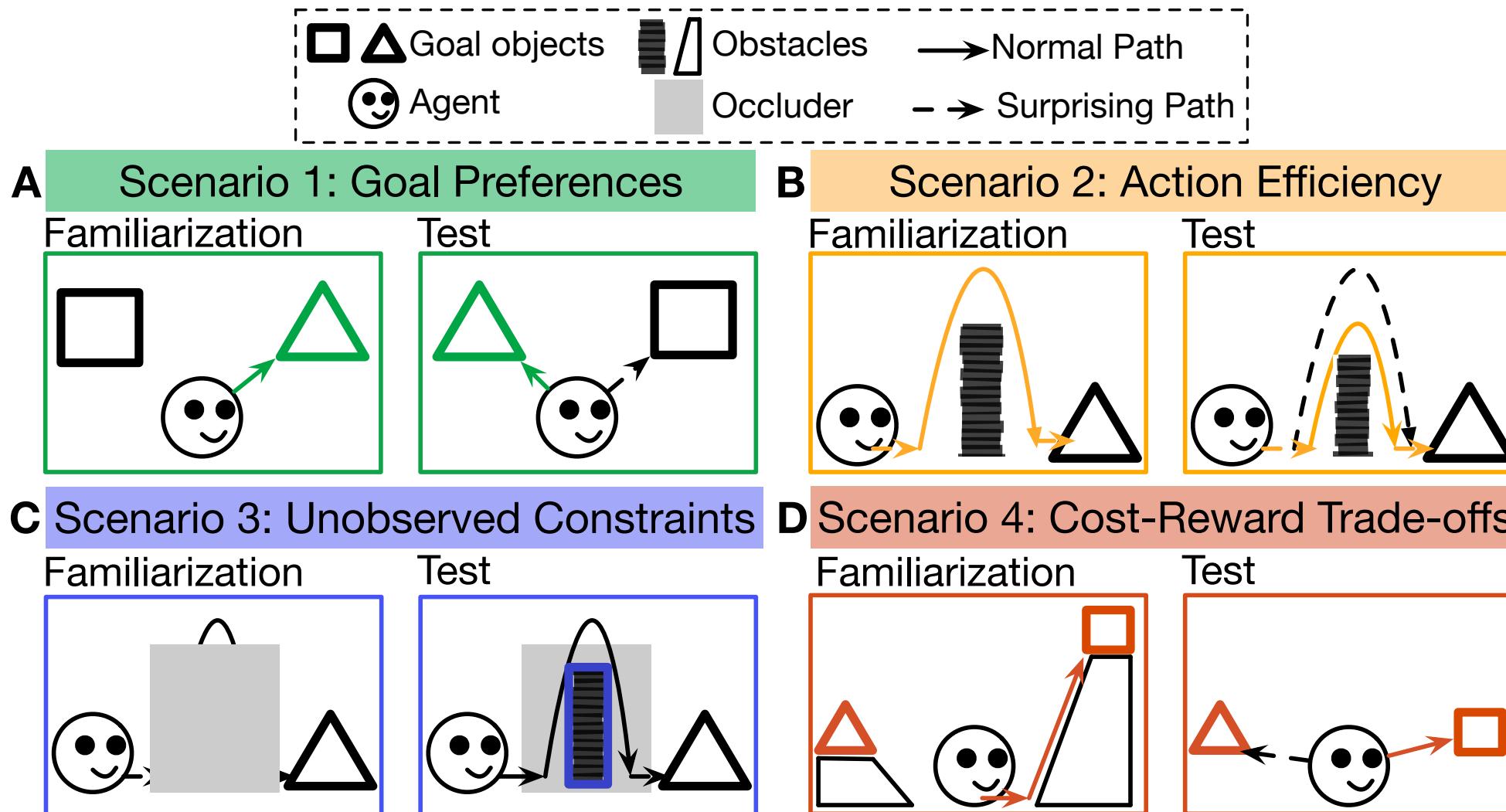
# Theory of Mind in humans

Theory of Mind: Reasoning about **hidden mental** variables  
that drive **observable actions**

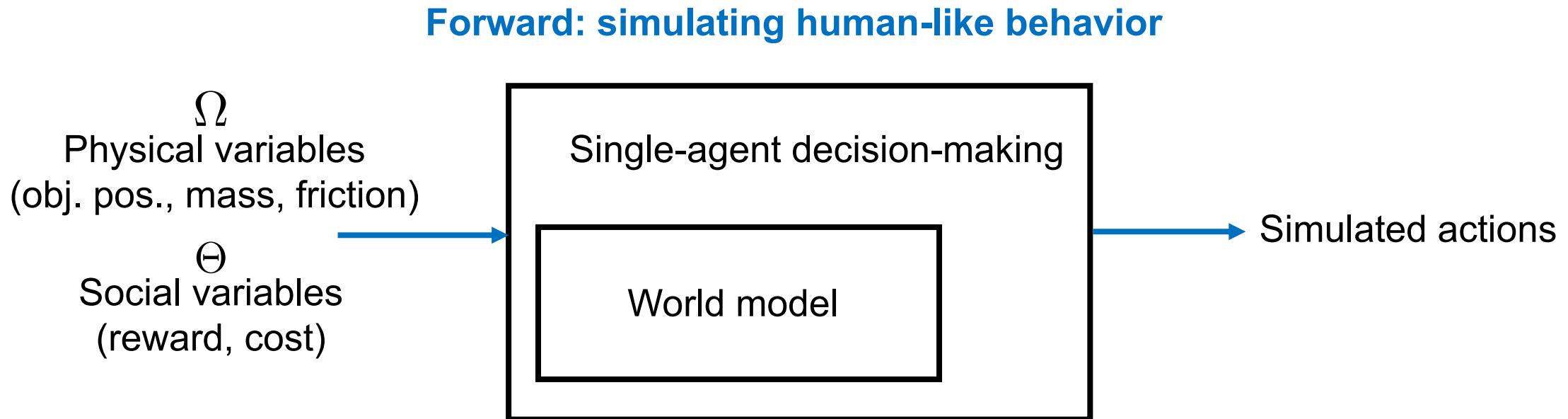


The agent takes the most **efficient action** to reach its **goal** (blue diamond), maximizing its **reward** and minimizing the **cost**, subject to **constraints**

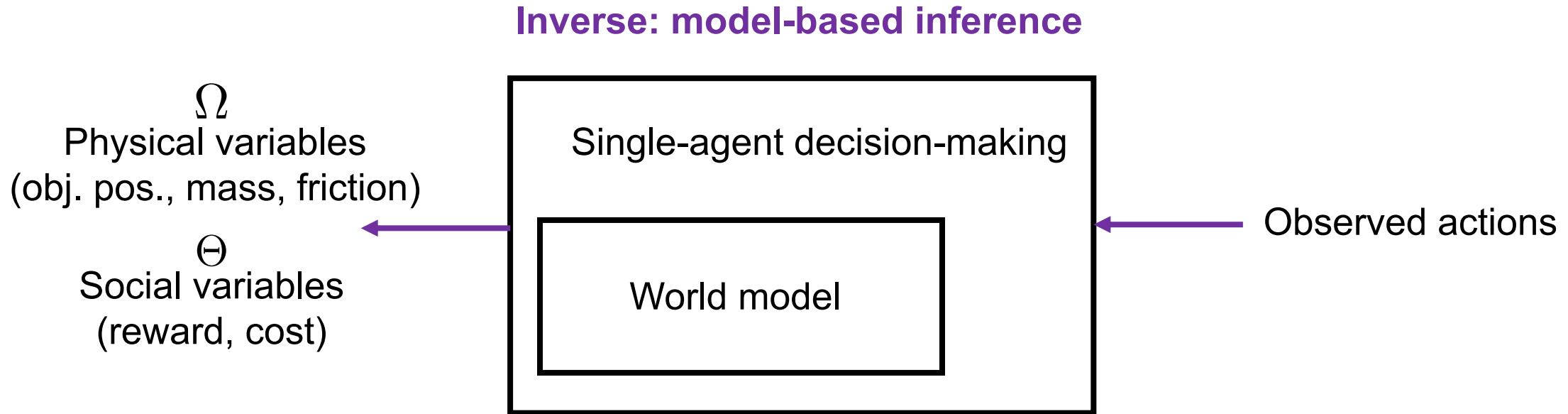
# AGENT: Action, Goal, Efficiency, coNstraint, uTility



# A generative model of single-agent behavior



# A generative model of single-agent behavior



# Theory of Mind

- Basic concepts and evaluation
- Single-agent decision making
- Inverse single-agent decision making

# **Single-agent decision making**

- Markov decision process (MDP)
- Value, policy
- Model-based approaches for solving MDPs
- Model-free RL (Q-learning)

# Model-based approaches

- Assume that we know the state transition probabilities (i.e., the world model)
- Policy improvement
- Policy iteration
- Value iteration

# Policy improvement theorem

- $V(s)$  tells us how good it is to following current policy from  $s$ . Can we do better with a new policy? I.e., take a different action at  $s$ .

$$\begin{aligned} q^\pi(s, a) &= E[r_{t+1} + \gamma v^\pi(s_{t+1}) | s_t = s, a_t = a] \\ &= \sum_{s', r} p(s', r | s, a)[r + \gamma v^\pi(s')] \end{aligned}$$

- Policy improvement theorem:
- Given any pair of deterministic policies such that, for all states  $s \in \mathcal{S}$

$$q^\pi(s, \pi'(s)) \geq v^\pi(s)$$

- Then the policy  $\pi'$  must be as good as, or better than,  $\pi$ . That is, it must obtain greater or equal expected return from all states  $s \in \mathcal{S}$ :

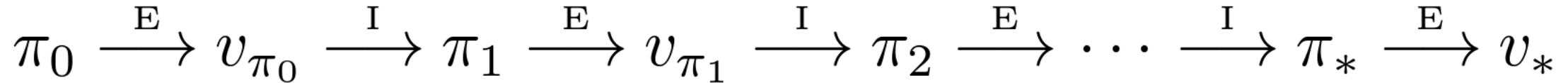
$$v^{\pi'}(s) \geq v^\pi(s)$$

# Policy improvement

- Constructing a greedy policy, which meets the condition of the policy improvement theorem → as good as or better than the original policy

$$\begin{aligned}\pi'(s) &= \arg \max_a q^\pi(s, a) \\ &= \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v^\pi(s')]\end{aligned}$$

# Policy iteration



Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$ ;  $V(\text{terminal}) \doteq 0$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each  $s \in \mathcal{S}$ :

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until  $\Delta < \theta$  (a small positive number determining the accuracy of estimation)

**Drawback:**

Each iteration involves policy evaluation.

