# Offline Reinforcement Learning
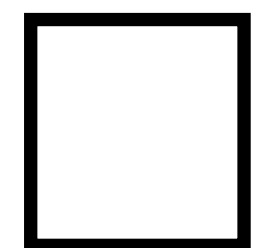
Sanjiban Choudhury

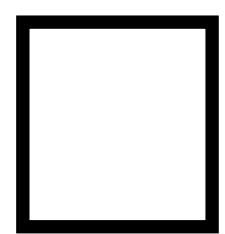# The story thus far ...

✅ Decision-making

✅ Perception
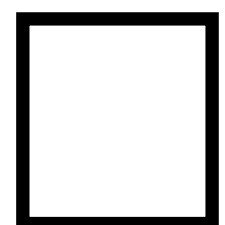
✅ Models of humans

☐ Practical Robot Learning

Today-> ☐ Offline RL

☐ Sim-to-Real

# Today's class

□ What is offline RL? Why do we need it for robots?

□ Paradigm 1: Offline RL via Pessimism
  □ Problem with Q-learning
  □ Pessimism to the rescue

□ Paradigm 2: RL via Supervised Learning
  □ Return-conditioned Supervised Learning
  □ Problem in Stochastic MDPs

# Why do we need offline RL for robots?

# Robots today still only work in CLOSED world



The Dream

Reality

# Generalize to variations of the OPEN world?

# Why can't we do RL with robots in the real world?

# Machine learning's answer!



Big Data                                    Big Models

# Efforts underway to scale up robotics data!

1M trajectories, 22 robots, 21 different institutions



(a) # Datasets per Robot Embodiment

(b) # Scenes per Embodiment

(c) # Trajectories per Embodiment

(d) Common Dataset Skills

(e) Common Dataset Objects

**Open-X Embodiment Dataset**

# Hope: Data grows logarithmically with tasks



On the quest for shared priors
w/ machine learning

Data

\# Tasks

Interact with the **physical** world to learn bottom-up commonsense

↑

i.e. "how the world works"

Credit: **Andy Zeng**

# Reality: Data grows linearly with tasks



On the quest for shared priors
w/ machine learning

Data

Reality

Expectation

Complexity in environment, embodiment, contact, etc.
↓

\# Tasks

Interact with the **physical** world to learn bottom-up commonsense
↑
i.e. "how the world works"

But for today, let's pretend we can collect a
ton of data
that "covers" tasks we care about

How can we learn *optimal* from large data collected by *any* policy?

# Goal: Offline Reinforcement Learning

big datasets
from past
interaction

train for
**many** epochs

occasionally
get more data

Credit: **Sergey Levine** "Offline RL lecture"

# Different paradigms of RL

## on-policy RL

# Different paragms of RL

## off-policy RL

# offline reinforcement learning



## on-policy RL

## off-policy RL

data collected **once** with **any** policy

training phase

$\{(s_i, a_i, s_i', r_i)\}$

rollout(s)   buffer $\mathcal{D}$   learn $\pi$   deployment

Offline RL enables robots to learn:
from pre-collected datasets
without real-time interaction,
enabling safer training
and leveraging diverse experiences.

# Today's class

☑ What is offline RL? Why do we need it for robots?

(Enables safer training, leverages diverse experience)

☐ Paradigm 1: Offline RL via Pessimism
  ☐ Problem with Q-learning
  ☐ Pessimism to the rescue

☐ Paradigm 2: RL via Supervised Learning
  ☐ Return-conditioned Supervised Learning
  ☐ Problem in Stochastic MDPs

Let's begin with a simple "offline" RL algorithm

We have already covered
a fundamental algorithm
in class that can learn
from offline data.

What is it?

# Q-learning



dataset of transitions ("replay buffer")

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

off-policy Q-learning

$\pi(\mathbf{a}|\mathbf{s})$ (e.g., $\epsilon$-greedy)

For every $\left(s_t, a_t, r_t, s_{t+1}\right)$

Can learn from any data!

$$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(r(s_t, a_t) + \gamma \max_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$$

# Fitted Q-Iteration

*Fitted* Q-iteration



**Given** $\{s_i, a_i, r_i, s'_i\}_{i=1}^N$

**Init** $Q_\theta(s, a) \leftarrow 0$

**while** *not converged* **do**

$\quad D \leftarrow \varnothing$

$\quad$ **for** $i \in 1, \ldots, N$ $\qquad$ *Use old copy of Q*

$\qquad$ input $\leftarrow \{s_i, a_i\}$ $\qquad$ *to set target*

$\qquad$ target $\leftarrow r_i + \gamma \max_{a'} Q_\theta(s'_i, a')$

$\qquad D \leftarrow D \cup \{\text{input, target}\}$

$\quad Q_\theta \leftarrow \textbf{Train}(D)$

**return** $Q_\theta$

*Training is a regression problem*

$$\ell(\theta) = \sum_{i=1}^N \left(Q_\theta(s_i, a_i) - target\right)^2$$

23

Activity!

