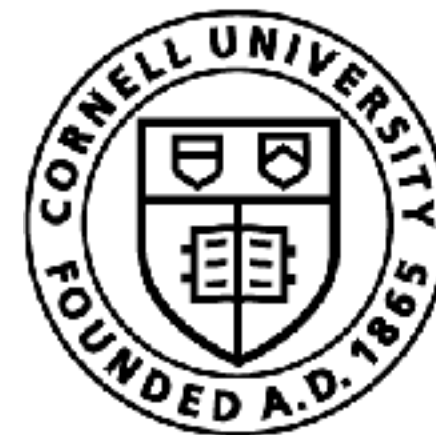


Frontiers: Offline Reinforcement Learning

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science



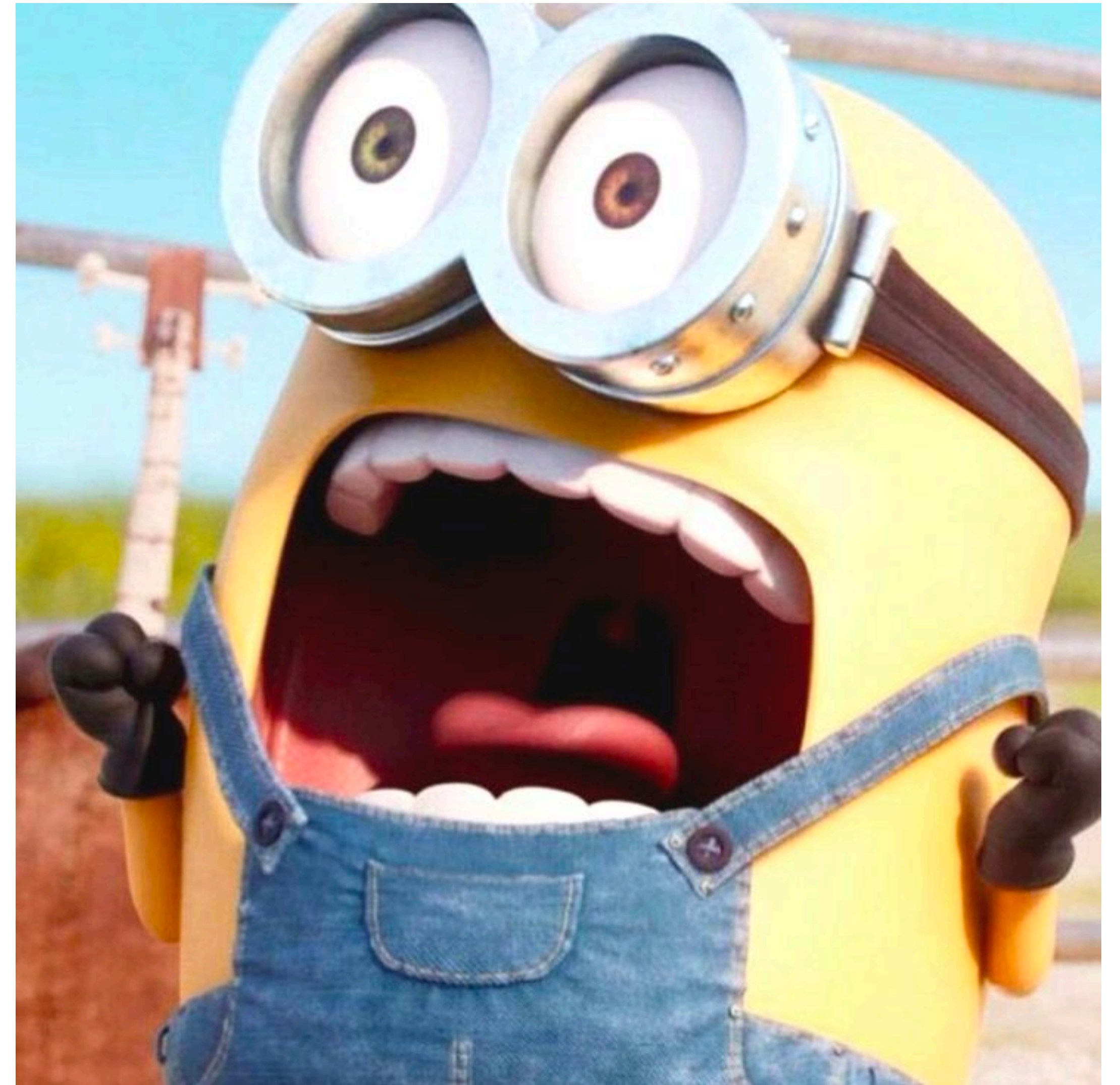
Frontiers

Problem: *Insane* number of papers out there!!

Impossible for outsiders to find any sort of scaffolding

Many of these papers recycle old ideas while butchering the insight

Hope: Sparse set of papers that give you reach



Our Strategy

Goal: Engage with various frontiers of research
on robot decision making

Strategy: Equip you with a sparse “support vector” of papers
that gives you maximal reach on the problem

Expectation: For details and concrete implementation,
you should be able to look that up

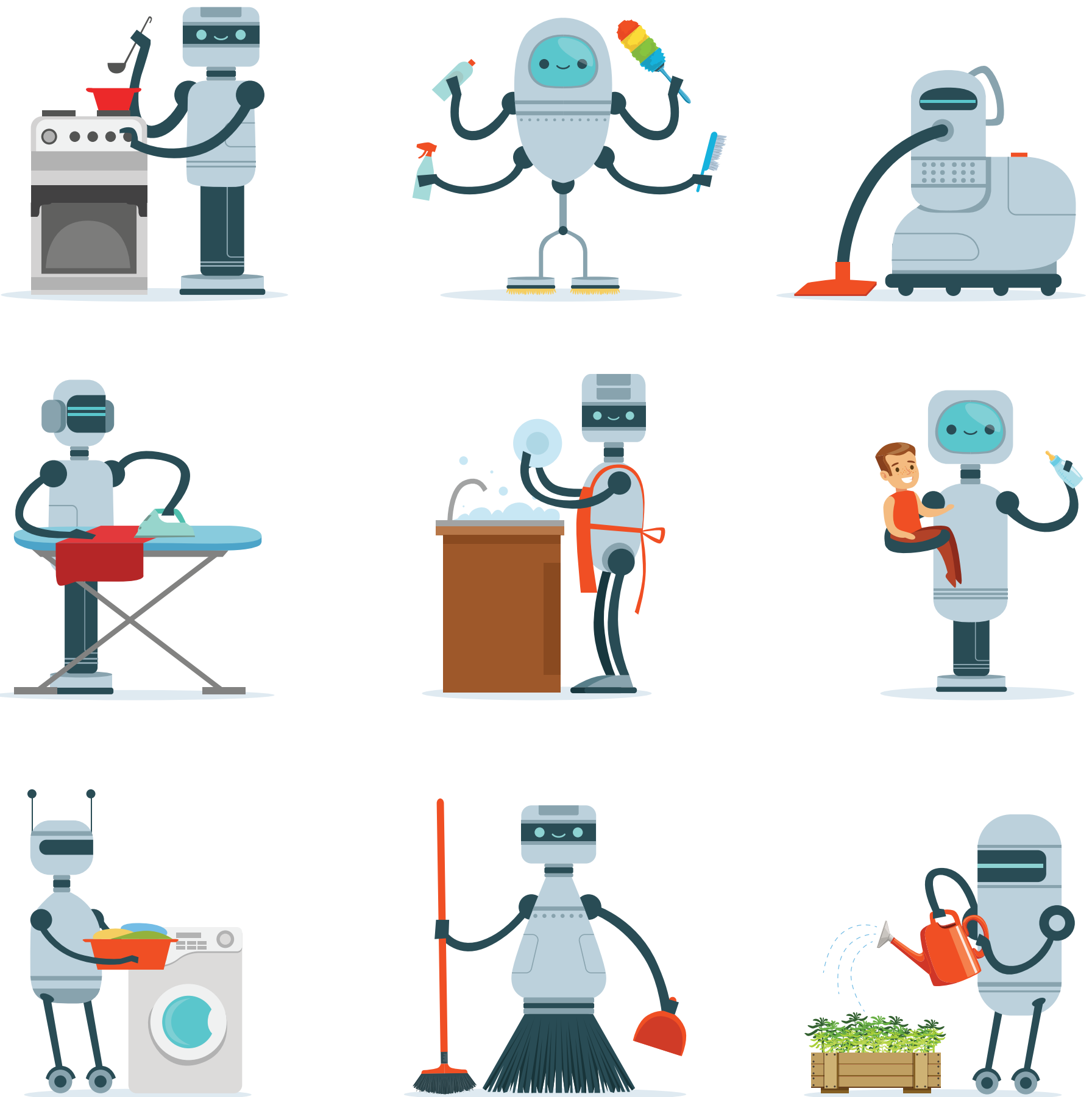
The Problem

Real World, Real Problems

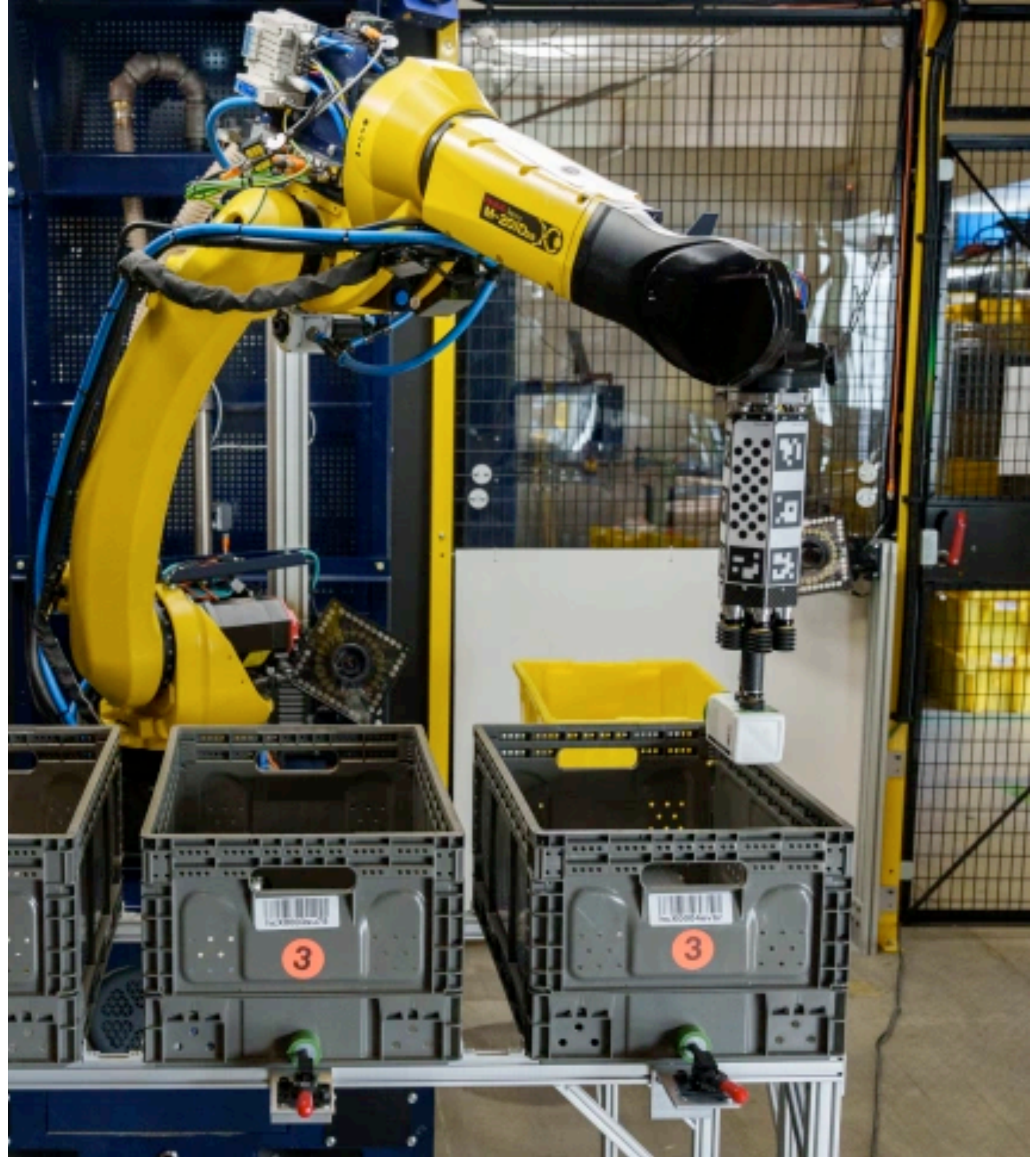


Robots can augment human capabilities to tackle these problems

Robots only really work in the CLOSED world



The Dream

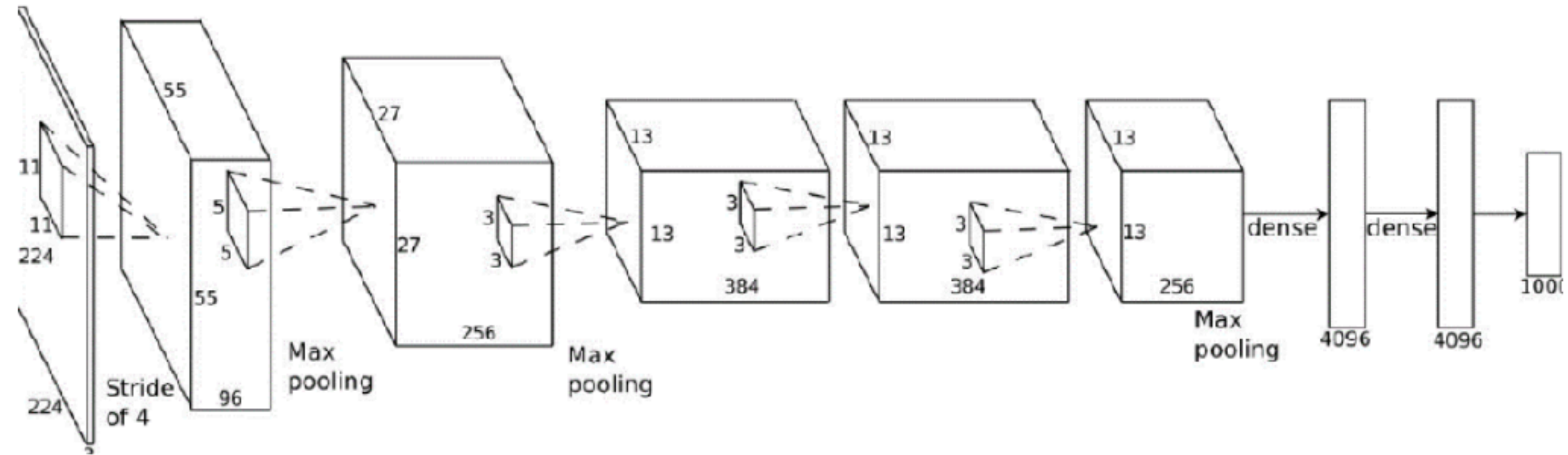
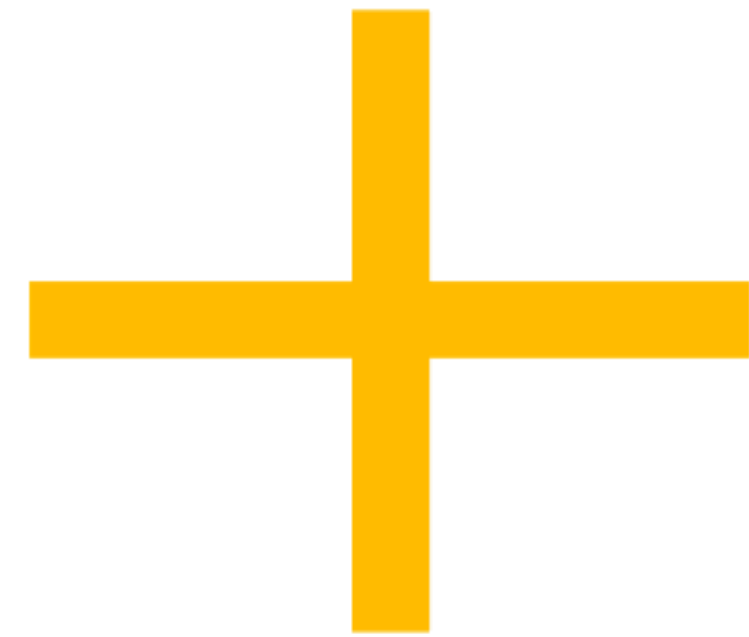
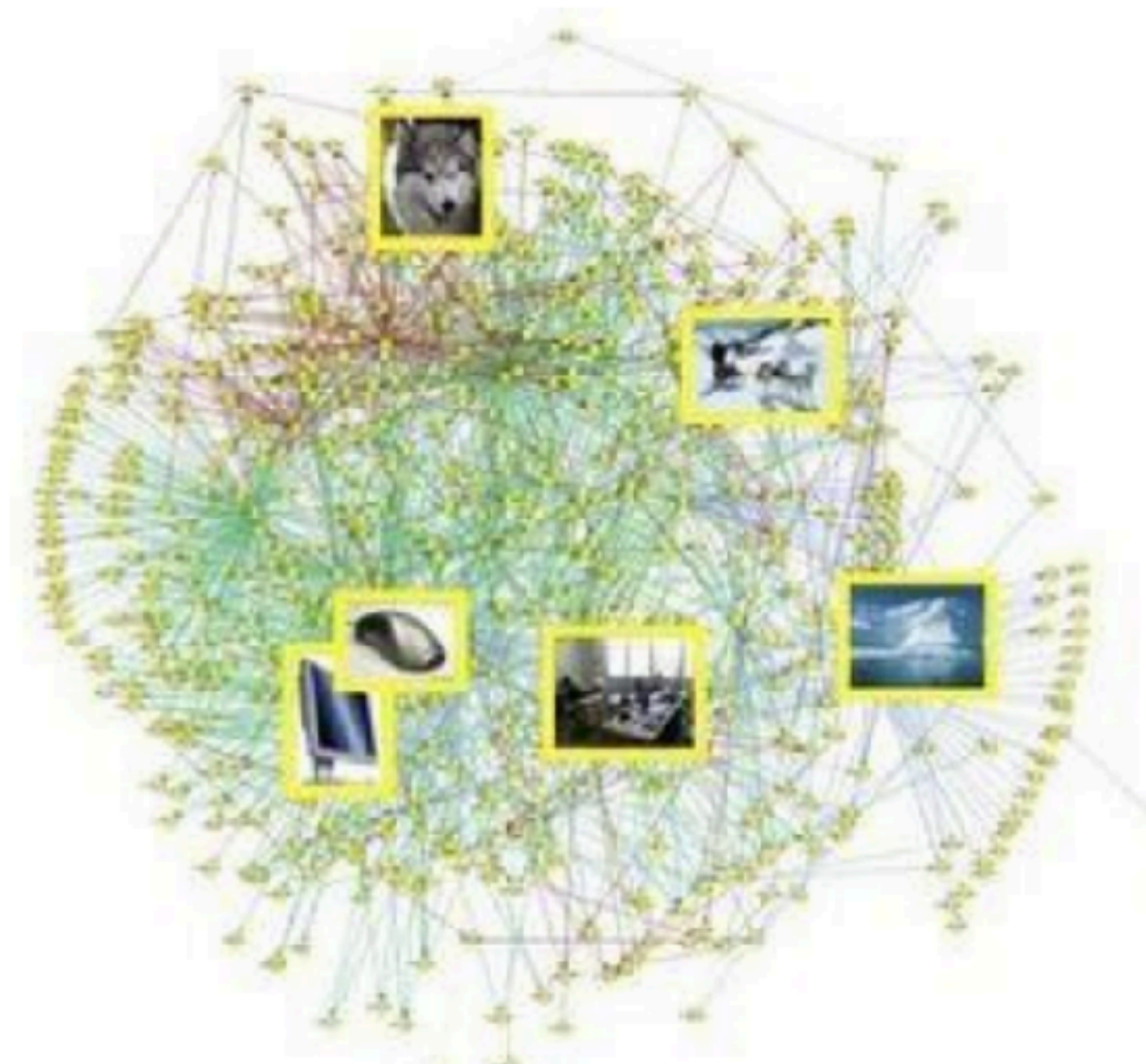


Reality

Generalize to variations of the OPEN world?



Machine learning's answer!

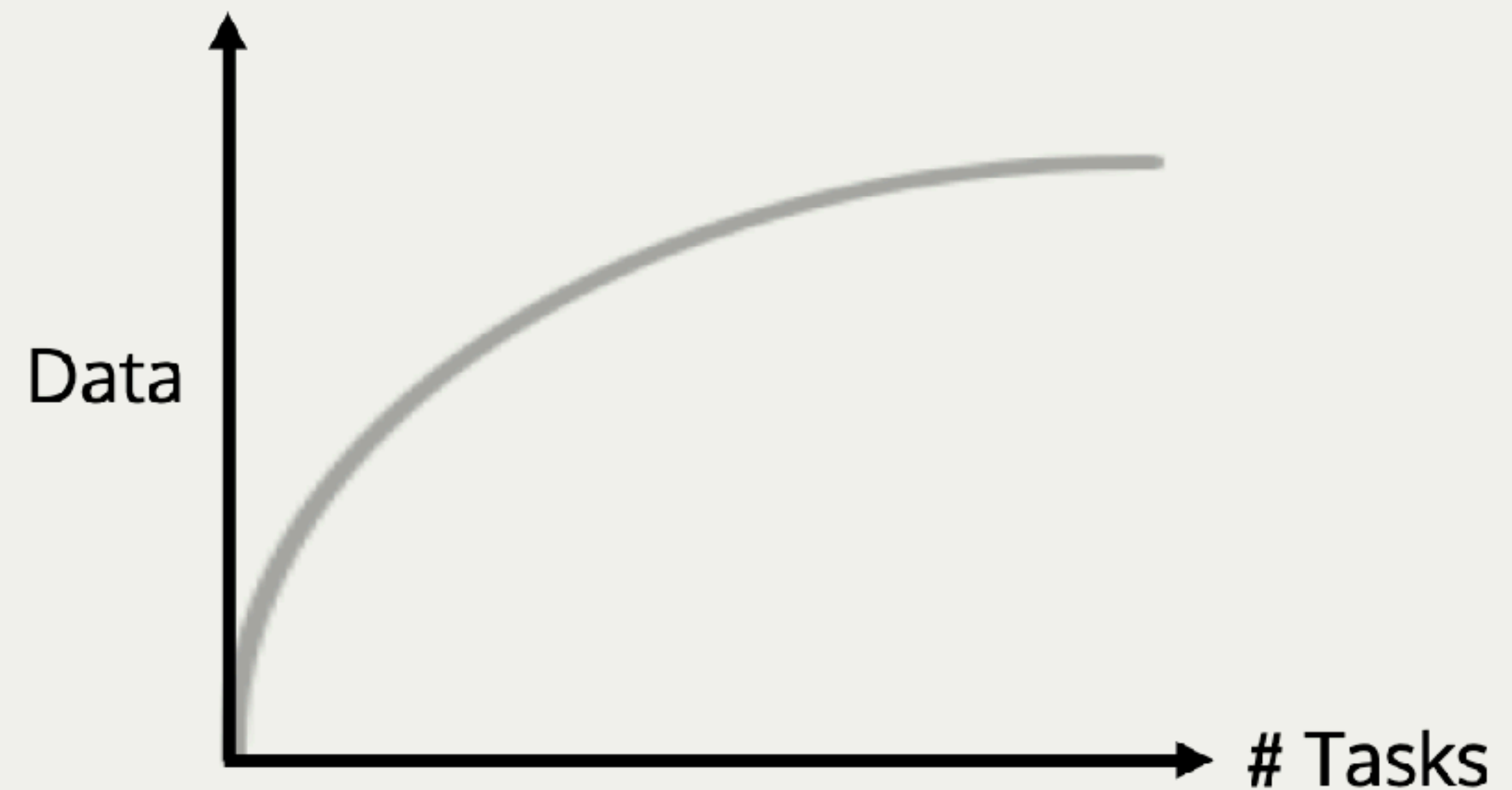


Big Data

Big Models

Hasn't quite been true so far robotics ...

