

# **Reinforcement Learning & Markov Decision Processes**

**Lucas Janson and Sham Kakade**

**CS/Stat 184: Introduction to Reinforcement Learning  
Fall 2022**

# Today

- HW 1 due this thurs
- Today: what is Reinforcement Learning?
  - examples/concepts
  - definition of Markov Decision Processes

# Today:

## Intro to Markov Decision Processes

# Four main themes we will cover in this course:

1. Bandits (horizon  $H = 1$ )
2. Two models, with horizon  $H > 1$ :
  - Markov Decision Process: Dynamic Programming & planning
  - Continuous Control

(technically, this is still an MDP, but with special structure)
3. Learning in “Large” Markov Decision Process
4. Advanced Topics



# Supplementary Reading Materials: Reinforcement Learning: Theory & Algorithms

<https://rltheorybook.github.io/>

This is an advanced RL book.  
We will pick **specific subsections**, to further your knowledge.

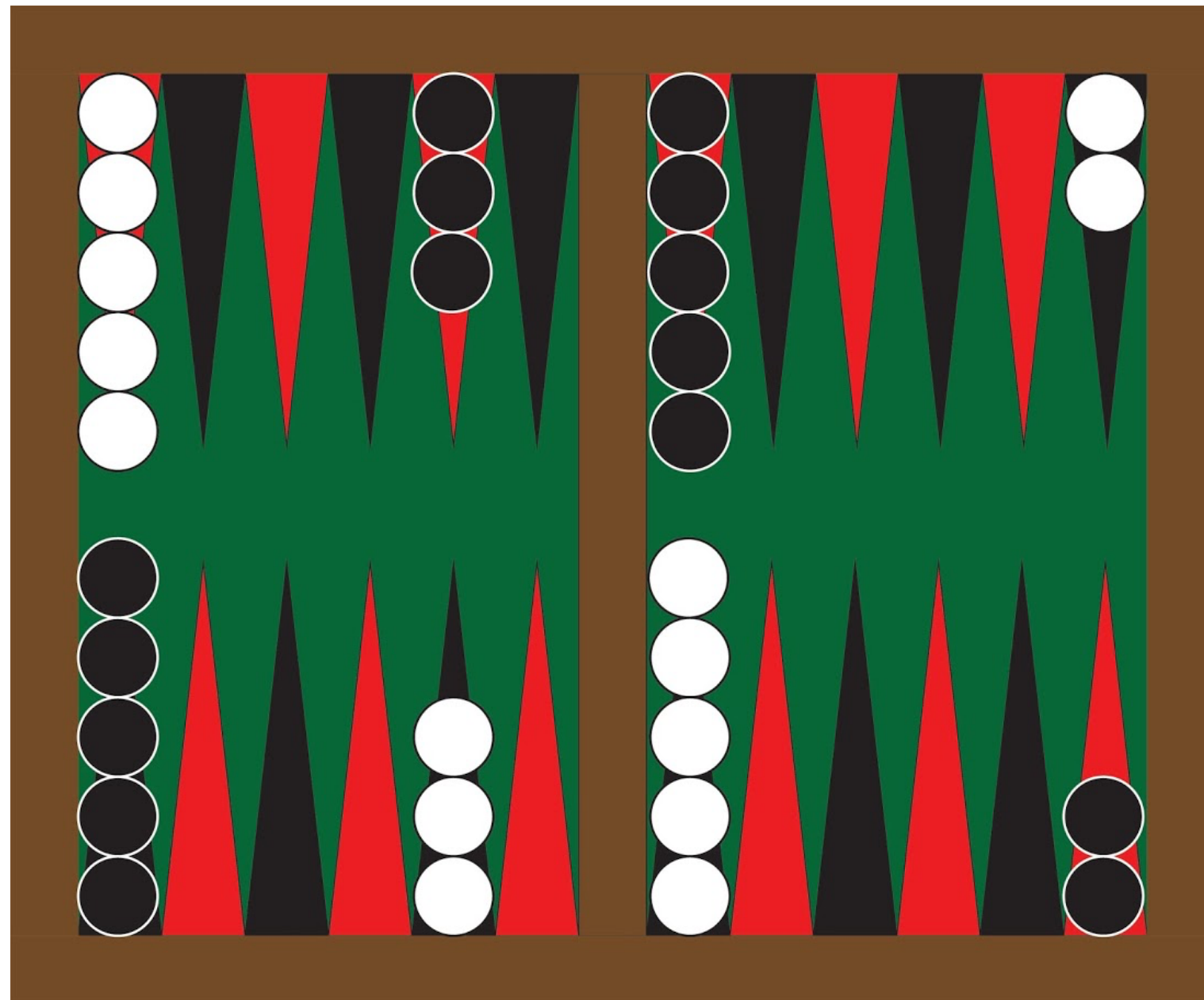
Please let us know if you find any typos or mistakes in the book

