# Reinforcement Learning & Markov Decision Processes

## Lucas Janson

**CS/Stat 184(0): Introduction to Reinforcement Learning**
**Fall 2024**

# Today

- Logistics (Welcome!)

- Overview of RL

- Markov Decision Processes

  - Problem statement

  - Policy Evaluation

# Course staff introductions

# Course staff introductions

- **Instructor:** Lucas Janson

# Course staff introductions

- **Instructor:** Lucas Janson

- **TFs:** Anvit Garg, Nowell Closser

# Course staff introductions

- **Instructor:** Lucas Janson

- **TFs:** Anvit Garg, Nowell Closser

- **CAs:** Jayden Personnat, Sibi Raja, Alex Cai, Ethan Tan, Neil Shah, Jason Wang, Russell Li, Sid Bharthulwar, Andrew Gu, Ian Moore

# Course staff introductions

- **Instructor:** Lucas Janson

- **TFs:** Anvit Garg, Nowell Closser

- **CAs:** Jayden Personnat, Sibi Raja, Alex Cai, Ethan Tan, Neil Shah, Jason Wang, Russell Li, Sid Bharthulwar, Andrew Gu, Ian Moore

- Homework 0 is posted!

  - This is "review" homework for material you should be familiar with to take the course.

# Course Overview

**All policies are stated on the course website:**
**http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

# Course Overview

**All policies are stated on the course website:**
**http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.

# Course Overview

**All policies are stated on the course website: http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.
- **Grades: Participation; HW0 +HW1-HW4; Midterm; Project**

# Course Overview

**All policies are stated on the course website:
http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.
- **Grades: Participation; HW0 +HW1-HW4; Midterm; Project**
- Participation (5%): not meant to be onerous (see website)
  - Just attending regularly will suffice
  - If you can't, then increase your participation in Ed/section.
  - Let us know if you have some hard conflict, let us know via Ed.

# Course Overview

**All policies are stated on the course website:**
**http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.
- **Grades: Participation; HW0 +HW1-HW4; Midterm; Project**
- Participation (5%): not meant to be onerous (see website)
  - Just attending regularly will suffice
  - If you can't, then increase your participation in Ed/section.
  - Let us know if you have some hard conflict, let us know via Ed.
- HWs (45%): will have math and programming components.
  - We will have an "embedded ethics lecture" + assignment

# Course Overview

**All policies are stated on the course website:**
**http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.

- **Grades: Participation; HW0 +HW1-HW4; Midterm; Project**

- Participation (5%): not meant to be onerous (see website)

  - Just attending regularly will suffice

  - If you can't, then increase your participation in Ed/section.

  - Let us know if you have some hard conflict, let us know via Ed.

- HWs (45%): will have math and programming components.

  - We will have an "embedded ethics lecture" + assignment

- Midterm (20%): this will be in class.

# Course Overview

**All policies are stated on the course website:
http://lucasjanson.fas.harvard.edu/CS_Stat_184_0.html**

- We want you to obtain fundamental and practical knowledge of RL.
- **Grades: Participation; HW0 +HW1-HW4; Midterm; Project**
- Participation (5%): not meant to be onerous (see website)
  - Just attending regularly will suffice
  - If you can't, then increase your participation in Ed/section.
  - Let us know if you have some hard conflict, let us know via Ed.
- HWs (45%): will have math and programming components.
  - We will have an "embedded ethics lecture" + assignment
- Midterm (20%): this will be in class.
- Project (30%): 2-3 people per project. Will be empirical.

# Other Points

# Other Points

- Our policies aim for consistency among all the students.

# Other Points

- Our policies aim for consistency among all the students.
- Participation: we will have a web-based attendance form

# Other Points

- Our policies aim for consistency among all the students.

- Participation: we will have a web-based attendance form

- Communication: please only use Ed to contact us

# Other Points

- Our policies aim for consistency among all the students.

- Participation: we will have a web-based attendance form

- Communication: please only use Ed to contact us

- Late policy (basically): you have 96 cumulative hours of late time.

# Other Points

- Our policies aim for consistency among all the students.

- Participation: we will have a web-based attendance form

- Communication: please only use Ed to contact us

- Late policy (basically): you have 96 cumulative hours of late time.

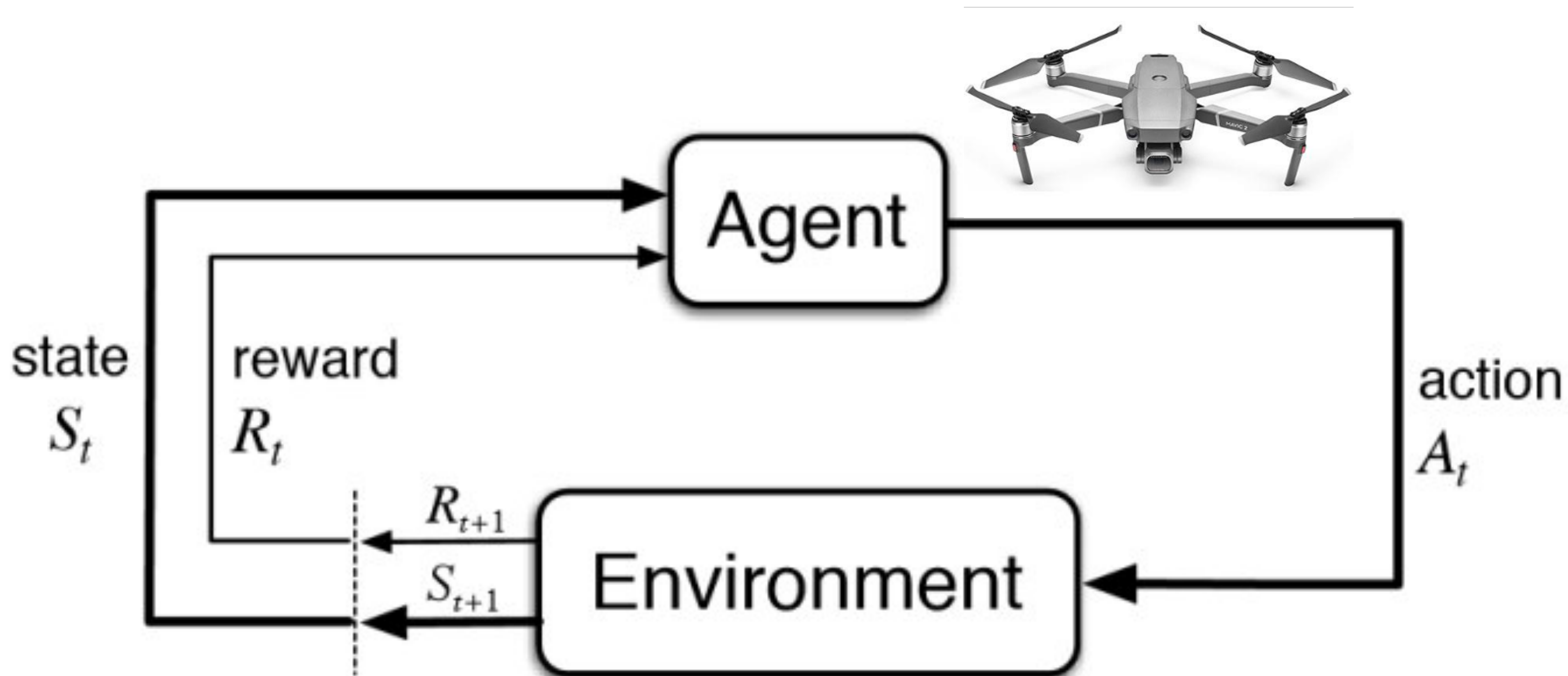  - *Please use this to plan for unforeseen circumstances.*

# Other Points

- Our policies aim for consistency among all the students.

- Participation: we will have a web-based attendance form

- Communication: please only use Ed to contact us

- Late policy (basically): you have  96 cumulative hours of late time.

  - *Please use this to plan for unforeseen circumstances.*

- Regrading: ask us in writing on Ed within a week

# Today

✓ • Logistics (Welcome!)

• Overview of RL

• Markov Decision Processes

  • Problem statement

  • Policy Evaluation

# The RL Setting, basically

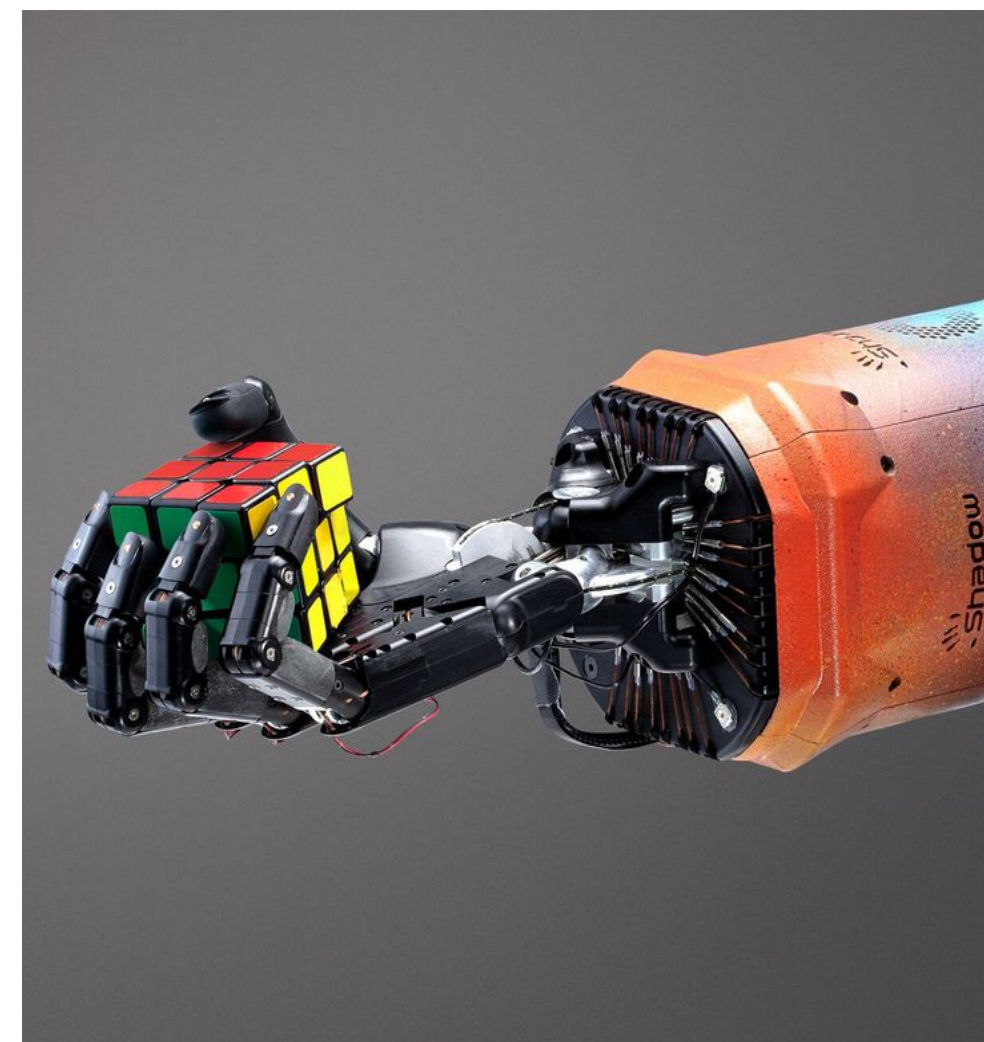# Many RL Successes
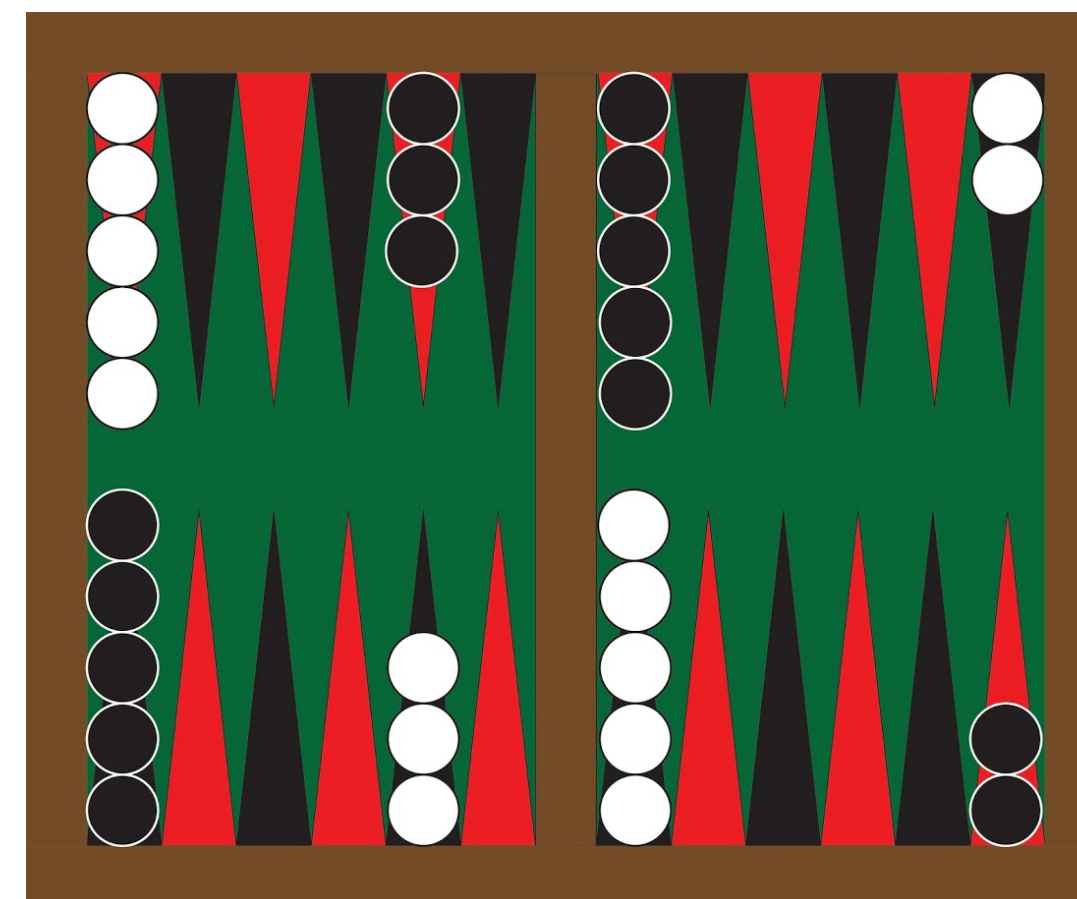