# Consider the following MDP

# Let's say I collected some data from the MDP

# What policy would Q-learning pick?



S0

L          R

S1                    S2

L        R        L        R

50%   -1   -10   -10   -2   50%

Assume we are in tabular case

Initialize Q values with 0's

For every $(s_t, a_t, r_t, s_{t+1})$

$Q^*(s_t, a_t) = Q^*(s_t, a_t) + \alpha(r(s_t, a_t) + \gamma \max_{a'} Q^*(s_{t+1}, a') - Q^*(s_t, a_t))$

# Think-Pair-Share!

Think (30 sec): What policy would Q-learning
pick in the tabular setting? Why? Ideas to fix it?

Pair: Find a partner

Share (45 sec): Partners exchange
                ideas

# The Problem with Q-learning

**Fundamental problem:** counterfactual queries

**Training data**

**What the policy wants to do**

Is this good? Bad?

How do we know if we didn't see it in the data?

<span style="color:red">Q-learning can be incorrectly optimistic about actions it has not see in the data</span>

# Today's class

☑ What is offline RL? Why do we need it for robots?

<span style="color:red">(Enables safer training, leverages diverse experience)</span>

☐ Paradigm 1: Offline RL via Pessimism

☑ Problem with Q-learning <span style="color:red">(Incorrectly optimistic about unseen actions)</span>

☐ Pessimism to the rescue

☐ Paradigm 2: RL via Supervised Learning

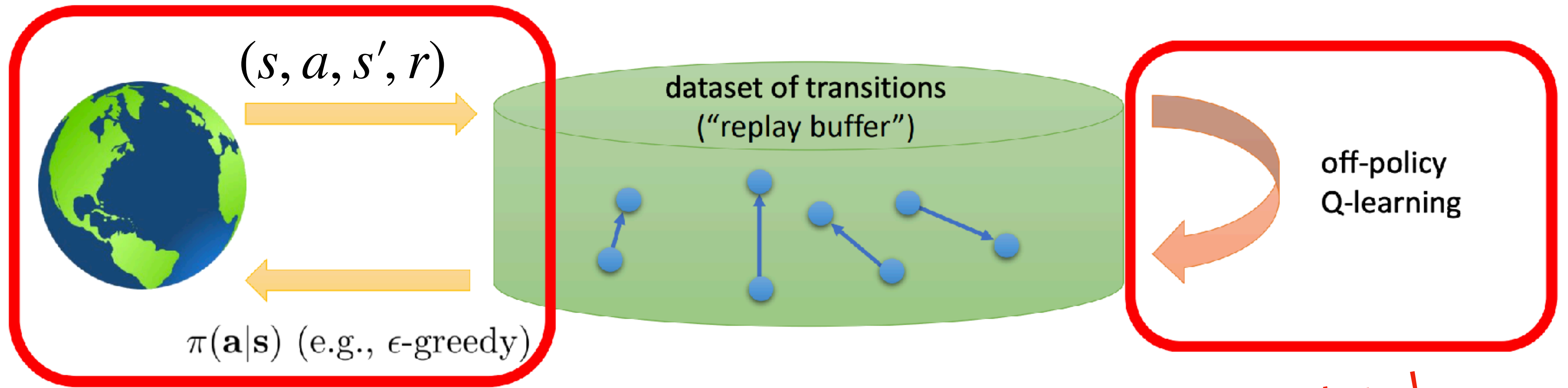☐ Return-conditioned Supervised Learning

☐ Problem in Stochastic MDPs

# Pessimism

# Pessimism as a policy constraint

Don't deviate too much from the data collecting policy

# Pessimism as a policy constraint

Don't deviate too much from the data collecting policy

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\mathrm{new}}}[Q(\mathbf{s}', \mathbf{a}')]$$

$$\pi_{\mathrm{new}}(\mathbf{a}|\mathbf{s}) = \arg\max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})]$$

Typical Q-learning

# Pessimism as a policy constraint

Don't deviate too much from the data collecting policy

$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\mathrm{new}}}[Q(\mathbf{s}', \mathbf{a}')]$$

$$\pi_{\mathrm{new}}(\mathbf{a}|\mathbf{s}) = \arg\max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})] \text{ s.t. } \boxed{D_{\mathrm{KL}}(\pi \| \pi_\beta) \leq \epsilon}$$

Typical Q-learning

Add a constraint
on policy

$$\pi_\beta(\mathbf{a}|\mathbf{s})$$

# TD3+BC: Most simple and effective offline RL!

## A Minimalist Approach to Offline Reinforcement Learning

Scott Fujimoto[1,2]    Shixiang Shane Gu[2]