# Outlines:

1. Introduction: Applications of RL, RL versus Supervised Learning
2. Basics of Markov Decision Process (MDP): model, example,  $V$  &  $Q$  functions



# Big Successful Stories of Reinforcement Learning



TD GAMMON [Tesauro 95]



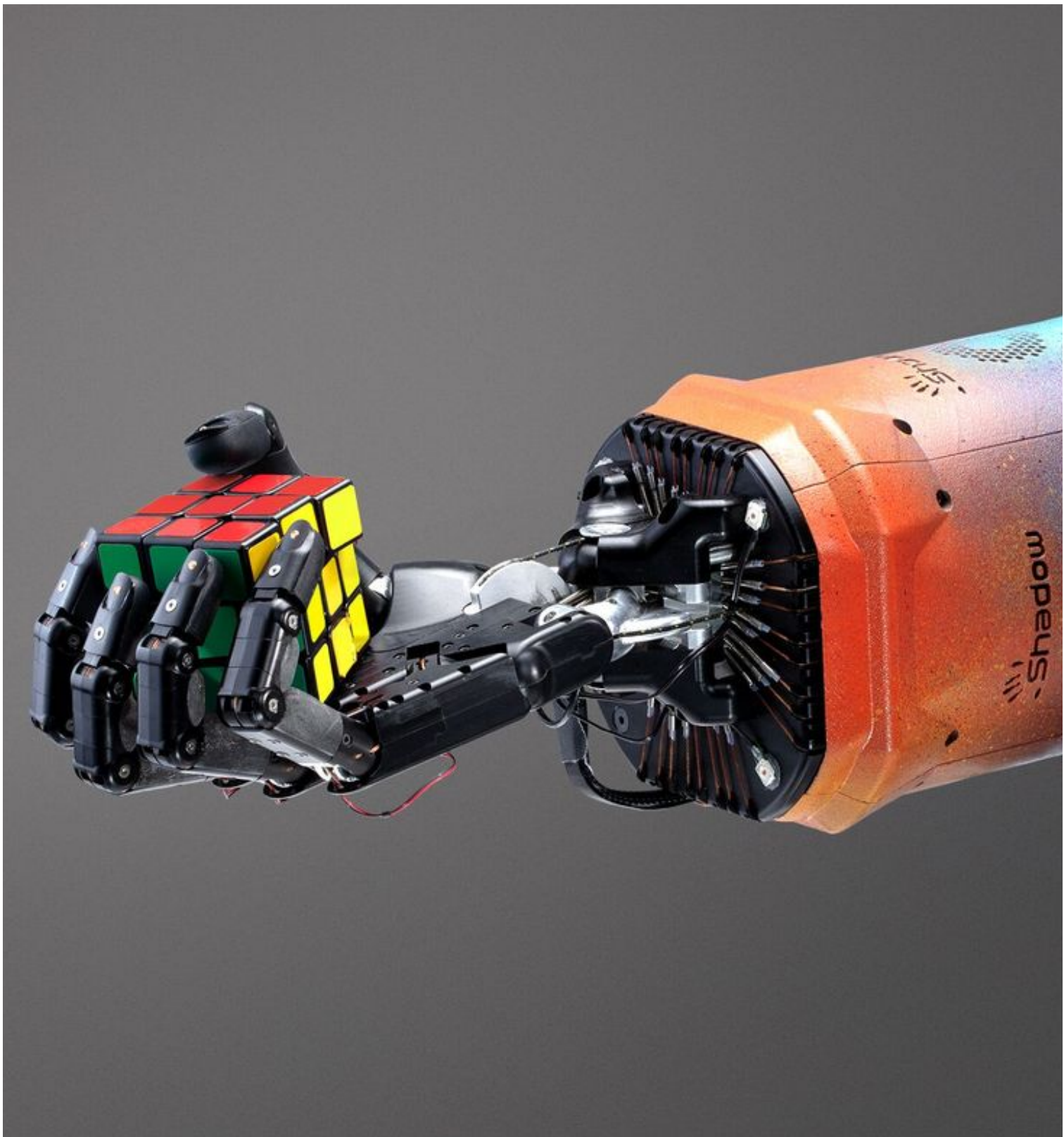
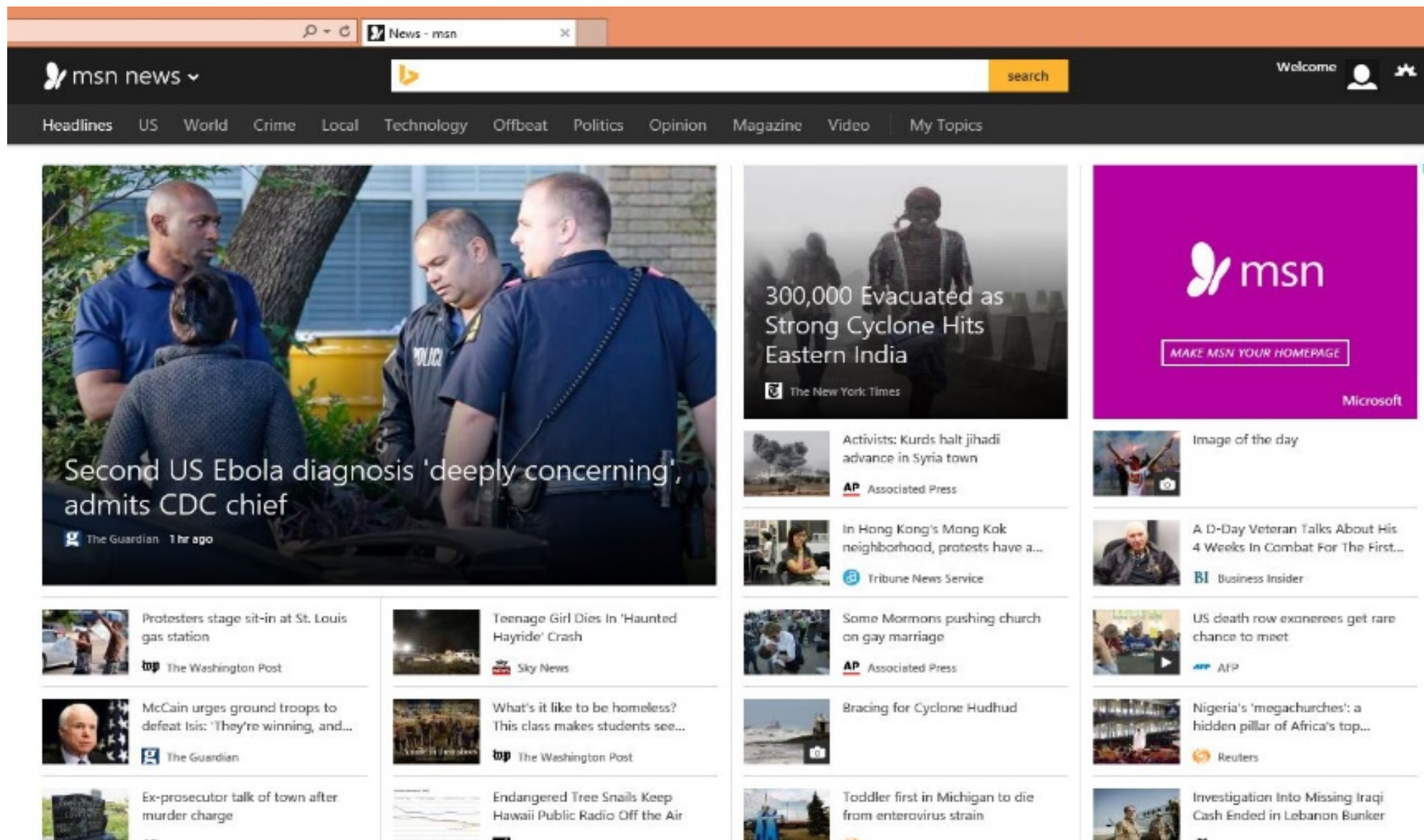
[AlphaZero, Silver et.al, 17]