On the quest for shared priors
w/ machine learning

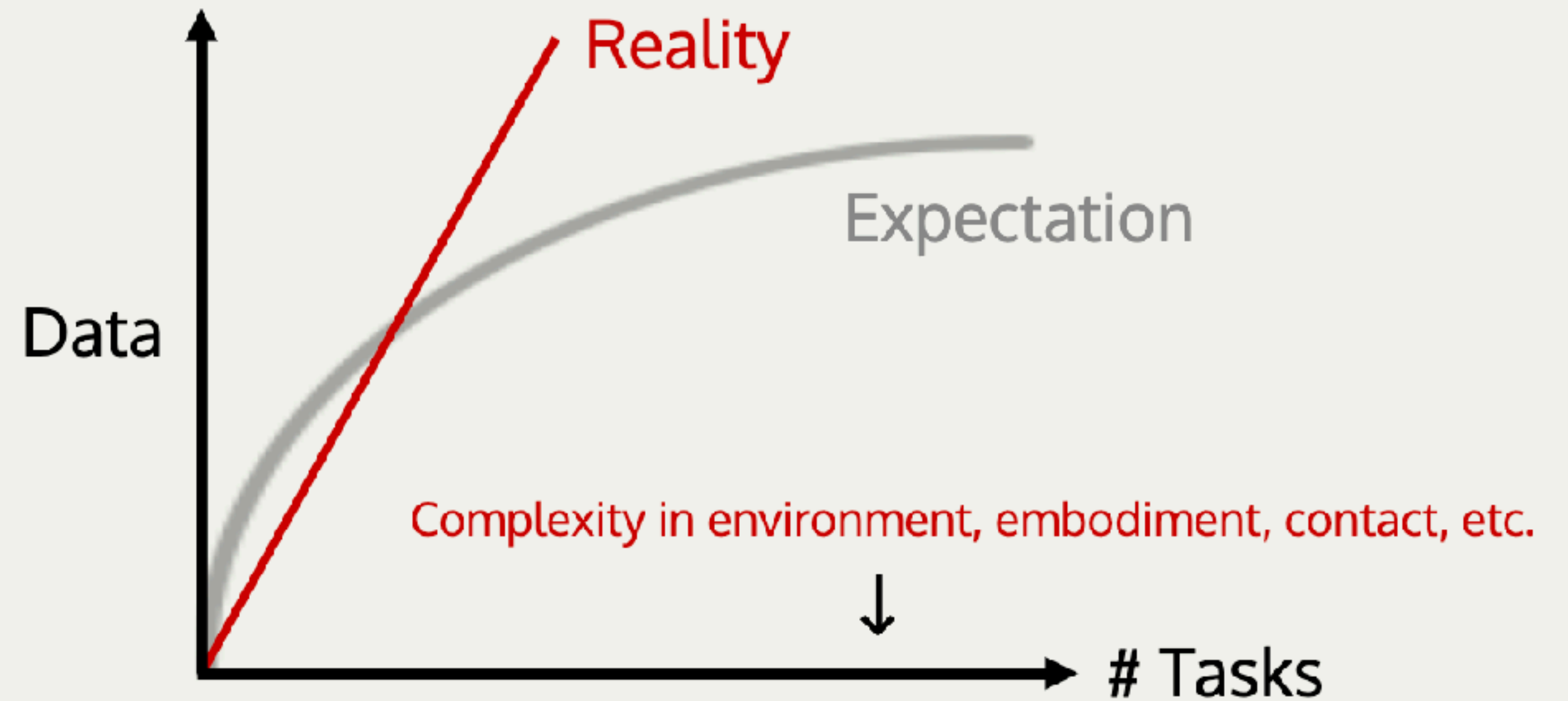


Interact with the **physical** world to learn **bottom-up commonsense**

↑
i.e. "how the world works"

Hasn't quite been true so far robotics ...

On the quest for shared priors
w/ machine learning



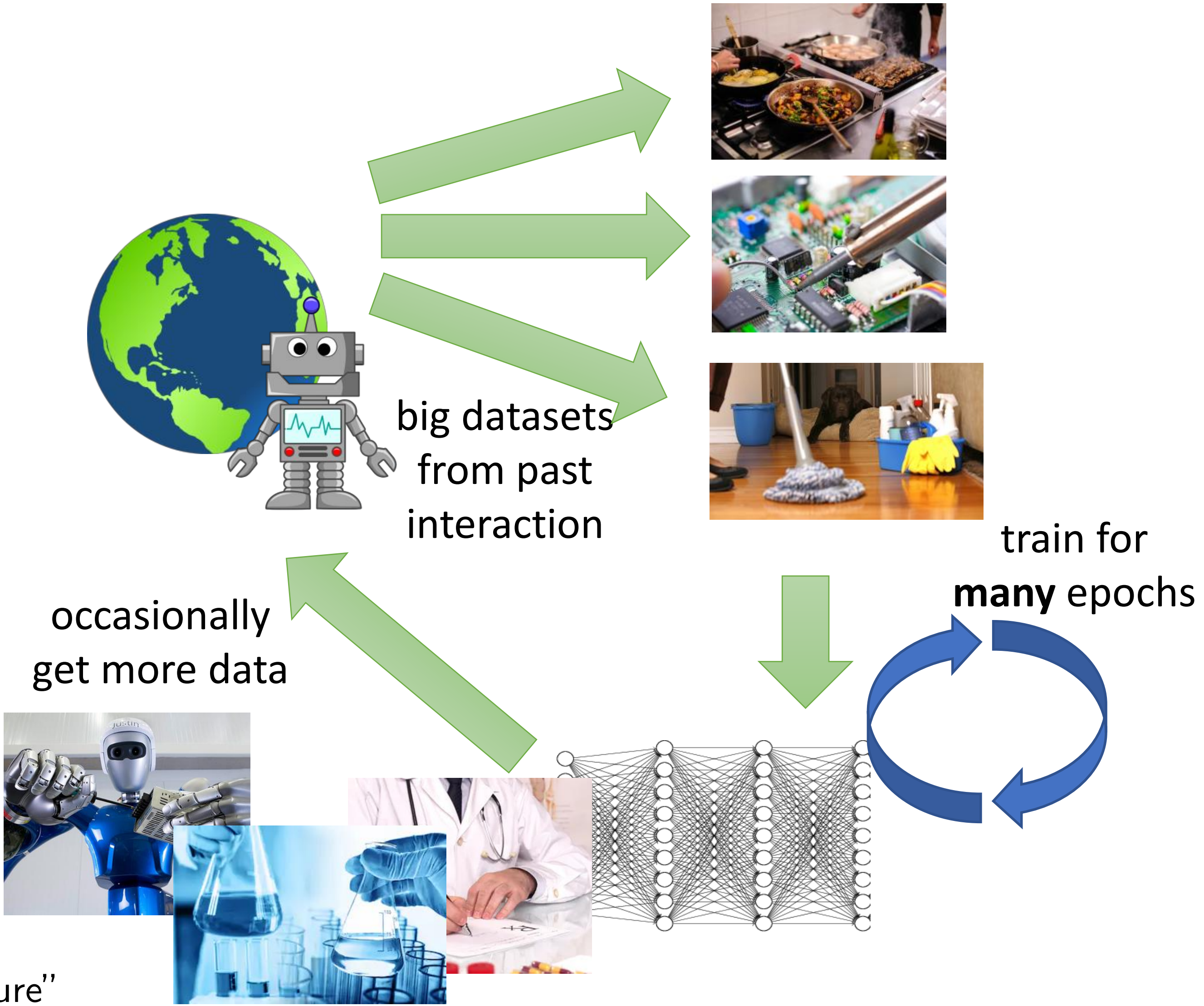
Interact with the **physical** world to learn **bottom-up commonsense**

↑
i.e. "how the world works"

But for today, let's pretend we can collect a ton of data

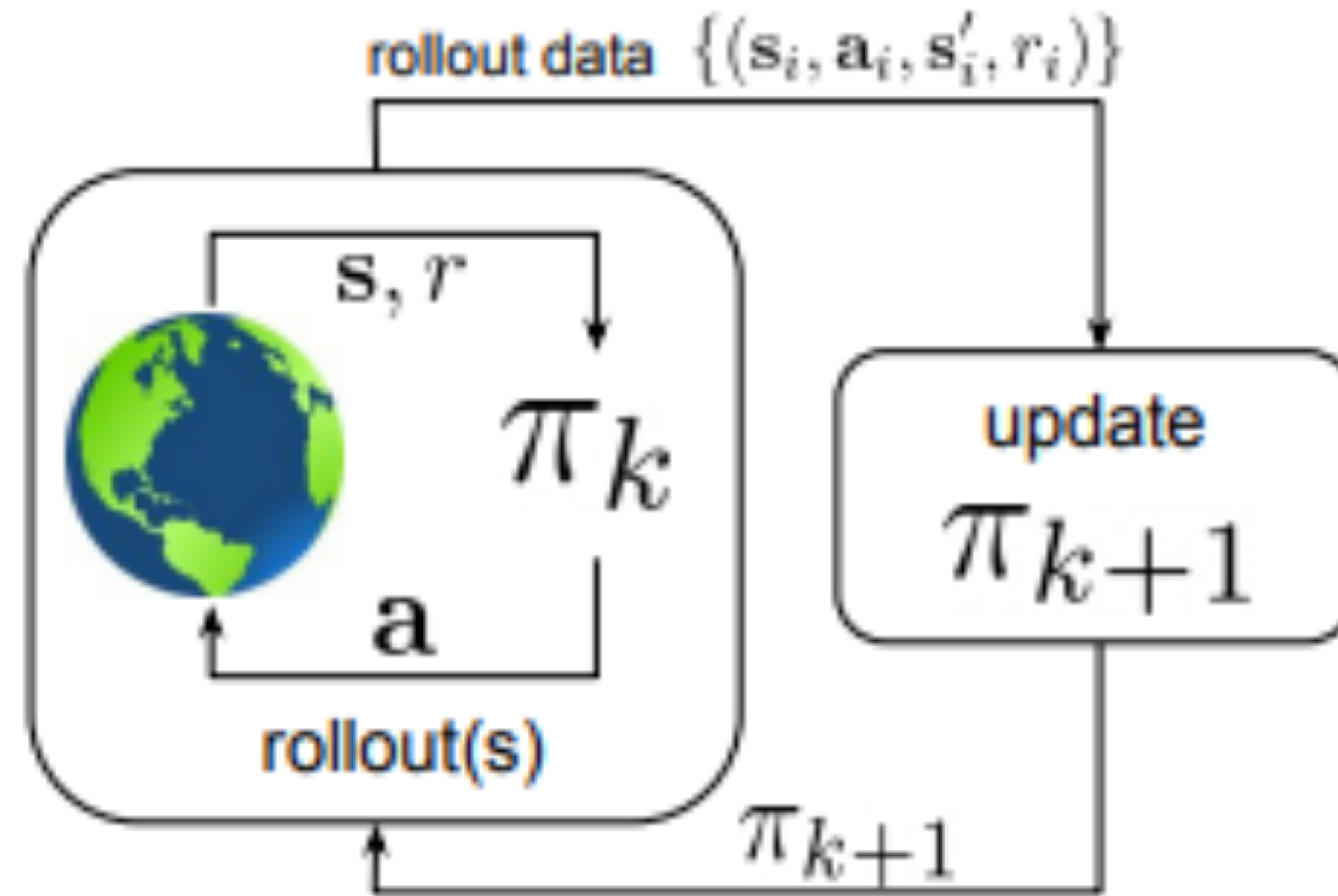
How can we learn “optimal” from large data collected by *any* policy?

Goal: Offline Reinforcement Learning



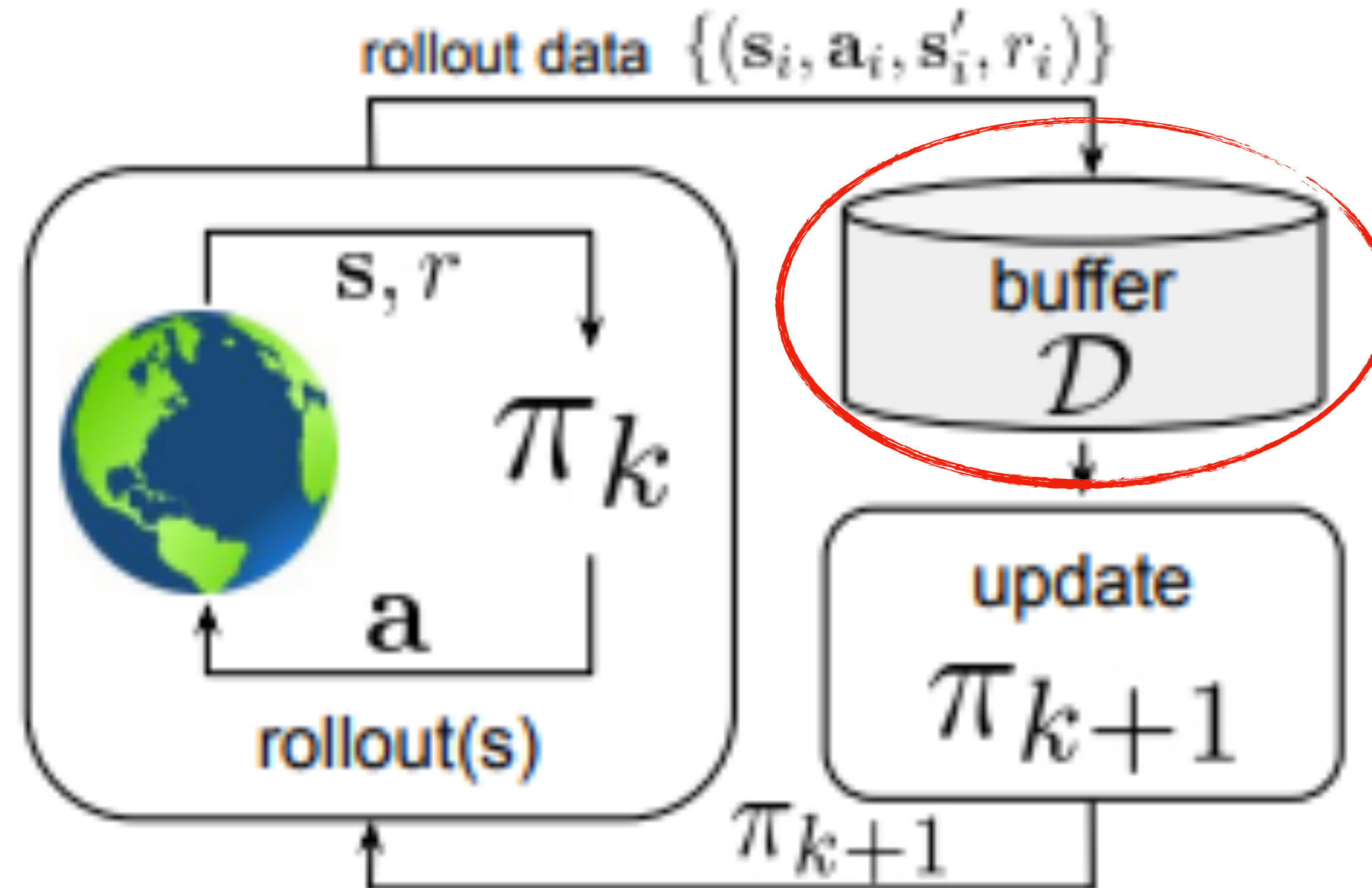
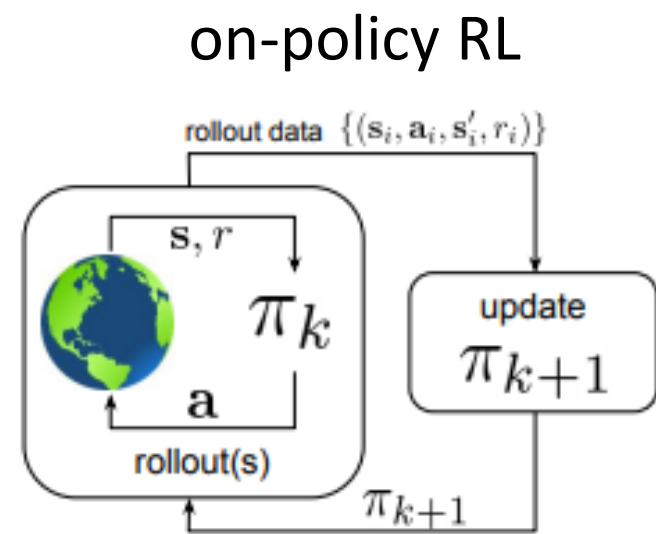
Different paradigms of RL

on-policy RL



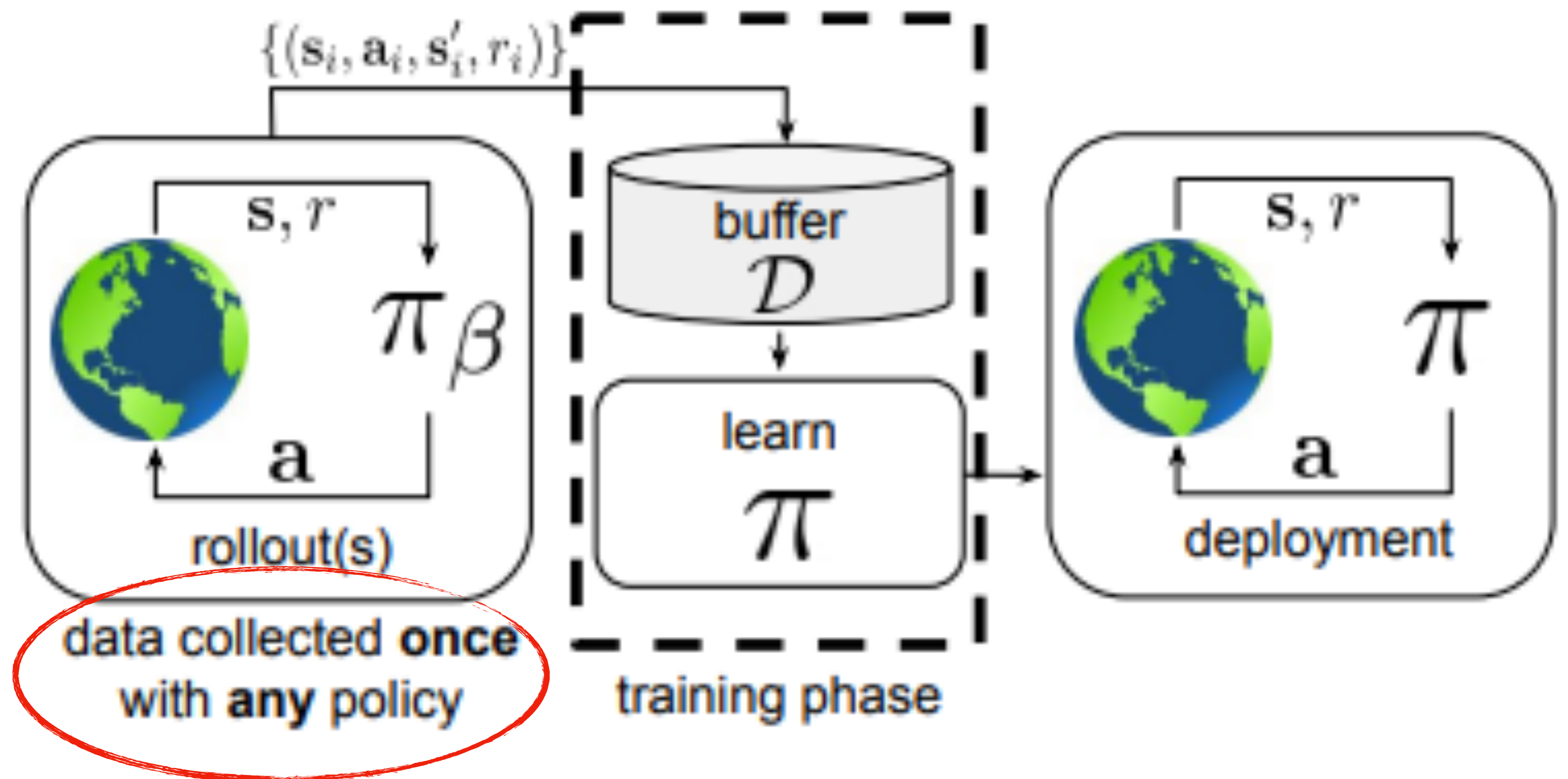
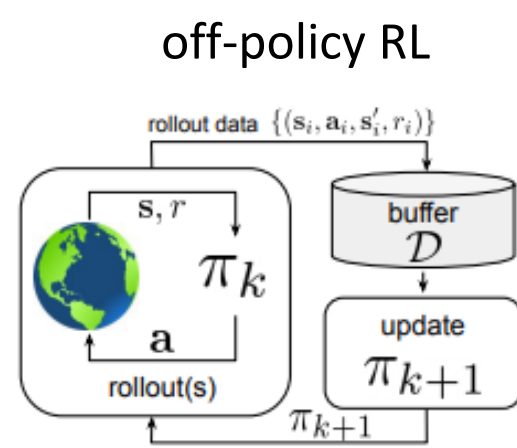
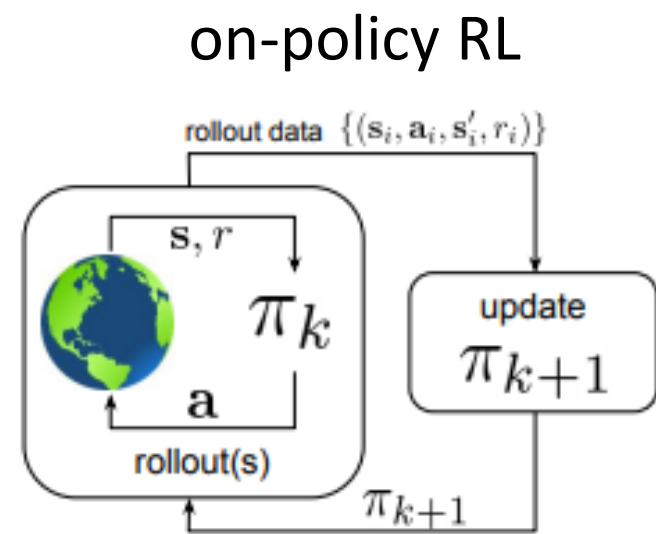
Different paradigms of RL