[1]Mila, McGill University

[2]Google Research, Brain Team

scott.fujimoto@mail.mcgill.ca

$$\pi = \arg\max_{\pi} \mathbb{E}_{(s,a)\sim\mathcal{D}}[Q(s,\pi(s))].$$

$$\pi = \arg\max_{\pi} \mathbb{E}_{(s,a)\sim\mathcal{D}}\left[\lambda Q(s,\pi(s)) - (\pi(s) - a)^2\right],$$

| | | BC | BRAC-p | AWAC | CQL | Fisher-BRC | TD3+BC |
|---|---|---|---|---|---|---|---|
| Random | HalfCheetah | 2.0 ±0.1 | 23.5 | 2.2 | 21.7 ±0.9 | 32.2 ±2.2 | 10.2 ±1.3 |
| | Hopper | 9.5 ±0.1 | 11.1 | 9.6 | 10.7 ±0.1 | 11.4 ±0.2 | 11.0 ±0.1 |
| | Walker2d | 1.2 ±0.2 | 0.8 | 5.1 | 2.7 ±1.2 | 0.6 ±0.6 | 1.4 ±1.6 |
| Medium | HalfCheetah | 36.6 ±0.6 | 44.0 | 37.4 | 37.2 ±0.3 | 41.3 ±0.5 | 42.8 ±0.3 |
| | Hopper | 30.0 ±0.5 | 31.2 | 72.0 | 44.2 ±10.8 | 99.4 ±0.4 | 99.5 ±1.0 |
| | Walker2d | 11.4 ±6.3 | 72.7 | 30.1 | 57.5 ±8.3 | 79.5 ±1.0 | 79.7 ±1.8 |
| Medium Replay | HalfCheetah | 34.7 ±1.8 | 45.6 | - | 41.9 ±1.1 | 43.3 ±0.9 | 43.3 ±0.5 |
| | Hopper | 19.7 ±5.9 | 0.7 | - | 28.6 ±0.9 | 35.6 ±2.5 | 31.4 ±3.0 |
| | Walker2d | 8.3 ±1.5 | -0.3 | - | 15.8 ±2.6 | 42.6 ±7.0 | 25.2 ±5.1 |
| Medium Expert | HalfCheetah | 67.6 ±13.2 | 43.8 | 36.8 | 27.1 ±3.9 | 96.1 ±9.5 | 97.9 ±4.4 |
| | Hopper | 89.6 ±27.6 | 1.1 | 80.9 | 111.4 ±1.2 | 90.6 ±43.3 | 112.2 ±0.2 |
| | Walker2d | 12.0 ±5.8 | -0.3 | 42.7 | 68.1 ±13.1 | 103.6 ±4.6 | 101.1 ±9.3 |
| Expert | HalfCheetah | 105.2 ±1.7 | 3.8 | 78.5 | 82.4 ±7.4 | 106.8 ±3.0 | 105.7 ±1.9 |
| | Hopper | 111.5 ±1.3 | 6.6 | 85.2 | 111.2 ±2.1 | 112.3 ±0.2 | 112.2 ±0.2 |
| | Walker2d | 56.0 ±24.9 | -0.2 | 57.0 | 103.8 ±7.6 | 79.9 ±32.4 | 105.7 ±2.7 |
| | Total | 595.3 ±91.5 | 284.1 | - | 764.3 ±61.5 | 974.6 ±108.3 | 979.3 ±33.4 |

36

# Works on real self-driving problems!

**Imitation Is Not Enough: Robustifying Imitation with Reinforcement Learning for Challenging Driving Scenarios**

Yiren Lu[1], Justin Fu[1], George Tucker[2], Xinlei Pan[1], Eli Bronstein[1], Rebecca Roelofs[2], Benjamin Sapp[1], Brandyn White[1], Aleksandra Faust[2], Shimon Whiteson[1], Dragomir Anguelov[1], Sergey Levine[2,3]

# Works on real self-driving problems!

**Imitation Is Not Enough: Robustifying Imitation with Reinforcement Learning for Challenging Driving Scenarios**

Yiren Lu[1], Justin Fu[1], George Tucker[2], Xinlei Pan[1], Eli Bronstein[1], Rebecca Roelofs[2], Benjamin Sapp[1], Brandyn White[1], Aleksandra Faust[2], Shimon Whiteson[1], Dragomir Anguelov[1], Sergey Levine[2,3]

# Many more sophisticated offline RL methods

**Conservative Q-Learning
for Offline Reinforcement Learning**

Aviral Kumar[1], Aurick Zhou[1], George Tucker[2], Sergey Levine[1,2]
[1]UC Berkeley, [2]Google Research, Brain Team
aviralk@berkeley.edu

**Adversarially Trained Actor Critic for Offline Reinforcement Learning**

Ching-An Cheng[*1]  Tengyang Xie[*2]  Nan Jiang[2]  Alekh Agarwal[3]

Instead of constraining policy, compute pessimistic Q values

Optimize the best worst case performance

# Today's class

☑ What is offline RL? Why do we need it for robots?