Do we must wait for exact convergence?

3. Policy Improvement

*policy-stable*  $\leftarrow$  true

For each  $s \in \mathcal{S}$ :

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

If  $\text{old-action} \neq \pi(s)$ , then *policy-stable*  $\leftarrow$  false

If *policy-stable*, then stop and return  $V \approx v_*$  and  $\pi \approx \pi_*$ ; else go to 2

# Value iteration

- Key idea: Policy evaluation is stopped after just one sweep (one update of each state) instead of wait until convergence

## Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold  $\theta > 0$  determining accuracy of estimation

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

$$|\quad \Delta \leftarrow 0$$

  | Loop for each  $s \in \mathcal{S}$ :

$$|\quad \quad v \leftarrow V(s)$$

$$|\quad \quad V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$|\quad \quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$$

  | until  $\Delta < \theta$

Combines the policy improvement and truncated policy evaluation steps

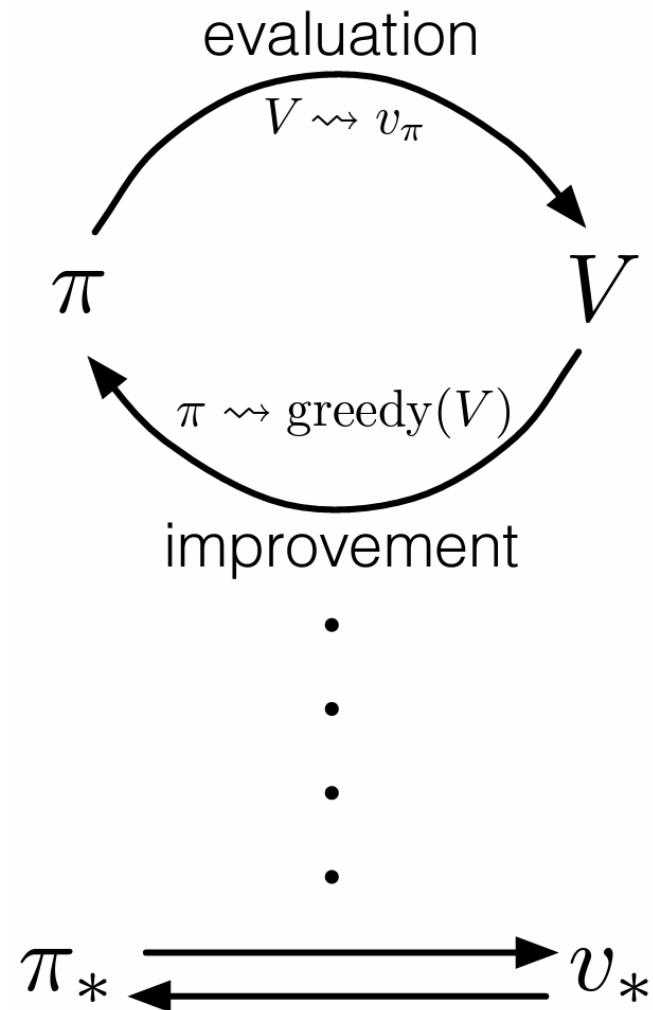
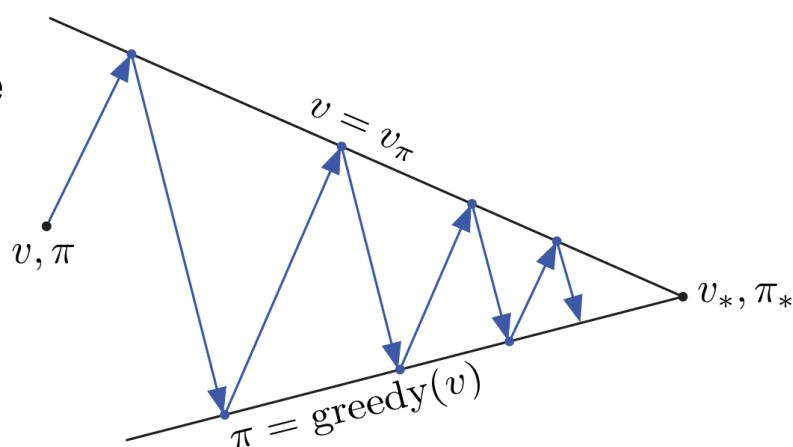
Output a deterministic policy,  $\pi \approx \pi_*$ , such that

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

# Generalization policy iteration (GPI)

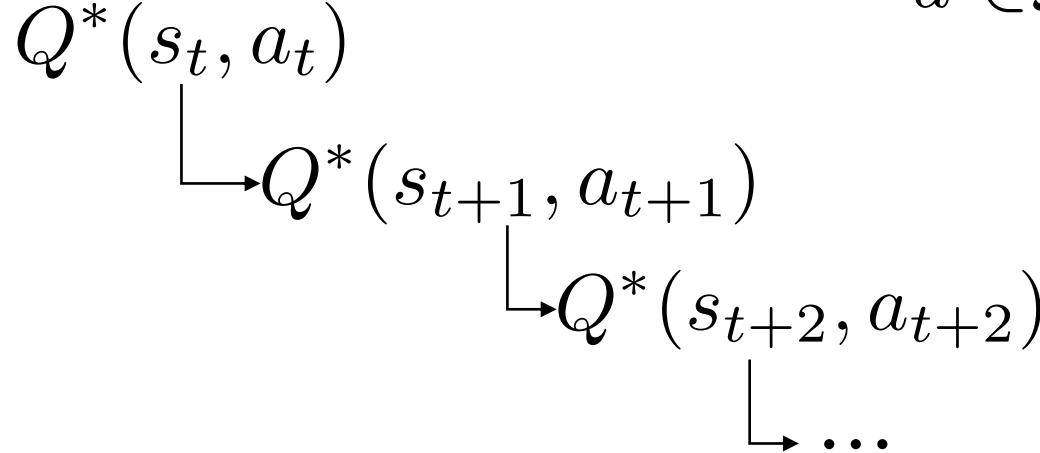
- The general idea of letting policy-evaluation and policy-improvement process interact
- Almost all reinforcement learning methods are well described as GPI
- All have identifiable policies and value functions, with the policy always being improved with respect to value function and the value function always being driven toward the value function for the policy

Simplified 2D space



# How to solve the Bellman equation for $Q^*(s, a)$

$$Q^*(s, a) = R(s, a) + \gamma \max_{a' \in \mathcal{A}} Q^*(s', a')$$

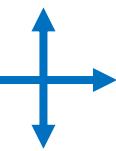


Many iterative methods:

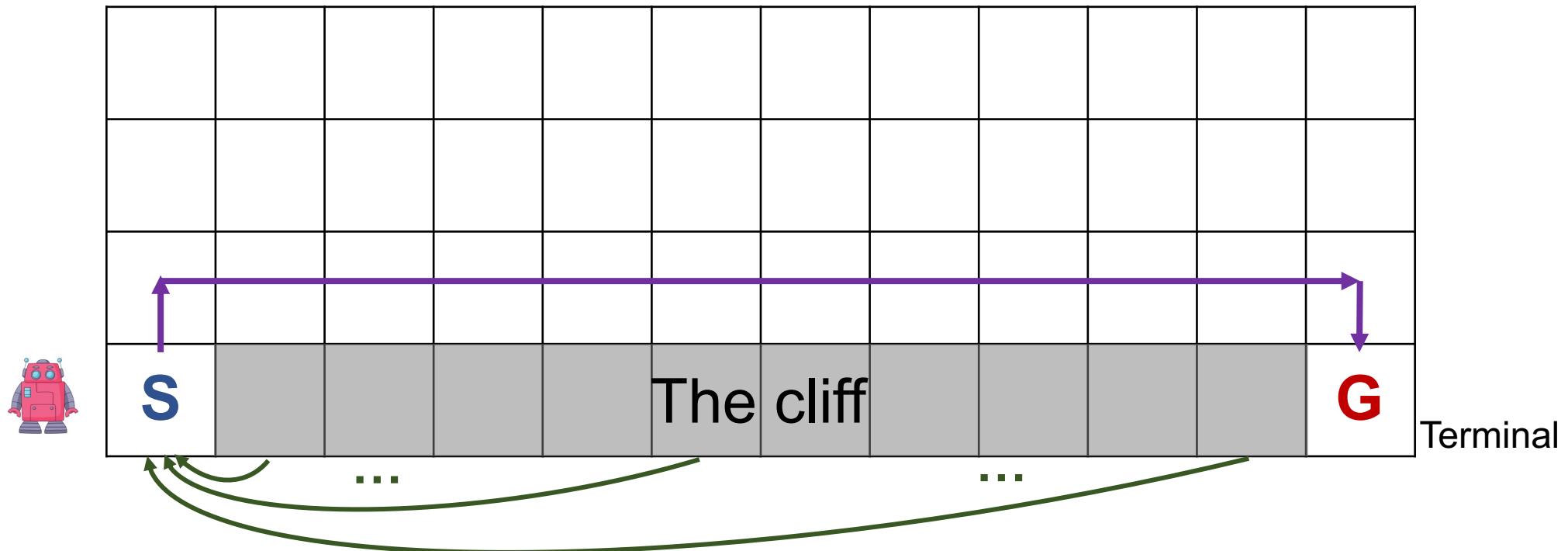
- Policy iteration
- Value iteration
- Sarsa
- • Q-learning