[OpenAI Five, 18]



# Reinforcement Learning in Real World:

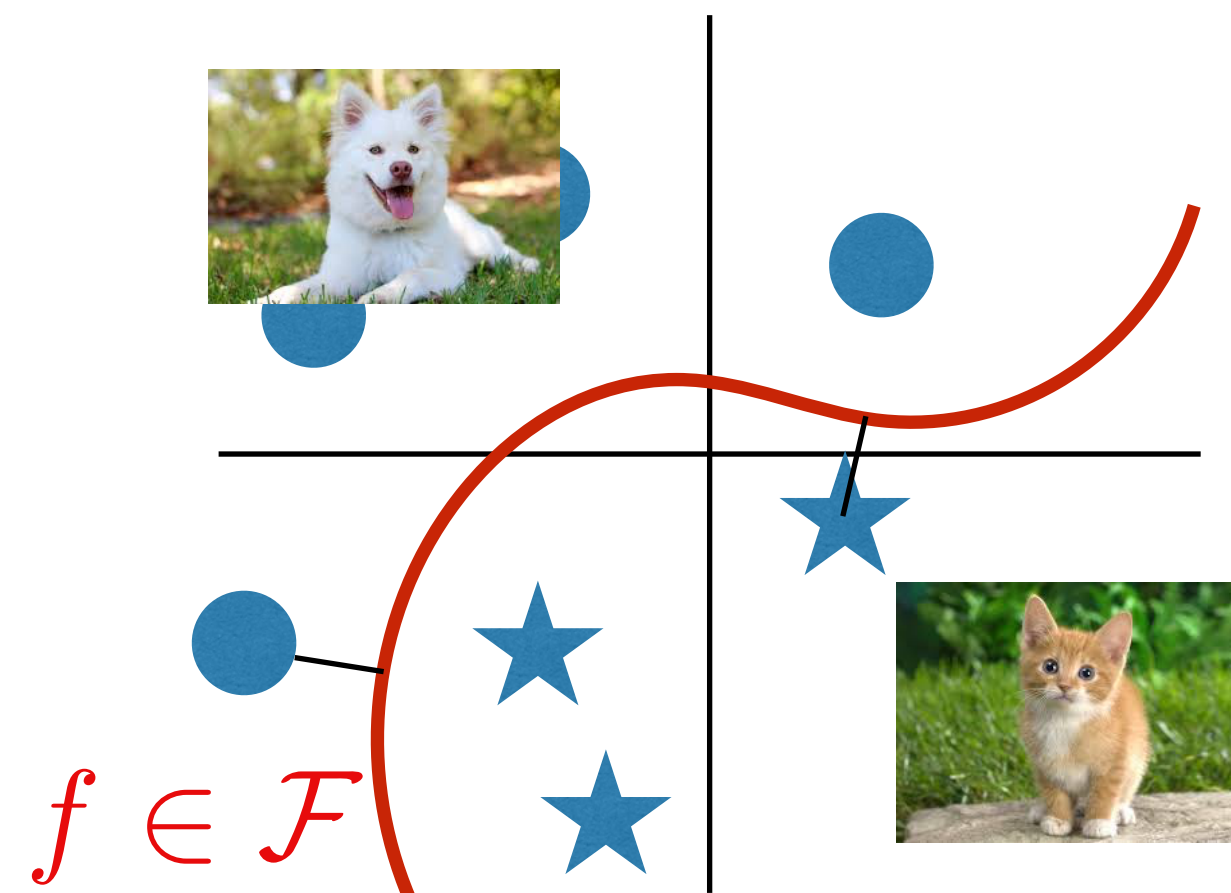
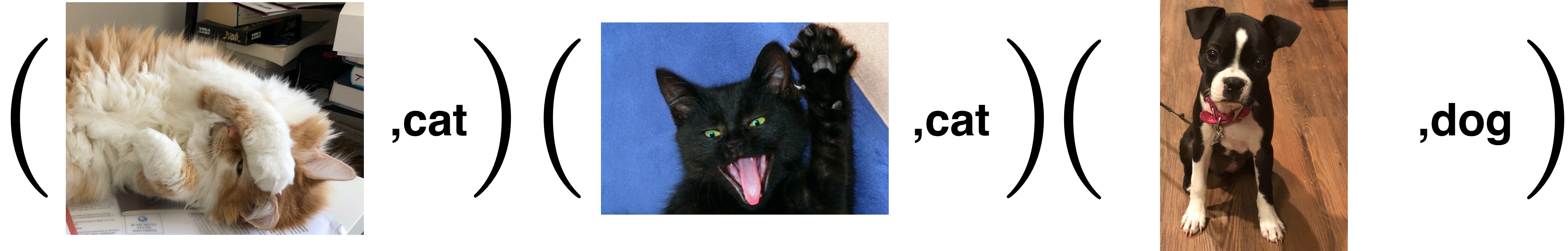


**To better understand RL,  
let's summarize “supervised learning”**



# Recap: Supervised Learning

Given i.i.d examples at training:



Passive:

Prediction



Data Distribution

AgentLinear  
Selected Actions:

RIGHT

SPEED

Active: Decisions → Data Distribution

# Summary so far:

1. In RL, we often start from zero data

2. In RL, **decisions/predictions have consequences:**

Future data is determined by our past historical decisions/predictions

3. To solve the task, we often need to make a **long sequence of decisions**

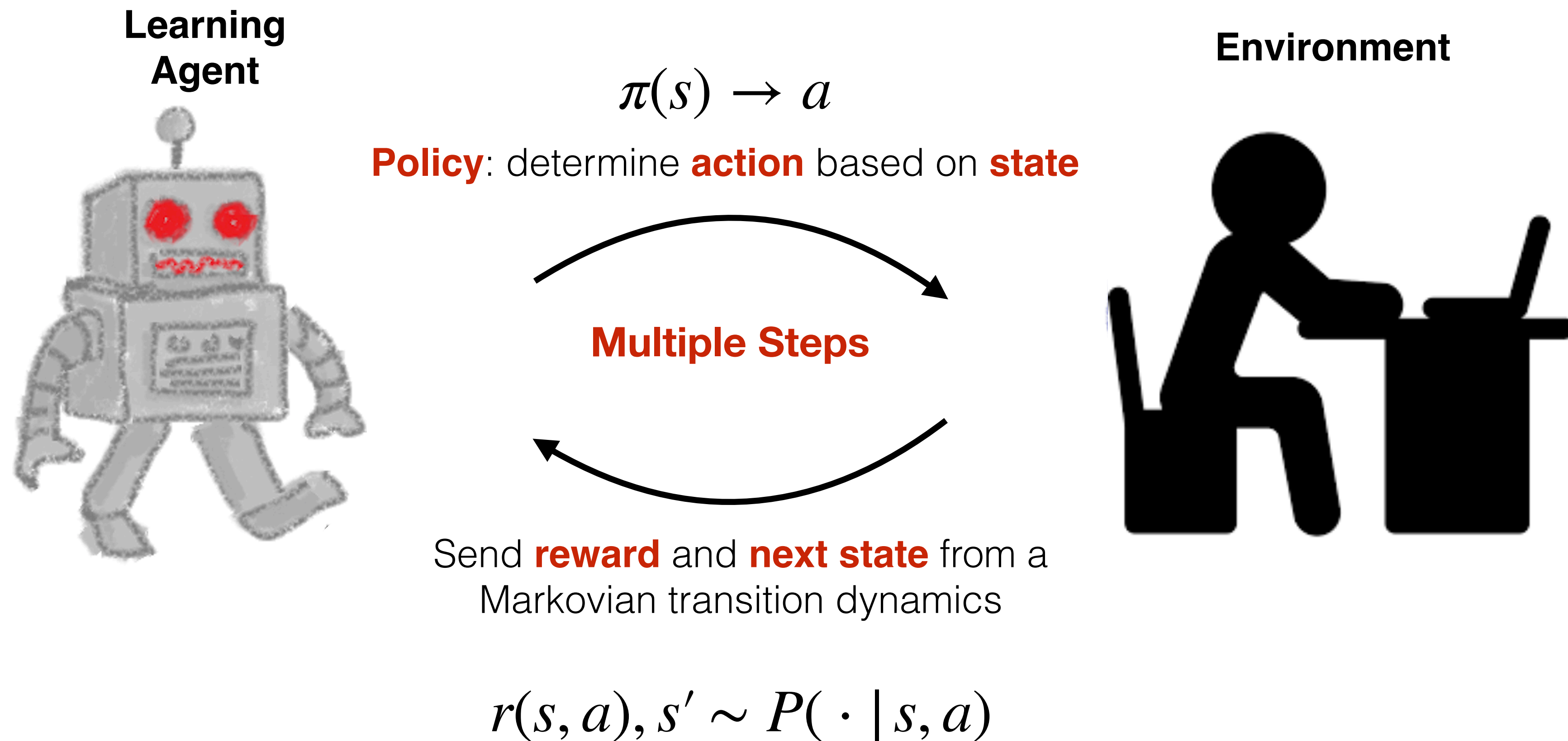


# Outlines:

 1. Introduction: Applications of RL, RL versus Supervised Learning

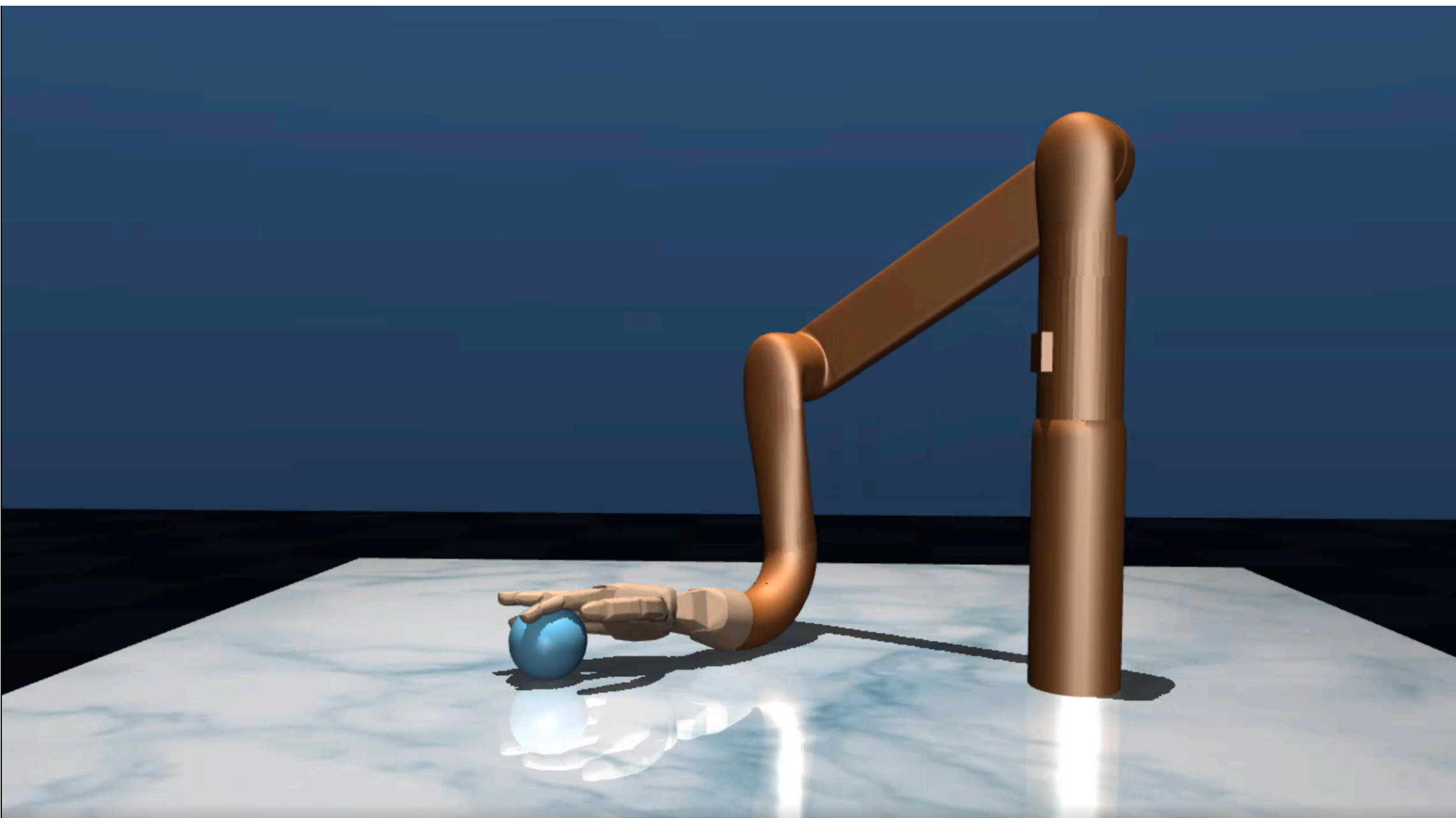
2. Basics of Markov Decision Process (MDP): model, example, V & Q functions