(Enables safer training, leverages diverse experience)

☐ Paradigm 1: Offline RL via Pessimism

  ☑ Problem with Q-learning (Incorrectly optimistic about unseen actions)

  ☑ Pessimism to the rescue (Constrain policy to not deviate from data)

☐ Paradigm 2: RL via Supervised Learning

  ☐ Return-conditioned Supervised Learning

  ☐ Problem in Stochastic MDPs

# Reinforcement Learning is Hard ...

# Many horror stories of RL!



Nightmares of
Policy Optimization

# Need many tricks to make Q-learning work in practice!

**Rainbow: Combining Improvements in Deep Reinforcement Learning**

| Matteo Hessel | Joseph Modayil | Hado van Hasselt | Tom Schaul | Georg Ostrovski |
|---|---|---|---|---|
| DeepMind | DeepMind | DeepMind | DeepMind | DeepMind |

| Will Dabney | Dan Horgan | Bilal Piot | Mohammad Azar | David Silver |
|---|---|---|---|---|
| DeepMind | DeepMind | DeepMind | DeepMind | DeepMind |

Double Q Learning

Prioritized Replay

Dueling Networks

Multi-step Learning

Distributional RL

Noisy Nets

Can we just go back to good old supervised learning?

# Supervised Learning success stories

# What if I did supervised learning (BC) here?

+2

+3

+2

-3

-7    -5

# What if I did supervised learning (BC) only on the top % rollouts?

+2

+3

+2

$\pi(a \mid s)$

-3

-7      -5

# An embarrassingly simply algorithm: BC%

1. Collect offline dataset using whatever behavior policy

2. Get the top % trajectories based on returns

3. Do BC on just that!

# Does this even work ?!?

| Dataset | Environment | 10%BC | 25%BC | 40%BC | 100%BC | CQL |
|---------|-------------|-------|-------|-------|--------|-----|
| Medium | HalfCheetah | 42.9 | 43.0 | 43.1 | 43.1 | **44.4** |
| Medium | Hopper | 65.9 | 65.2 | 65.3 | 63.9 | 58.0 |
| Medium | Walker | 78.8 | **80.9** | 78.8 | 77.3 | 79.2 |
| Medium | Reacher | 51.0 | 48.9 | 58.2 | **58.4** | 26.0 |
| Medium-Replay | HalfCheetah | 40.8 | 40.9 | 41.1 | 4.3 | **46.2** |
| Medium-Replay | Hopper | 70.6 | 58.6 | 31.0 | 27.6 | 48.6 |
| Medium-Replay | Walker | **70.4** | 67.8 | 67.2 | 36.9 | 26.7 |
| Medium-Replay | Reacher | **33.1** | 16.2 | 10.7 | 5.4 | 19.0 |
| **Average** | | **56.7** | 52.7 | 49.4 | 39.5 | 43.5 |

49

# An embarrassingly simply algorithm: BC%

1. Collect offline dataset using whatever behavior policy

2. Get the top % trajectories based on returns

3. Do BC on just that!

Challenge with BC%:
What happens as I vary % from small to high values?

Can we have a more principled approach?

# Idea: Train a policy *conditioned* on the returns



$$\pi(a \,|\, s, R)$$

+2

+3

+2

-3

-7

-5

# RvS: What is Essential for Offline RL via Supervised Learning?

**Scott Emmons**[1], **Benjamin Eysenbach**[2], **Ilya Kostrikov**[1], **Sergey Levine**[1]

[1]UC Berkeley, [2]Carnegie Mellon University

emmons@berkeley.edu

# The Idea



(a) replay buffer

# The Idea



$(s_1, \quad \omega_1, \quad a_1)$

$(s_2, \quad \omega_2, \quad a_2)$

$(s_3, \quad \omega_3, \quad a_3)$

(b) training dataset

# The Idea



(c) network architecture

# The Algorithm

For all trajectories:  For all timesteps in that trajectory:  For all achieved outcomes:

$$\max_{\theta} \sum_{\tau \in \mathcal{D}} \sum_{1 \leq t \leq |\tau|} \mathbb{E}_{\omega \sim f(\omega | \tau_{t:H})}[\log \pi_\theta(a_t \mid s_t, \omega)].$$