Model-free RL: without access of the world model

# An example: cliff walk in a 2D grid world

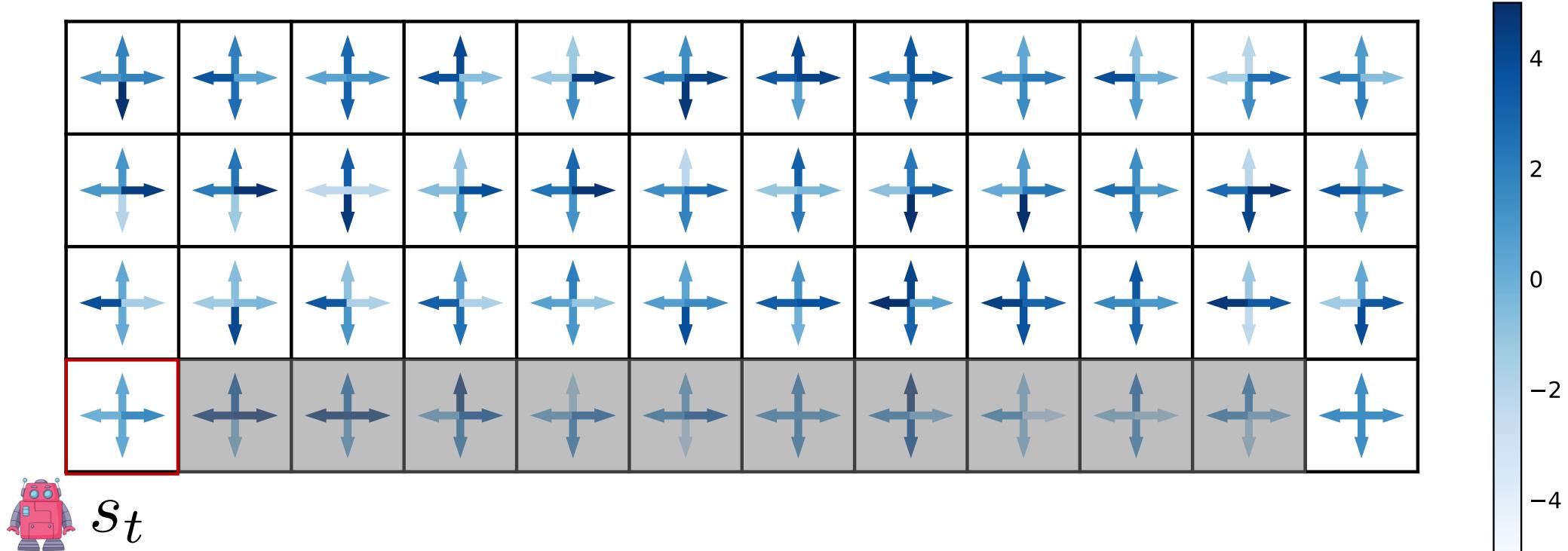
State  $s$ :  $(x, y)$     Action  $a$ :     Reward  $R(s, a) = \begin{cases} -10 & \text{if step into the cliff} \\ -1 & \text{otherwise} \end{cases}$

$$\gamma = 0.95$$

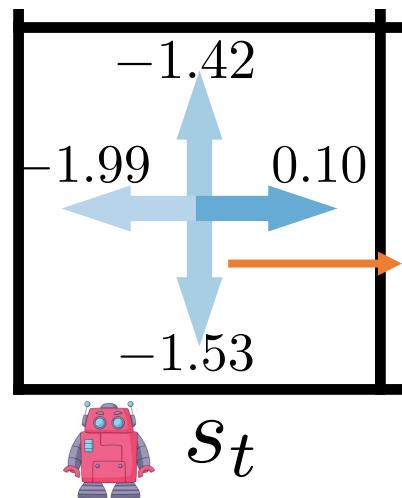


# Q-learning

Start from a random estimation  $Q(s, a), \forall s \in \mathcal{S}, \forall a \in \mathcal{A}$



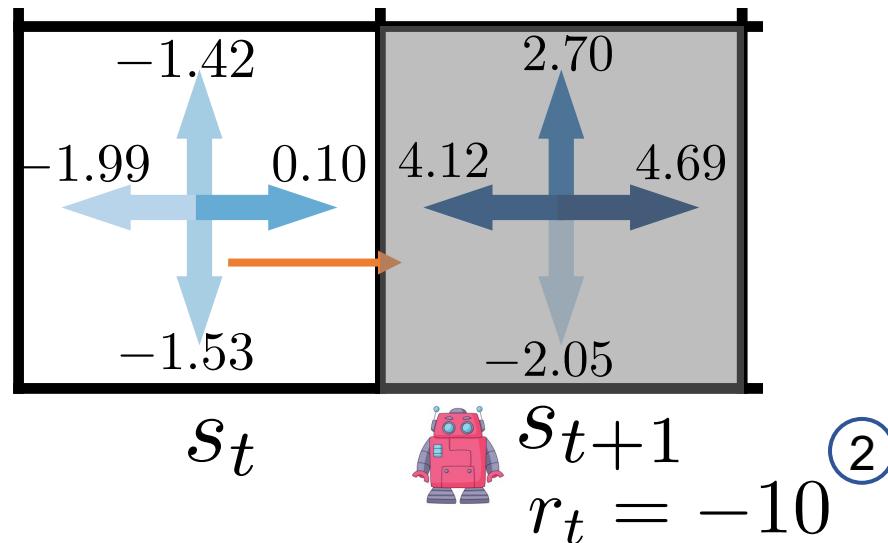
# Q-learning



$$a_t = \operatorname{argmax}_{a_t} Q(s_t, a_t)$$

$$Q(s_t, a_t) = 0.10$$

# Q-learning



$$\textcircled{1} \quad a_t = \underset{a_t}{\operatorname{argmax}} Q(s_t, a_t)$$

$$Q(s_t, a_t) = \boxed{0.10}$$

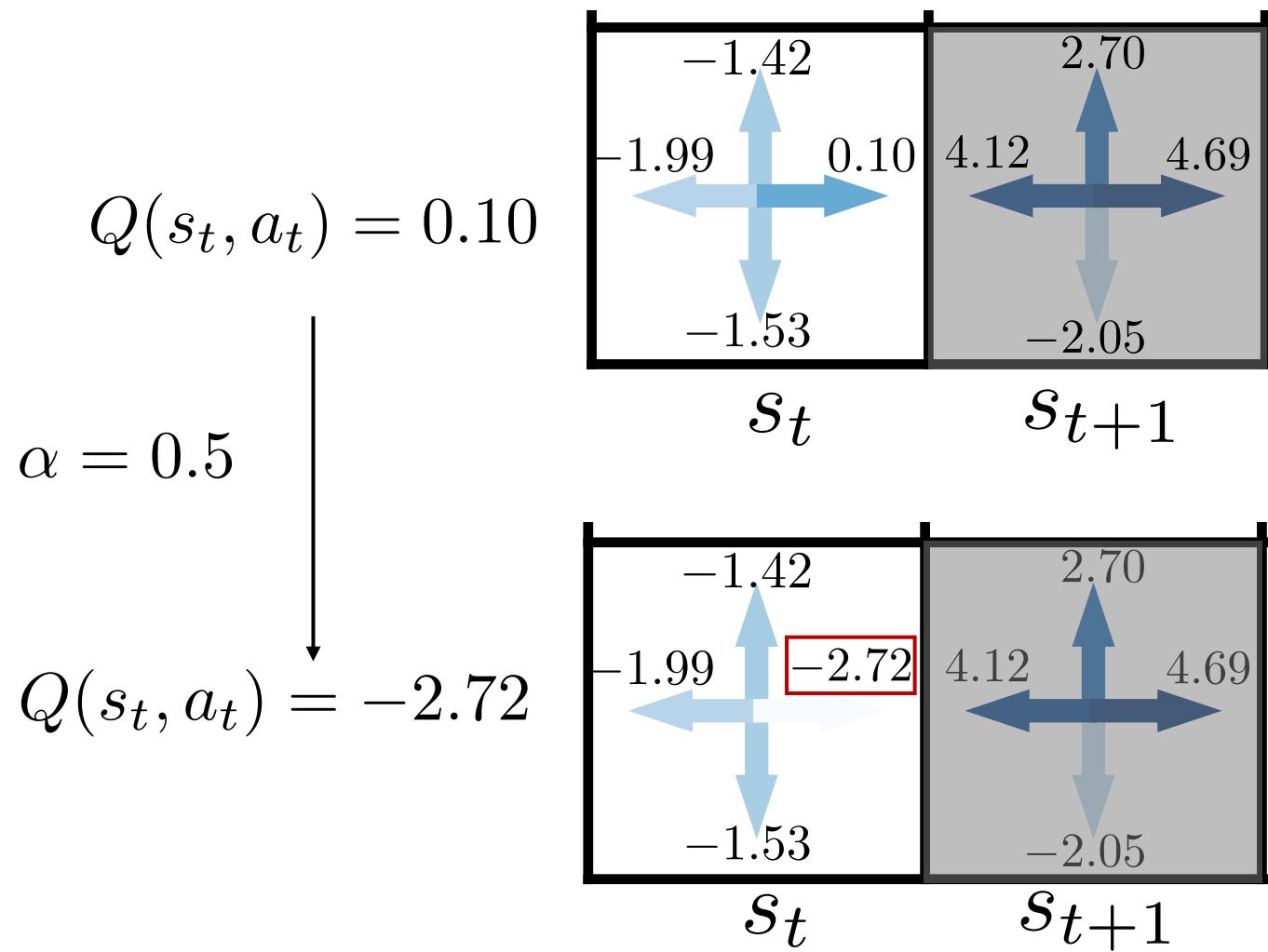
$$Q^{\text{new}}(s_t, a_t) = r_t + \gamma \max_{a_{t+1} \in \mathcal{A}} Q(s_{t+1}, a_{t+1})$$