off-policy RL



Different paradigms of RL

offline reinforcement learning



Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems

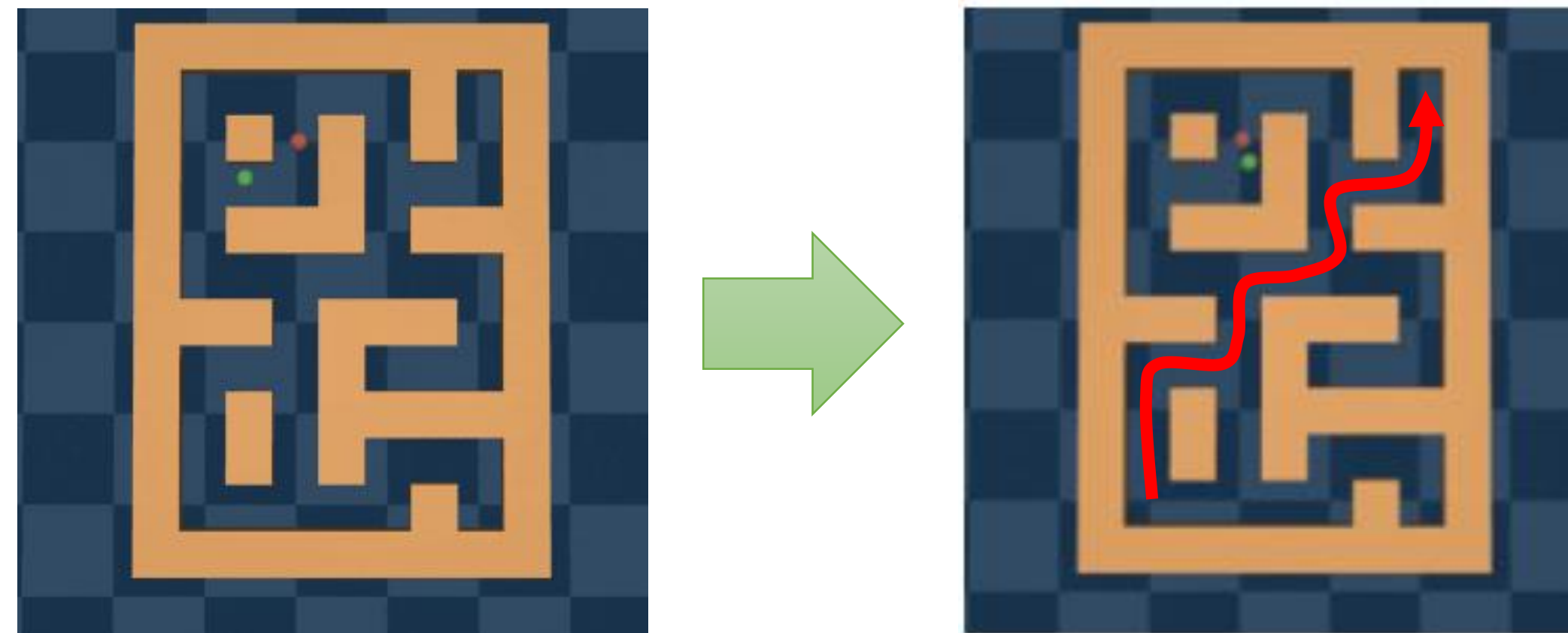
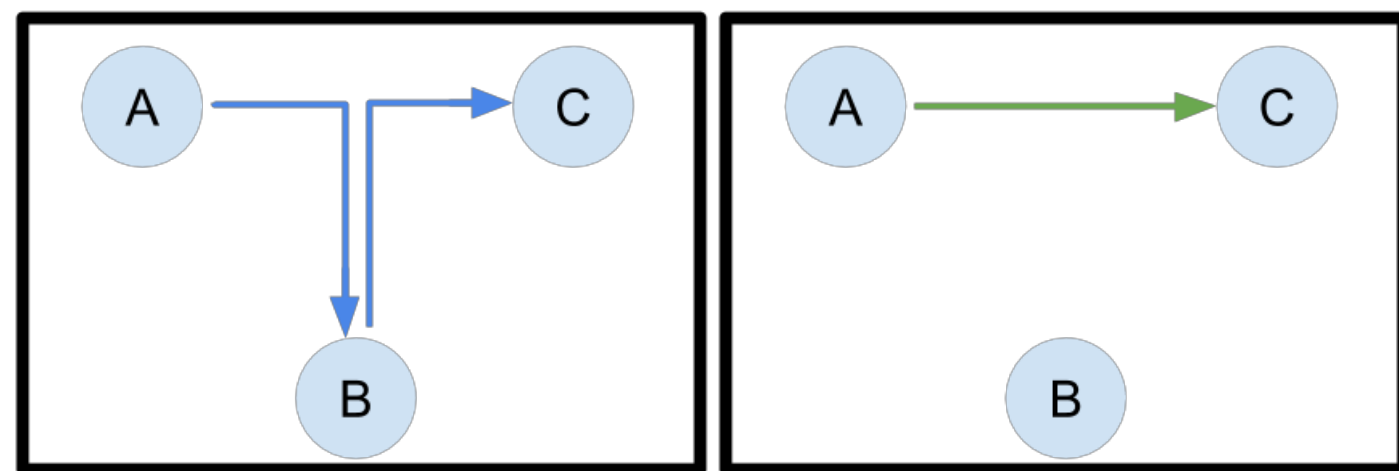
Sergey Levine^{1,2}, Aviral Kumar¹, George Tucker², Justin Fu¹
¹UC Berkeley, ²Google Research, Brain Team



Fun collab tutorial: <https://colab.research.google.com/drive/1oJOYIAIOI9d1JjIutPY66KmfPkwPCgEE?usp=sharing>

How is this even possible?

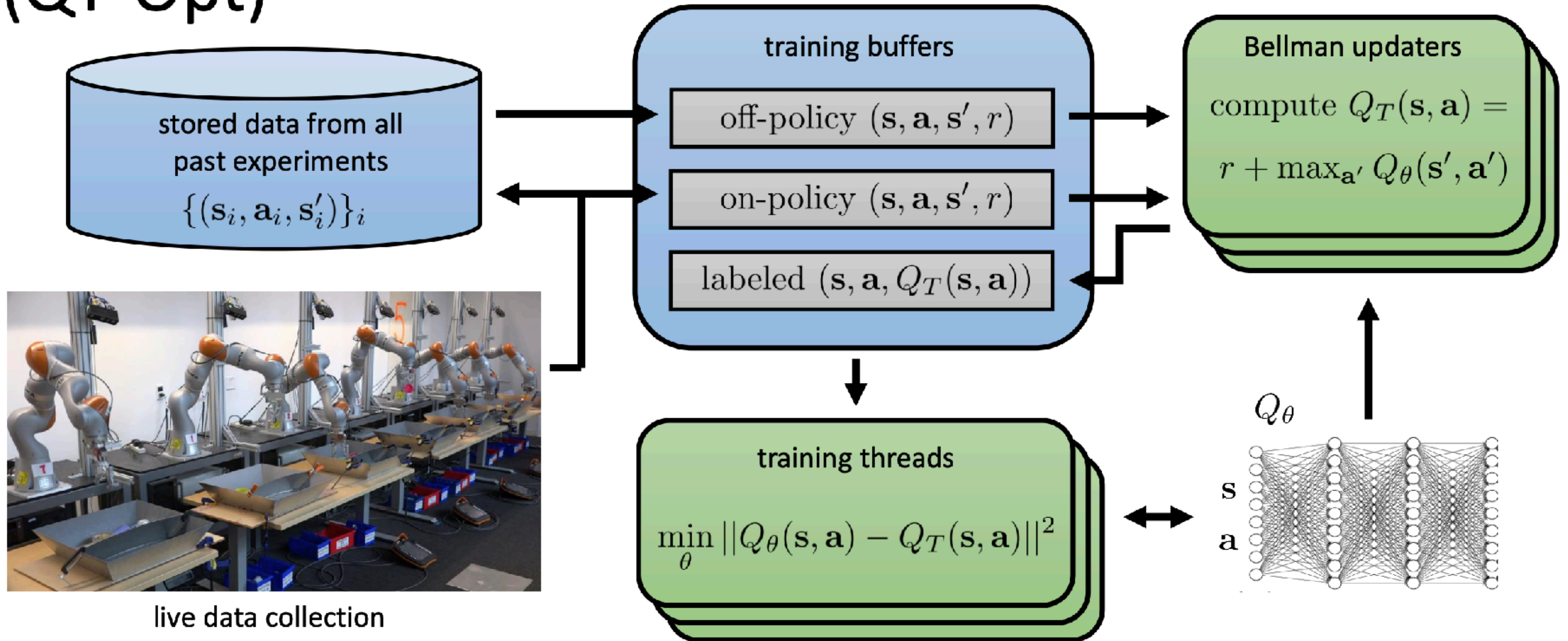
1. Find the “good stuff” in a dataset full of good and bad behaviors
2. Generalization: good behavior in one place may suggest good behavior in another place
3. “Stitching”: parts of good behaviors can be recombined



Does it work?

Sometimes*

Large-scale Q-learning with continuous actions (QT-Opt)



Optimal Insulin Dose

Offline reinforcement learning for safer blood glucose control in people with type 1 diabetes

Harry Emerson ^{a,*}, Matthew Guy ^{a,b}, Ryan McConville ^a

^a University of Bristol, 1 Cathedral Square, Bristol, BS1 5TS, United Kingdom

^b University Hospital Southampton, Tremona Road, Southampton, SO16 6YD, Hampshire, United Kingdom

Table 2

The mean performance of the offline RL algorithms: BCQ, CQL and TD3-BC against the online RL approach SAC-RNN and the control baseline PID. TD3-BC can be seen to significantly improve the proportion of TIR when compared to the PID and the SAC-RNN algorithms. This is done so without any associated increase in risk (reward) or TBR. Statistical significance was confirmed via a Friedman rank test for all glucose metrics ($p < 0.05$). †, ‡ and § indicate an offline RL, online RL and classical control algorithm respectively, with the best performing algorithm highlighted in bold.

Algorithm	Reward	TIR (%)	TBR (%)	CV (%)	Failure (%)
BCQ [†]	-41,034 ± 1,060	65.8 ± 0.6	1.0 ± 0.1	35.1 ± 0.4	0.00
CQL [†]	-45,259 ± 1,071	56.2 ± 0.5	0.1 ± 0.1	30.3 ± 0.3	0.00
TD3-BC[†]	-37,955 ± 547	65.3 ± 0.5	0.2 ± 0.1	33.3 ± 0.2	0.00
SAC-RNN [‡]	-93,480 ± 71,826	34.9 ± 3.1	4.1 ± 0.7	29.6 ± 1.3	13.3
PID [§]	-49,077 ± 556	61.6 ± 0.3	0.4 ± 0.1	33.5 ± 0.2	0.00

Combustion control in power stations

DeepThermal: Combustion Optimization for Thermal Power Generating Units Using Offline Reinforcement Learning

Xianyuan Zhan^{1*}, Haoran Xu^{2,3,4*}, Yue Zhang^{2,3}, Xiangyu Zhu^{2,3}, Honglei Yin^{2,3}, Yu Zheng^{2,3,4}

¹ Institute for AI Industry Research (AIR), Tsinghua University, Beijing, China

² JD iCity, JD Technology, Beijing, China

³ JD Intelligent Cities Research, Beijing, China

⁴ Xidian University, Xi'an, China

{zhanxianyuan, ryanxhr, zhangyuezjx, zackxiangyu, yinhonglei93}@gmail.com, msyuzheng@outlook.com

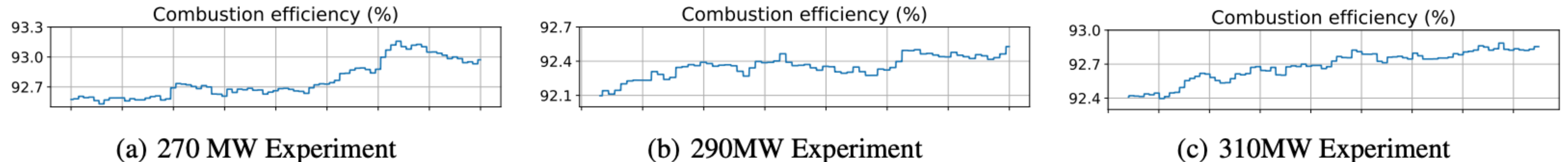


Figure 4: Real-world experiments at CHN Energy Nanning Power Station

Goal-directed conversation

CHAI: A Chatbot AI for Task-Oriented Dialogue with Offline Reinforcement Learning

[Siddharth Verma](#)
UC Berkeley

[Justin Fu](#)
UC Berkeley

[Mengjiao Yang](#)
UC Berkeley

[Sergey Levine](#)
UC Berkeley

Metric	Fluency	Coherency	On-Topic	Human-Likeness	Total
CHAI-prop	4.31 ± 0.97	3.91 ± 1.17	4.16 ± 0.99	3.47 ± 1.27	15.84 ± 3.86
He et al. (2018) (Utility)	3.56 ± 1.34	2.47 ± 1.39	3.09 ± 1.40	2.13 ± 1.13	11.25 ± 4.50
Lang. Model	4.06 ± 1.11	2.66 ± 1.36	3.63 ± 1.18	2.50 ± 1.10	12.84 ± 3.66

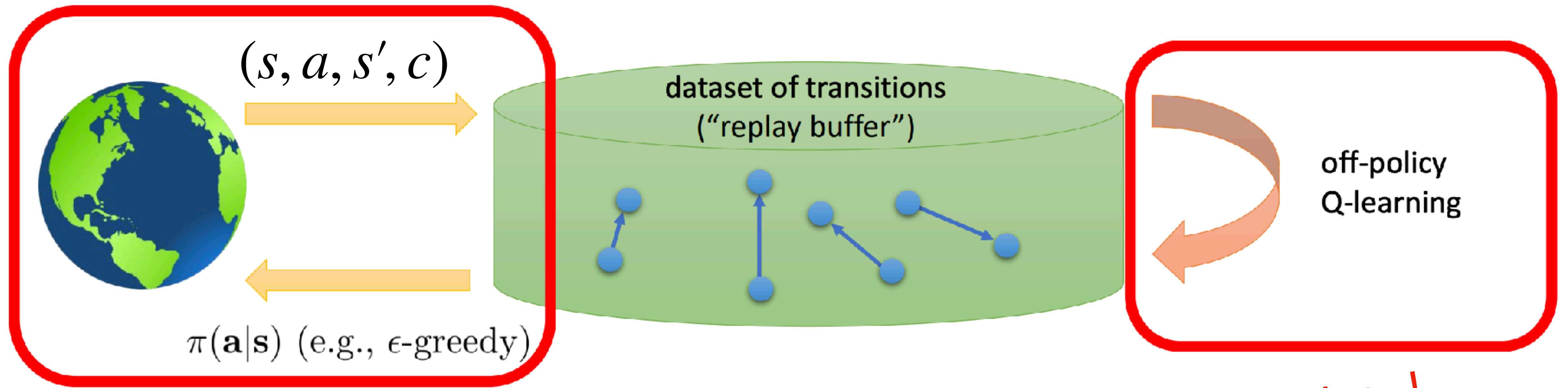
Table 2: Human evaluation scores comparing CHAI, [He et al. \(2018\)](#), and language model (higher is better). Numbers are reported as means and standard deviations over 32 trials. CHAI scores the highest across all metrics.

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

What is it?



Q-learning



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

Can learn from any data!

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

Q-learning

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

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

Notice we are *not* approximating $Q^\pi(s_t, a_t)$

We don't even care about π

We can learn from any data!

Q-learning

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

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

Conditions for
convergence

1. Each state-action pair is visited infinite times
2. $\lim_{k \rightarrow \infty} \sum_{k=0}^{\infty} \alpha_k = \infty$
3. $\lim_{k \rightarrow \infty} \sum_{k=0}^{\infty} \alpha_k^2 < \infty,$

Q-learning

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

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

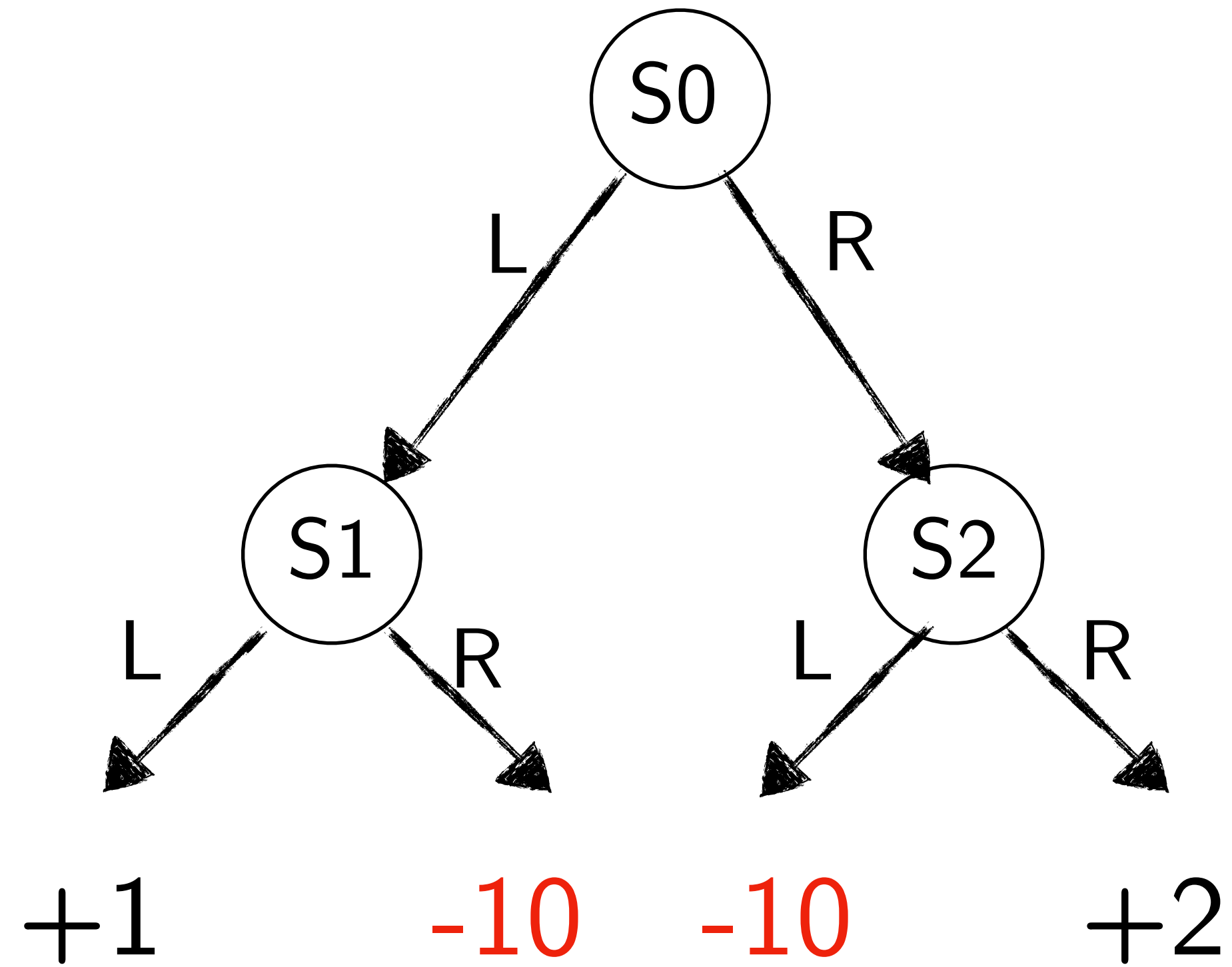
What happens
when this is
not met?

1. Each state-action pair is visited infinite times
2. $\lim_{k \rightarrow \infty} \sum_{k=0}^{\infty} \alpha_k = \infty$
3. $\lim_{k \rightarrow \infty} \sum_{k=0}^{\infty} \alpha_k^2 < \infty,$

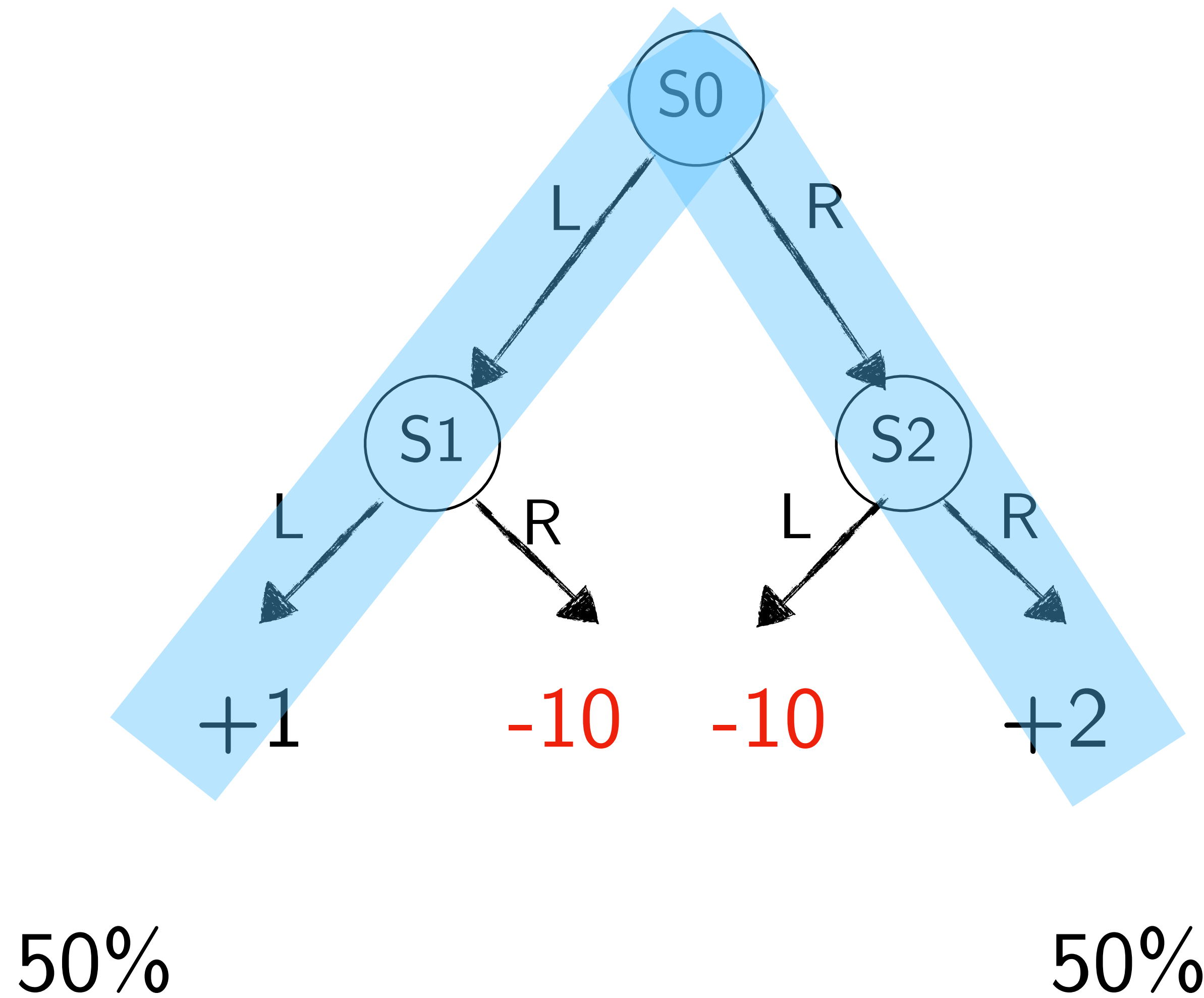
Activity!