---

**Algorithm 1** RvS-Learning

---

1: **Input**: Dataset of trajectories, $\mathcal{D} = \{\tau\}$
2: Initialize policy $\pi_\theta(a \mid s, \omega)$.
3: **while** not converged **do**
4:     Randomly sample trajectories: $\tau \sim \mathcal{D}$.
5:     Sample time index for each trajetory, $t \sim [1, H]$, and
   sample a corresponding outcome: $\omega \sim f(\omega \mid \tau_{t:H})$.
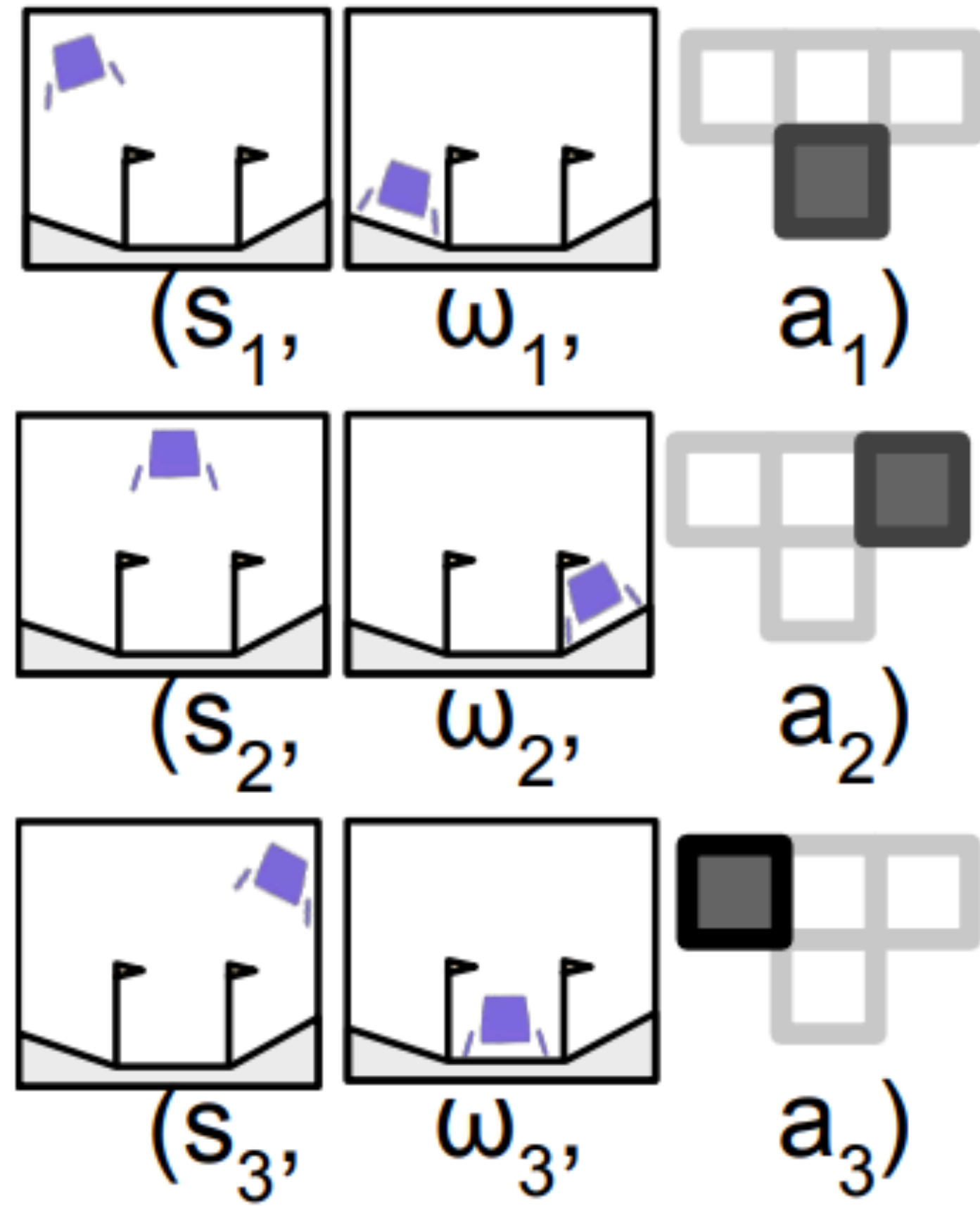6:     Compute loss: $\mathcal{L}(\theta) \leftarrow \sum_{(s_t, a_t, \omega)} \log \pi_\theta(a_t \mid s_t, \omega)$
7:     Update policy parameters: $\theta \leftarrow \theta + \eta \nabla_\theta \mathcal{L}(\theta)$
8: **end while**
9: **return** Conditional policy $\pi_\theta(a \mid s, \omega)$

---

# What are some choices for "outcomes"?

Option 1: What is the future state the agent ended up at?

RvS-G (Goal conditioned)

Option 2: What is the total return that the agent got?

RvS-R (Return conditioned)

# A very *popular* idea

Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. **Decision transformer: Reinforcement learning via sequence modeling**

Felipe Codevilla, Matthias Muller, Antonio Lopez, Vladlen Koltun, and Alexey Dosovitskiy. **End-to-end driving via conditional imitation learning**

Yiming Ding, Carlos Florensa, Pieter Abbeel, and Mariano Phielipp. **Goal-conditioned imitation learning.**

Michael Janner, Qiyang Li, and Sergey Levine. **Offline reinforcement learning as one big sequence modeling problem**

Aviral Kumar, Xue Bin Peng, and Sergey Levine. **Reward-conditioned policies**

Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. **Advantage-weighted regression: Simple and scalable off-policy reinforcement learning**

Rupesh Kumar Srivastava, Pranav Shyam, Filipe Mutz, Wojciech Jaskowski, and Jurgen Schmidhuber. ¨ **Training agents using upside-down reinforcement learning**

# Many popular algorithm, e.g. Decision Transformer

· · · $\widehat{R}_{t\text{-}1}$

# Today's class

☑ What is offline RL? Why do we need it for robots?