Online advertising



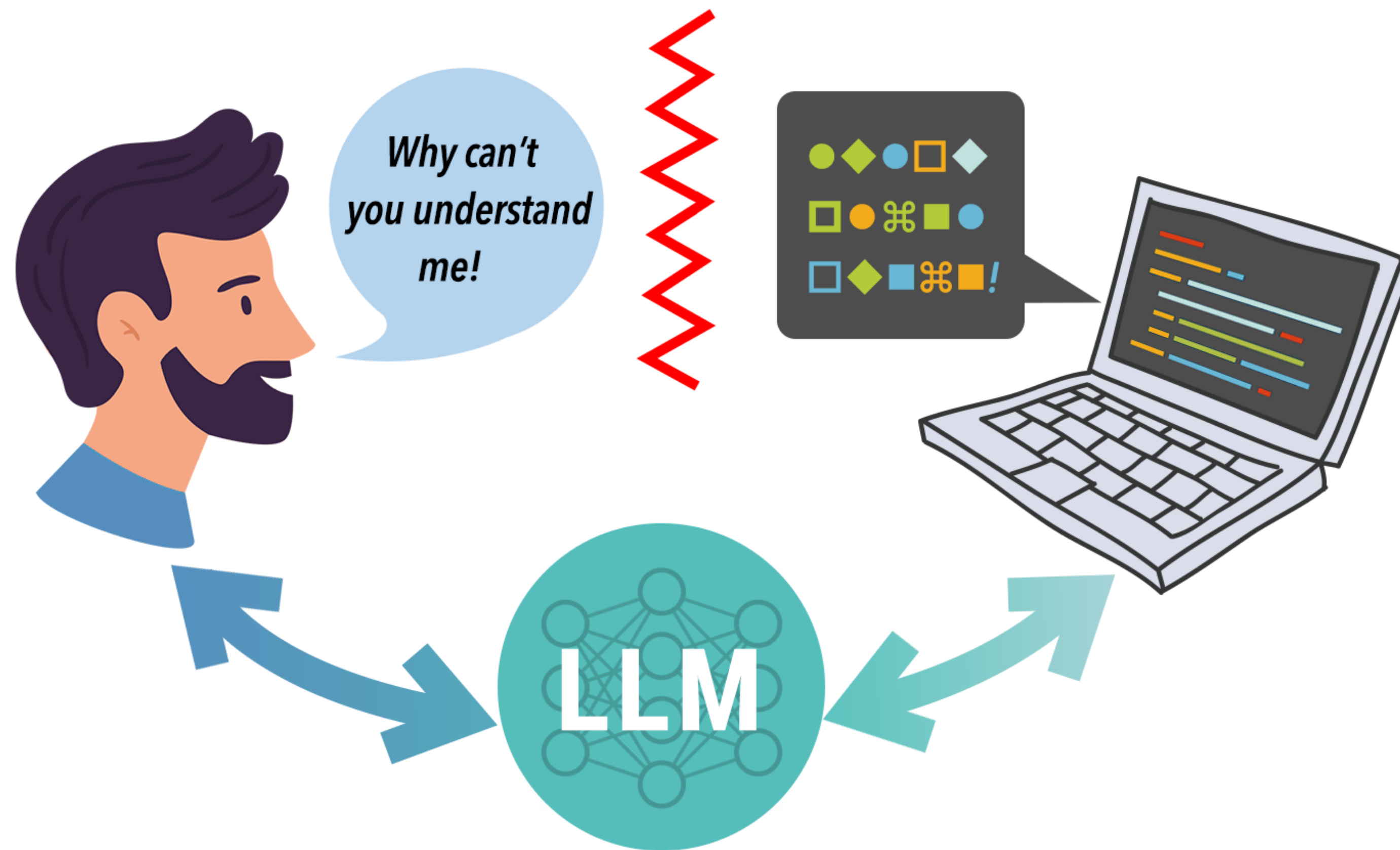[AlphaZero, Silver et.al, 17]



[OpenAI Five, 18]



[OpenAI,19]



TD GAMMON [Tesauro 95]



Supply Chains [Madeka et al '23]

8

# Many Future RL Challenges

# Vs Other Settings

| | Learn from Experience | Generalize | Interactive | Exploration | Credit assignment |
|---|---|---|---|---|---|
| Supervised Learning | ✔ | ✔ | | | |
| Bandits ("horizon 1"-RL) | ✔ | ✔ | ✔ | ✔ | |
| "Full" Reinforcement Learning | ✔ | ✔ | ✔ | ✔ | ✔ |

# Vs Other Settings

| | Learn from Experience | Generalize | Interactive | Exploration | Credit assignment |
|---|---|---|---|---|---|
| **Supervised Learning** | ✔ | ✔ | | | |
| **Bandits ("horizon 1"-RL)** | ✔ | ✔ | ✔ | ✔ | |
| **"Full" Reinforcement Learning** | ✔ | ✔ | ✔ | ✔ | ✔ |

Dog

Dog

Not Dog

→ Supervised Learning → Predictive Model

# Vs Other Settings

| | Learn from Experience | Generalize | Interactive | Exploration | Credit assignment |
|---|---|---|---|---|---|
| **Supervised Learning** | ✔ | ✔ | | | |
| **Bandits ("horizon 1"-RL)** | ✔ | ✔ | ✔ | ✔ | |
| **"Full" Reinforcement Learning** | ✔ | ✔ | ✔ | ✔ | ✔ |

Dog

Dog

Not Dog

→ Supervised Learning → Predictive Model
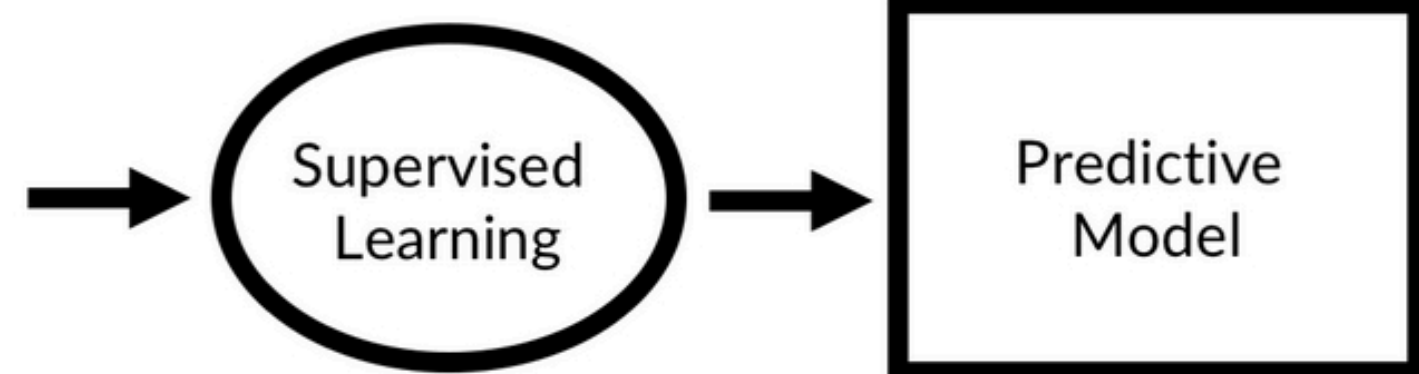
Online Advertising

# Why study RL?

# Why study RL?

- Applications to many important domains

# Why study RL?

- Applications to many important domains

- Very general and intuitive formulation—could be seen as a model for basically anything anyone does in the world

# Why study RL?

- Applications to many important domains

- Very general and intuitive formulation—could be seen as a model for basically anything anyone does in the world

- To me: a more natural way (than supervised learning) to think about "learning" as I do it in my life, where I'm not just predicting but also acting in the world, and *inter*acting with the world through the data I choose to collect

# Why study RL?

- Applications to many important domains

- Very general and intuitive formulation—could be seen as a model for basically anything anyone does in the world

- To me: a more natural way (than supervised learning) to think about "learning" as I do it in my life, where I'm not just predicting but also acting in the world, and *inter*acting with the world through the data I choose to collect

  - I think every human is some sort of reinforcement learner (not as clear that we're supervised learners in the same way, IMO)

# Why study RL?

- Applications to many important domains

- Very general and intuitive formulation—could be seen as a model for basically anything anyone does in the world

- To me: a more natural way (than supervised learning) to think about "learning" as I do it in my life, where I'm not just predicting but also acting in the world, and *inter*acting with the world through the data I choose to collect

    - I think every human is some sort of reinforcement learner (not as clear that we're supervised learners in the same way, IMO)

- Surprising how much you can learn *without* any knowledge of supervised learning
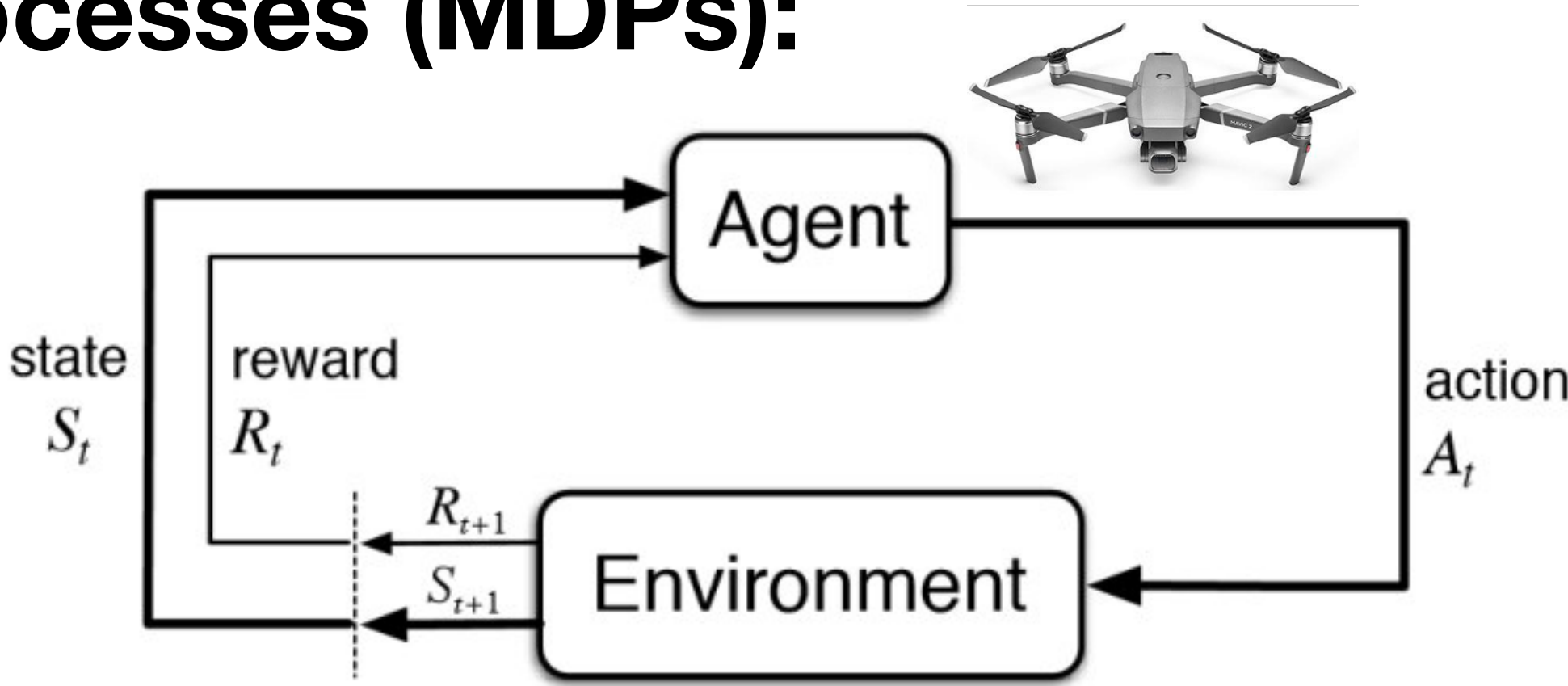
# Why study RL?