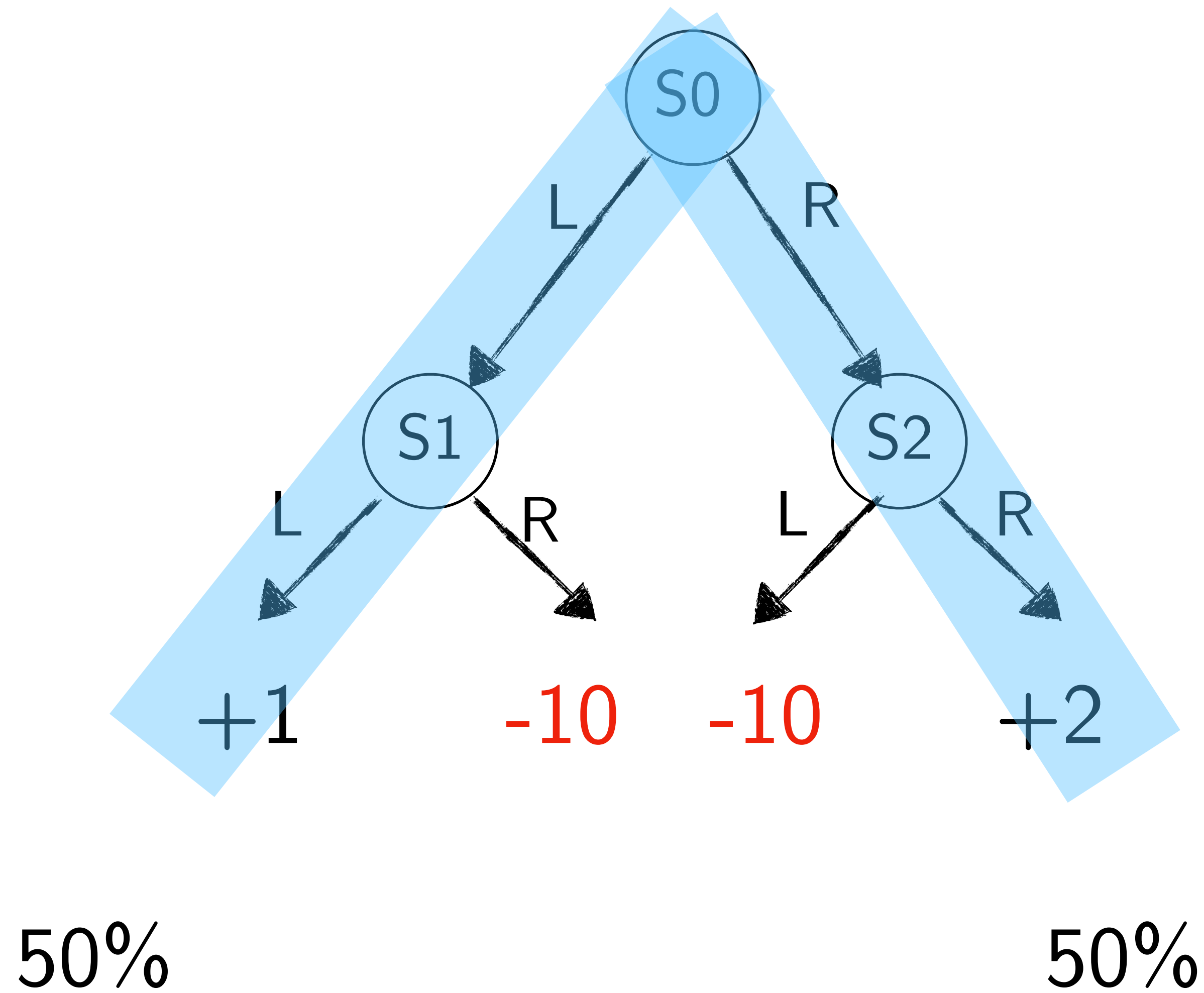
Consider the following MDP



Let's say I collected some data from the MDP



What would happen if I did Q-learning with this data?

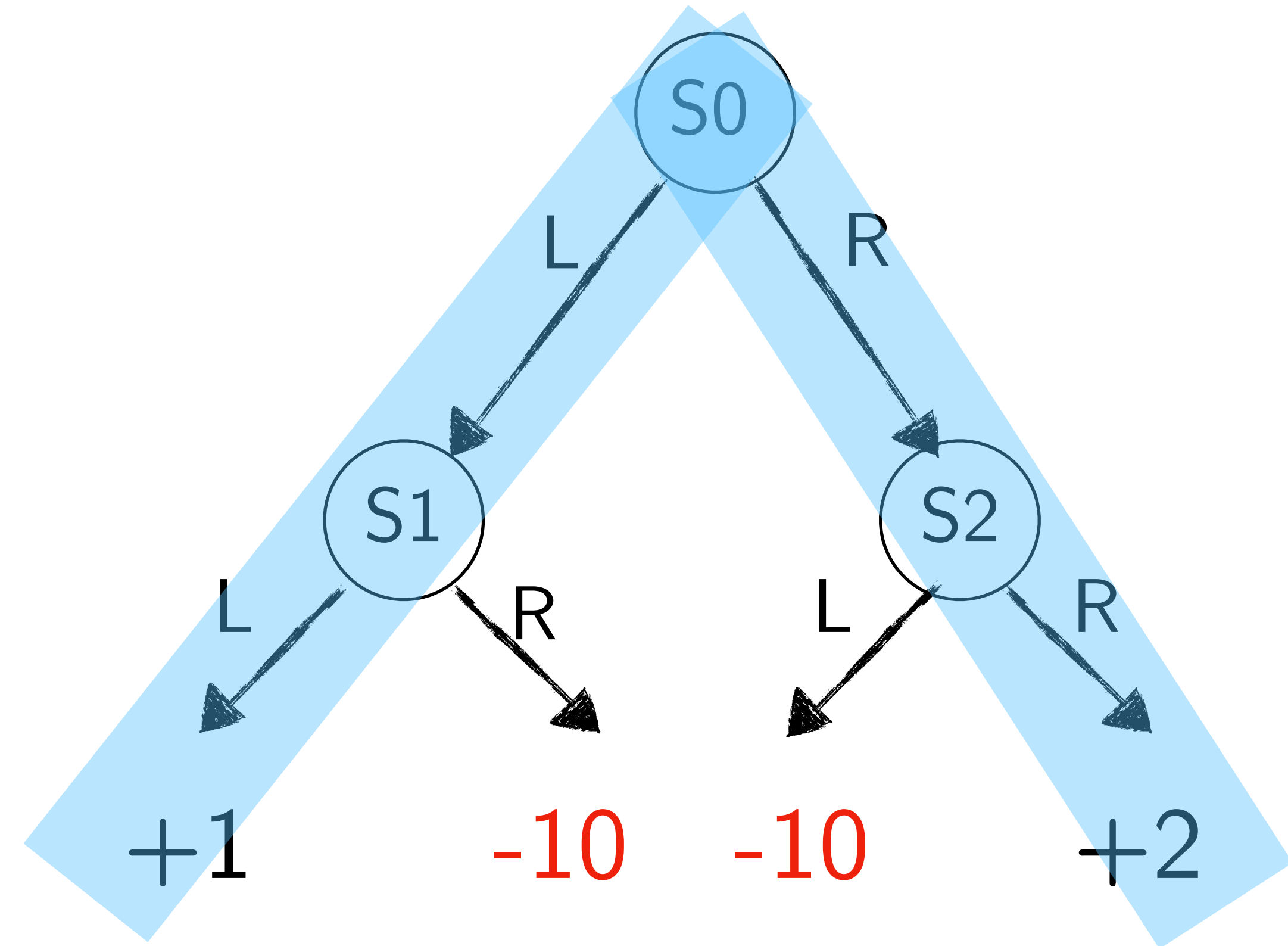


Think-Pair-Share!

Think (30 sec): What would happen if we did Q-learning with this data? Ideas on how to fix it.

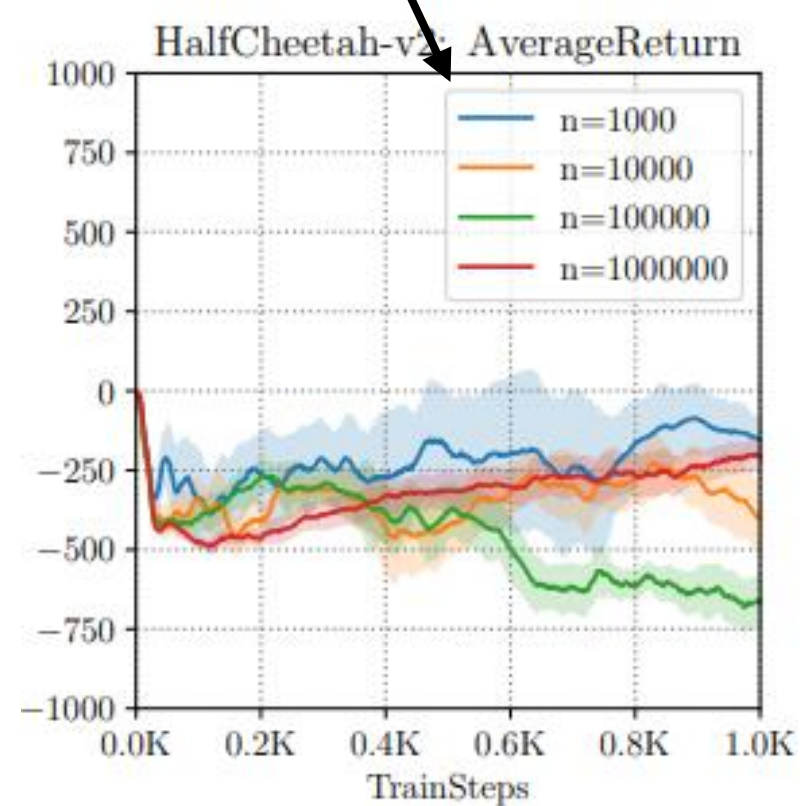
Pair: Find a partner

Share (45 sec): Partners exchange ideas



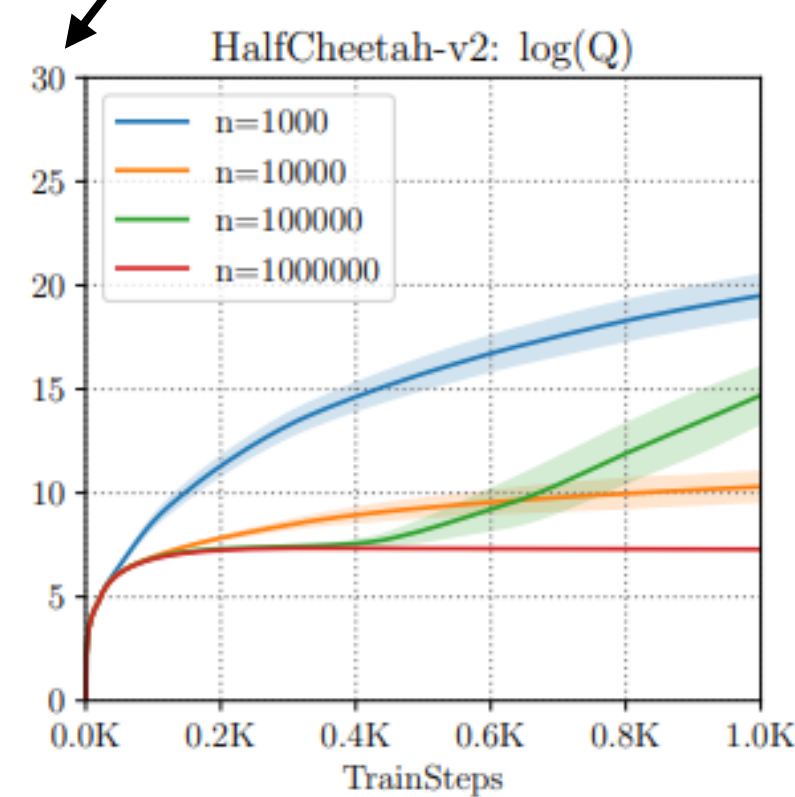
Why is offline RL hard?

amount of data



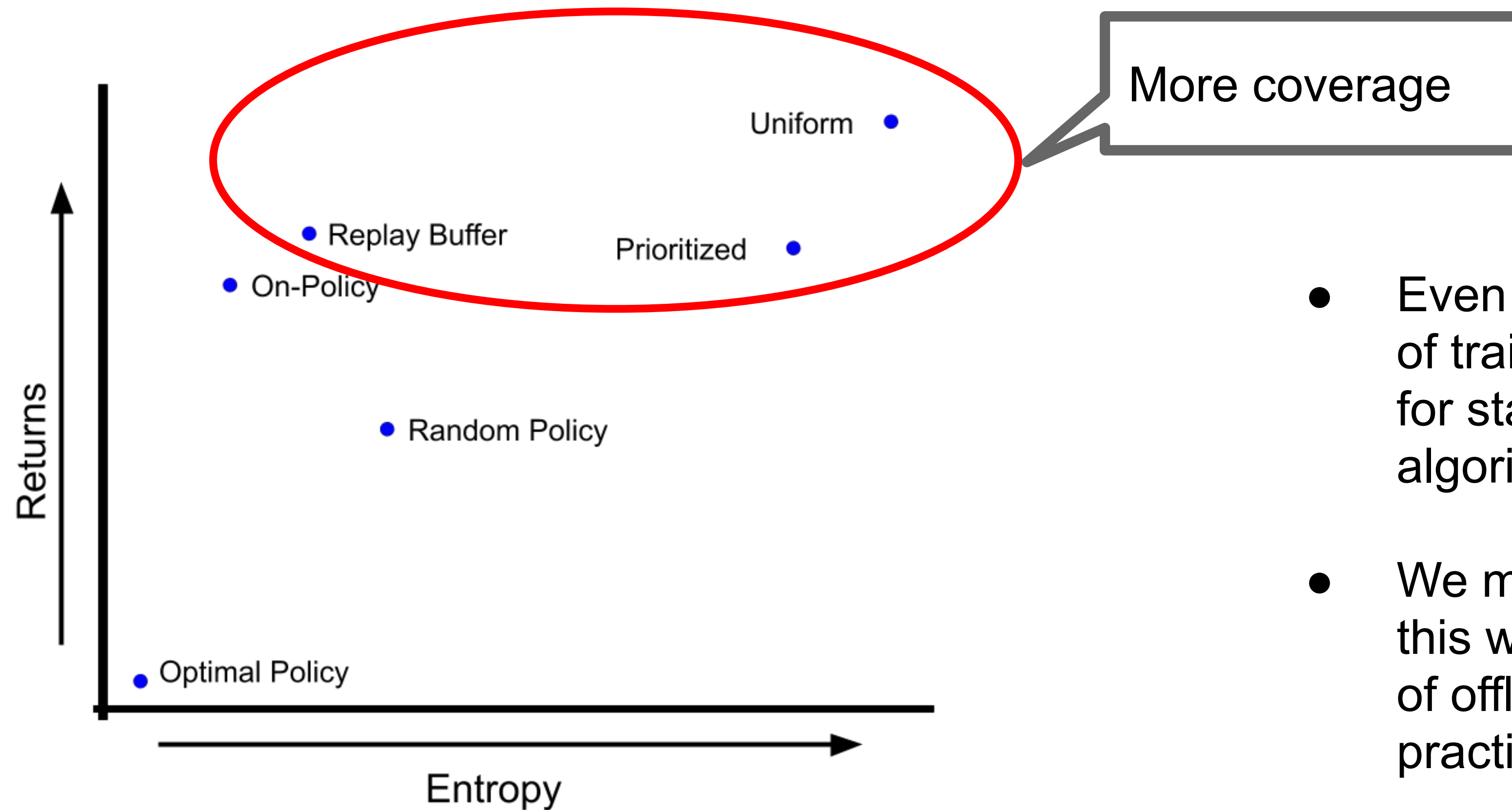
how well it does

log scale (massive overestimation)



how well it *thinks*
it does (Q-values)

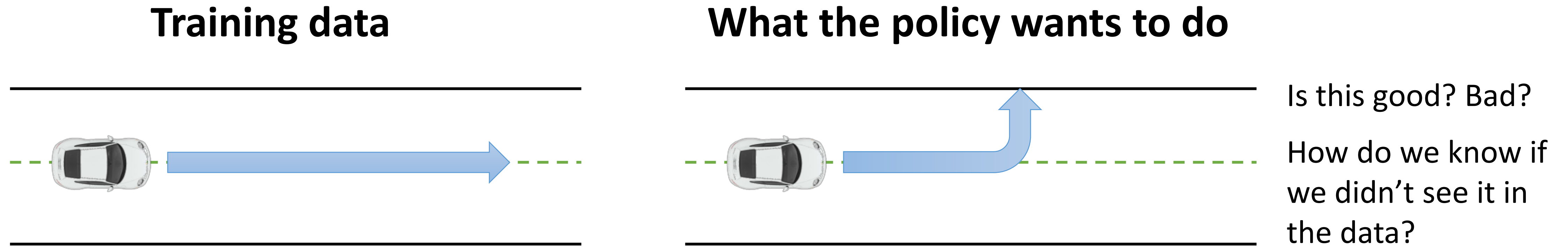
Illustration of Distribution Shift



- Even on gridworlds, the choice of training distribution matters for standard Q-learning algorithms.
- We might only wonder how this will affect the performance of offline RL algorithms in practice.

Why is offline RL hard?

Fundamental problem: counterfactual queries

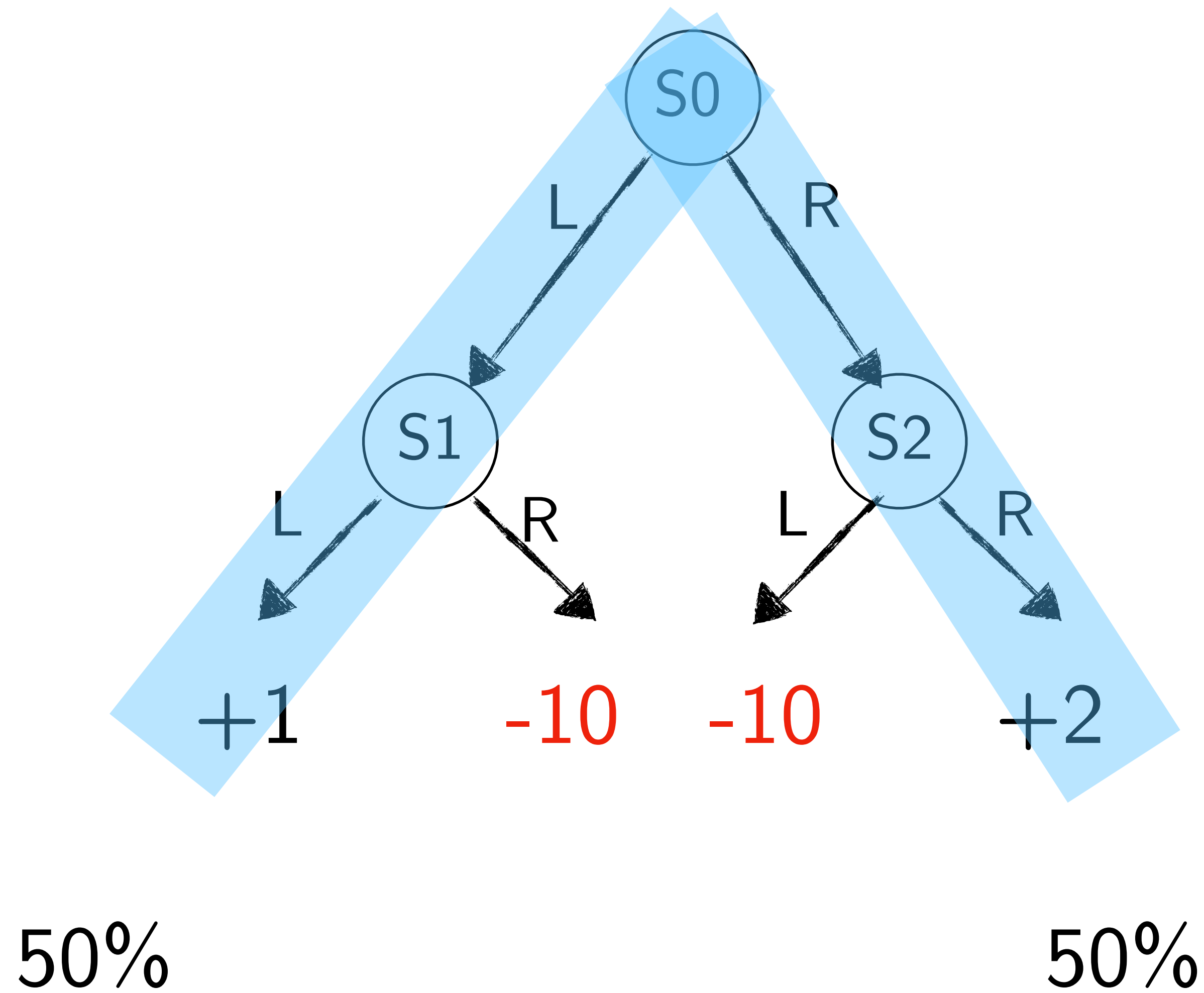


Online RL algorithms don't have to handle this, because they can simply **try** this action and see what happens

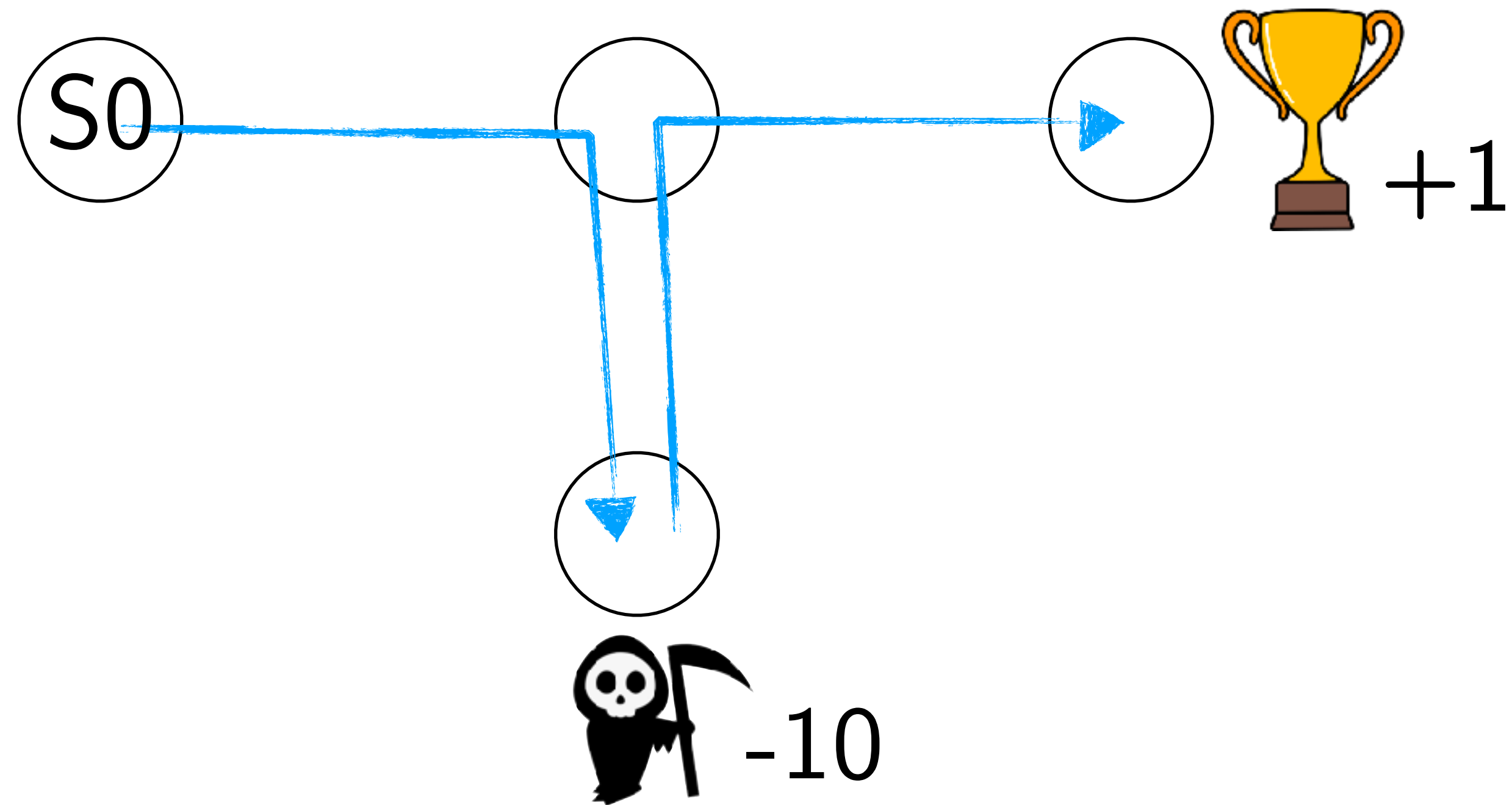
Offline RL methods must somehow account for these unseen ("out-of-distribution") actions, ideally in a safe way

...while still making use of generalization to come up with behaviors that are better than the best thing seen in the data!