(Enables safer training, leverages diverse experience)

☐ Paradigm 1: Offline RL via Pessimism

　☑ Problem with Q-learning (Incorrectly optimistic about unseen actions)

　☑ Pessimism to the rescue (Constrain policy to not deviate from data)

☐ Paradigm 2: RL via Supervised Learning

　☑ Return-conditioned Supervised Learning

　☐ Problem in Stochastic MDPs

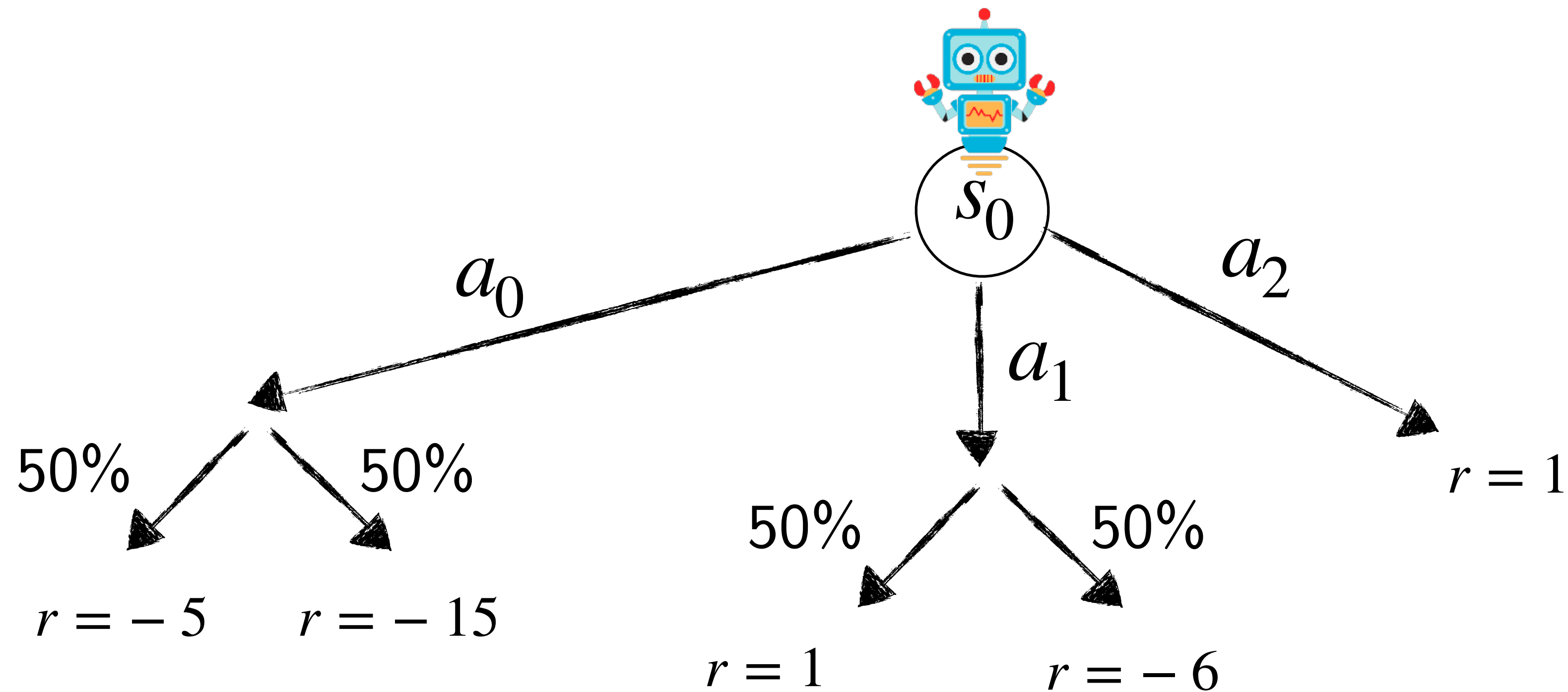# Activity!