- Applications to many important domains

- Very general and intuitive formulation—could be seen as a model for basically anything anyone does in the world

- To me: a more natural way (than supervised learning) to think about "learning" as I do it in my life, where I'm not just predicting but also acting in the world, and *inter*acting with the world through the data I choose to collect

  - I think every human is some sort of reinforcement learner (not as clear that we're supervised learners in the same way, IMO)

- Surprising how much you can learn *without* any knowledge of supervised learning

  - In some sense, the fundamentals of RL are orthogonal to supervised learning
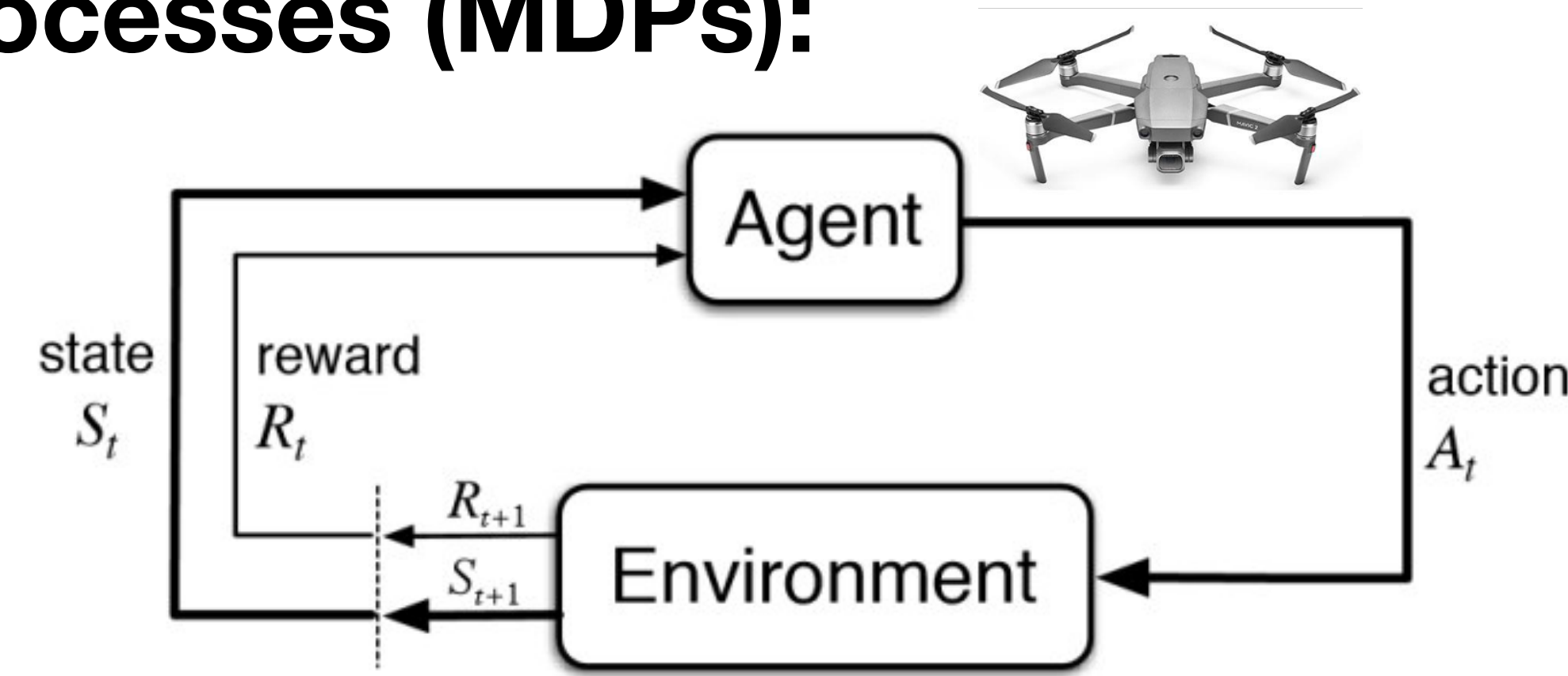
# Today

✓ • Logistics (Welcome!)

✓ • Overview of RL

• Markov Decision Processes

  • Problem statement

  • Policy Evaluation

# Finite Horizon Markov Decision Processes (MDPs):
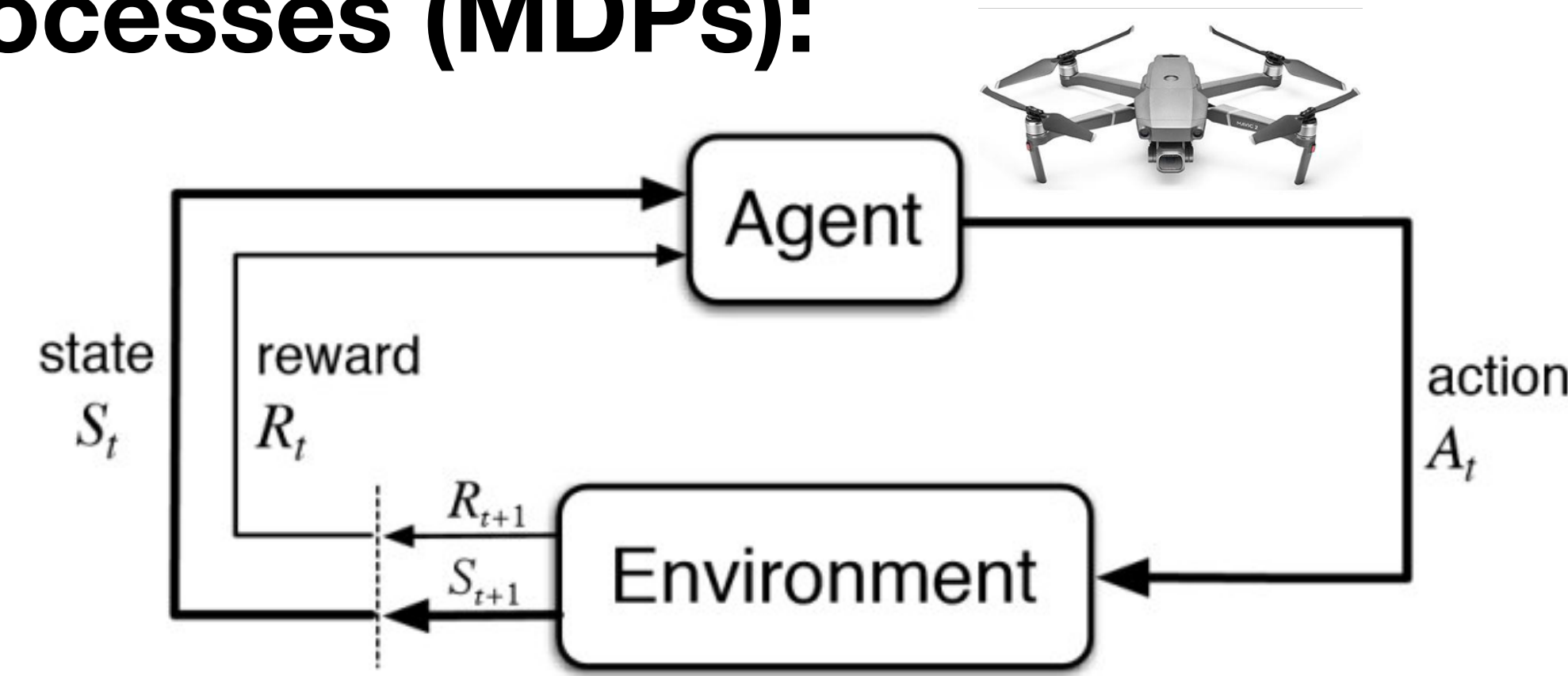
# Finite Horizon Markov Decision Processes (MDPs):

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$



state $S_t$ | reward $R_t$ | action $A_t$

Agent

Environment

$R_{t+1}$

$S_{t+1}$

# Finite Horizon Markov Decision Processes (MDPs):



- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)

# Finite Horizon Markov Decision Processes (MDPs):
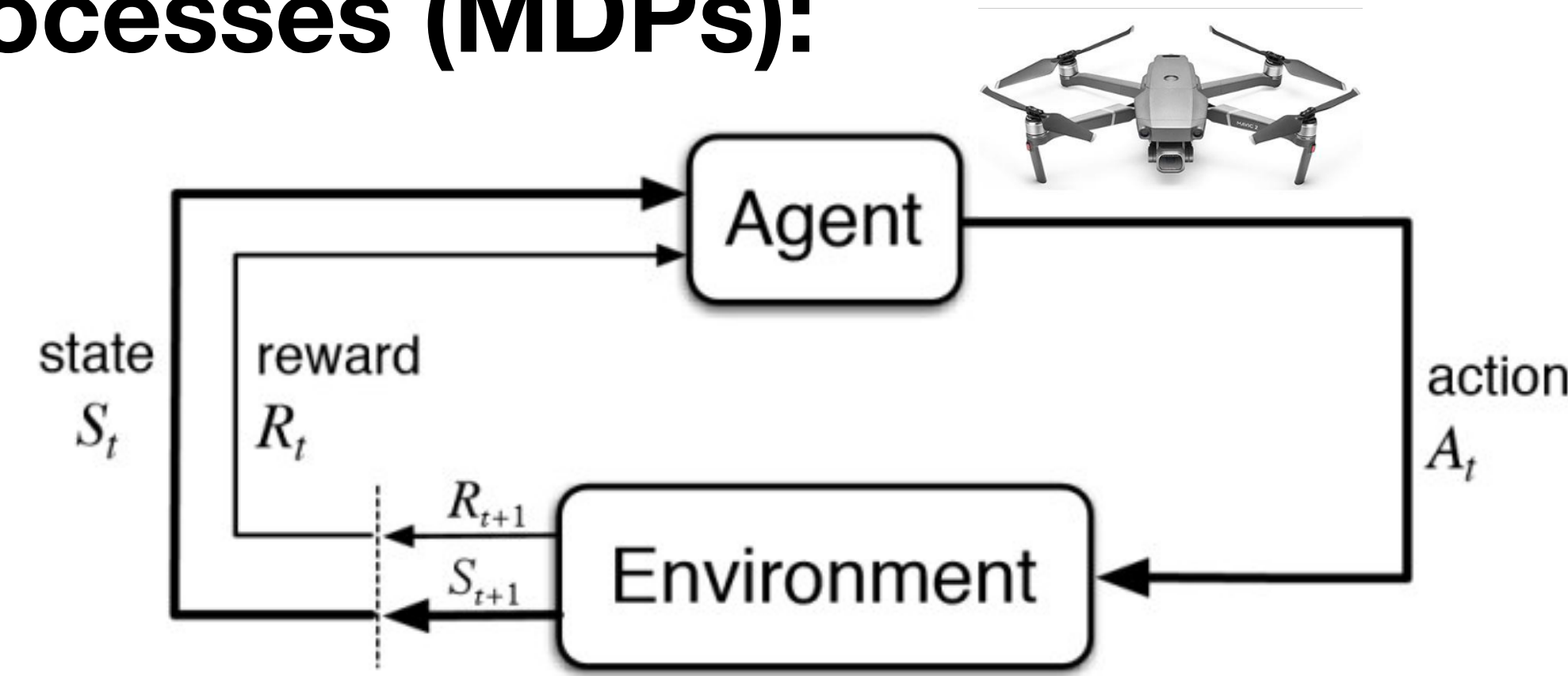


- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)
  - $S$ a set of states

state
$S_t$

reward
$R_t$

action
$A_t$

$R_{t+1}$

$S_{t+1}$

Agent

Environment

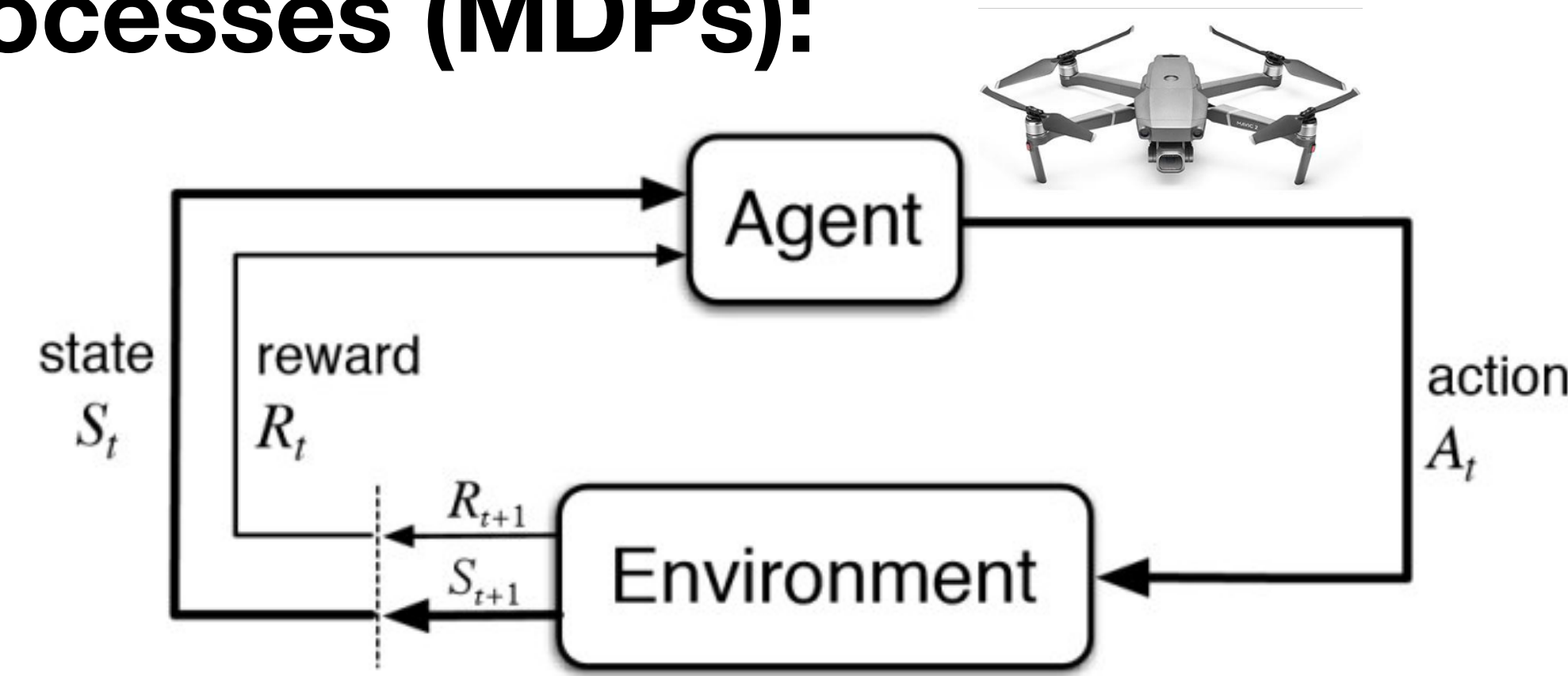## Finite Horizon Markov Decision Processes (MDPs):

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)
  - $S$ a set of states
  - $A$ a set of actions

# Finite Horizon Markov Decision Processes (MDPs):



- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)
  - $S$ 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 state $s$ via action $a$