Why not just do imitation learning?



Now consider this MDP

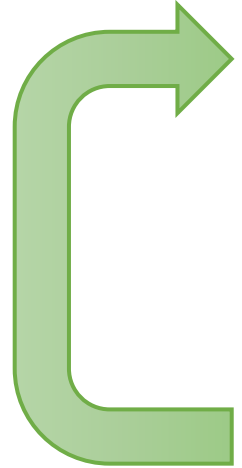


What is the optimal policy? What would imitation learning do?

Pessimism

Pessimism as a policy constraint

Don't deviate too much from the data collecting policy


$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\text{new}}} [Q(\mathbf{s}', \mathbf{a}')] \\ \pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \text{ s.t. } D_{\text{KL}}(\pi \parallel \pi_{\beta}) \leq \epsilon$$

Choose any divergence, e.g. KL!

TD3+BC: Most simple and effective offline RL!

A Minimalist Approach to Offline Reinforcement Learning

Scott Fujimoto^{1,2} Shixiang Shane Gu²

¹Mila, McGill University

²Google Research, Brain Team

scott.fujimoto@mail.mcgill.ca

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

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

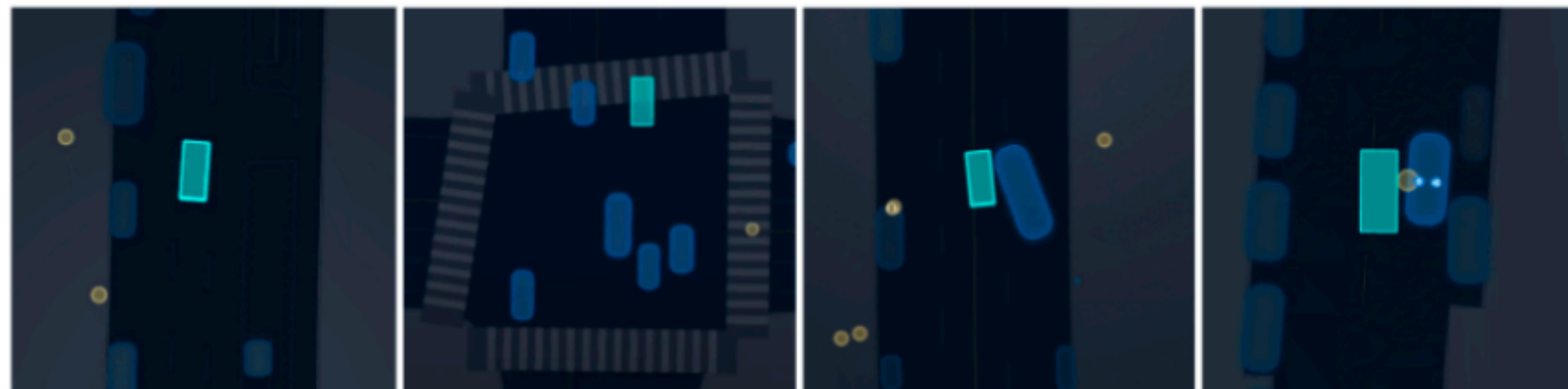
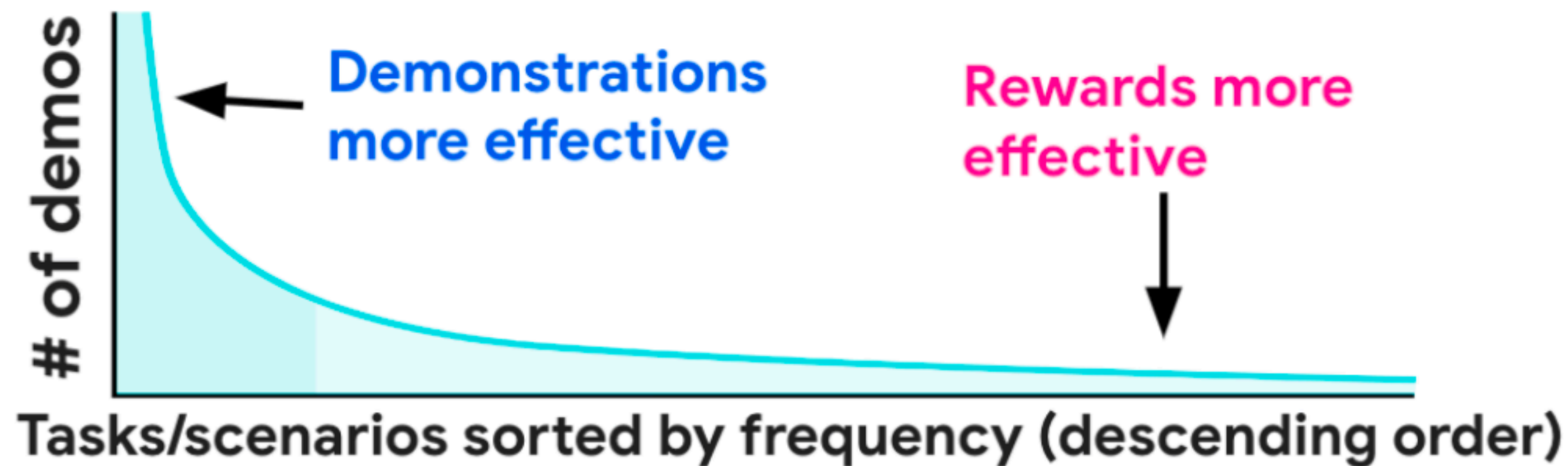
		BC	BRAC-p	AWAC	CQL	Fisher-BRC	TD3+BC
Random	HalfCheetah	2.0 \pm 0.1	23.5	2.2	21.7 \pm 0.9	32.2 \pm 2.2	10.2 \pm 1.3
	Hopper	9.5 \pm 0.1	11.1	9.6	10.7 \pm 0.1	11.4 \pm 0.2	11.0 \pm 0.1
	Walker2d	1.2 \pm 0.2	0.8	5.1	2.7 \pm 1.2	0.6 \pm 0.6	1.4 \pm 1.6
Medium	HalfCheetah	36.6 \pm 0.6	44.0	37.4	37.2 \pm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

Table 2: Average normalized score over the final 10 evaluations and 5 seeds. The highest performing scores are highlighted. CQL and Fisher-BRC are re-run using author-provided implementations to ensure an identical evaluation process, while BRAC and AWAC use previously reported results. \pm captures the standard deviation over seeds. TD3+BC achieves effectively the same performances as the state-of-the-art Fisher-BRC, despite being much simpler to implement and tune and more than halving the computation cost.

Works on real self-driving problems!

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

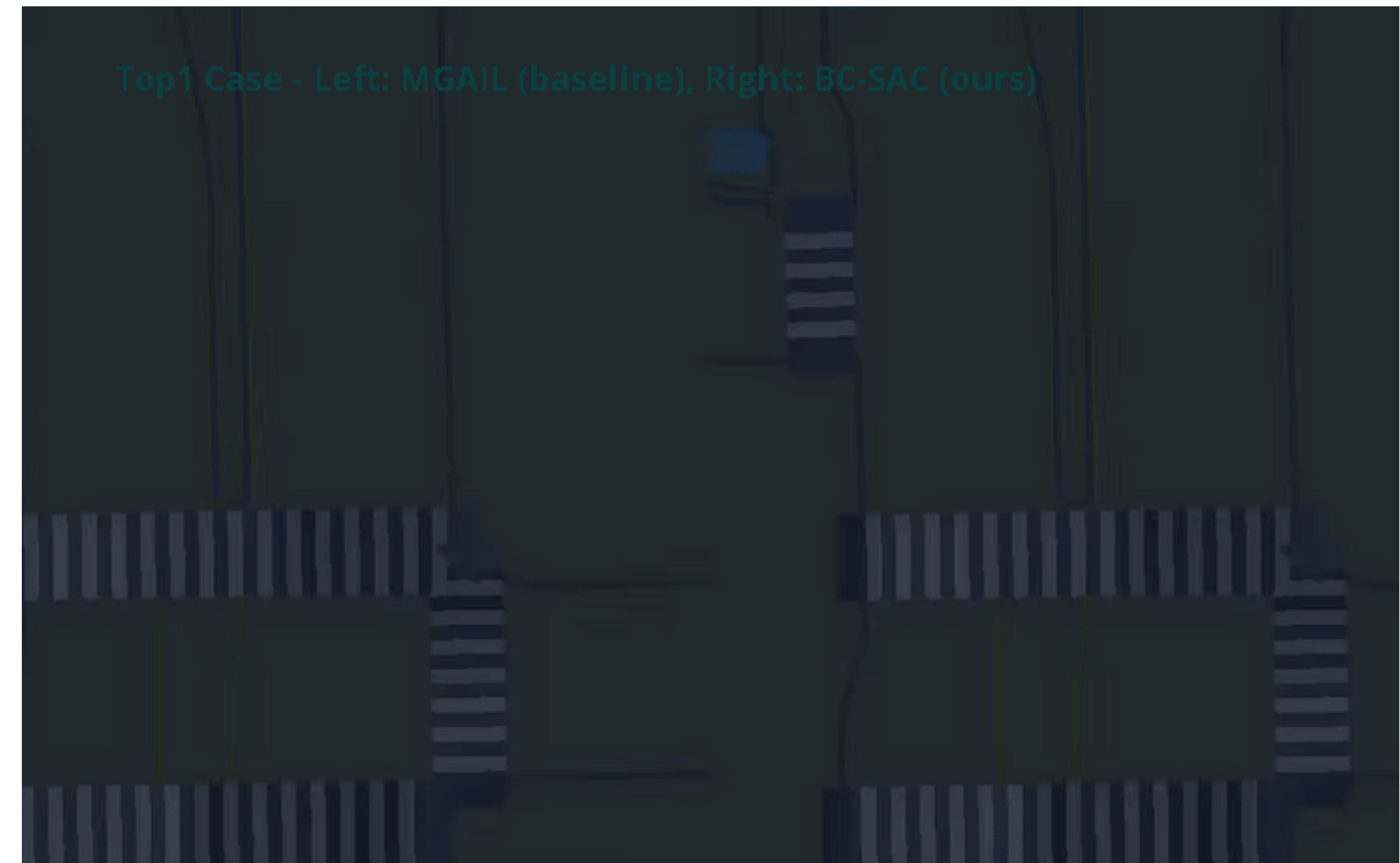
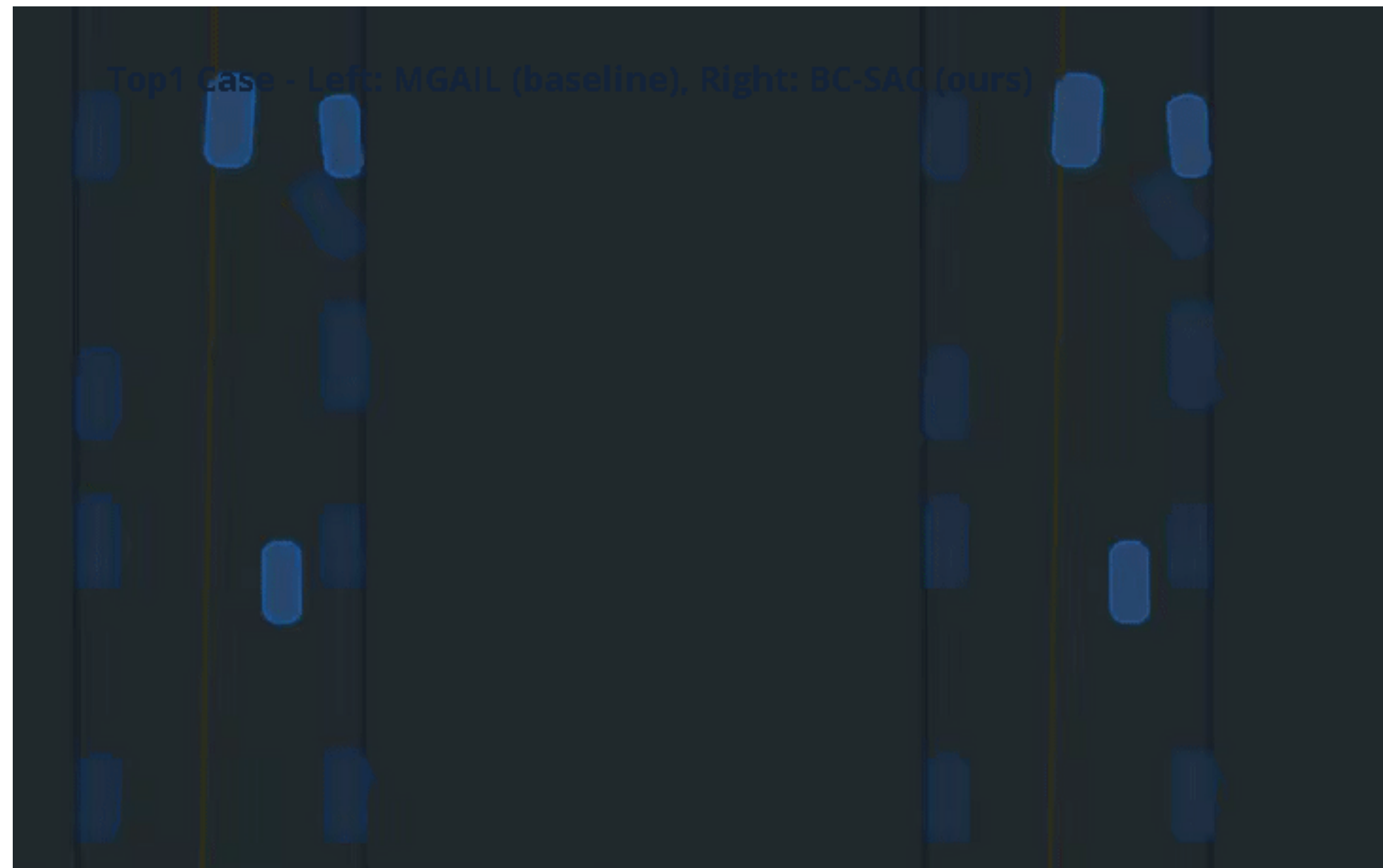
Yiren Lu¹, Justin Fu¹, George Tucker², Xinlei Pan¹, Eli Bronstein¹, Rebecca Roelofs², Benjamin Sapp¹,
Brandyn White¹, Aleksandra Faust², Shimon Whiteson¹, Dragomir Anguelov¹, Sergey Levine^{2,3}



Works on real self-driving problems!

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

Yiren Lu¹, Justin Fu¹, George Tucker², Xinlei Pan¹, Eli Bronstein¹, Rebecca Roelofs², Benjamin Sapp¹,
Brandyn White¹, Aleksandra Faust², Shimon Whiteson¹, Dragomir Anguelov¹, Sergey Levine^{2,3}



But choosing a divergence seems arbitrary?



↻ $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + E_{\mathbf{a}' \sim \pi_{\text{new}}} [Q(\mathbf{s}', \mathbf{a}')]$
 $\pi_{\text{new}}(\mathbf{a}|\mathbf{s}) = \arg \max_{\pi} E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})} [Q(\mathbf{s}, \mathbf{a})] \text{ s.t. } D_{\text{KL}}(\pi || \pi_{\beta}) \leq \epsilon$

$$D_{\text{KL}}(\pi || \pi_{\beta}) \leq \epsilon$$

Another notion of pessimism

Can we make the Q-value itself pessimistic
on actions it has not seen?

Conservative Q-Learning (CQL)



Conservative Q-Learning for Offline Reinforcement Learning

Aviral Kumar¹, Aurick Zhou¹, George Tucker², Sergey Levine^{1,2}

¹UC Berkeley, ²Google Research, Brain Team

aviralk@berkeley.edu

Conservative Q-Learning (CQL)



1. Learn \hat{Q}^π using offline data \mathcal{D} .
2. Optimize policy w.r.t. $\hat{Q}^\pi : \pi \leftarrow \arg \max_{\pi} \mathbb{E}_{\pi}[\hat{Q}^\pi] - \alpha D(\pi_{\phi}, \pi_{\beta})$

Approach 2: Directly modify the Q-function to be pessimistic

- Key idea behind CQL: Learn lower-bounds on Q-values

CQL Algorithm:

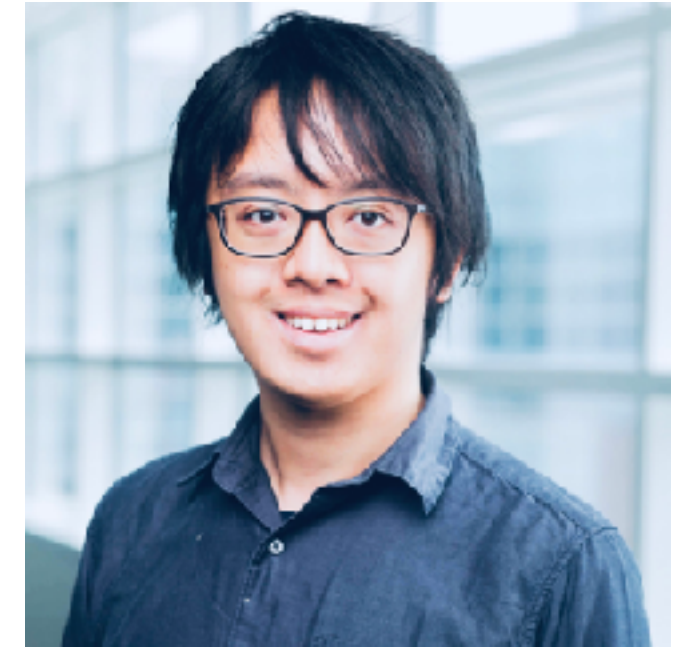
1. Learn \hat{Q}_{CQL}^π using offline data \mathcal{D} .
2. Optimize policy w.r.t. $\hat{Q}_{\text{CQL}}^\pi : \pi \leftarrow \arg \max_{\pi} \mathbb{E}_{\pi}[\hat{Q}_{\text{CQL}}^\pi]$.