# Consider the following MDP



$s_0$

$a_0$

$a_1$

$a_2$

$r = 1$

50%    50%

$r = -5$    $r = -15$

50%    50%

$r = 1$    $r = -6$

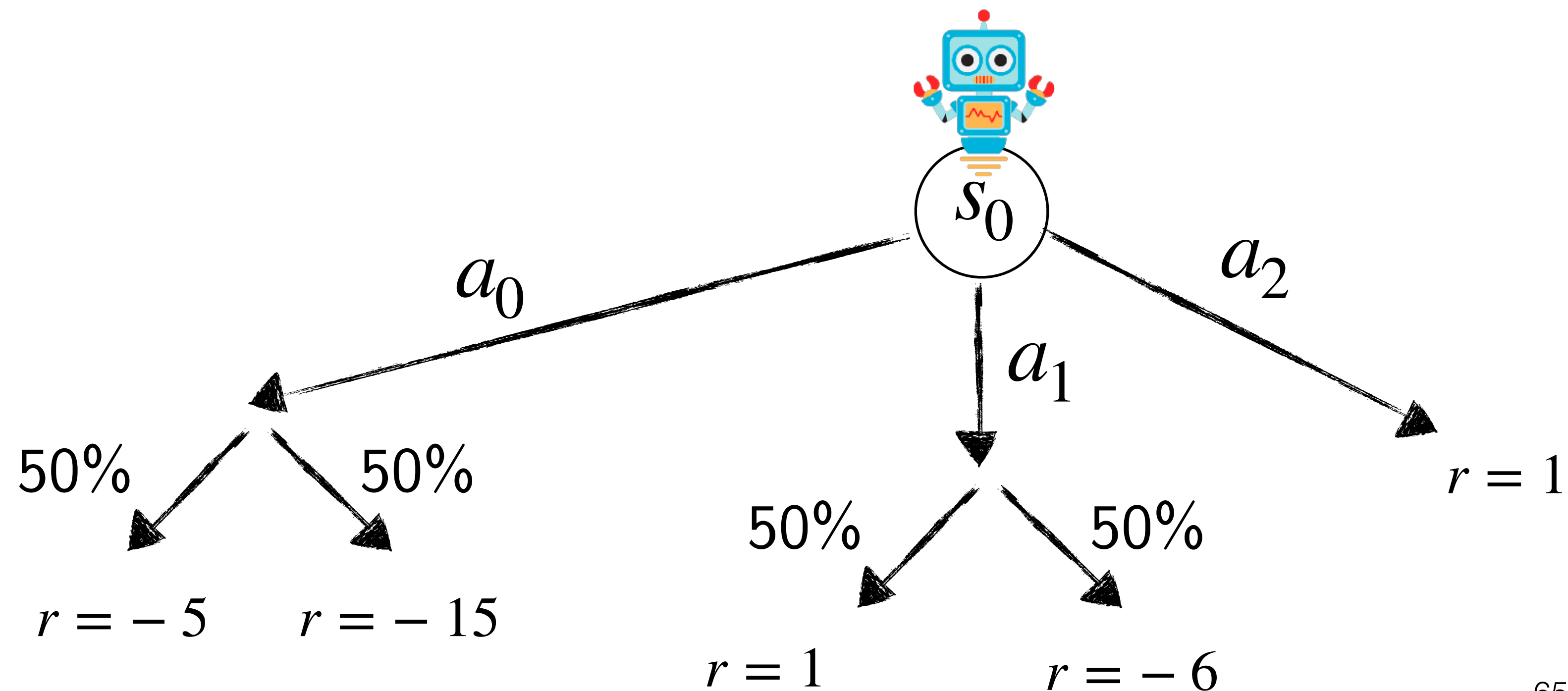# Consider the following MDP



What is the optimal action? What will Decision Transformer play?

# Think-Pair-Share!
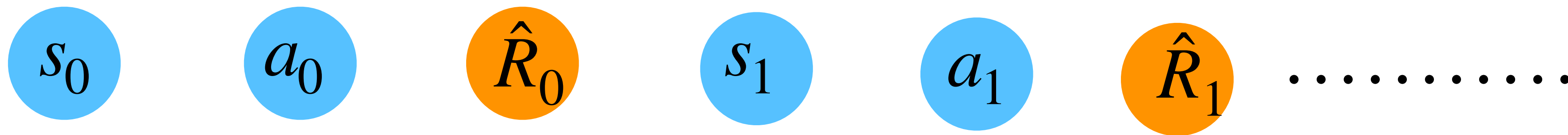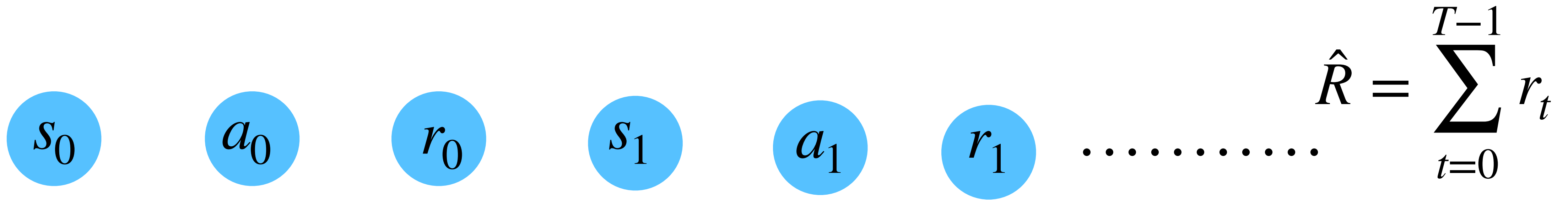
Think (30 sec): What is the optimal action? What would decision transformers play?