# Finite Horizon Markov Decision Processes (MDPs):

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
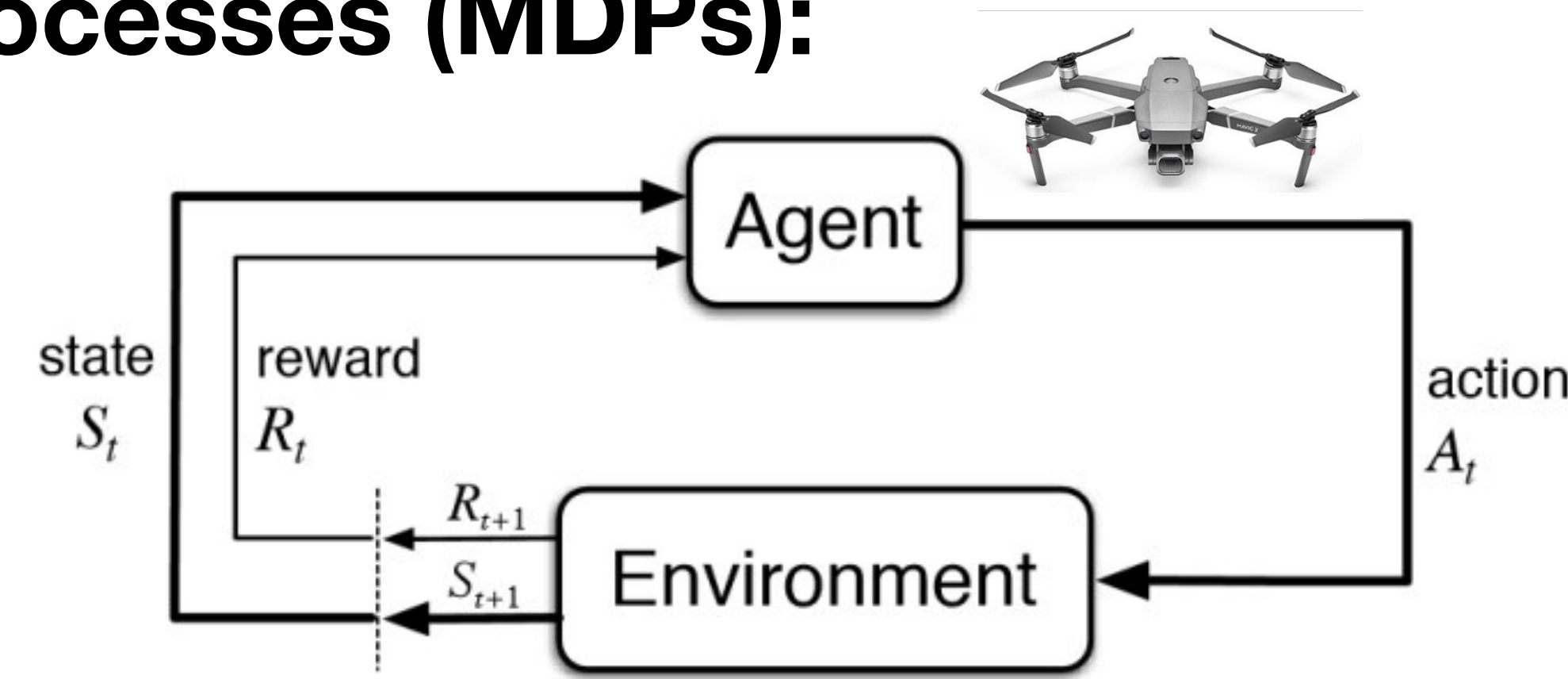  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)
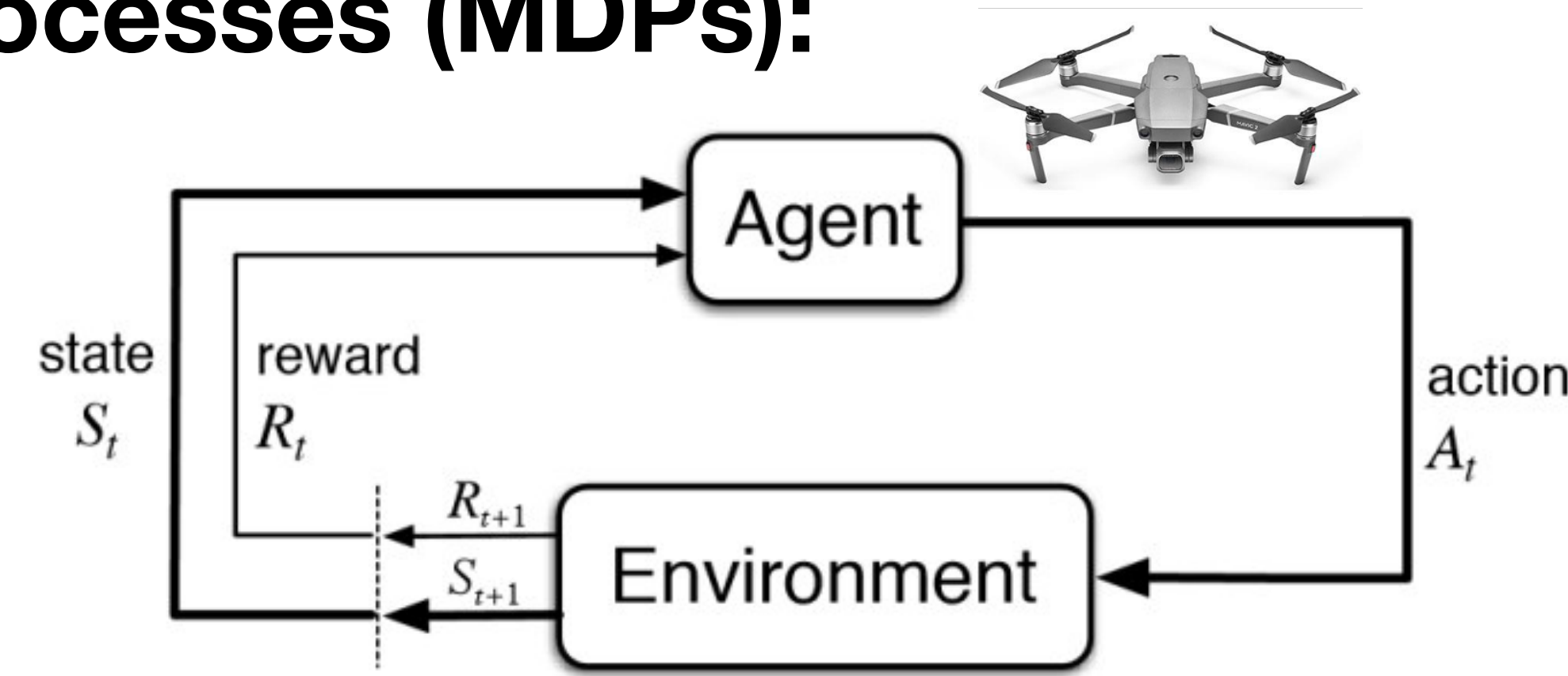  - $S$ a set of states
  - $A$ a set of actions
  - $P : S \times A \mapsto \Delta(S)$ specifies the dynamics model,
    i.e. $P(s'|s, a)$ is the probability of transitioning to $s'$ from state $s$ via action $a$
  - $r : S \times A \to [0,1]$
    - For now, let's assume this is a deterministic function
    - (sometimes we use a cost $c : S \times A \to [0,1]$)

state
$S_t$

reward
$R_t$

action
$A_t$

$R_{t+1}$
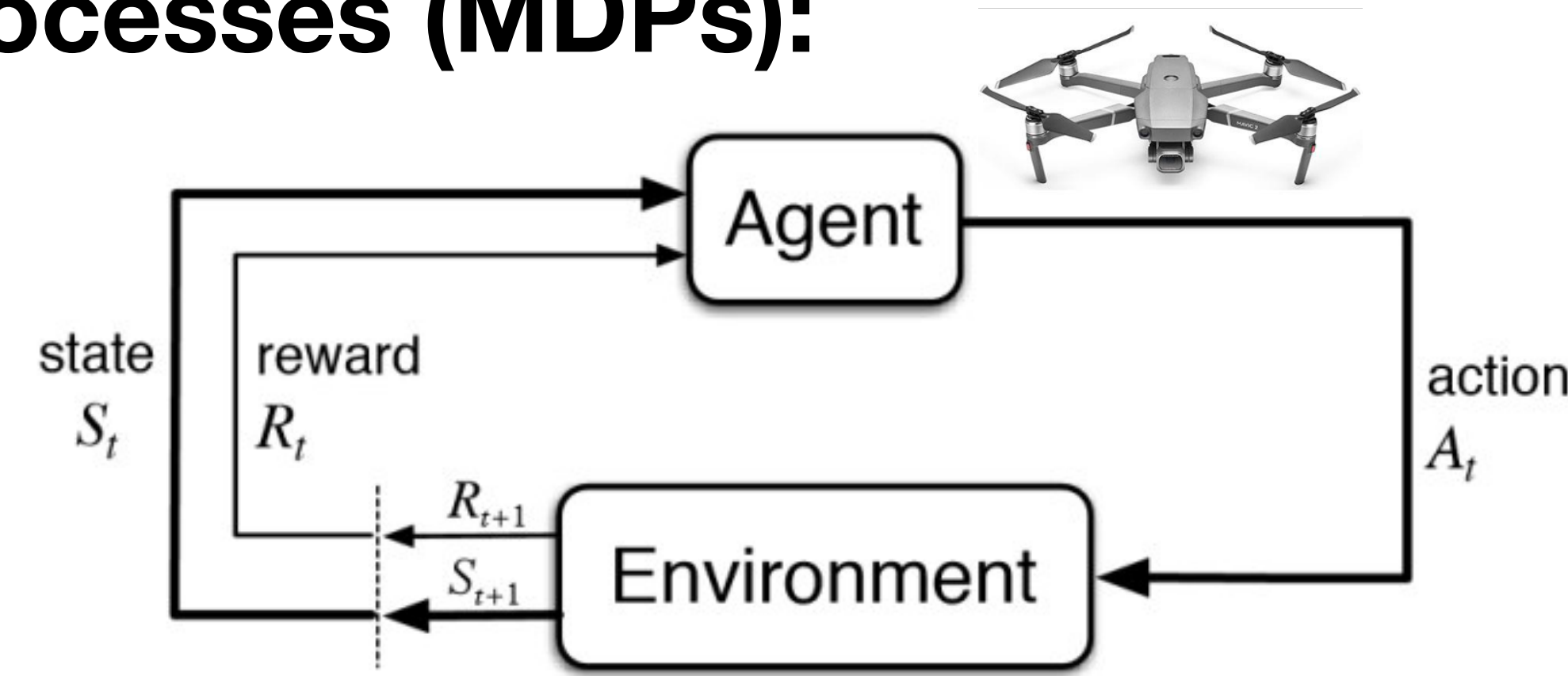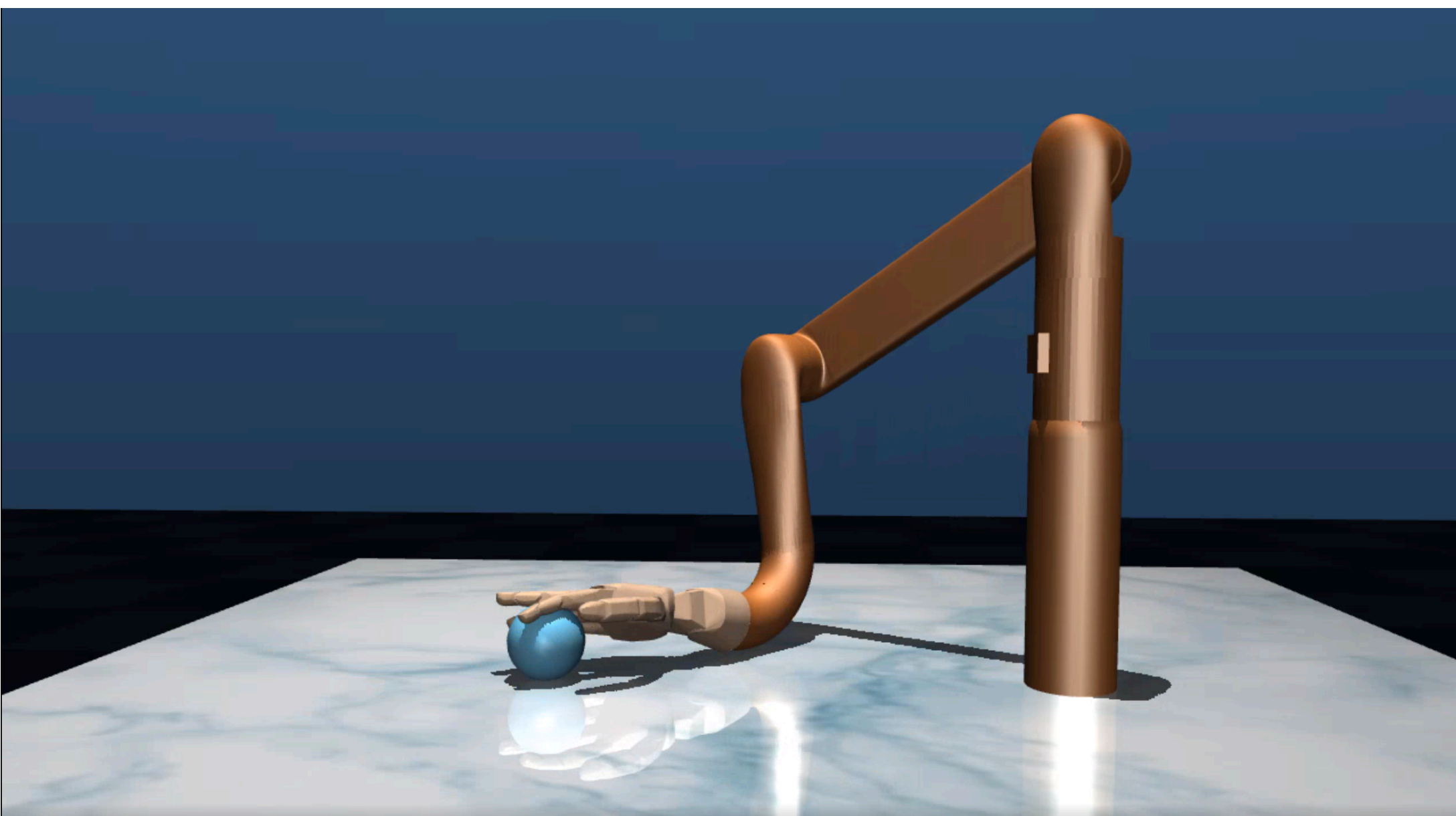
$S_{t+1}$

Agent

Environment

# Finite Horizon Markov Decision Processes (MDPs):



- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
  - $\mu$ is a distribution over initial states
    (sometimes we assume we start a given state $s_0$)
  - $S$ 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 state $s$ via action $a$
  - $r : S \times A \to [0,1]$
    - For now, let's assume this is a deterministic function
    - (sometimes we use a cost $c : S \times A \to [0,1]$)
  - A time horizon $H \in \mathbb{N}$