# The Mathematical framework: Markov Decision Process



## Example:

robot hand needs to pick the ball and hold it in a goal (x,y,z) position



**State**  $s$ : robot configuration (e.g., joint angles) and the ball's position

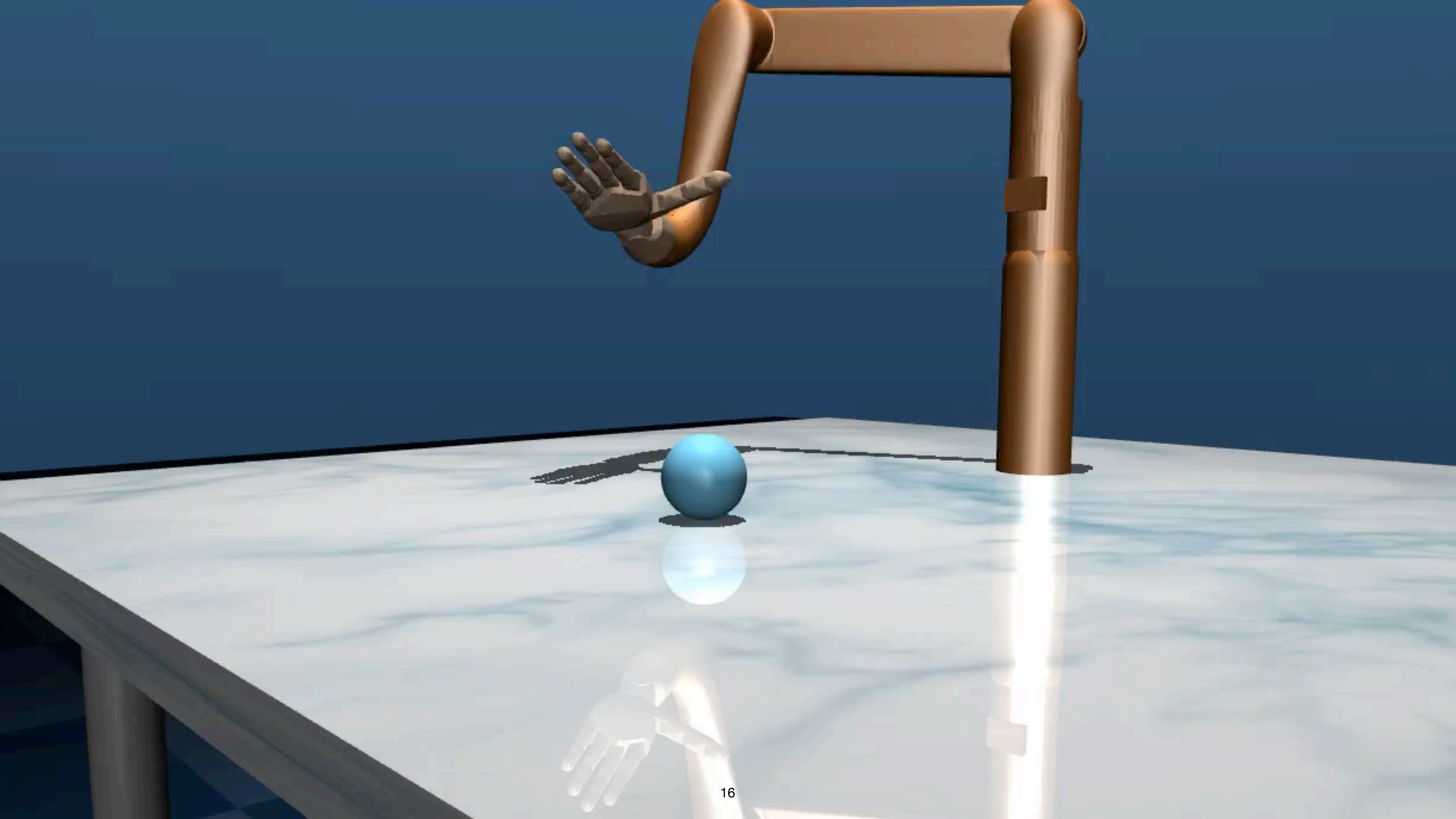
**Action**  $a$ : Torque on joints in arm & fingers