$$= -10 + 0.95 * 4.69 \\ = \boxed{-5.54}$$

$$\textcircled{3} \quad Q(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha Q^{\text{new}}(s_t, a_t) \quad \text{A better estimation than } Q(s_t, a_t)!$$

$\alpha \in (0, 1]$  Why weighted average?  $Q^{\text{new}}$  is not yet the optimal value

## After one-step update



# Pseudocode for Q-learning

Set the algorithm parameter:  $\alpha \in (0, 1]$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ , randomly except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Agent starts from the initial state  $S_0$

Loop for each step  $t$  of episode:

Choose  $a_t$  from  $\mathcal{S}_t$  using policy derived from  $Q(s_t, a_t)$

Take action  $a_t$ , receive  $r_t$  and  $S_{t+1}$

$$Q^{\text{new}}(s_t, a_t) = r_t + \gamma \max_{a_{t+1} \in \mathcal{A}} Q(s_{t+1}, a_{t+1})$$

$$Q(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha Q^{\text{new}}(s_t, a_t)$$

Until reach the terminal state

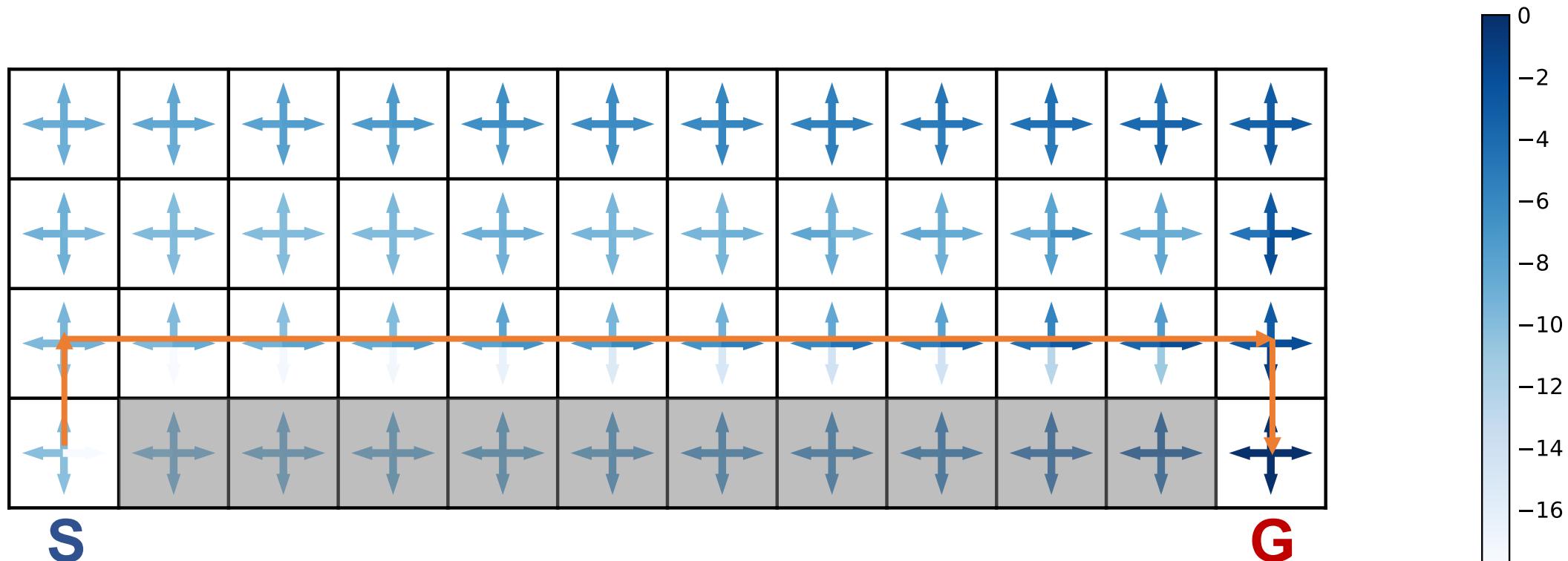
Temporal difference (TD) error

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha(Q^{\text{new}}(s_t, a_t) - Q(s_t, a_t))$$