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

# 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

# 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

# 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

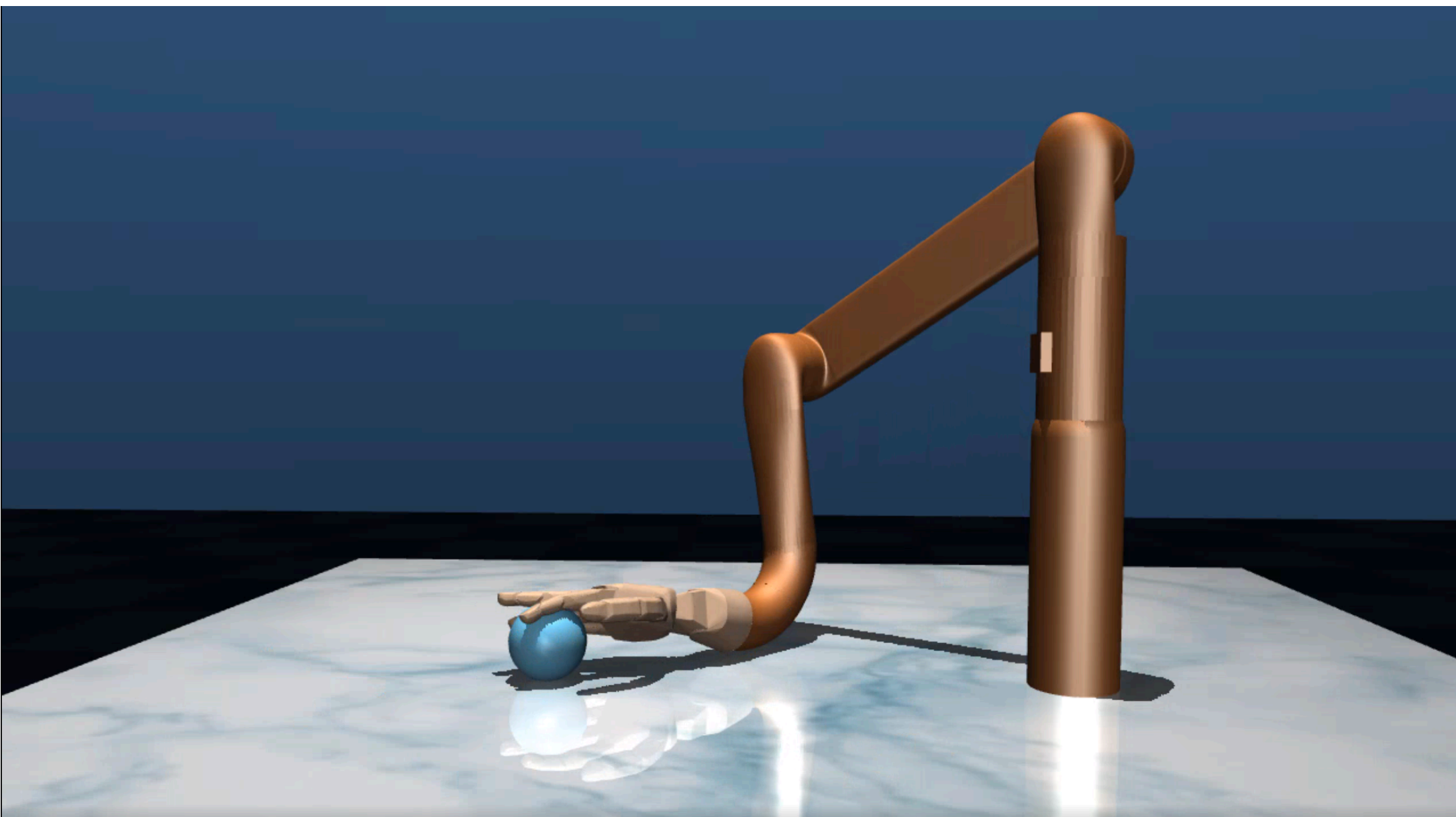**Transition** $s' \sim P(\cdot \mid s, a)$: physics + some noise

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



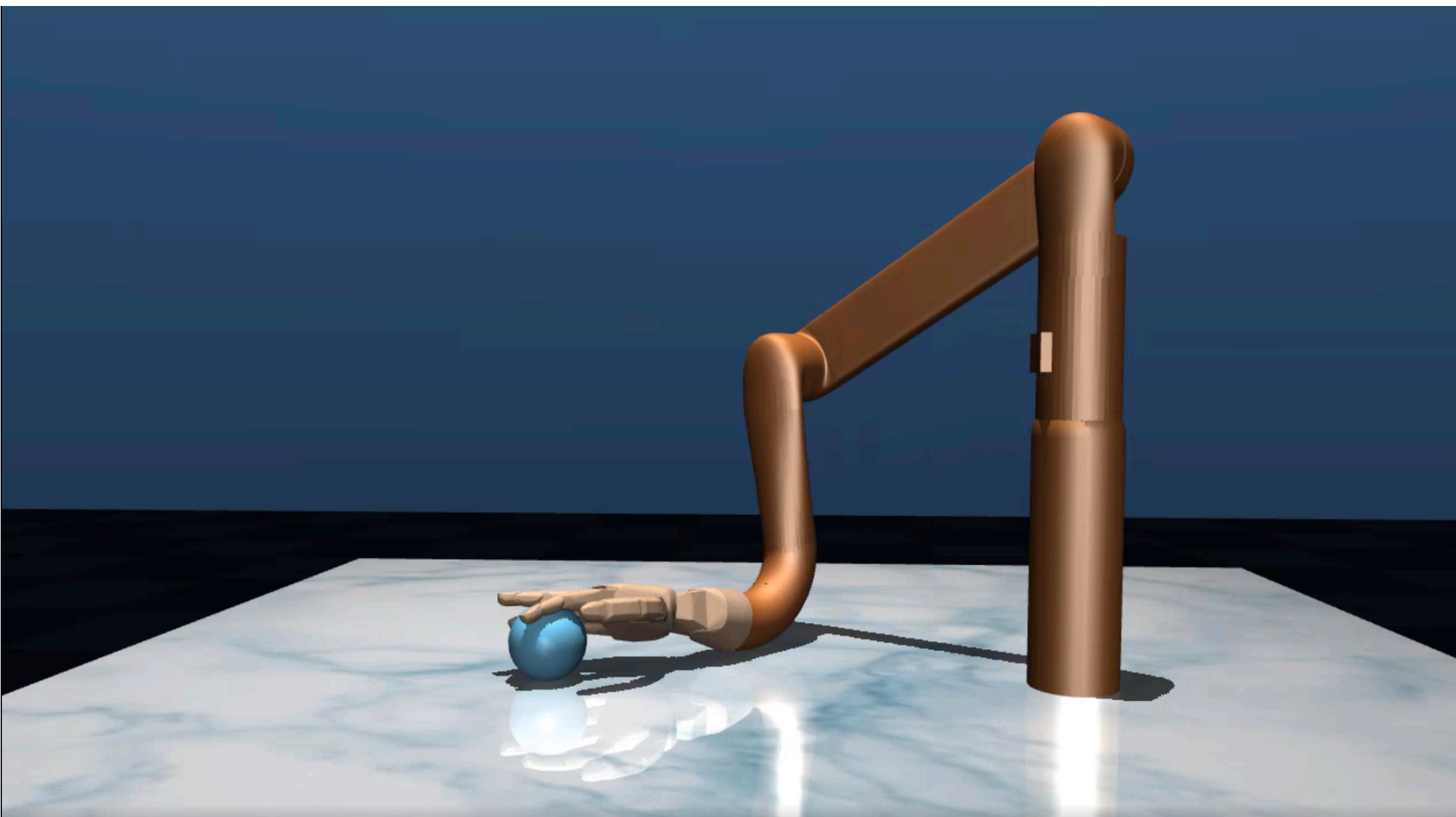**State** $s$: robot configuration (e.g., joint angles) and the ball's position

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

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

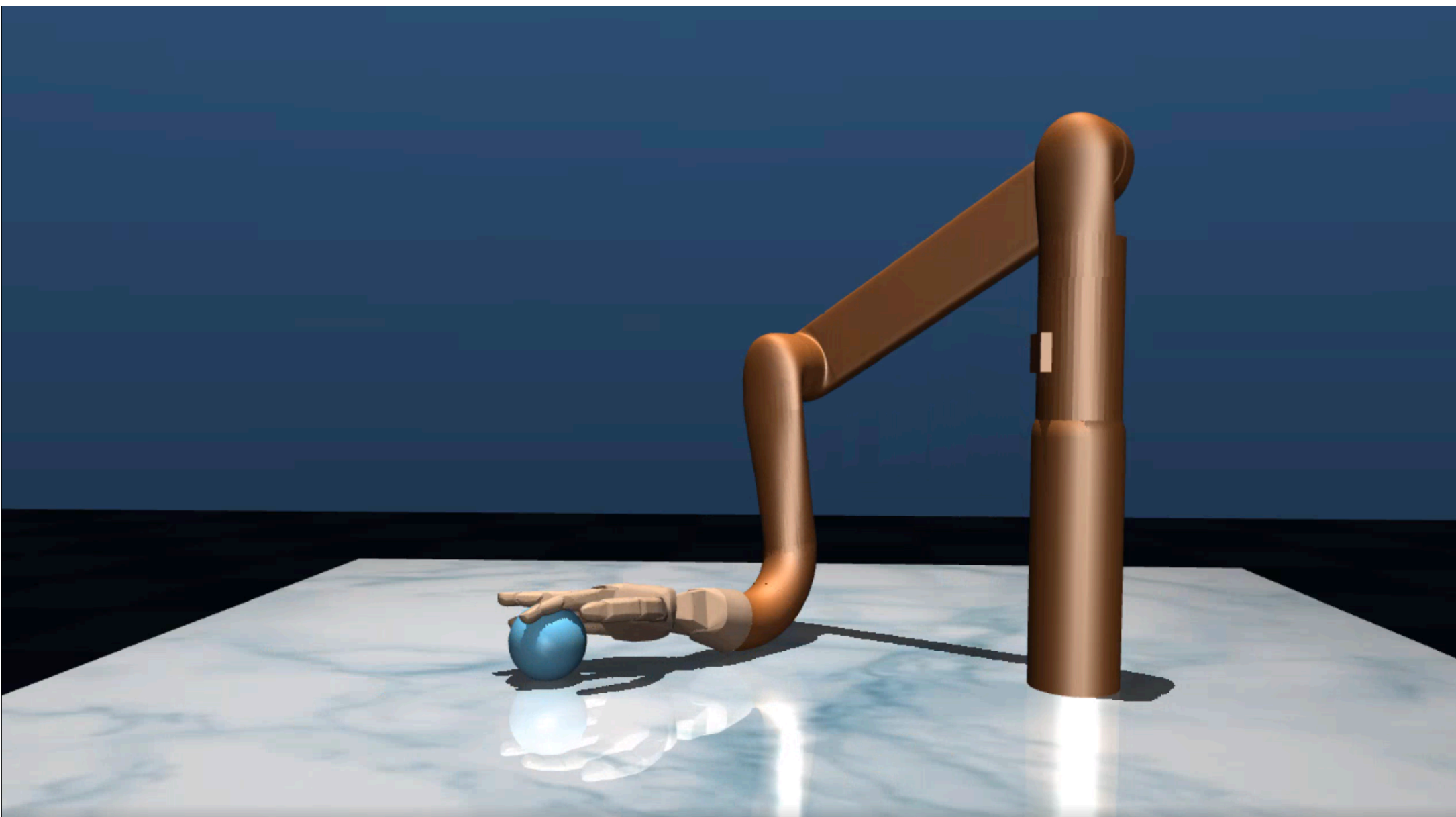**Policy** $\pi(s)$: a function mapping from robot state to action (i.e., torque)

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

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

**Reward/Cost:**

  $r(s, a)$: immediate reward at state $(s, a)$, or

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

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



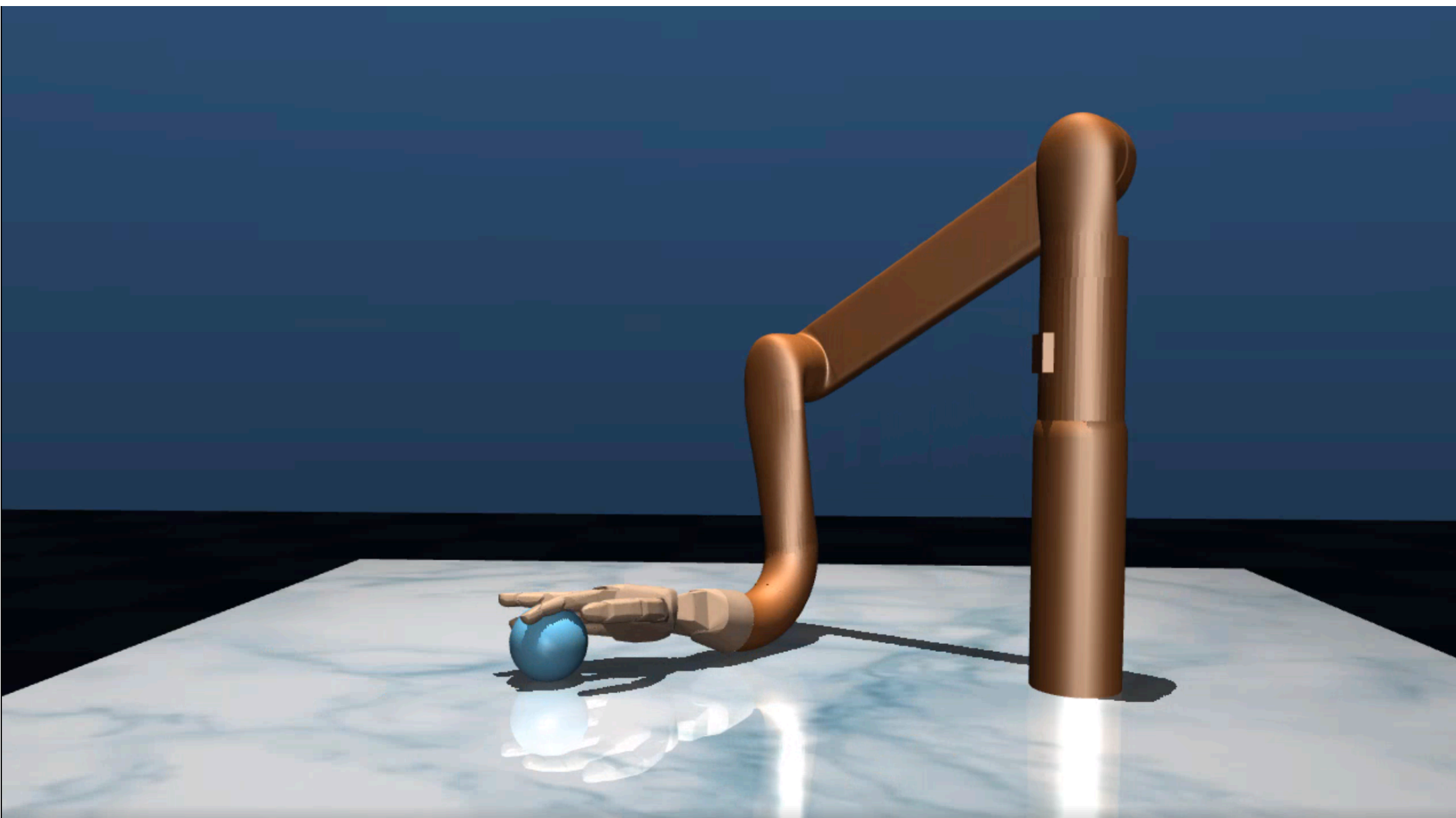**State** $s$: robot configuration (e.g., joint angles) and the ball's position