Pair: Find a partner

Share (45 sec): Partners exchange ideas



$s_0$

$a_0$

$a_1$

$a_2$

50%    50%

$r = -5$    $r = -15$

50%    50%

$r = 1$    $r = -6$

$r = 1$

$$\hat{R} = \sum_{t=0}^{T-1} r_t$$

$$\hat{R}_0 = \sum_{t=0}^{T-1} r_t \qquad \hat{R}_1 = \sum_{t=1}^{T-1} r_t$$

$\cdots$ $\widehat{R}_{t-1}$

linear decoder

causal transformer

emb. + pos. enc.

$\hat{R}_{t-1}$  $s_{t-1}$  $a_{t-1}$  $\hat{R}_t$  $s_t$  $a_t$

$a_{t-1}$  $a_t$

21

return    state    action

**Introducing Decision Transformers on Hugging Face** 🤗

Published March 28, 2022

Update on GitHub

edbeeching
Edward Beeching

ThomasSimonini
Thomas Simonini

# Test Time

Start at initial state $s_0$

Specify the desired target return $R_0$

$a_0 = \text{Transformer}(R_0, s_0)$

Execute action, observe reward and next state $(r_0, s_1)$
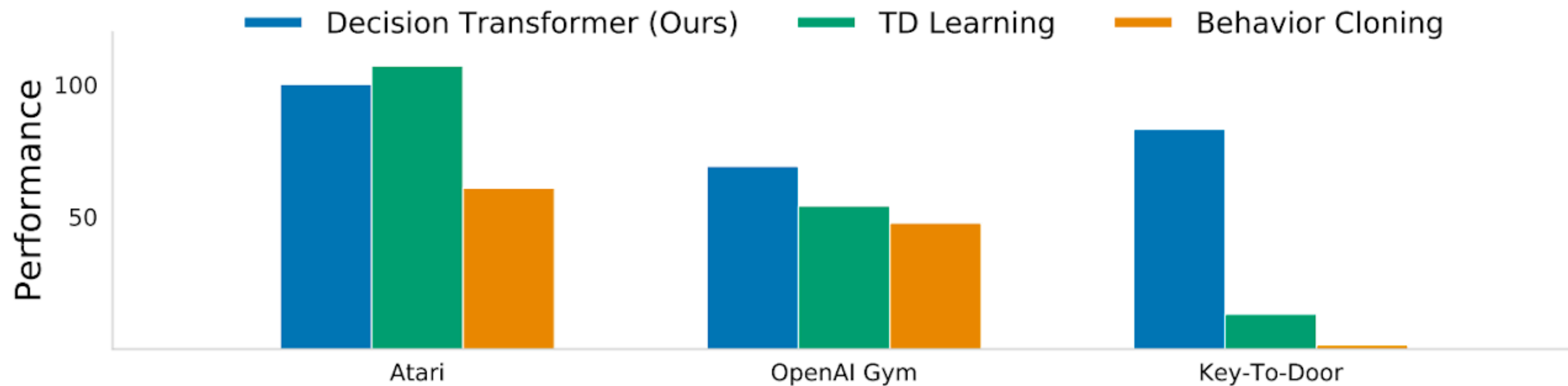
Decrement the target return $R_1 = R_0 - r_0$

$a_1 = \text{Transformer}(R_0, s_0, a_0, R_1, s_1)$

# Seems to work!

# Seems to work!

| Game | DT (Ours) | CQL | QR-DQN | REM | BC |
|------|-----------|-----|--------|-----|-----|
| Breakout | $\mathbf{267.5 \pm 97.5}$ | $211.1$ | $17.1$ | $8.9$ | $138.9 \pm 61.7$ |
| Qbert | $15.4 \pm 11.4$ | $\mathbf{104.2}$ | $0.0$ | $0.0$ | $17.3 \pm 14.7$ |
| Pong | $106.1 \pm 8.1$ | $\mathbf{111.9}$ | $18.0$ | $0.5$ | $85.2 \pm 20.0$ |
| Seaquest | $\mathbf{2.5 \pm 0.4}$ | $1.7$ | $0.4$ | $0.7$ | $2.1 \pm 0.3$ |

Atari

# Today's class

☑ What is offline RL? Why do we need it for robots?

<span style="color:red">(Enables safer training, leverages diverse experience)</span>

☐ Paradigm 1: Offline RL via Pessimism

    ☑ Problem with Q-learning <span style="color:red">(Incorrectly optimistic about unseen actions)</span>

    ☑ Pessimism to the rescue <span style="color:red">(Constrain policy to not deviate from data)</span>

☐ Paradigm 2: RL via Supervised Learning

    ☐ Return-conditioned Supervised Learning

    ☐ Problem in Stochastic MDPs