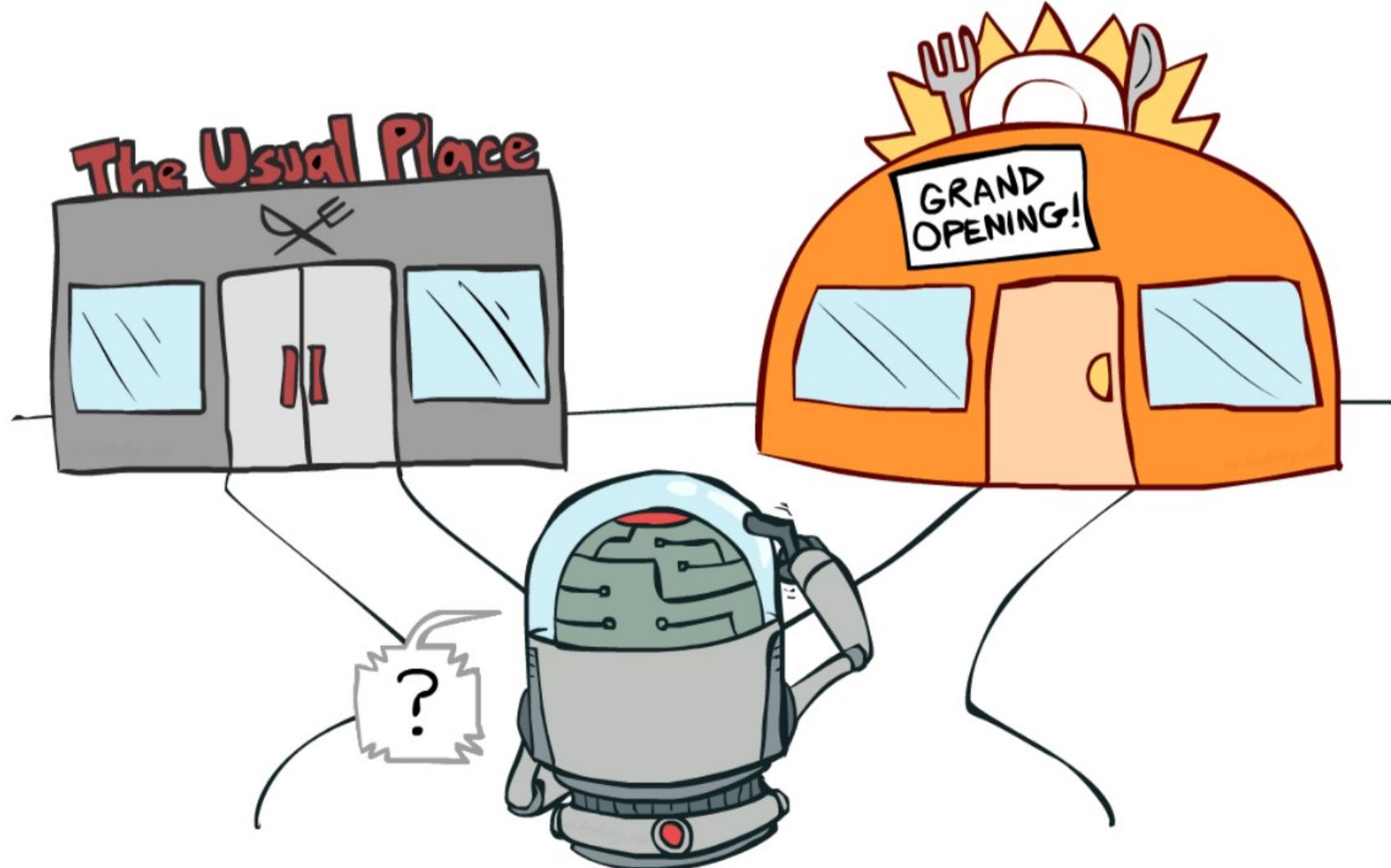
# Training results

$$\alpha = 0.5$$

Learned value function after 200 episodes



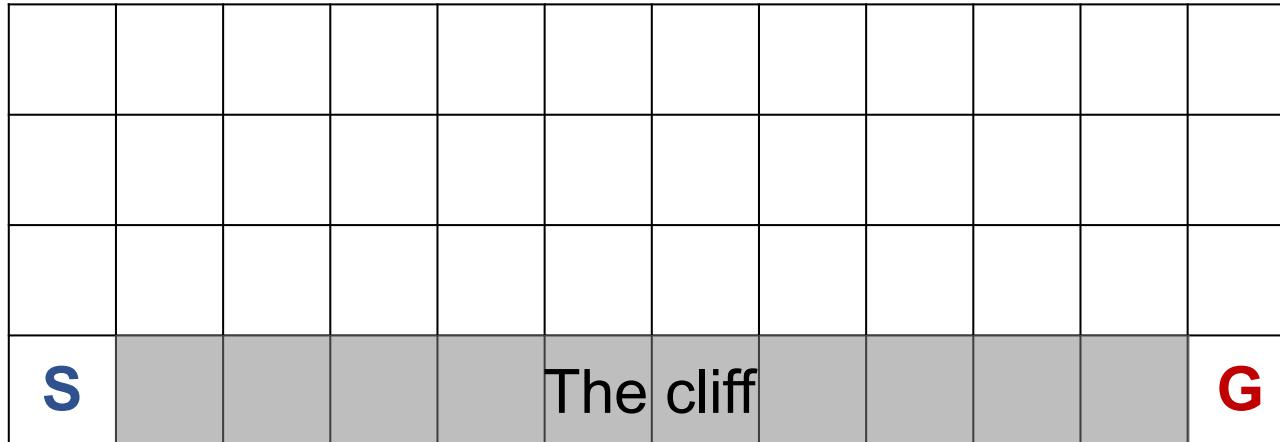
# Exploration vs Exploitation Dilemma



## $\epsilon$ -greedy

- Given  $\epsilon \in (0,1)$
- At each step, with a probability of  $\epsilon$ , sample a random action instead
- There are more advanced approaches
- E.g., soft Q-learning (Haarnoja et al., 2017)

# Beyond 2D grid words?

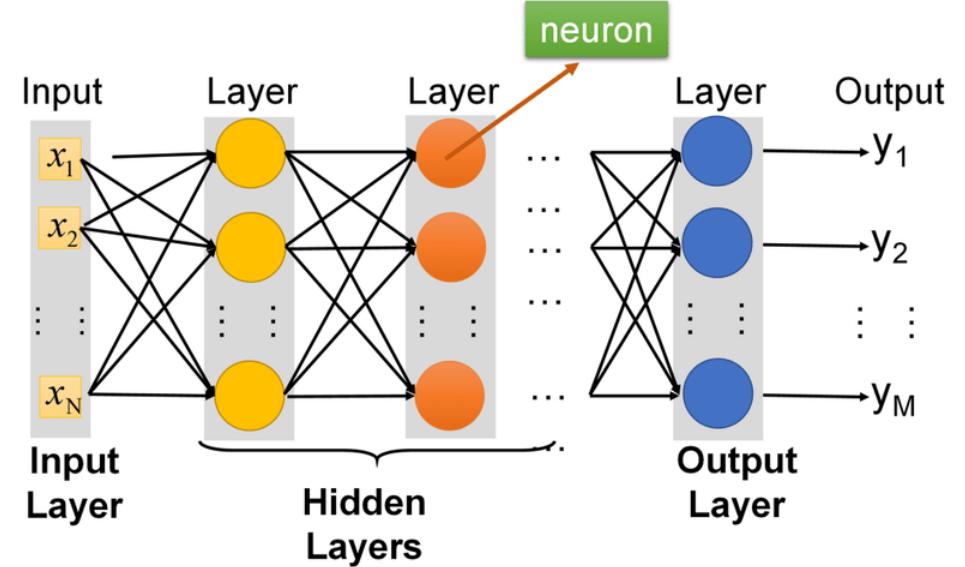


$Q($  (



,  $a)$

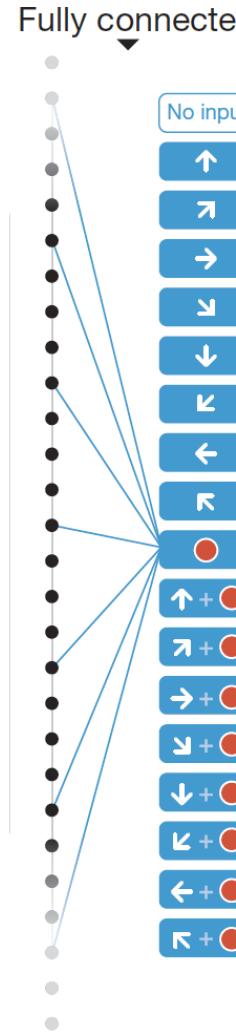
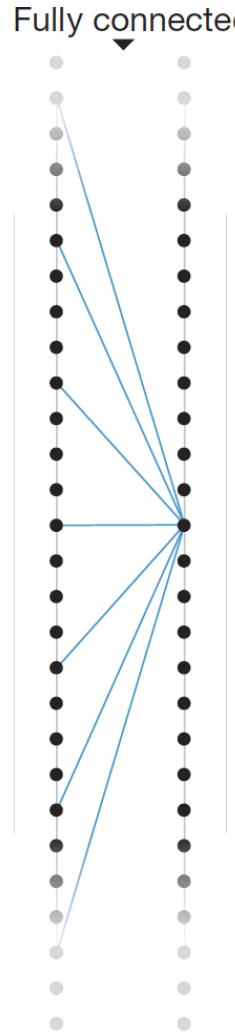
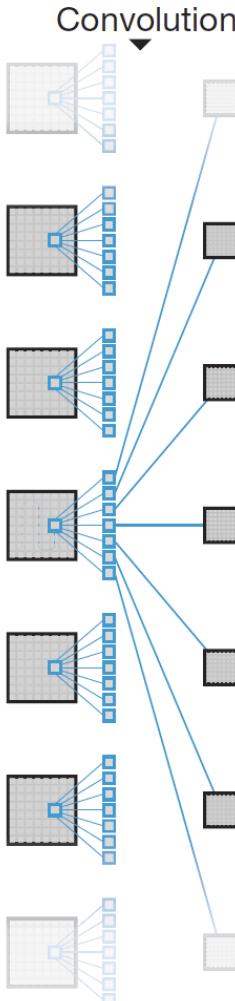
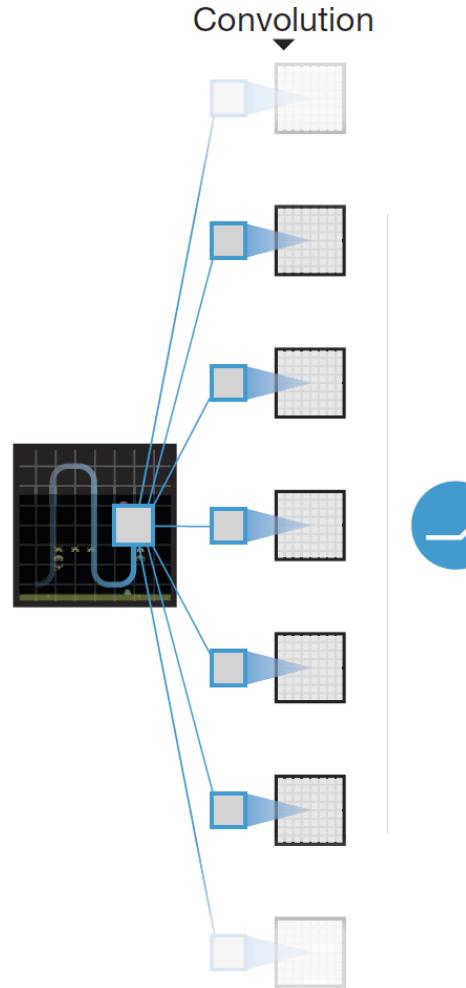
**Deep neural networks**  
Can approximate any function



# Deep neural networks to approximate the value function

## Deep Q-network (DQN)

$S$



$Q(s, a)$

No input



# Further readings

For learning about how to combine Q-learning with deep neural networks

**nature**

---

---

## Human-level control through deep reinforcement learning

Volodymyr Mnih<sup>1\*</sup>, Koray Kavukcuoglu<sup>1\*</sup>, David Silver<sup>1\*</sup>, Andrei A. Rusu<sup>1</sup>, Joel Veness<sup>1</sup>, Marc G. Bellemare<sup>1</sup>, Alex Graves<sup>1</sup>, Martin Riedmiller<sup>1</sup>, Andreas K. Fidjeland<sup>1</sup>, Georg Ostrovski<sup>1</sup>, Stig Petersen<sup>1</sup>, Charles Beattie<sup>1</sup>, Amir Sadik<sup>1</sup>, Ioannis Antonoglou<sup>1</sup>, Helen King<sup>1</sup>, Dharshan Kumaran<sup>1</sup>, Daan Wierstra<sup>1</sup>, Shane Legg<sup>1</sup> & Demis Hassabis<sup>1</sup>