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

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

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

**Reward/Cost:**

$r(s, a)$: immediate reward at state $(s, a)$, or

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

**Horizon:** timescale $H$

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



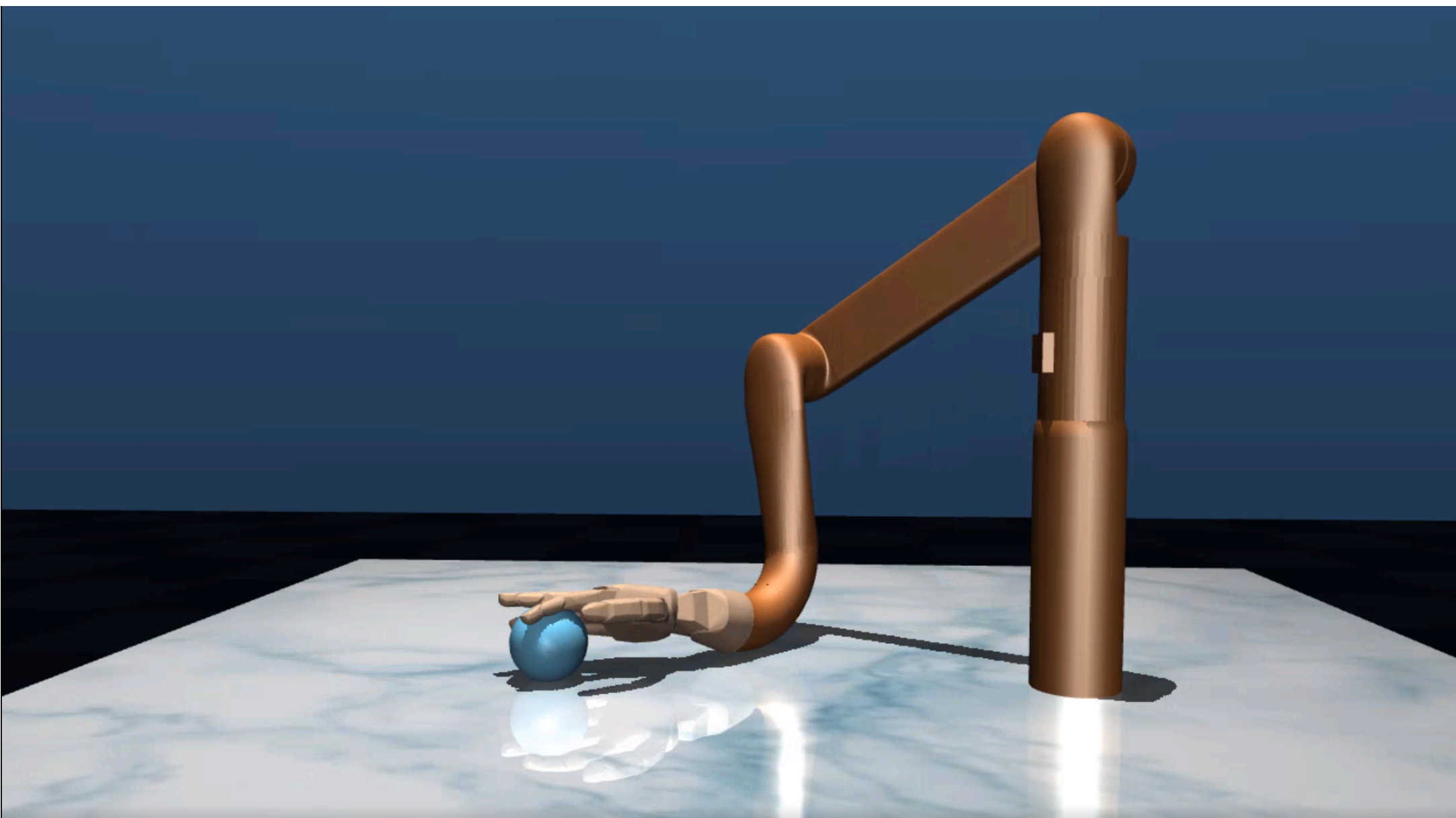**State** $s$: robot configuration (e.g., joint angles) and the ball's position

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

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

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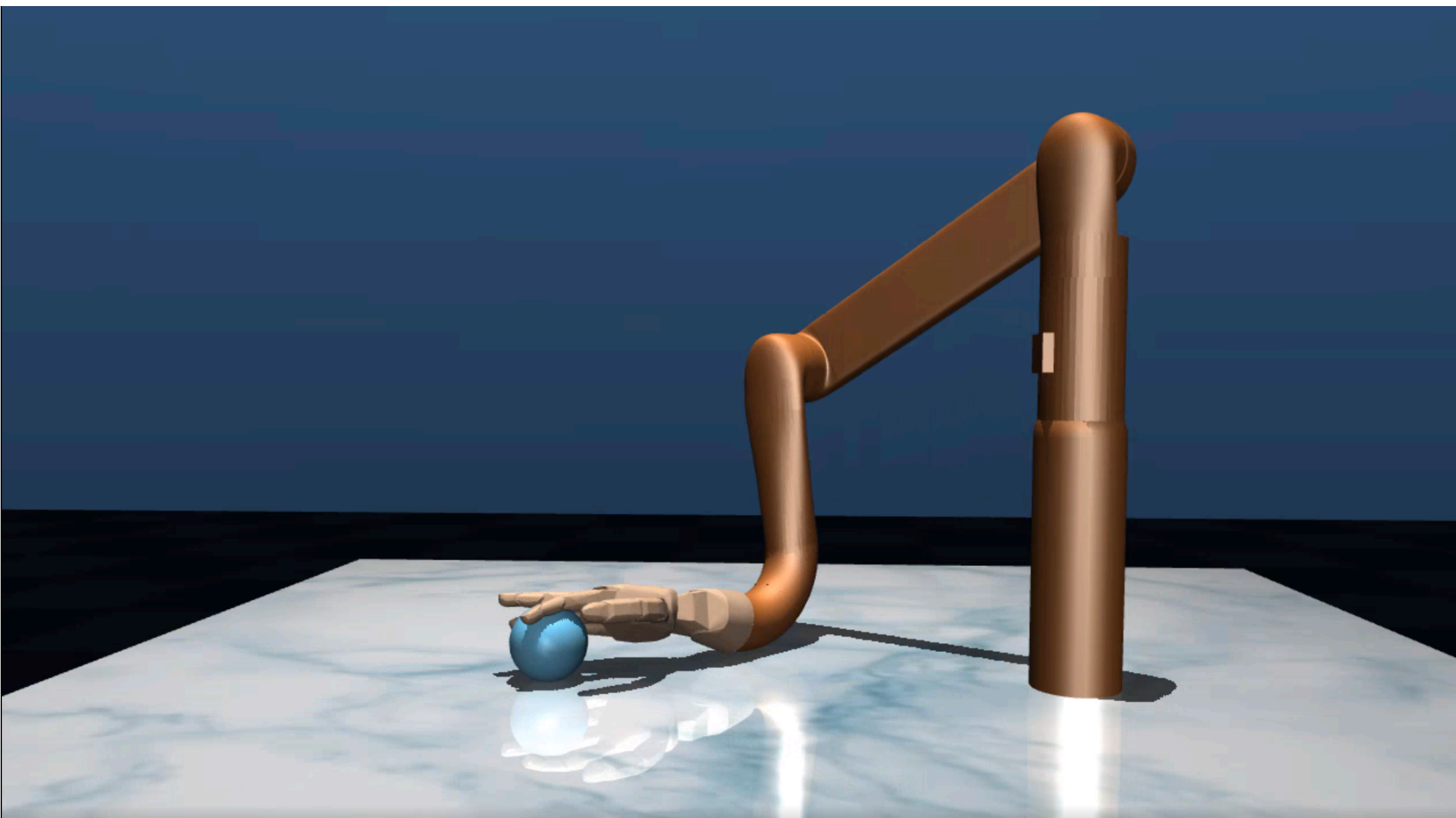
**Reward/Cost:**

$r(s, a)$: immediate reward at state $(s, a)$, or

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

**Horizon:** timescale $H$

$$\pi^\star = \arg\min_\pi \mathbb{E}\left[ c(s_0, a_0) + c(s_1, a_1) + c(s_2, a_2) + \ldots c(s_{H-1}, a_{H-1}) \,\middle|\, s_0, \pi \right]$$

# Today

✓ • Logistics (Welcome!)

✓ • Overview of RL

✓ • Markov Decision Processes

   • Problem statement

   • Policy Evaluation

# The Episodic Setting and Trajectories

# The Episodic Setting and Trajectories

- Policy $\pi := \{\pi_0, \pi_1, \ldots, \pi_{H-1}\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

# The Episodic Setting and Trajectories

- Policy $\pi := \left\{ \pi_0, \pi_1, \ldots, \pi_{H-1} \right\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

# The Episodic Setting and Trajectories

- Policy $\pi := \left\{ \pi_0, \pi_1, \ldots, \pi_{H-1} \right\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

  - Sample an initial state $s_0 \sim \mu$:

# The Episodic Setting and Trajectories

- Policy $\pi := \left\{ \pi_0, \pi_1, \ldots, \pi_{H-1} \right\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

  - Sample an initial state $s_0 \sim \mu$:

  - For $t = 0,1,2,\ldots H - 1$

    - Take action $a_t \sim \pi_t( \cdot \mid s_t)$

# The Episodic Setting and Trajectories

- Policy $\pi := \{\pi_0, \pi_1, \ldots, \pi_{H-1}\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

  - Sample an initial state $s_0 \sim \mu$:
  - For $t = 0,1,2,\ldots H-1$
    - Take action $a_t \sim \pi_t( \cdot \mid s_t)$
    - Observe reward $r_t = r(s_t, a_t)$

## The Episodic Setting and Trajectories

- Policy $\pi := \left\{ \pi_0, \pi_1, \ldots, \pi_{H-1} \right\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

  - Sample an initial state $s_0 \sim \mu$:

  - For $t = 0,1,2,\ldots H - 1$

    - Take action $a_t \sim \pi_t( \cdot \mid s_t)$
    - Observe reward $r_t = r(s_t, a_t)$
    - Transition to (and observe) $s_{t+1}$ where $s_{t+1} \sim P( \cdot \mid s_t, a_t)$

16

# The Episodic Setting and Trajectories

- Policy $\pi := \left\{ \pi_0, \pi_1, \ldots, \pi_{H-1} \right\}$

  - deterministic policies: $\pi_t : S \mapsto A$; stochastic policies: $\pi_t : S \mapsto \Delta(A)$
  - we also consider time-dependent policies (but not a function of the history)

- Sampling a trajectory $\tau$ on an episode: for a given policy $\pi$

  - Sample an initial state $s_0 \sim \mu$:
  - For $t = 0,1,2,\ldots H-1$
    - Take action $a_t \sim \pi_t(\,\cdot\,|\,s_t)$
    - Observe reward $r_t = r(s_t, a_t)$
    - Transition to (and observe) $s_{t+1}$ where $s_{t+1} \sim P(\,\cdot\,|\,s_t, a_t)$
  - The sampled trajectory is $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$

# The Probability of a Trajectory & The Objective

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory

$\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.
  - Shorthand: we sometimes write $\rho$ or $\rho_\pi$ when $\pi$ and/or $\mu$ are clear from context.

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.
  - Shorthand: we sometimes write $\rho$ or $\rho_\pi$ when $\pi$ and/or $\mu$ are clear from context.
  - The rewards in this trajectory must be $r_t = r(s_t, a_t)$ (else $\rho_\pi(\tau) = 0$).

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory

  $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.

  - Shorthand: we sometimes write $\rho$ or $\rho_\pi$ when $\pi$ and/or $\mu$ are clear from context.

  - The rewards in this trajectory must be $r_t = r(s_t, a_t)$ (else $\rho_\pi(\tau) = 0$).

  - For $\pi$ stochastic:

    $\rho_\pi(\tau) = \mu(s_0)\pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\ldots\pi(a_{H-2} \,|\, s_{H-2})P(s_{H-1} \,|\, s_{H-2}, a_{H-2})\pi(a_{H-1} \,|\, s_{H-1})$

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory
  $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.
  - Shorthand: we sometimes write $\rho$ or $\rho_\pi$ when $\pi$ and/or $\mu$ are clear from context.
  - The rewards in this trajectory must be $r_t = r(s_t, a_t)$ (else $\rho_\pi(\tau) = 0$).
  - For $\pi$ stochastic:
    $\rho_\pi(\tau) = \mu(s_0)\pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\ldots\pi(a_{H-2} \,|\, s_{H-2})P(s_{H-1} \,|\, s_{H-2}, a_{H-2})\pi(a_{H-1} \,|\, s_{H-1})$
  - For $\pi$ deterministic:
    $\rho_\pi(\tau) = \mu(s_0)\mathbf{1}\big(a_0 = \pi(s_0)\big)P(s_1 \,|\, s_0, a_0)\ldots P(s_{H-1} \,|\, s_{H-2}, a_{H-2})\mathbf{1}\big(a_{H-1} = \pi(s_{H-1})\big)$

# The Probability of a Trajectory & The Objective

- Probability of trajectory: let $\rho_{\pi,\mu}(\tau)$ denote the probability of observing trajectory

  $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \ldots, s_{H-1}, a_{H-1}, r_{H-1}\}$ when acting under $\pi$ with $s_0 \sim \mu$.