**Transition**  $s' \sim P(\cdot | s, a)$ : physics + some noise

**policy**  $\pi(s)$ : a function mapping from robot state to action (i.e., torque)

**Cost**  $c(s, a)$ : torque magnitude + dist to goal

$$\pi^{\star} = \arg \min_{\pi} \mathbb{E} \left[ c(s_0, a_0) + \gamma c(s_1, a_1) + \gamma^2 c(s_2, a_2) + \gamma^3 c(s_3, a_3) + \dots \mid s_0, \pi \right]$$





## MDPs, more formally:

- An MDP:  $\mathcal{M} = \{S, A, P, r, \gamma\}$ 
  - $S$  a set of states
  - $A$  a set of actions
  - $P : S \times A \mapsto \Delta(S)$  specifies the dynamics model,  
i.e.  $P(s' | s, a)$  is the probability of transitioning to  $s'$  from states  $s$  under action  $a$
  - $r : S \times A \rightarrow [0,1]$ 
    - let's assume this is deterministic
    - (sometimes we use a cost  $c : S \times A \rightarrow [0,1]$ )
  - A discount factor  $\gamma \in [0,1)$

# The Objective

- A “stationary” policy  $\pi : S \mapsto A$ 
  - “stationary” means not history dependent
  - we could also consider  $\pi$  to be random and a function of the history
- Sampling a trajectory: from a given policy  $\pi$  starting at state  $s_0$ :
  - For  $t = 0, 1, 2, \dots \infty$ 
    - Take action  $a_t = \pi(s_t)$
    - Observe reward  $r_t = r(s_t, a_t)$
    - Transition to (and observe)  $s_{t+1}$  where  $s_{t+1} \sim P(\cdot \mid s_t, a_t)$
- Objective: given state starting state  $s$ ,  
find a policy  $\pi$  that maximizes our expected, discounted future reward:

$$\max_{\pi} \mathbb{E} \left[ r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots \mid s_0 = s, \pi \right]$$

## Question:

**Assume we have  $|S|$  many states, and  $|A|$  many actions, how many different policies there are?**

(Hint: a policy is a mapping from  $s$  to  $a$ , we have  $A$  many choices per state  $s$ )

# Infinite horizon Discounted Setting

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto A$$

Quantities that allow us to reason policy's long-term effect:

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), \pi \right]$$



# Understanding Value function and Q functions

Value function  $V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right]$

Q function  $Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), \pi \right]$

# Bellman Consistency Equation for V-function:

$$V^\pi(s) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right]$$

We have that:

$$V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s')$$

# Proof: Bellman Consistency for V-function:

- By definition:

$$\begin{aligned} V^\pi(s) &= r(s, \pi(s)) + \mathbb{E} \left[ \sum_{h=1}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right] \\ &= r(s, \pi(s)) + \gamma \mathbb{E} \left[ r(s_1, a_1) + \gamma r(s_2, a_2) + \dots \mid s_0 = s, \pi \right] \end{aligned}$$

- By the “tower property” and due to that  $s_1 = s'$  with probability  $P(s' \mid s, \pi(s))$ ,

$$\begin{aligned} &= r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, \pi(s))} \left[ \mathbb{E} \left[ r(s_1, a_1) + \gamma r(s_2, a_2) + \dots \mid s_0 = s, s_1 = s', \pi \right] \right] \\ &= r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, \pi(s))} \left[ \mathbb{E} \left[ r(s_1, a_1) + \gamma r(s_2, a_2) + \dots \mid s_1 = s', \pi \right] \right] \\ &= r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, \pi(s))} \left[ V^\pi(s') \right] \end{aligned}$$