Many ways to construct a conservative Q value

Original CQL paper proposed one such way

Recent work has come up with more unified frameworks

Adversarially Trained Actor Critic (ATACL)



Adversarially Trained Actor Critic for Offline Reinforcement Learning

Ching-An Cheng^{*1} Tengyang Xie^{*2} Nan Jiang² Alekh Agarwal³

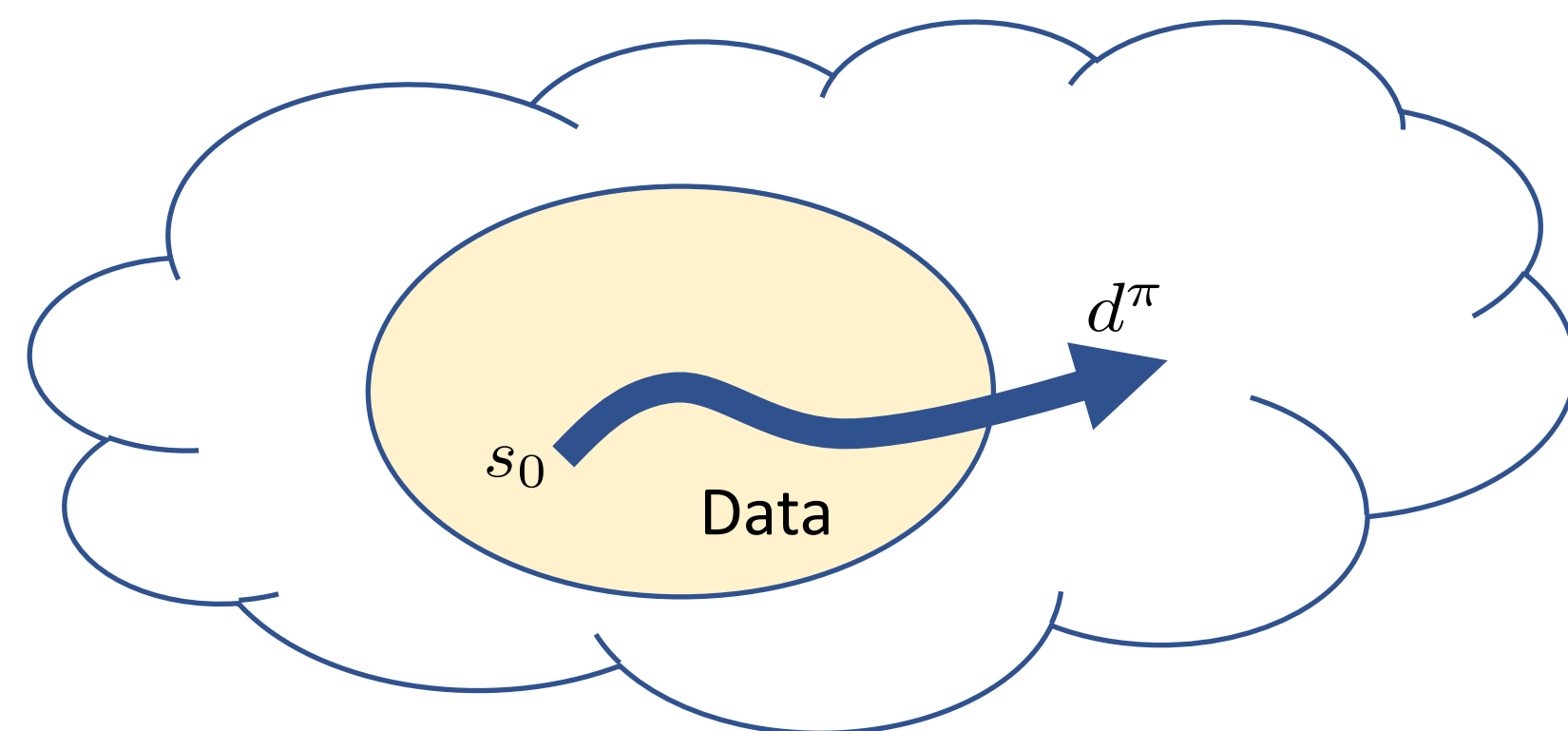


Key Idea: Relative Pessimism

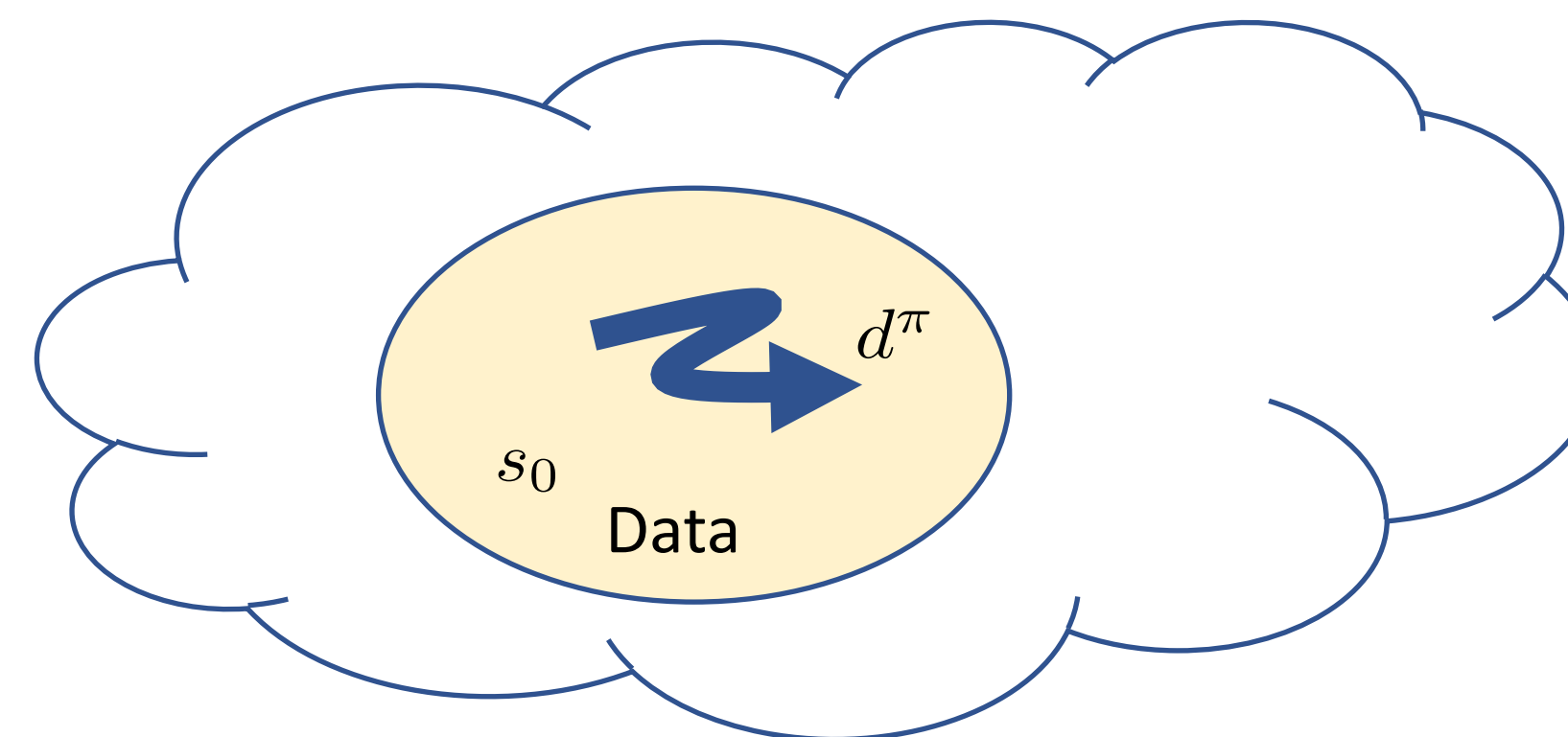
- Optimize for the **worst-case** performance compared with the **behavior policy μ** .

$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \text{Lower bound of } J(\pi) - J(\mu)$$

Lower bound $< J(\pi) - J(\mu)$



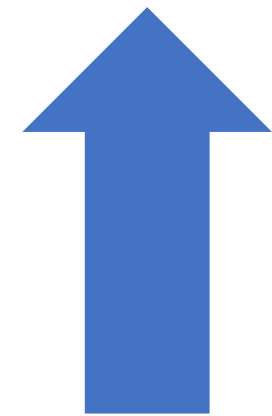
Lower bound $\approx J(\pi) - J(\mu)$



Key Idea: Relative Pessimism

- Optimize for the **best worst-case** performance compared with the **behavior policy μ** .

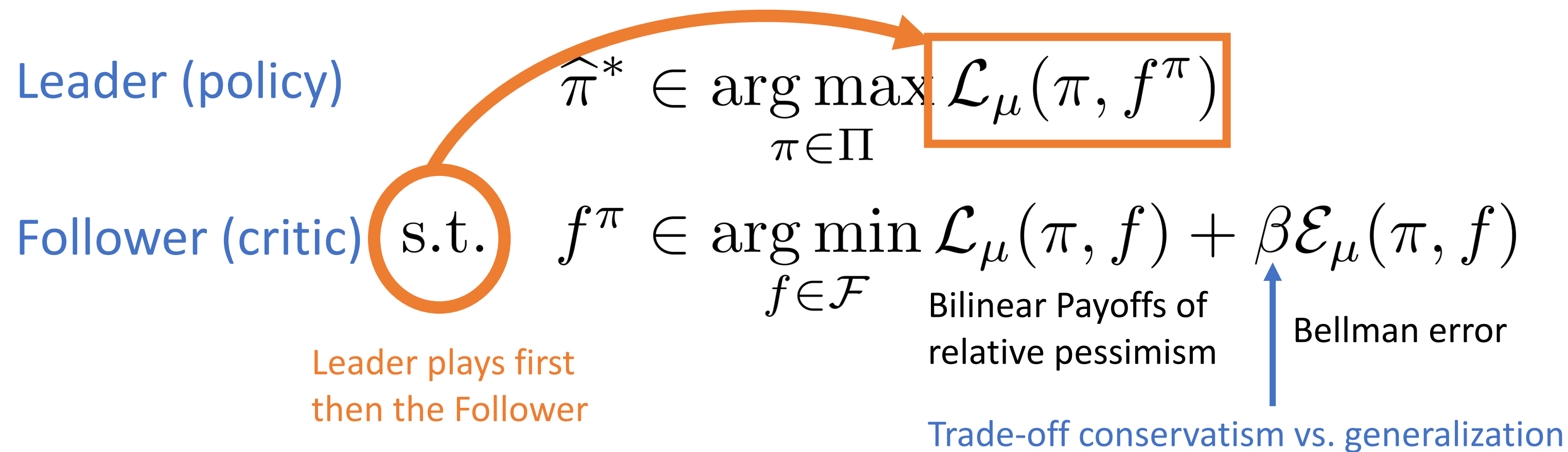
$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \text{Lower bound of } J(\pi) - J(\mu)$$



ATAC frames this problem as a Stackelberg game (i.e., bilevel optimization)

A Stackelberg Game for Offline RL

- ATAC optimizes for relative pessimism via solving a Stackelberg game



$$\mathcal{L}_\mu(\pi, f) := \mathbb{E}_\mu[f(s, \pi) - f(s, a)]$$

$$\mathcal{E}_\mu(\pi, f) := \mathbb{E}_\mu[((f - \mathcal{T}^\pi f)(s, a))^2].$$

A Stackelberg Game for Offline RL

- ATAC optimizes for relative pessimism via solving a Stackelberg game

Leader (policy)

$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi) \geq \mathcal{L}(\pi, Q^\pi) \equiv J(\pi) - J(\mu), \forall \beta \geq 0$$

Follower (critic)

s.t.

$$f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$$

Bilinear Payoffs of relative pessimism

Bellman error

Trade-off conservatism vs. generalization

Leader plays first then the Follower

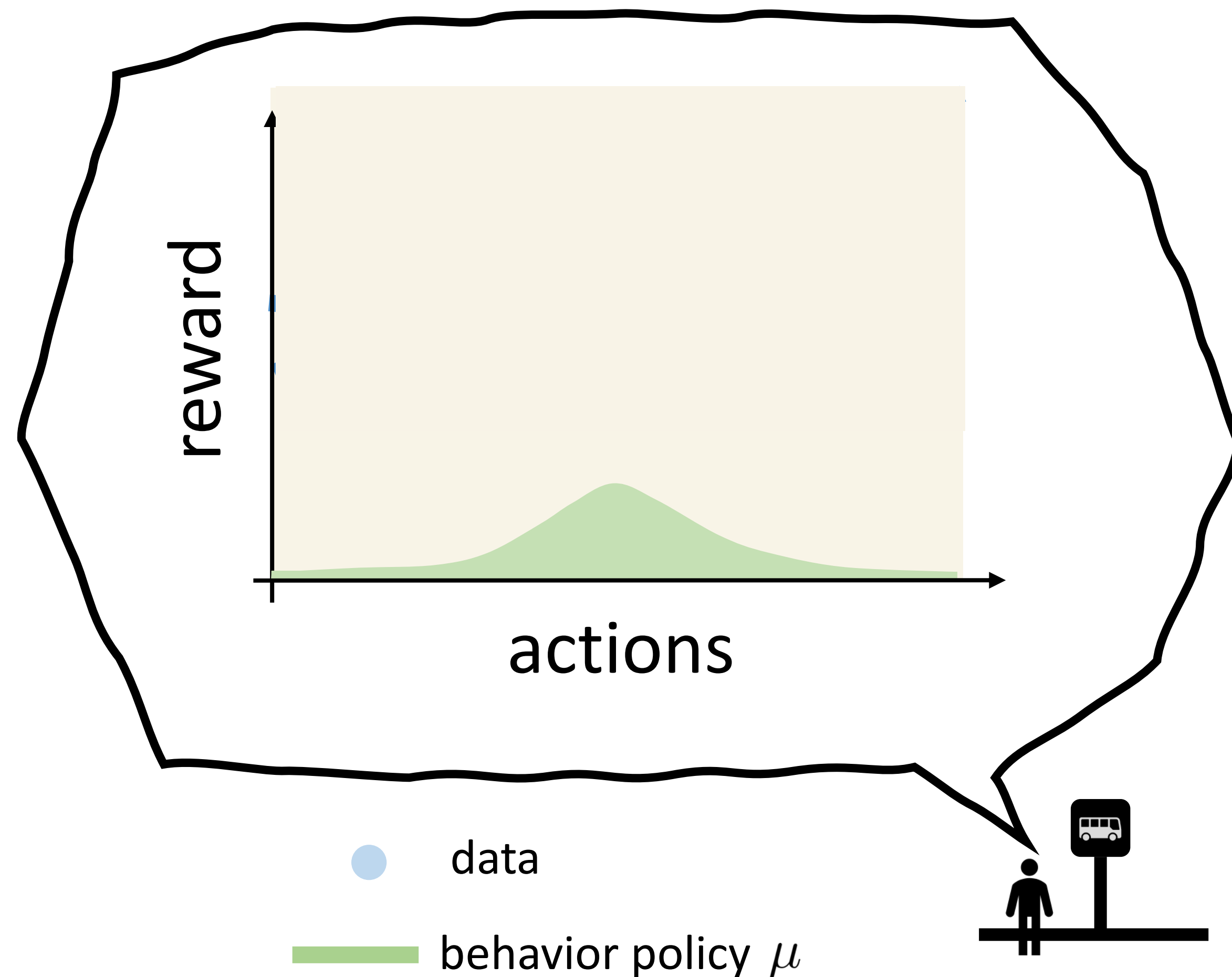
Robust Policy Improvement Property

For all $\beta \geq 0$, the ATAC policy is always no worse than the behavior policy that collected the data.

$$\begin{aligned} \mathcal{L}_\mu(\pi, f) &:= \mathbb{E}_\mu[f(s, \pi) - f(s, a)] \\ \mathcal{E}_\mu(\pi, f) &:= \mathbb{E}_\mu[((f - \mathcal{T}^\pi f)(s, a))^2]. \end{aligned}$$

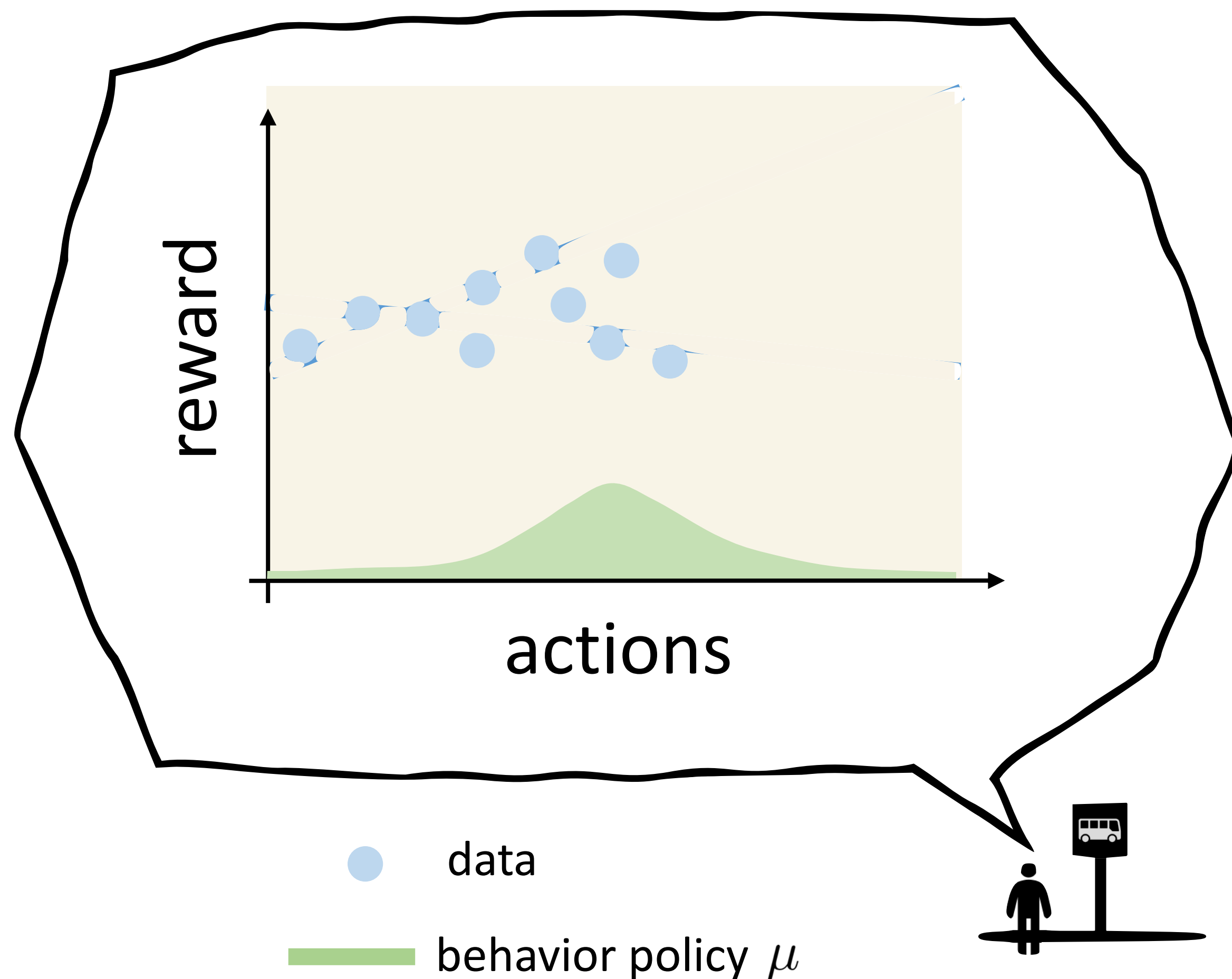
Let's look at a simple example

Let's say the time horizon $T=1$
(Multi-armed bandit!)

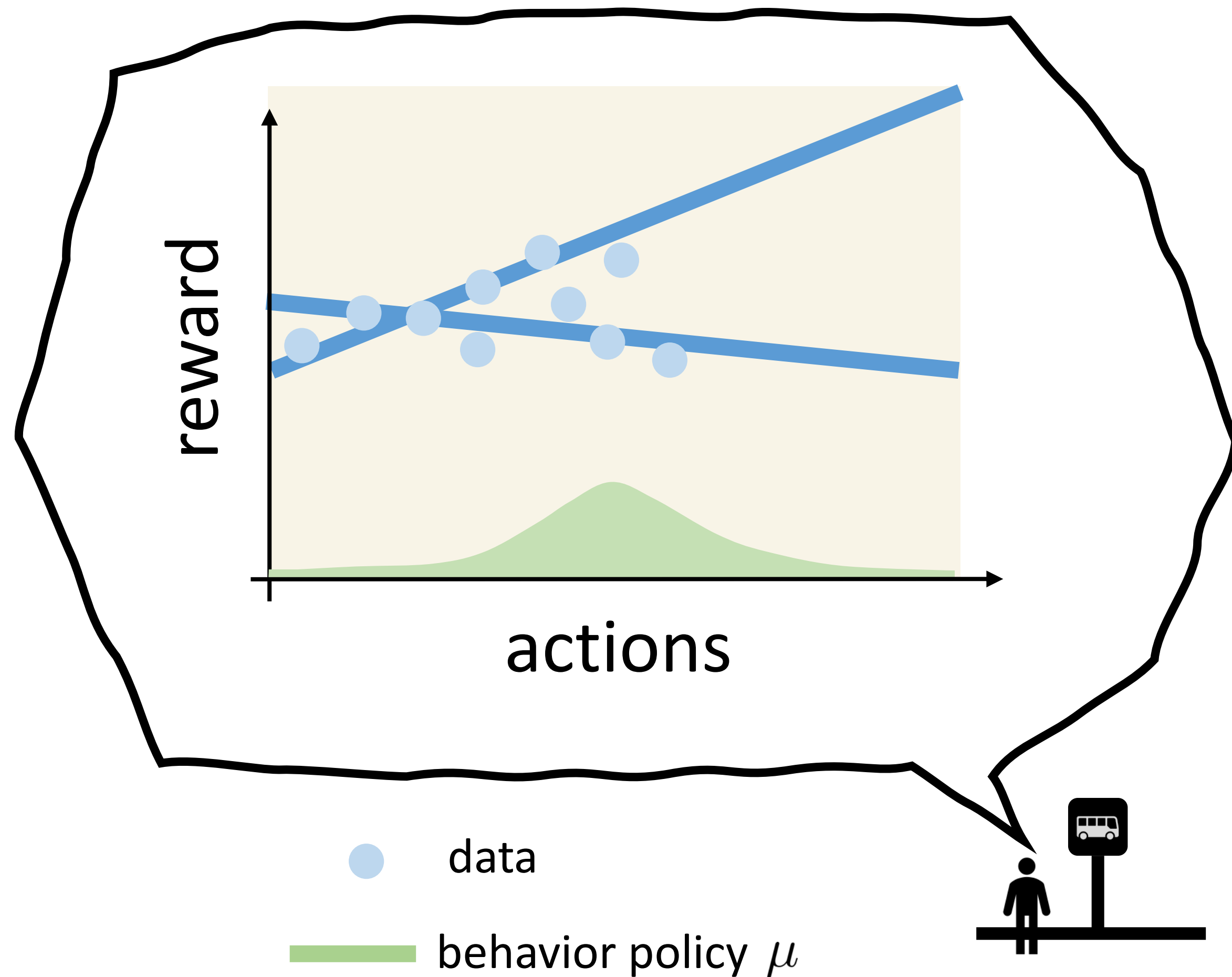


Let's look at a simple example

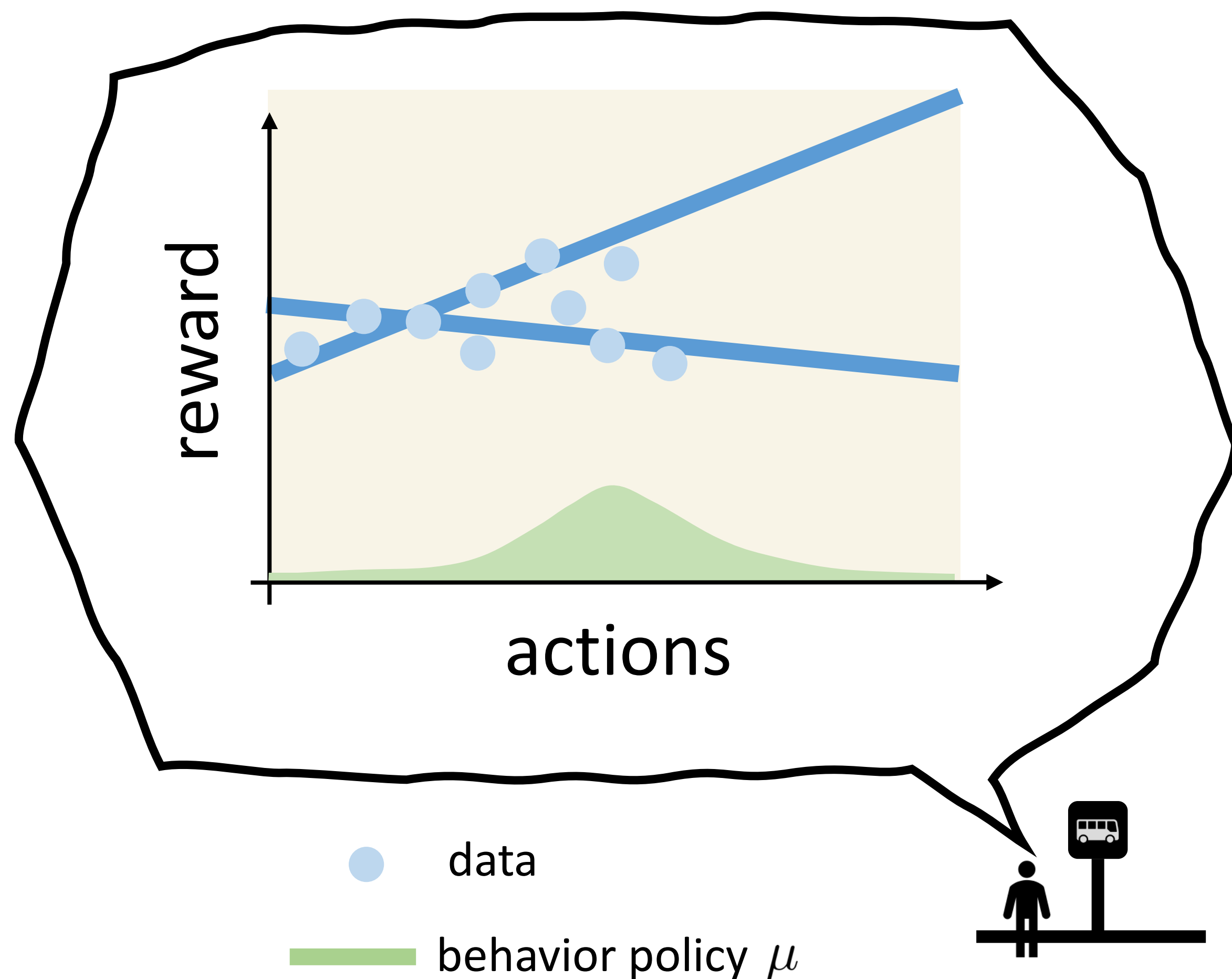
Let's say the time horizon $T=1$
(Multi-armed bandit!)