  - Shorthand: we sometimes write $\rho$ or $\rho_\pi$ when $\pi$ and/or $\mu$ are clear from context.
  - The rewards in this trajectory must be $r_t = r(s_t, a_t)$ (else $\rho_\pi(\tau) = 0$).
  - For $\pi$ stochastic:
    $\rho_\pi(\tau) = \mu(s_0)\pi(a_0 \,|\, s_0)P(s_1 \,|\, s_0, a_0)\ldots\pi(a_{H-2} \,|\, s_{H-2})P(s_{H-1} \,|\, s_{H-2}, a_{H-2})\pi(a_{H-1} \,|\, s_{H-1})$
  - For $\pi$ deterministic:
    $\rho_\pi(\tau) = \mu(s_0)\mathbf{1}\big(a_0 = \pi(s_0)\big)P(s_1 \,|\, s_0, a_0)\ldots P(s_{H-1} \,|\, s_{H-2}, a_{H-2})\mathbf{1}\big(a_{H-1} = \pi(s_{H-1})\big)$

- Objective: find policy $\pi$ that maximizes our expected cumulative episodic reward:
  $$\max_\pi \mathbb{E}_{\tau \sim \rho_\pi}\Big[r(s_0, a_0) + r(s_1, a_1) + \ldots + r(s_{H-1}, a_{H-1})\Big]$$

17

# Today

✓ • Logistics (Welcome!)

✓ • Overview of RL

✓ • Markov Decision Processes

✓     • Problem statement

    • Policy Evaluation

# Policy Evaluation = Computing Value function and/or Q function

# Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

# Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

- Value function $V_h^\pi(s) = \mathbb{E}\left[\sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, s_h = s\right]$

# Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

- Value function $V_h^\pi(s) = \mathbb{E}\left[\sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, s_h = s\right]$

- Q function $Q_h^\pi(s, a) = \mathbb{E}\left[\sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, (s_h, a_h) = (s, a)\right]$

# Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

- Value function $V_h^\pi(s) = \mathbb{E}\left[\sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, s_h = s\right]$

- Q function $Q_h^\pi(s, a) = \mathbb{E}\left[\sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, (s_h, a_h) = (s, a)\right]$

- At the last stage, what are:

$$Q_{H-1}^\pi(s, a) = \qquad\qquad\qquad V_{H-1}^\pi(s) =$$

# Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

- Value function $V_h^\pi(s) = \mathbb{E}\left[ \sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, s_h = s \right]$

- Q function $Q_h^\pi(s, a) = \mathbb{E}\left[ \sum_{t=h}^{H-1} r(s_t, a_t) \,\middle|\, (s_h, a_h) = (s, a) \right]$
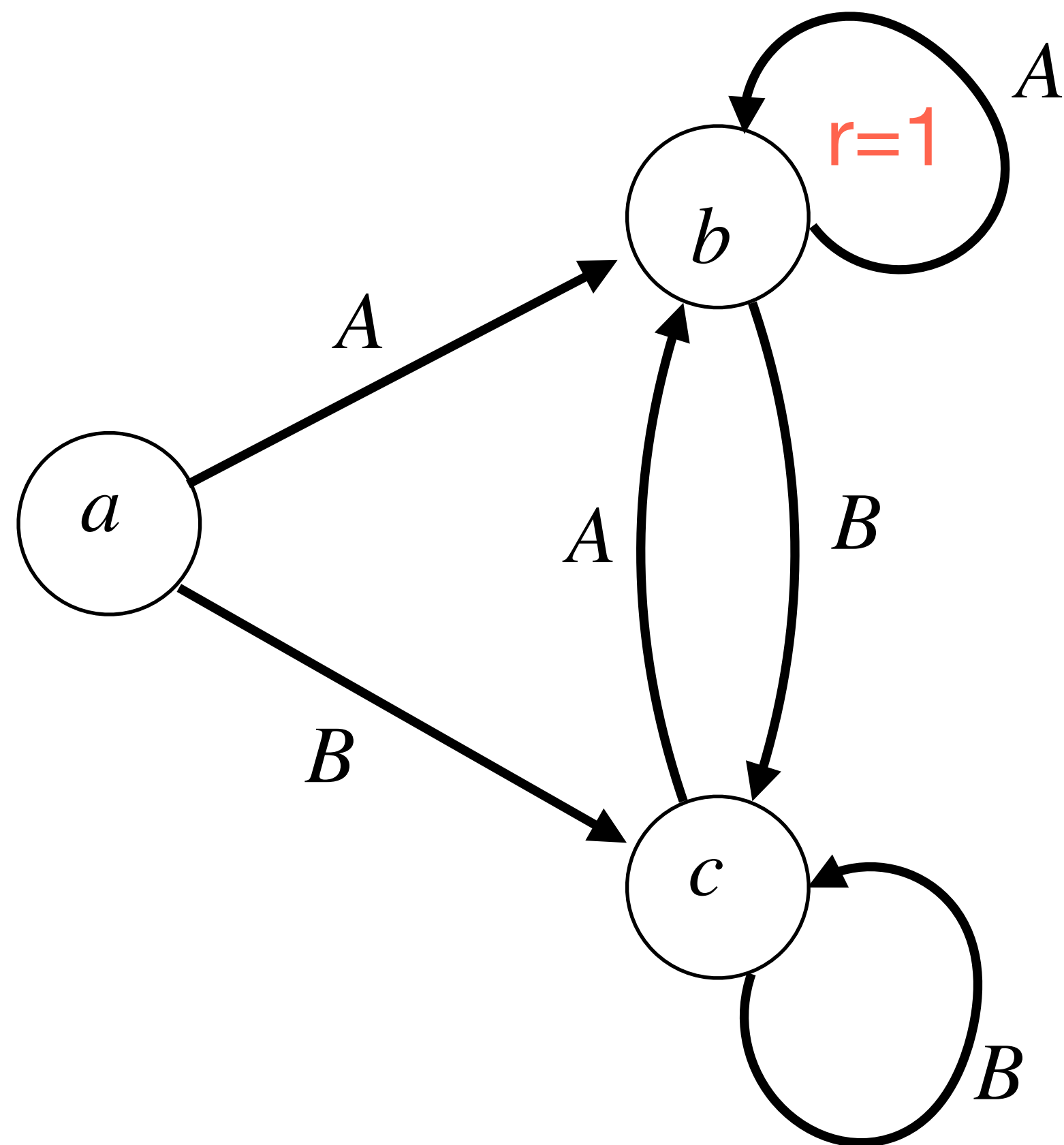
- At the last stage, for a stochastic policy,:

$$Q_{H-1}^\pi(s, a) = r(s, a) \qquad\qquad V_{H-1}^\pi(s) = \sum_a \pi_{H-1}(a \,|\, s) r(s, a)$$

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



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

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

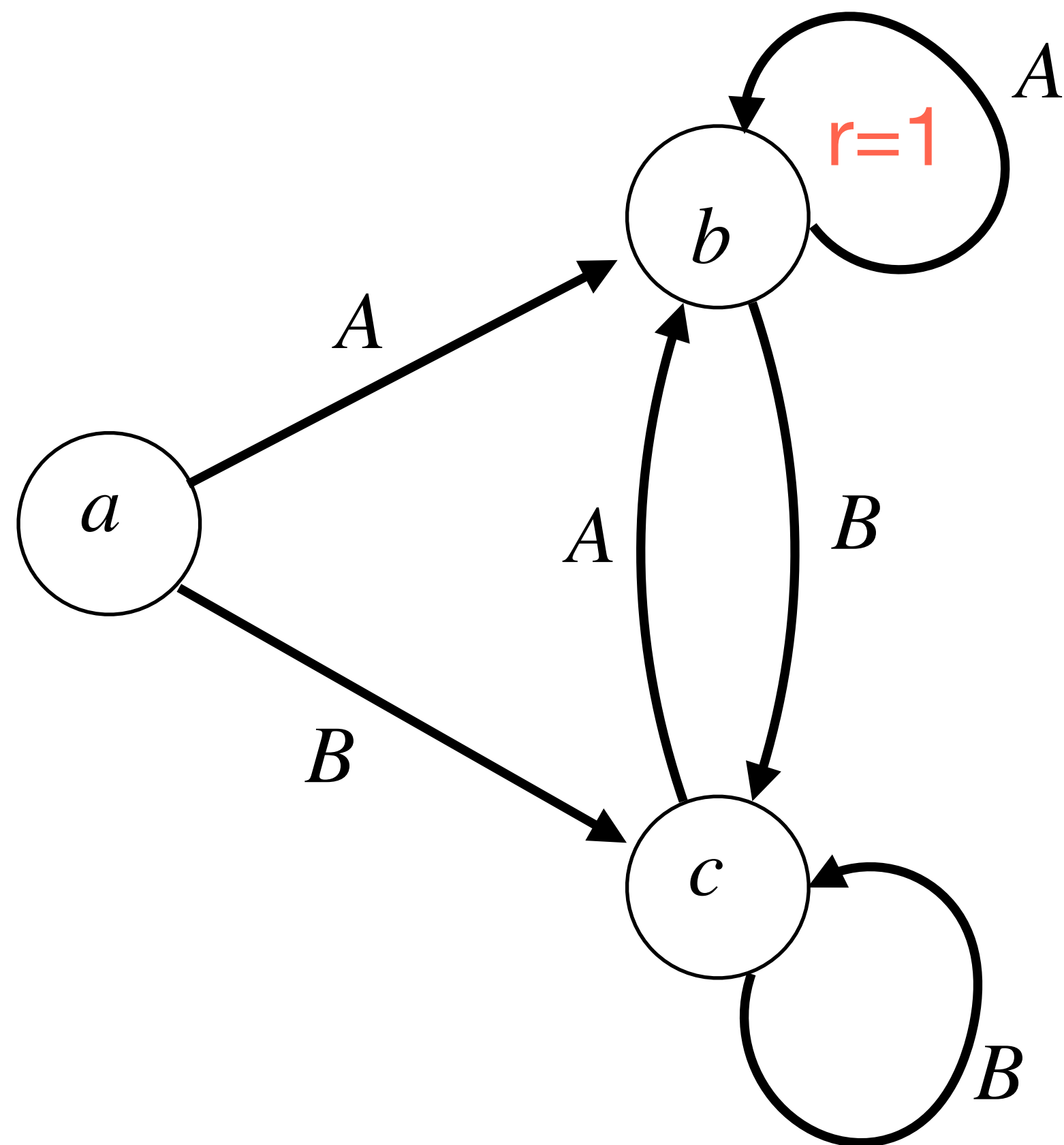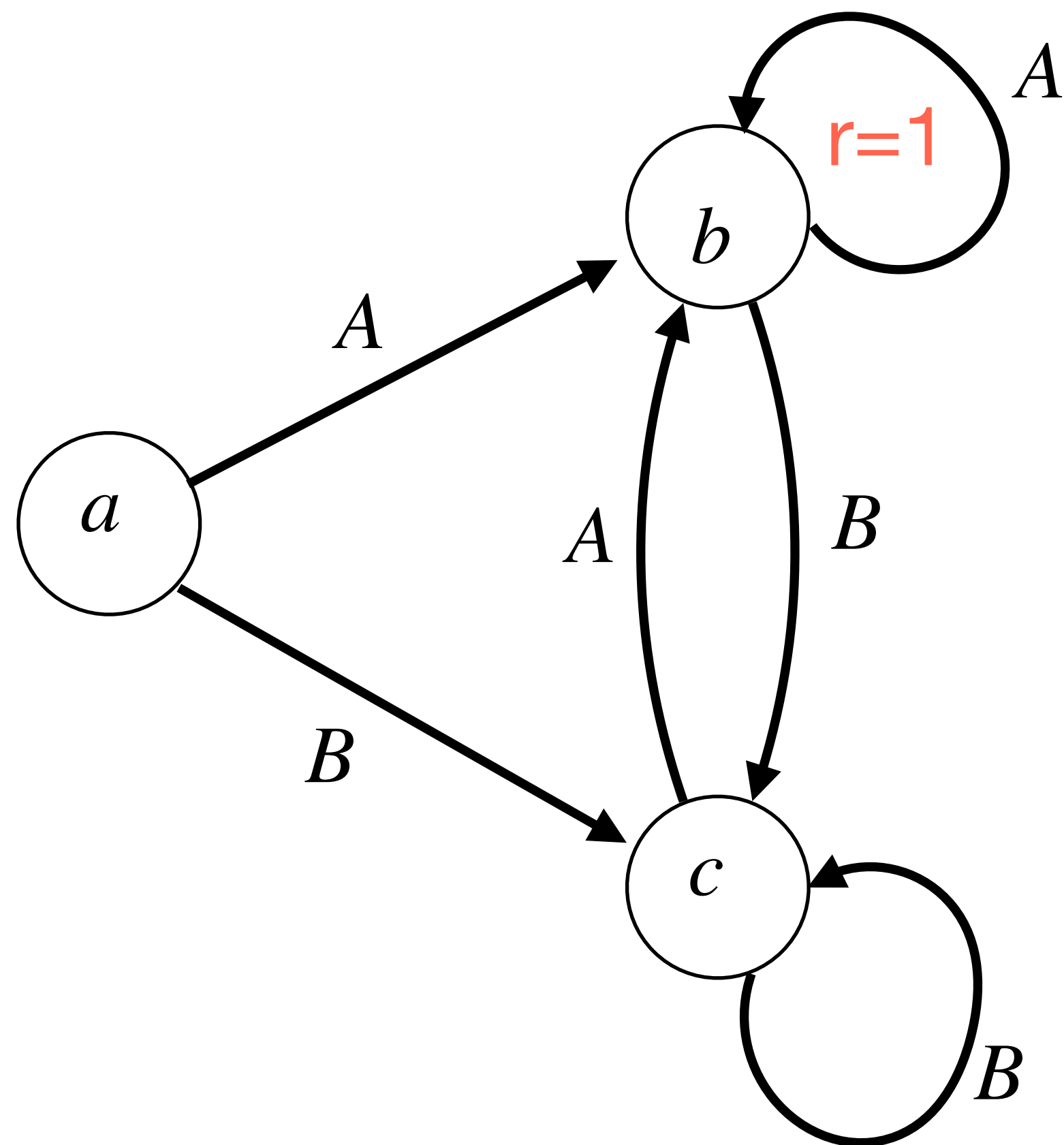Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