# Further readings

For learning about more advanced topics in RL

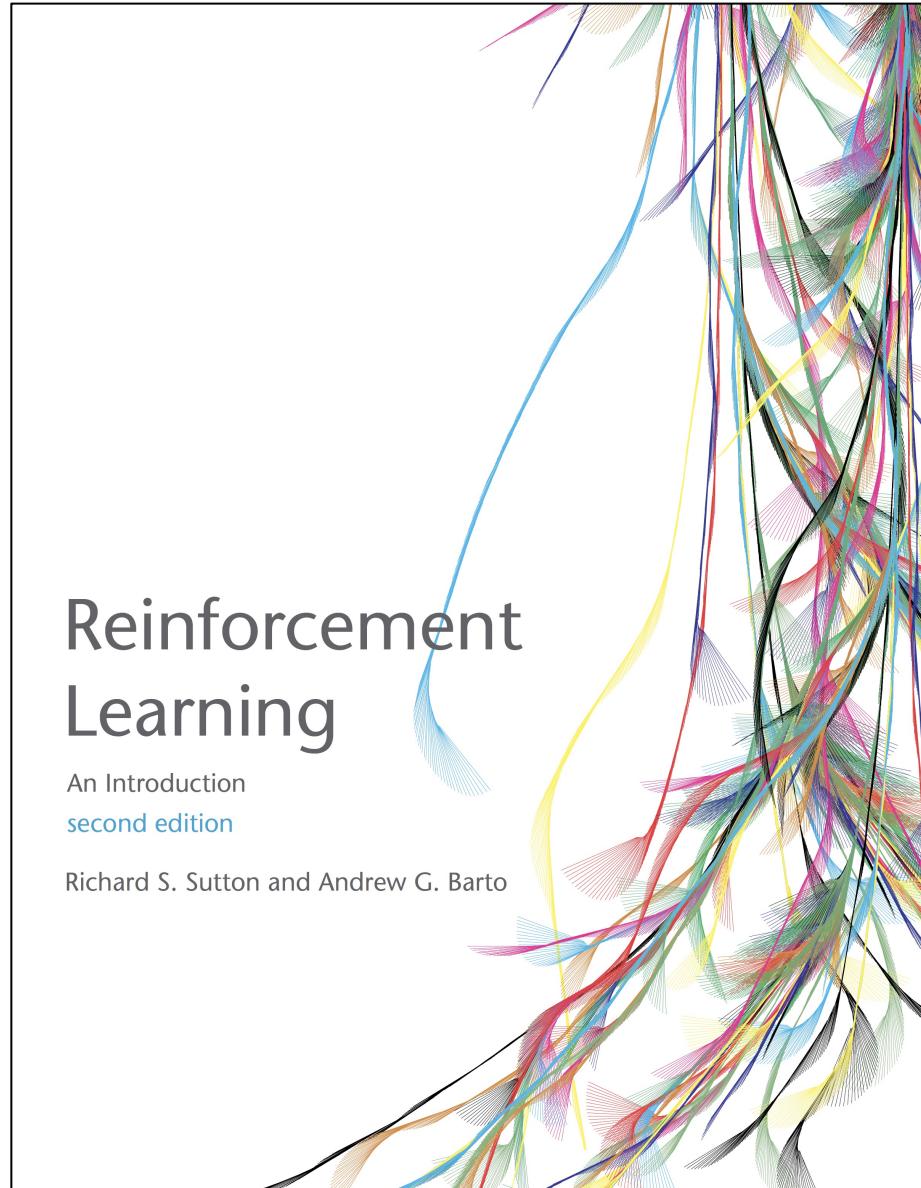
Pioneers in RL:



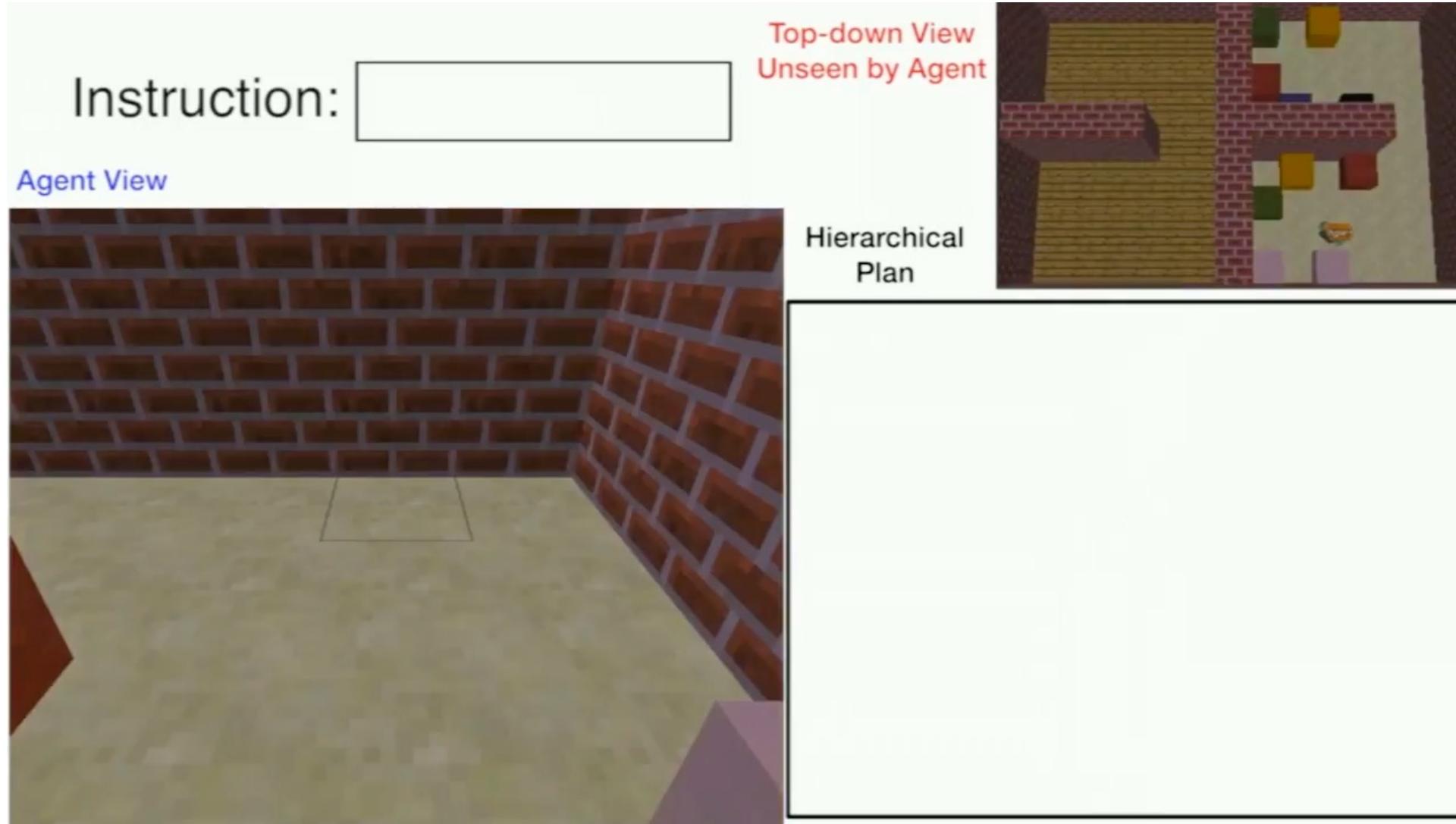
Richard S. Sutton



Andrew G. Barto



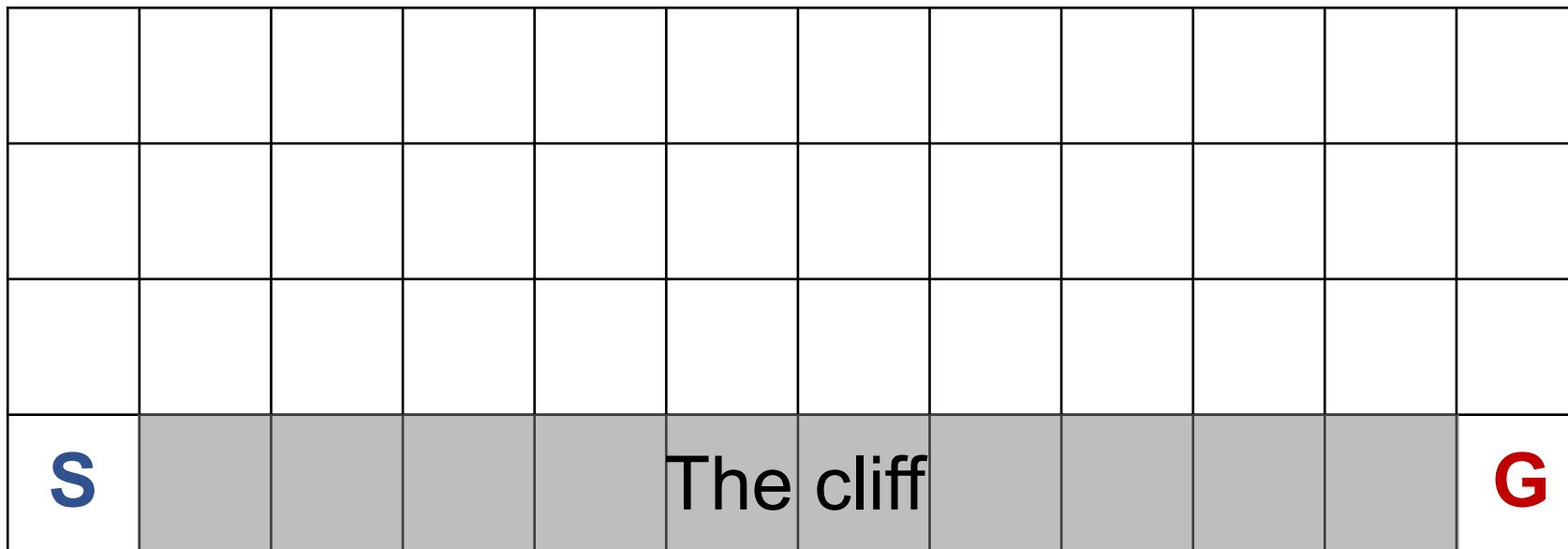
# Multimodal inputs and outputs



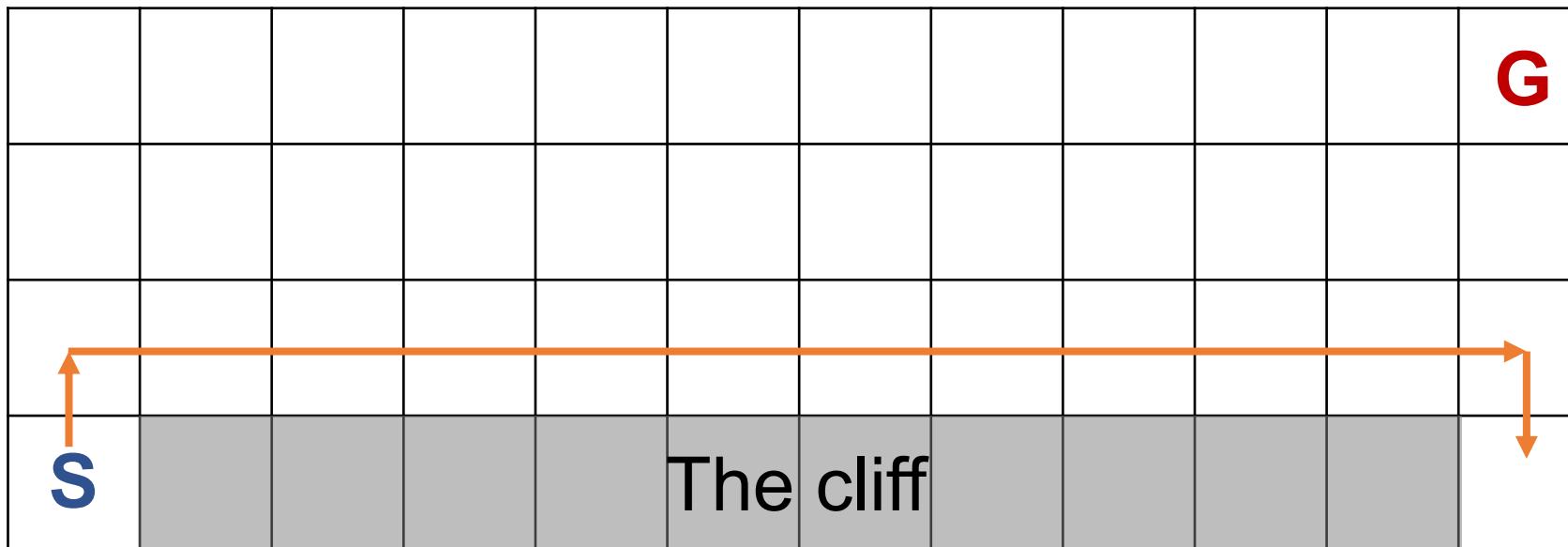
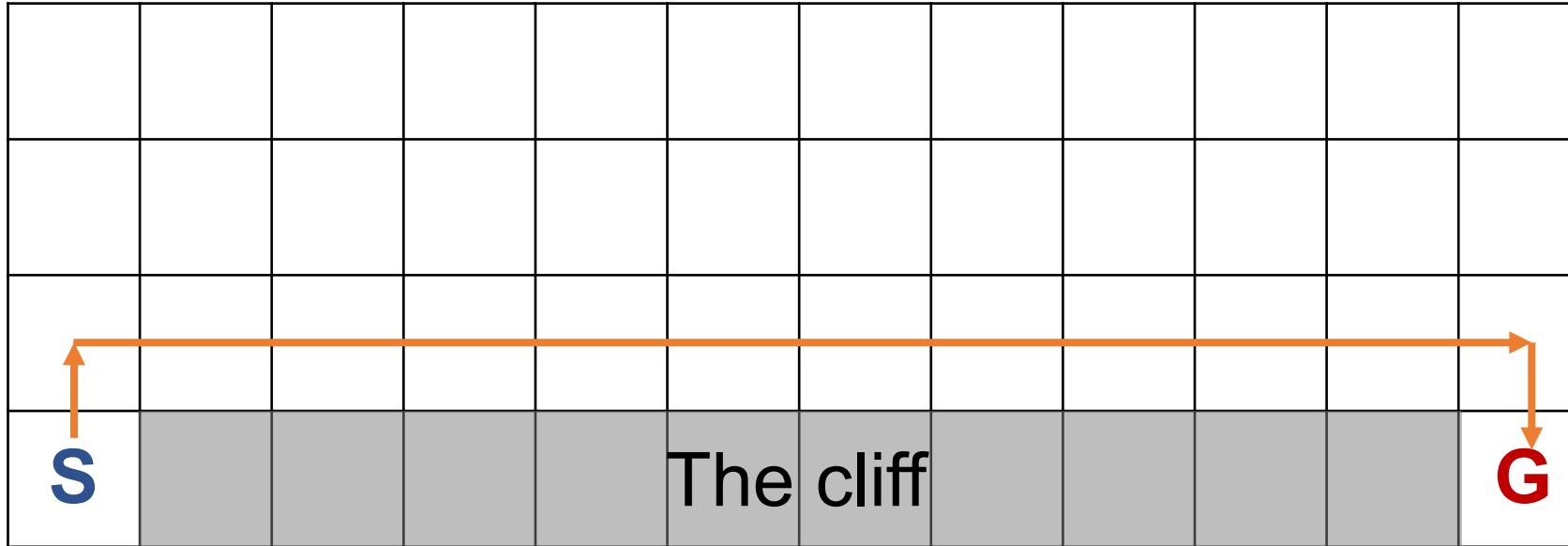
# Summary

- MDPs, policies and value functions
- Bellman equations
- Policy iteration, value iteration
- Q-learning