Let's look at a simple example



A Stackelberg Game for Offline RL



ATAC

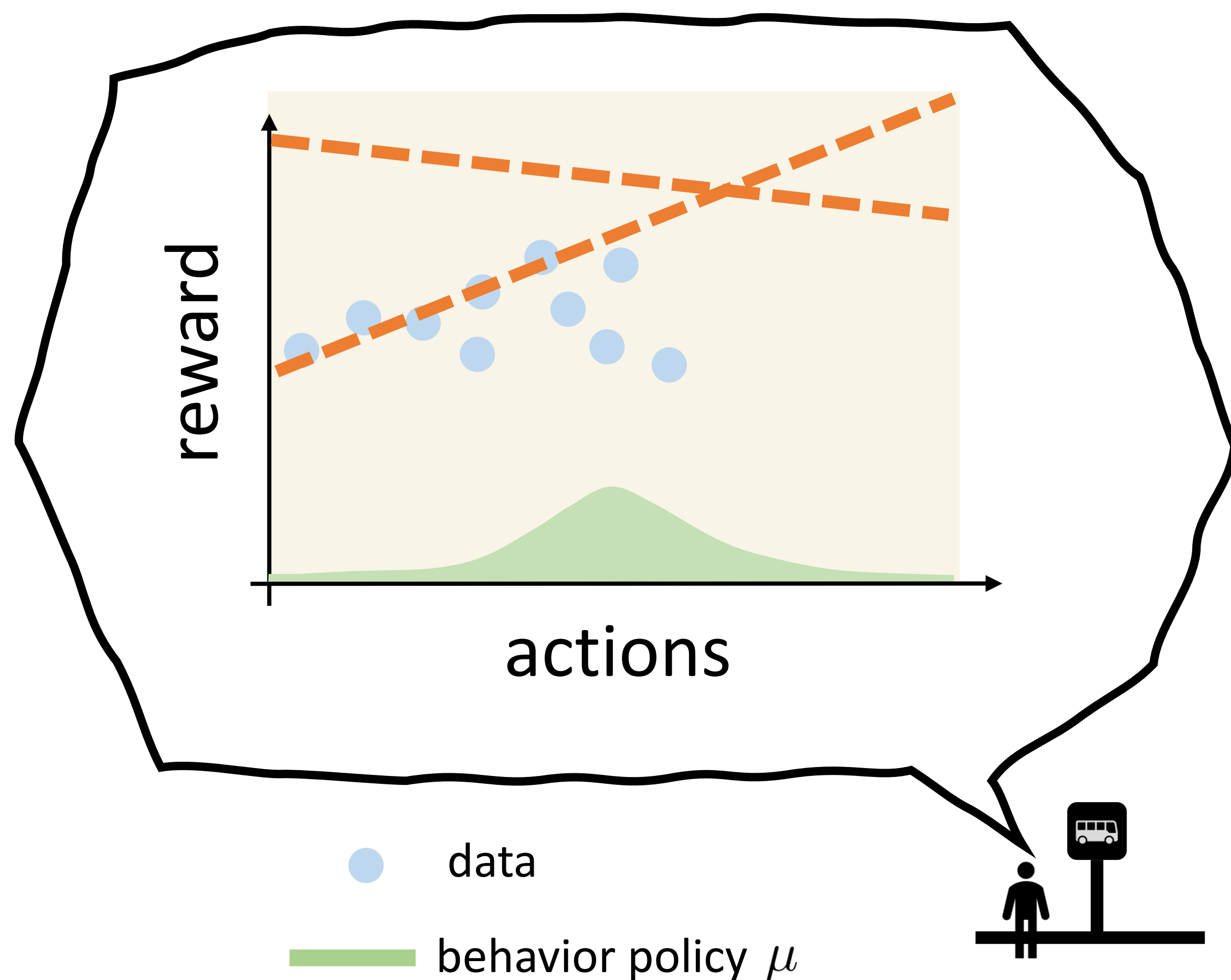
$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi)$$

s.t. $f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$

— hypothesis $f(s, \cdot)$ with small $\beta \mathcal{E}_\mu$

Functions that are consistent with the reward and the dynamics on the behavior data

A Stackelberg Game for Offline RL



ATAC

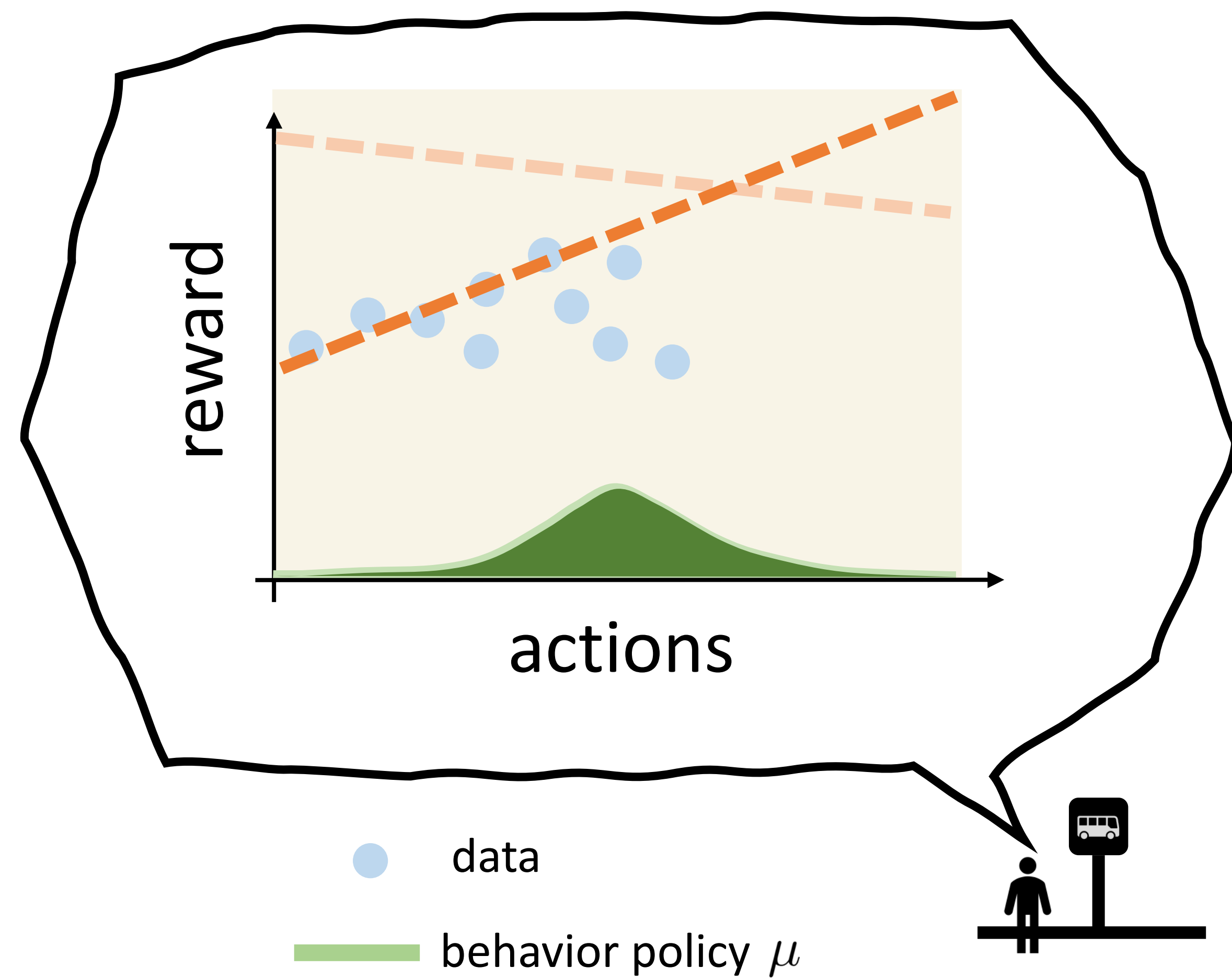
$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi)$$
$$\text{s.t. } f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$$

— value difference hypothesis $f(s, \cdot) - f(s, \mu)$

Functions shifted from the original hypotheses

What is the solution to the Stackelberg game?

A Stackelberg Game for Offline RL



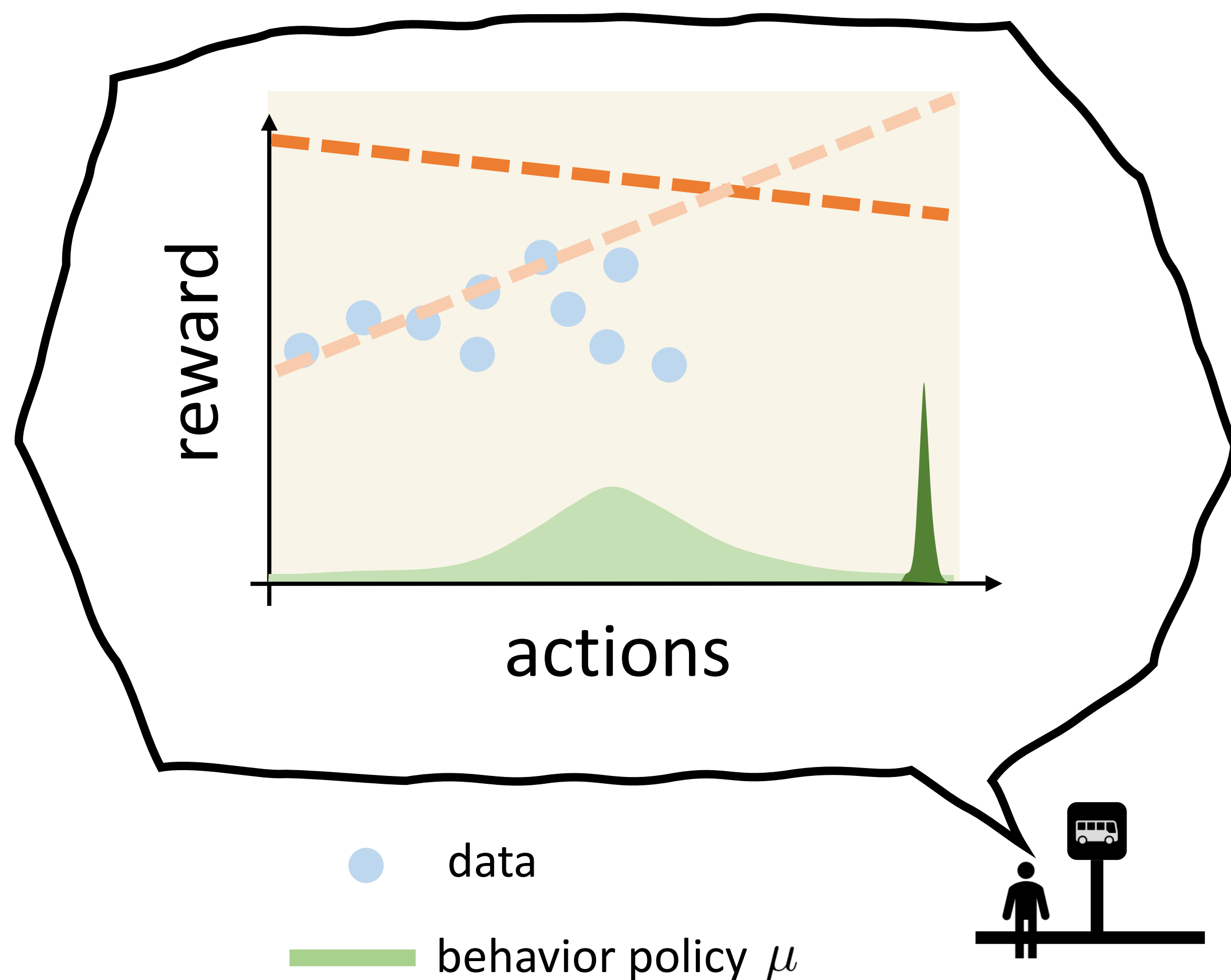
ATAC

$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi)$$
$$\text{s.t. } f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$$

- value difference hypothesis $f(s, \cdot) - f(s, \mu)$
- inactive value difference hypothesis
- decision policy π

Not the behavior policy in this case...

A Stackelberg Game for Offline RL



ATAC

$$\begin{aligned} \hat{\pi}^* &\in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi) \\ \text{s.t. } f^\pi &\in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f) \end{aligned}$$

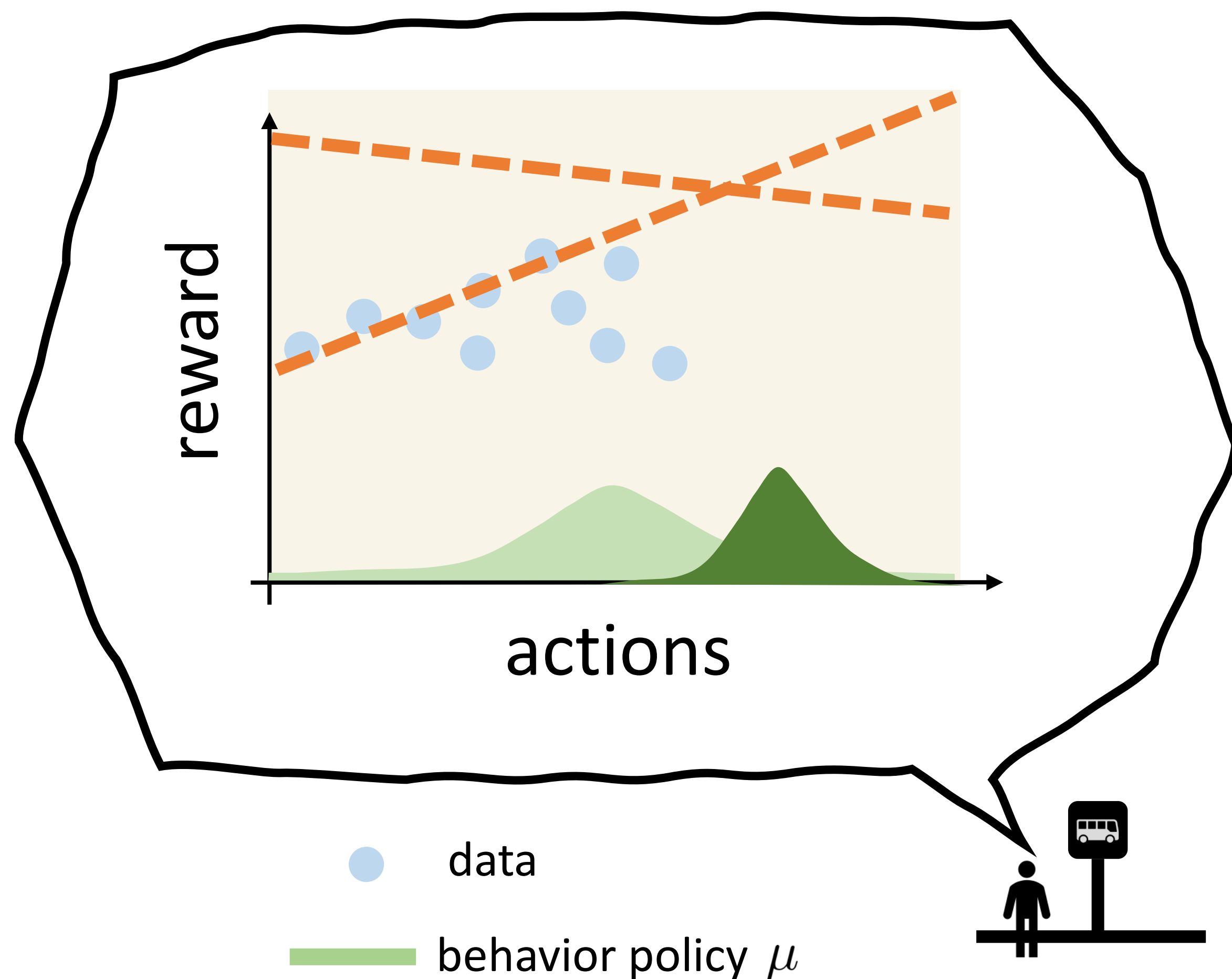
— value difference hypothesis $f(s, \cdot) - f(s, \mu)$

— inactive value difference hypothesis

— decision policy π




***Not the policy that maximizes
a single hypothesis...***

A Stackelberg Game for Offline RL



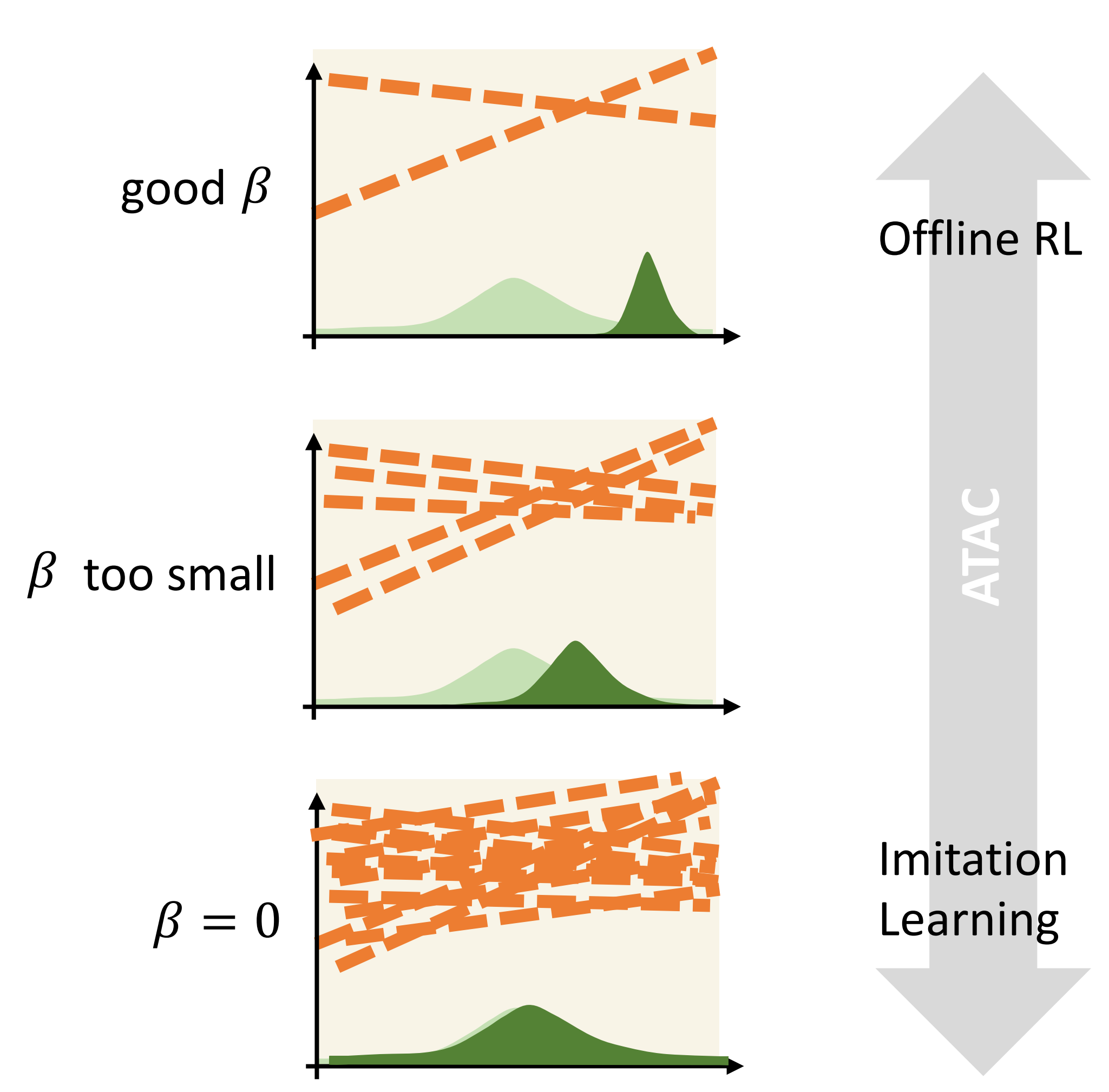
ATAC

$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi)$$
$$\text{s.t. } f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$$

-  value difference hypothesis $f(s, \cdot) - f(s, \mu)$
-  inactive value difference hypothesis
-  decision policy π

The optimal decision balances multiple hypotheses

A Stackelberg Game for Offline RL



Leader = Actor = Conditional generator
Follower = Critic = Discriminator

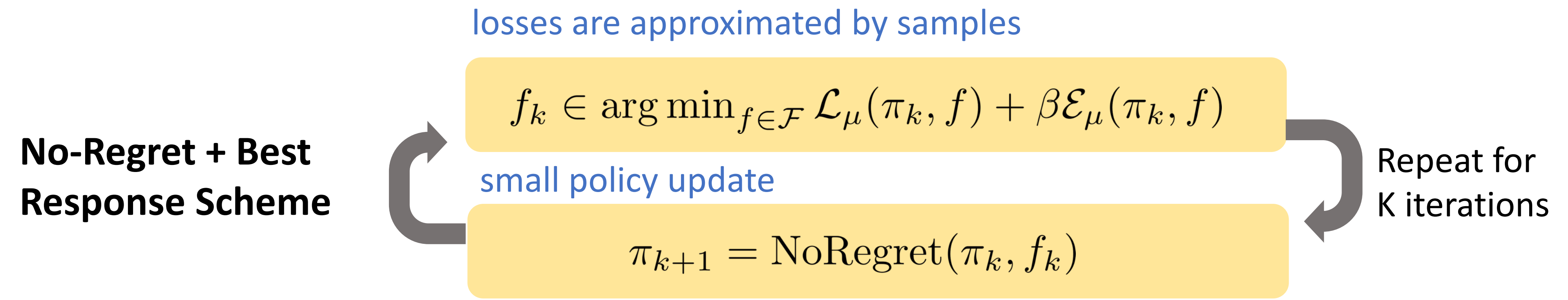
ATAC

$$\hat{\pi}^* \in \arg \max_{\pi \in \Pi} \mathcal{L}_\mu(\pi, f^\pi)$$
$$\text{s.t. } f^\pi \in \arg \min_{f \in \mathcal{F}} \mathcal{L}_\mu(\pi, f) + \beta \mathcal{E}_\mu(\pi, f)$$