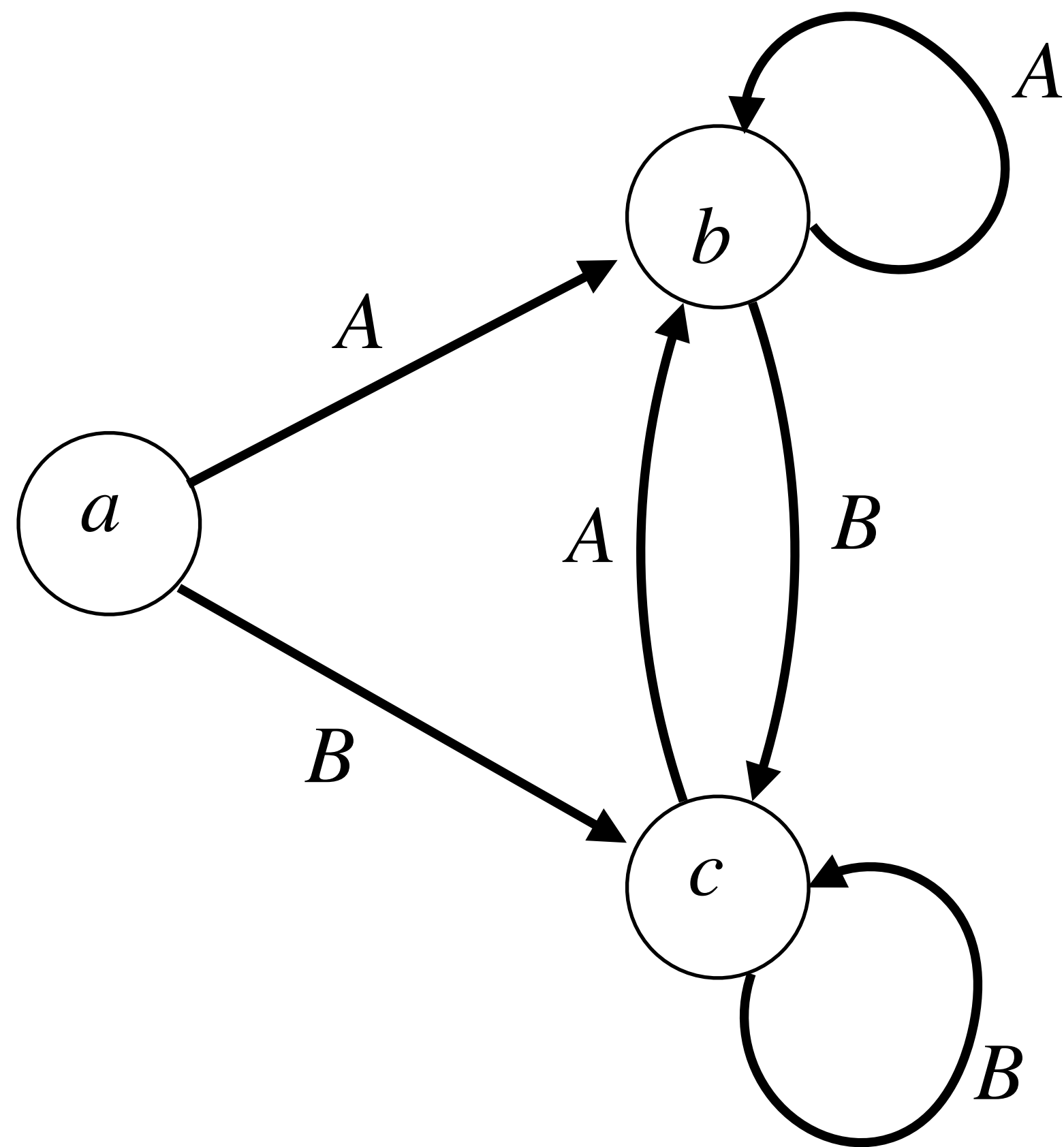
# Bellman Consistency Equation for Q-function:

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), \pi \right]$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\pi(s')$$

# Example of Optimal Policy $\pi^\star$

Consider the following **deterministic** MDP w/ 3 states & 2 actions



Let's say  $\gamma \in (0,1)$   
What's the optimal policy?

$$\pi^\star(s) = A, \forall s$$

$$V^\star(a) = \frac{\gamma}{1-\gamma}, V^\star(b) = \frac{1}{1-\gamma}, V^\star(c) = \frac{\gamma}{1-\gamma}$$

What about policy  $\pi(s) = B, \forall s$

$$V^\pi(a) = 0, V^\pi(b) = 0, V^\pi(c) = 0$$

Reward:  $r(b, A) = 1$ , & 0 everywhere else

# Summary:

- **RL is different from Supervised Learning:**
  - Our actions have consequences
  - Need to make sequence of decisions to complete the task
- **Discounted infinite horizon MDP:**
  - State, action, policy, transition, reward (or cost), discount factor
  - **V function and Q function**
  - Key concept: **Bellman consistency equations**

1-minute feedback form: <https://bit.ly/3RHtlxy>