# Challenges for RL in the real world: safety



# Challenges for RL in the real world: generalization

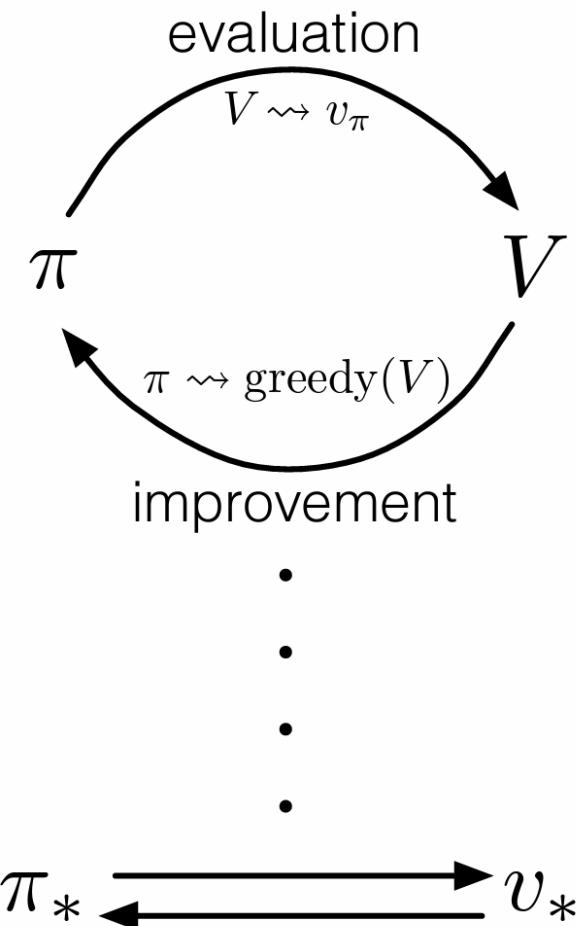


# Challenges for RL in the real world: generalization



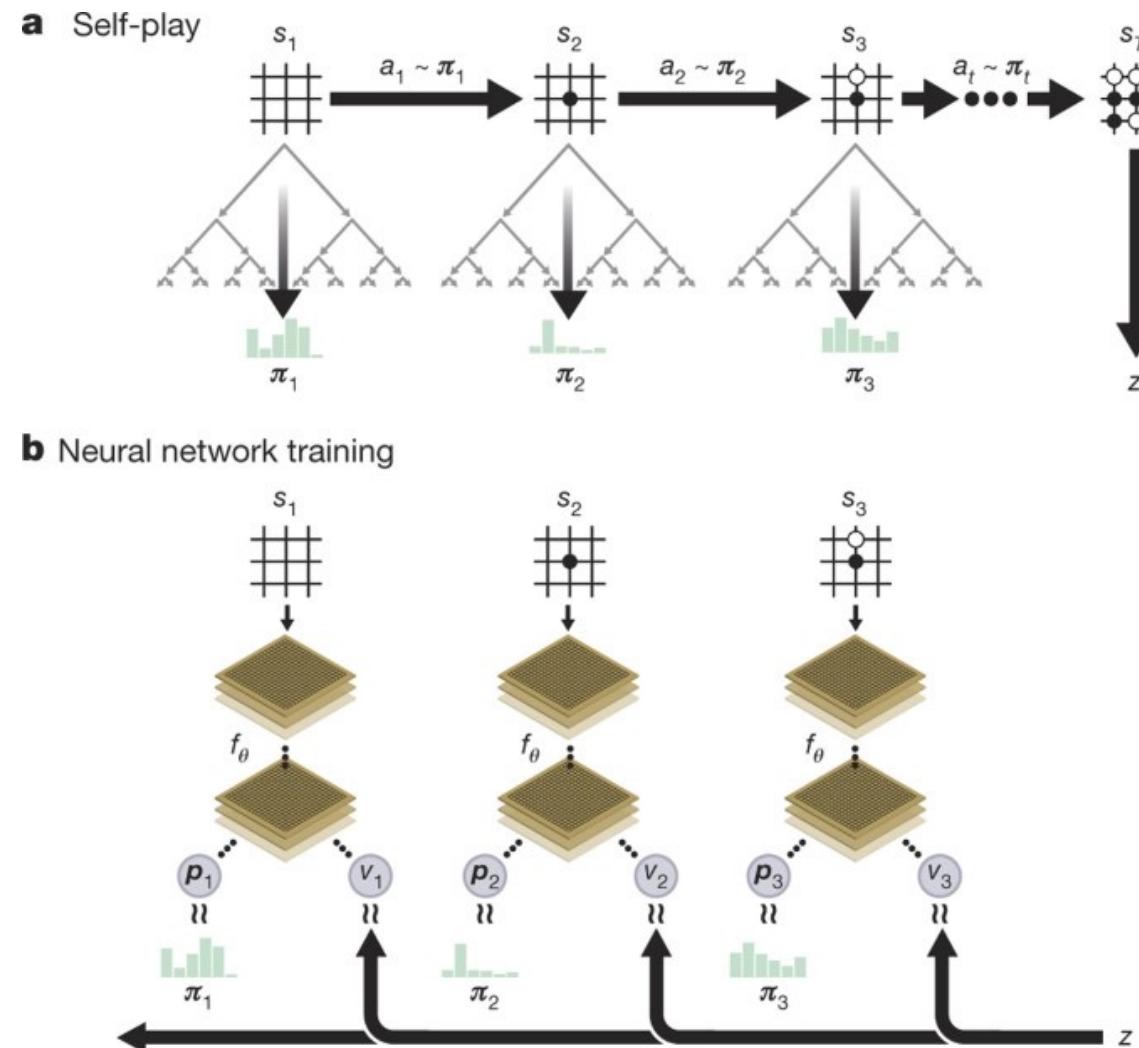
# Model-based learning

- Beyond traditional planning algorithms
- World model → model-based planning
- Policy and value function



# Model-based learning

- Beyond traditional planning algorithms
- AlphaGo



# Model-based learning

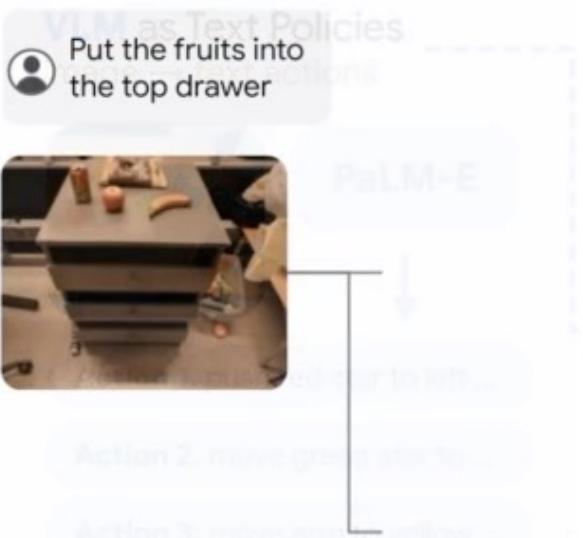
- Beyond traditional planning algorithms

## Video Language Planning

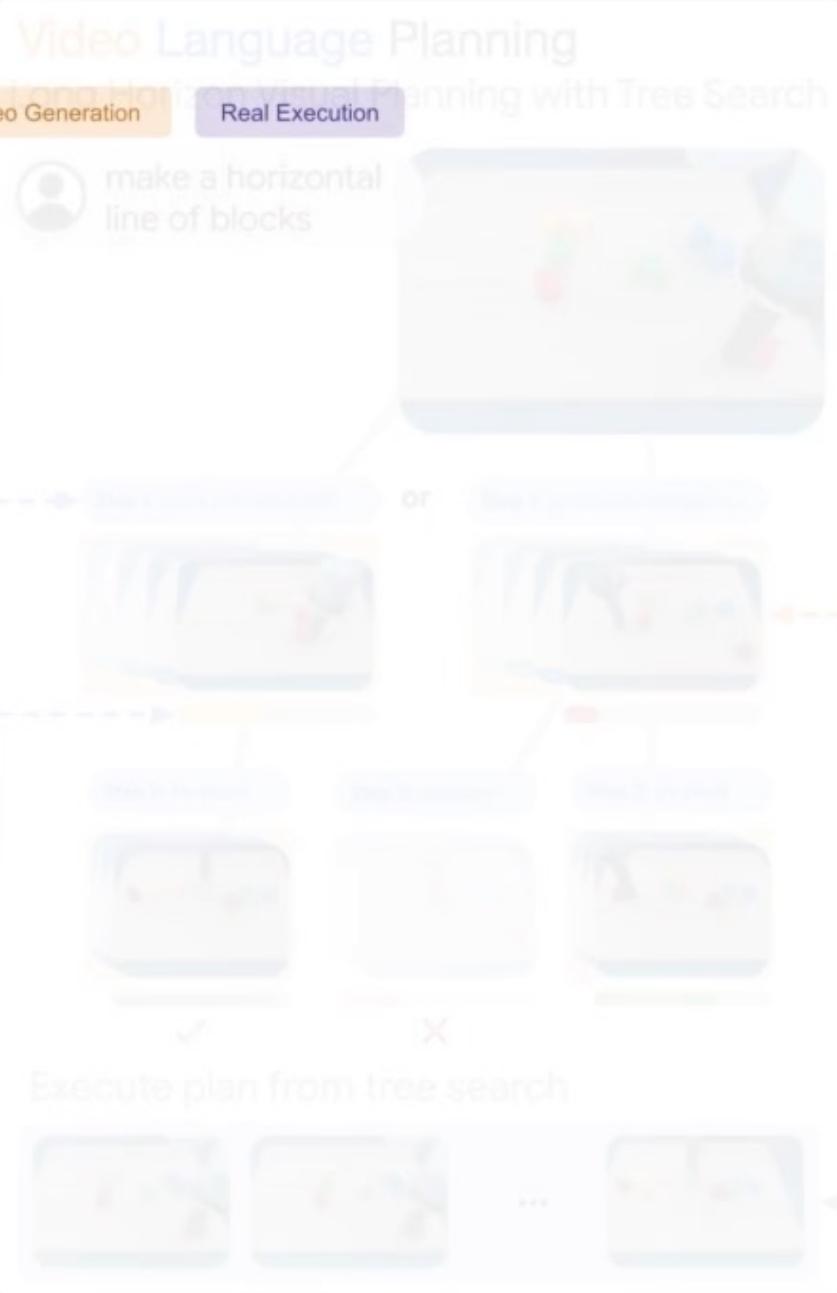
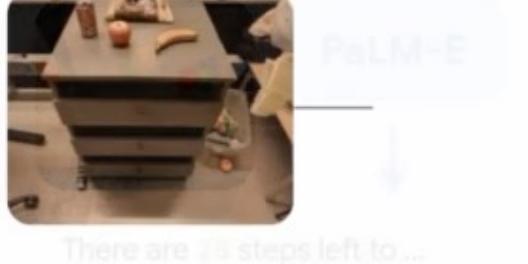
Yilun Du Sherry Yang Pete Florence Fei Xia Ayzaan Wahid Brian Ichter Pierre Sermanet Tianhe Yu  
Pieter Abbeel Joshua B. Tenenbaum Leslie Kaelbling Andy Zeng Jonathan Tompson

- Vision-language models to serve as both policies and value functions
- Text-to-video models as dynamics models.

## Video Language Planning



## VLM as Heuristic Functions



Video Model as Dynamics Model  
text + image → video rollouts of the future

push red star to green → Text-to-Video



Actions from Goal-Conditioned Policies



Policies

Experiments across Multiple Robots