- Consider the deterministic policy
  $\pi_0(s) = A, \pi_1(s) = A, \pi_2(s) = B, \forall s$

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

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

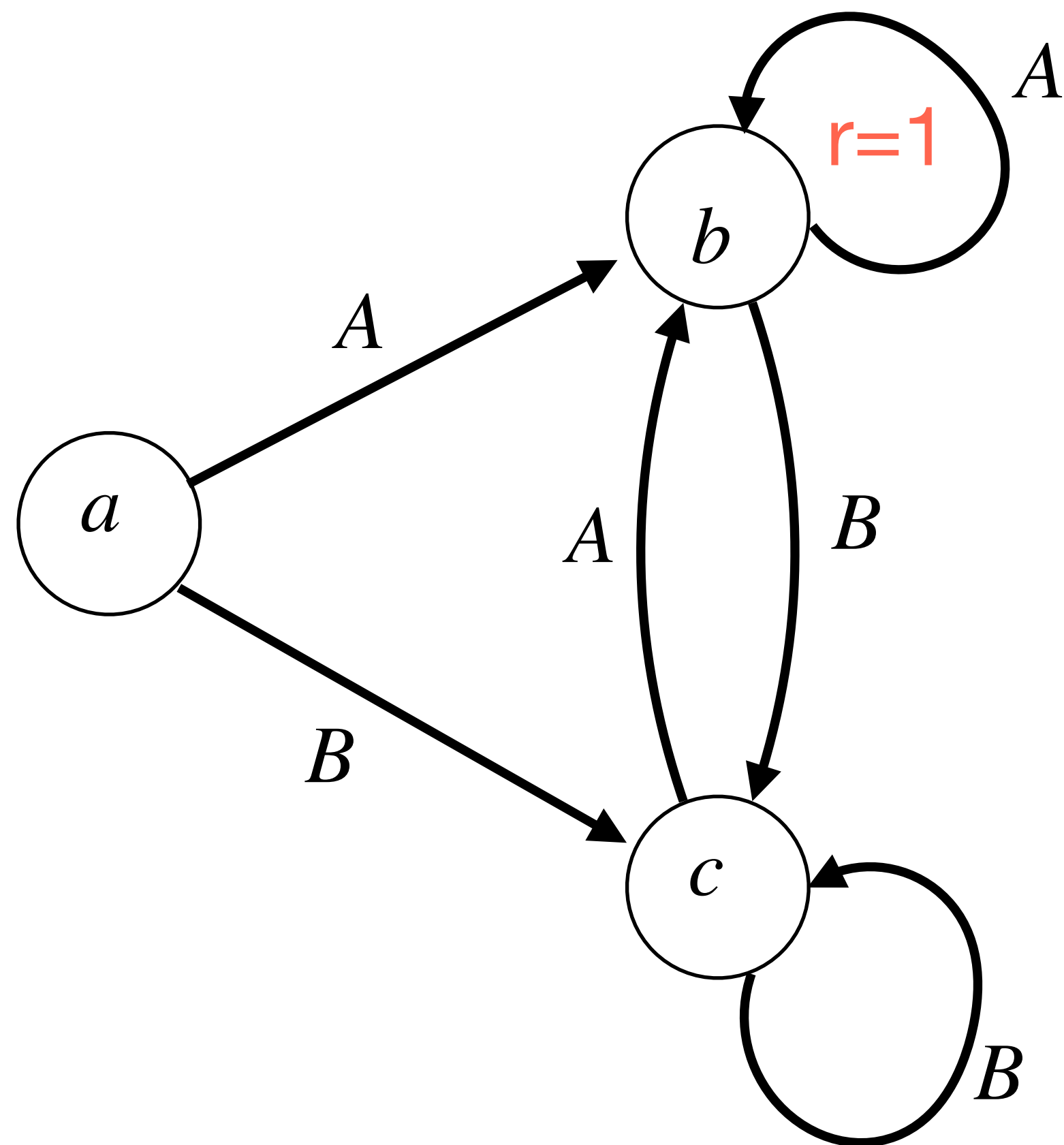Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



- Consider the deterministic policy
  $\pi_0(s) = A, \pi_1(s) = A, \pi_2(s) = B, \forall s$

- What is $V^\pi$?

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

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



r=1

- Consider the deterministic policy
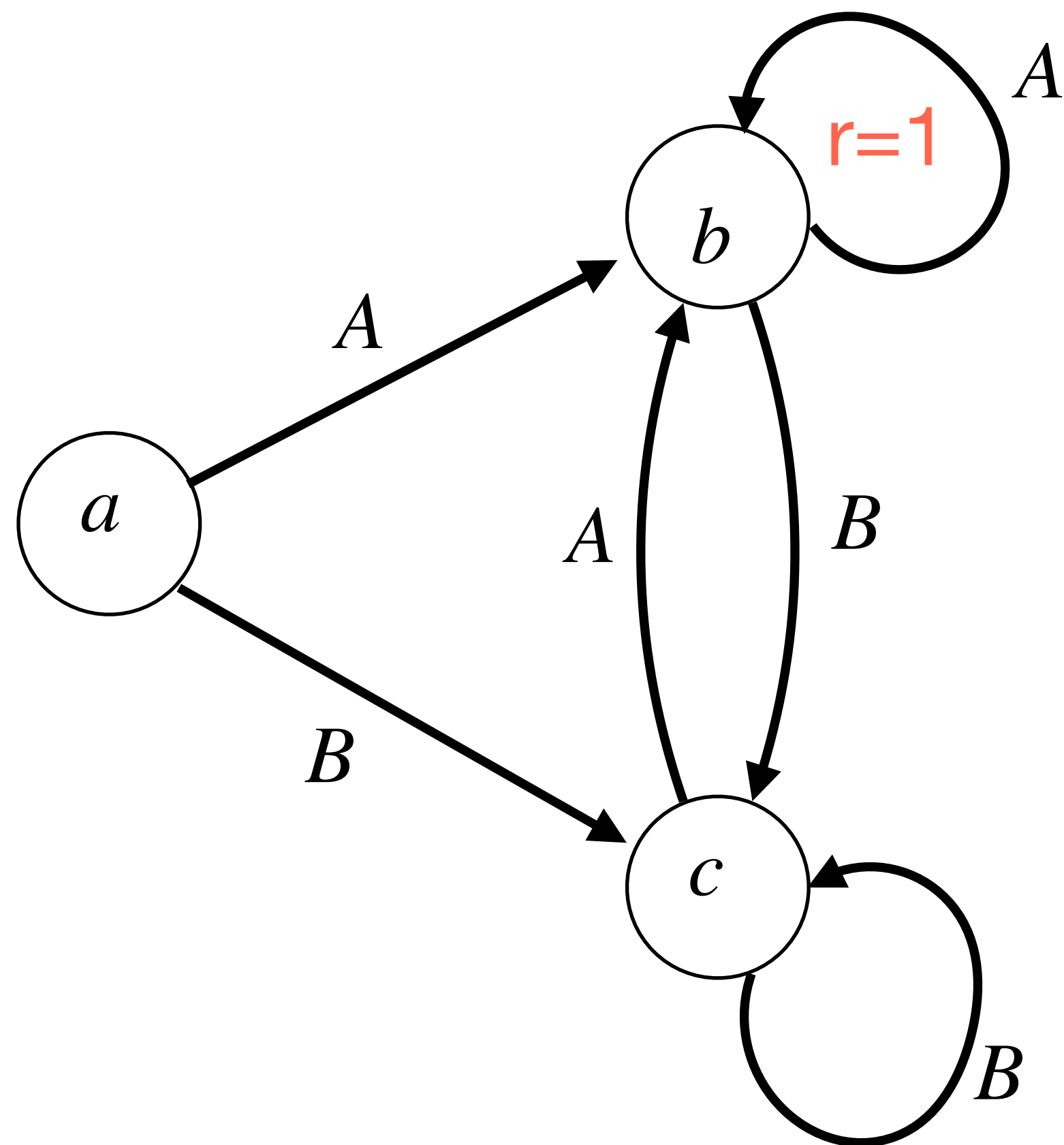  $$\pi_0(s) = A, \pi_1(s) = A, \pi_2(s) = B, \forall s$$

- What is $V^\pi$?
  $$V_2^\pi(a) = 0, V_2^\pi(b) = 0, V_2^\pi(c) = 0$$

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

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



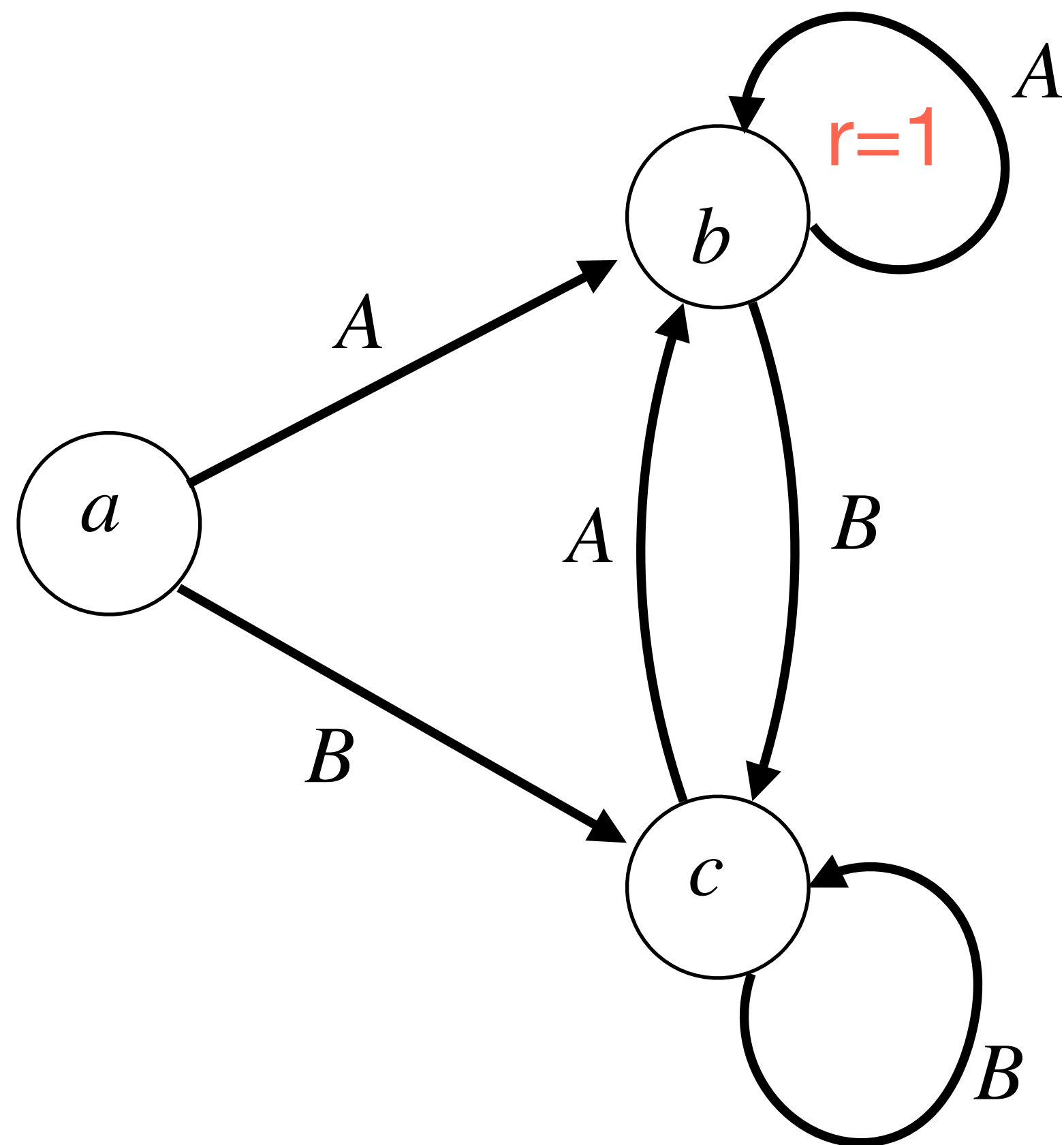- Consider the deterministic policy
$\pi_0(s) = A, \pi_1(s) = A, \pi_2(s) = B, \forall s$

- What is $V^\pi$?
$V_2^\pi(a) = 0, V_2^\pi(b) = 0, V_2^\pi(c) = 0$

$V_1^\pi(a) = 0, V_1^\pi(b) = 1, V_1^\pi(c) = 0$

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

# Example of Policy Evaluation (i.e. computing $V^\pi$ and $Q^\pi$)

Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



- Consider the deterministic policy
  $$\pi_0(s) = A, \pi_1(s) = A, \pi_2(s) = B, \forall s$$

- What is $V^\pi$?
  $$V_2^\pi(a) = 0, V_2^\pi(b) = 0, V_2^\pi(c) = 0$$

  $$V_1^\pi(a) = 0, V_1^\pi(b) = 1, V_1^\pi(c) = 0$$

  $$V_0^\pi(a) = 1, V_0^\pi(b) = 2, V_0^\pi(c) = 1$$

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

# Today

✓ • Logistics (Welcome!)

✓ • Overview of RL

✓ • Markov Decision Processes

   ✓ • Problem statement

   ✓ • Policy Evaluation

# Summary:

- **Finite horizon MDPs (a framework for RL):**
- Key concepts: **sampling a trajectory** $\rho_\pi(\tau)$, **V and Q functions**

Attendance:
bit.ly/3RcTC9T



Attendance Password:

Feedback:
bit.ly/3RHtlxy