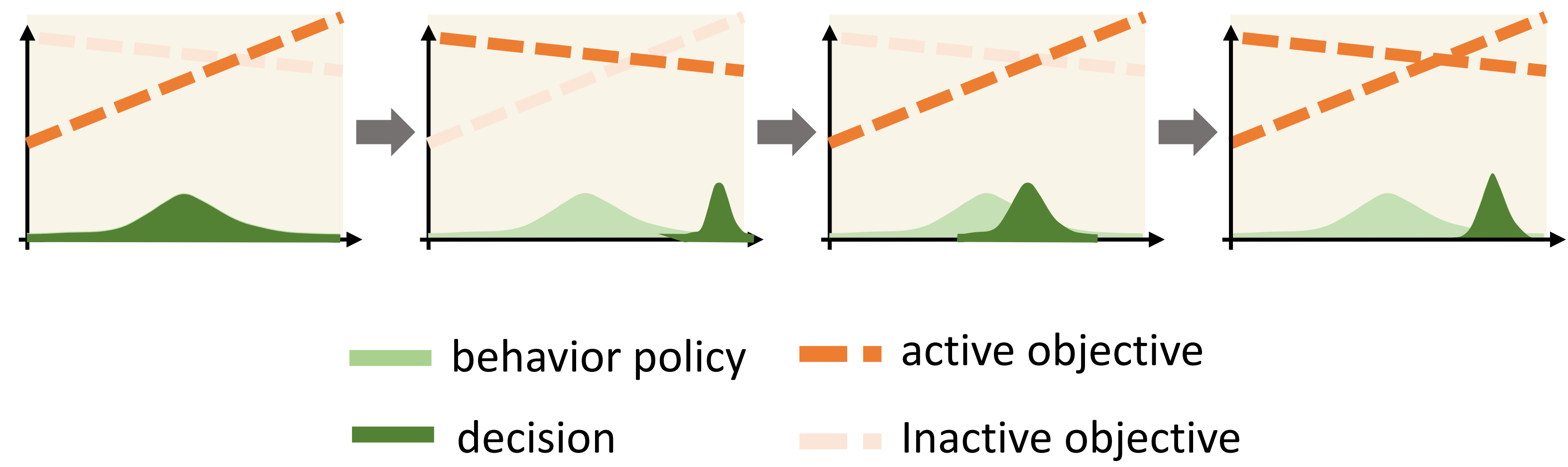
ATAC provides a bridge between offline RL and imitation learning with IPM via the lens of generative adversarial networks (GAN)

Offline RL + Relative Pessimism = IL + Bellman Regularization

Solving the Stackelberg Game



Output uniform mixture of policies (theory) or the last policy (practice)
In practice, the above is implemented by two-timescale SGD updates



ATAC Theory (Informal)

Learning Optimality

Assume \mathcal{F} satisfies realizability and completeness.

Given dataset \mathcal{D} s.t. $|\mathcal{D}| = N$. With $\beta = \Theta(N^{2/3})$. Then $\forall \pi \in \Pi$,

$$J(\pi) - J(\hat{\pi}) \leq \mathcal{O}\left(\frac{1}{(1-\gamma)N^{1/3}}\right) + \epsilon_{\text{generalization}}(\mathcal{F}, \pi, \mathcal{D})$$

average Bellman error of f_t on
the distribution of π

*With a **well tuned** β , ATAC can compete with any policy within the data coverage.*

ATAC Theory (Informal)

Robust Policy Improvement

Assume \mathcal{F} satisfies **realizability** without the need of completeness.

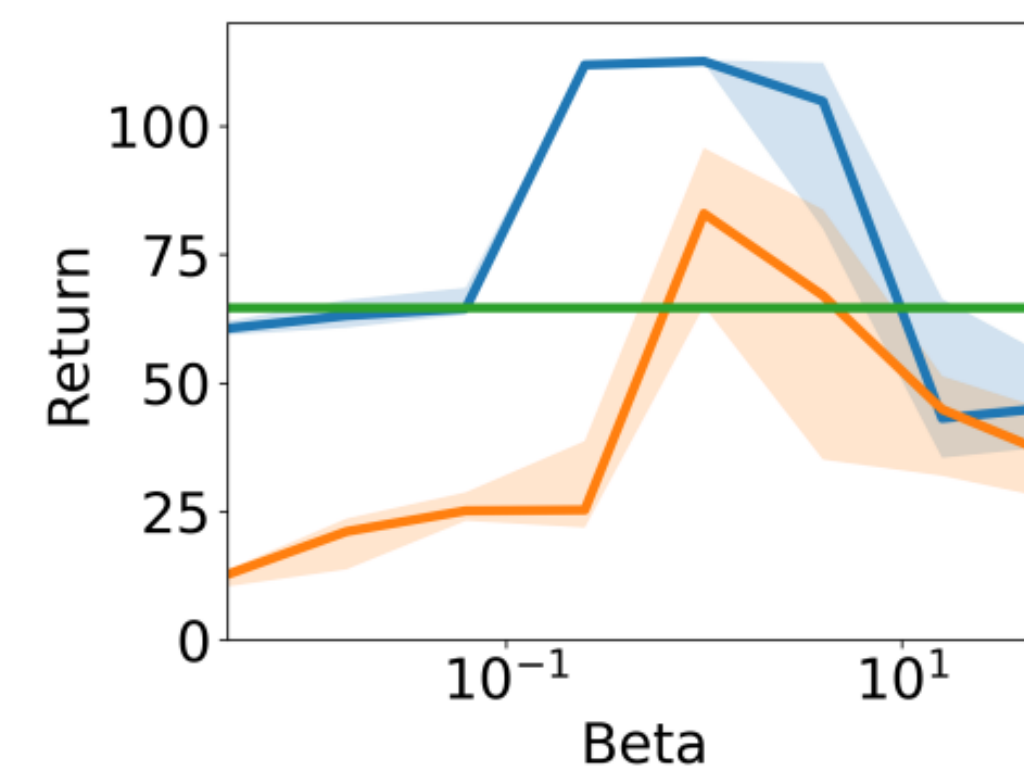
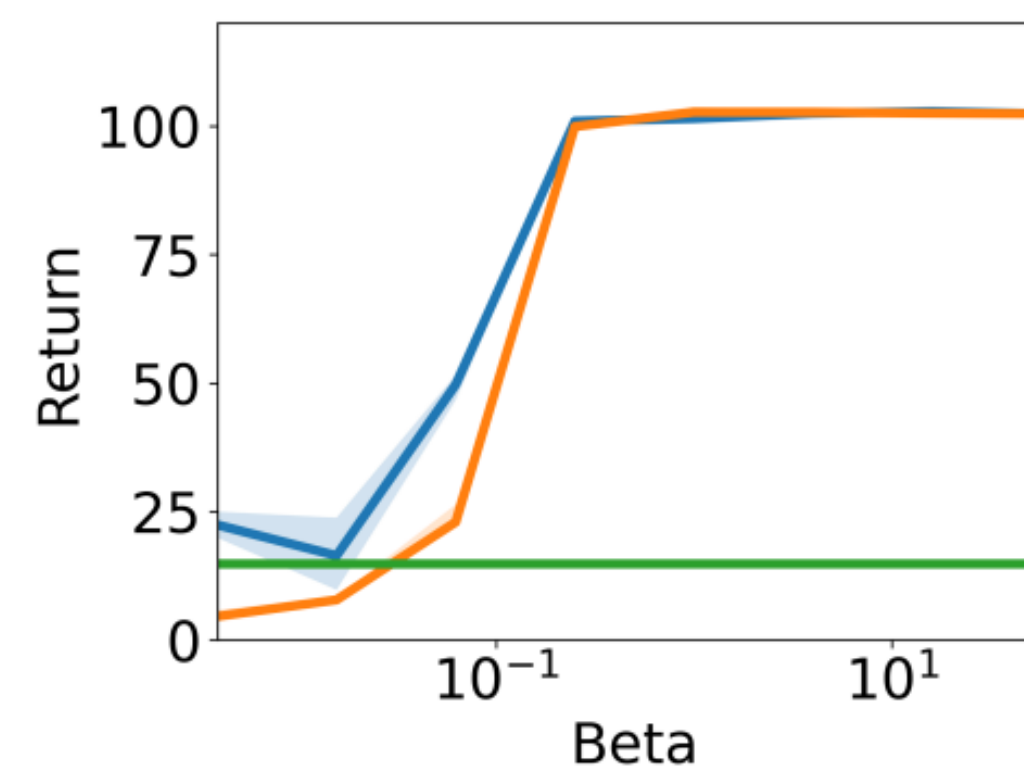
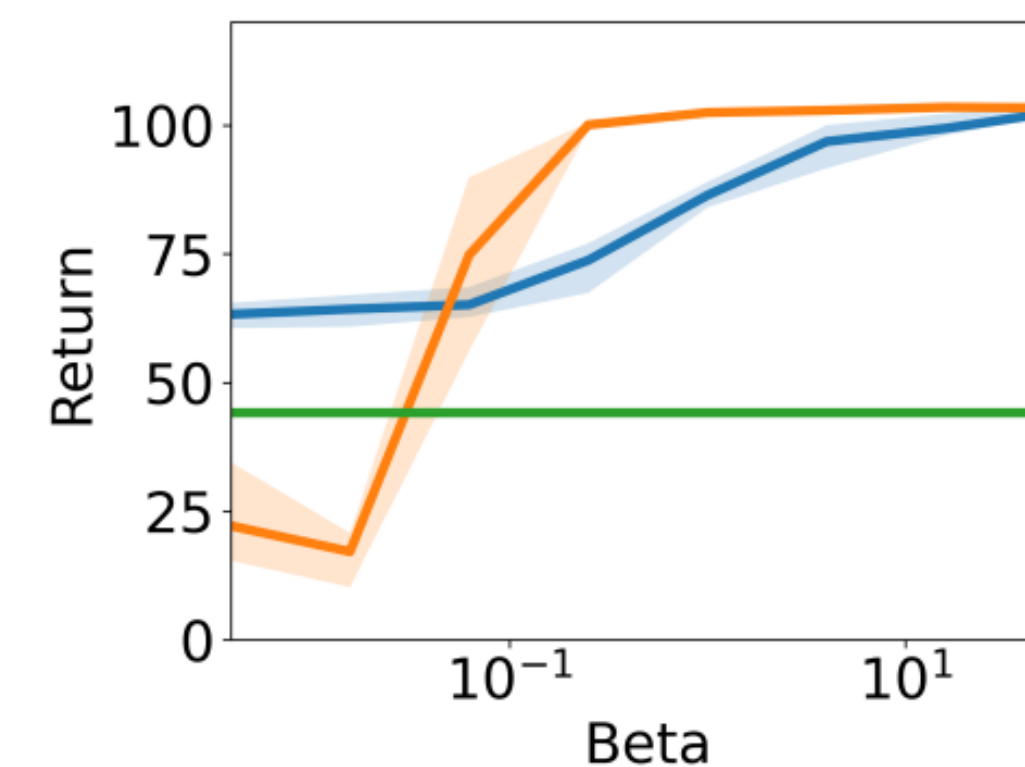
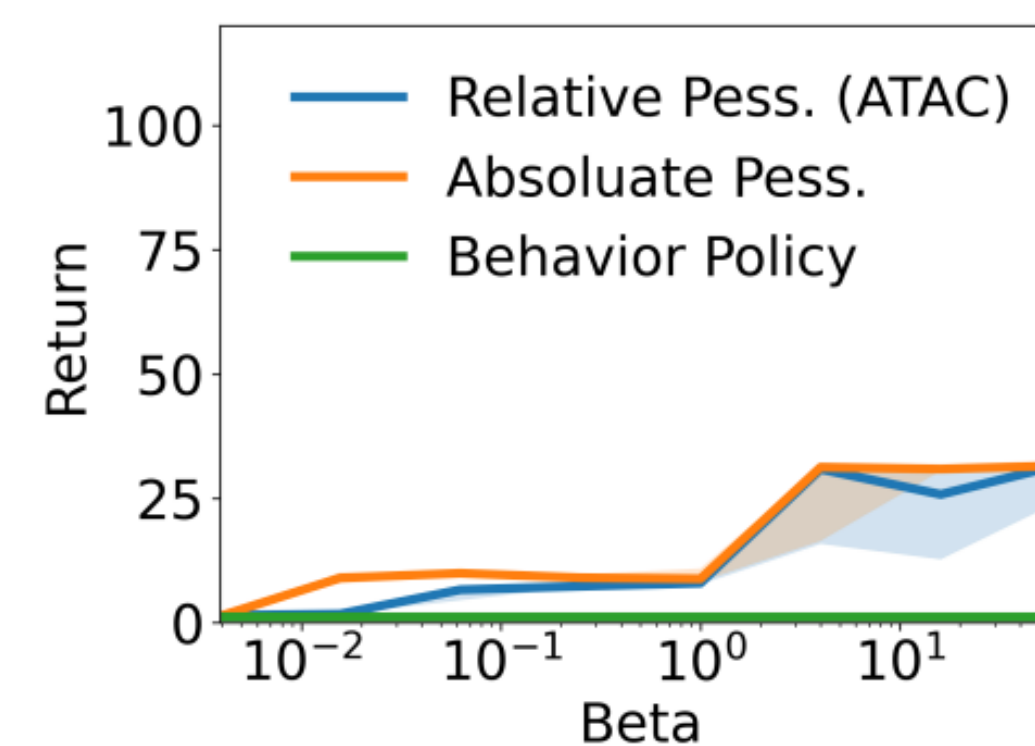
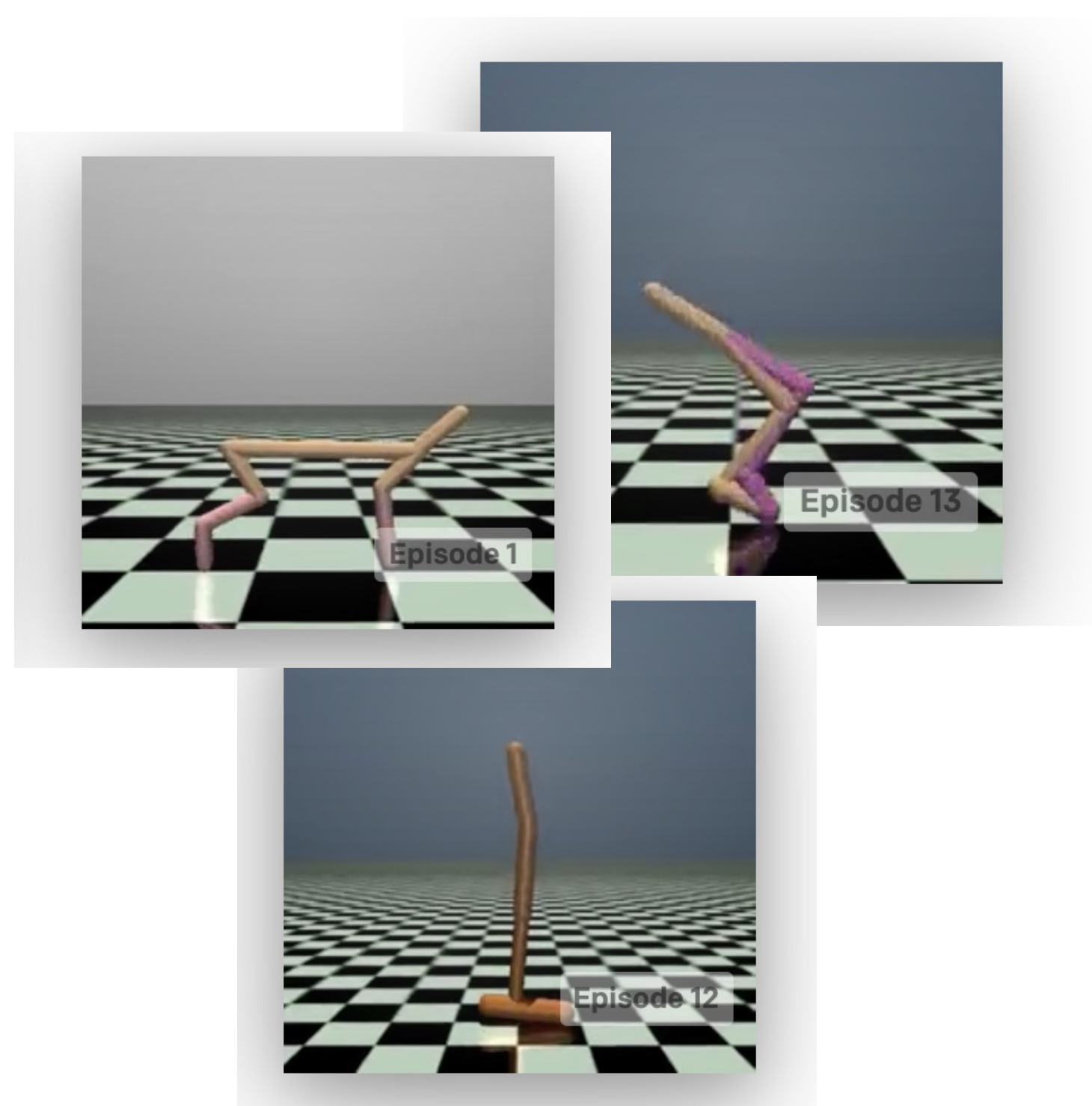
$$\text{If } \mu \in \Pi, \text{ then } J(\mu) - J(\bar{\pi}) \leq \mathcal{O} \left(\underbrace{\frac{1}{(1-\gamma)N^{1/2}}}_{\text{faster rate}} + \frac{\beta}{(1-\gamma)N} \right)$$

ATAC always improves over the behavior policy so long as $\beta = o(N)$.

Experimental Results

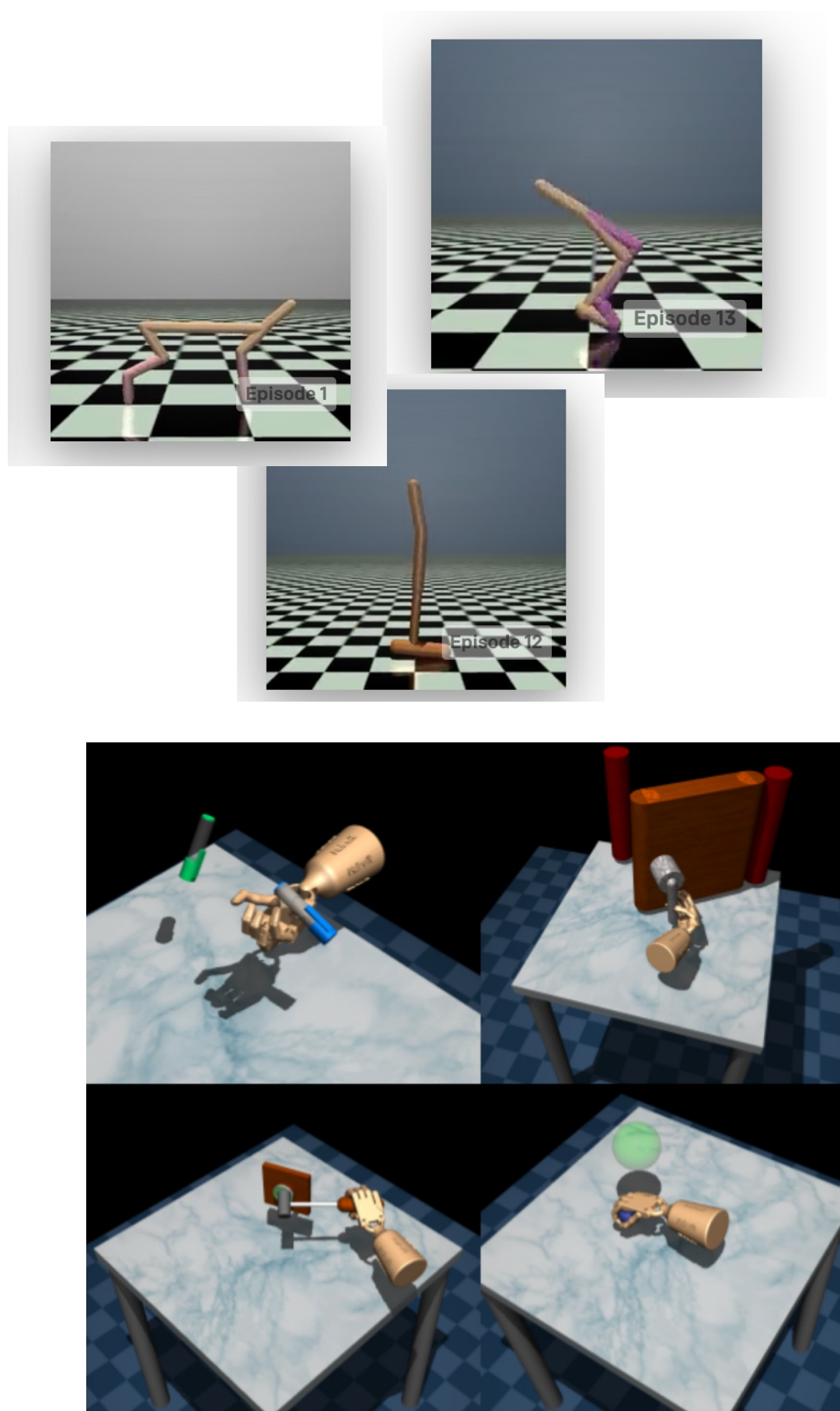
Robust Policy Improvement

ATAC's robustness property enables online HP selection. We can gradually increase β to tune its performance without breaking the baseline performance.



Experimental Results

ATAC achieves SOTA performance, outperforming baseline algorithms in most datasets



	Behavior	ATAC*	CQL	COMBO	TD3BC	IQL	BC
halfcheetah-rand	-0.1	4.8	35.4	38.8	10.2	-	2.1
walker2d-rand	0.0	8.0	7.0	7.0	1.4	-	1.6
hopper-rand	1.2	31.8	10.8	17.9	11.0	-	9.8
halfcheetah-med	40.6	54.3	44.4	54.2	42.8	47.4	36.1
walker2d-med	62.0	91.0	74.5	75.5	79.7	78.3	6.6
hopper-med	44.2	102.8	86.6	94.9	99.5	66.3	29.0
halfcheetah-med-replay	27.1	49.5	46.2	55.1	43.3	44.2	38.4
walker2d-med-replay	14.8	94.1	32.6	56.0	25.2	73.9	11.3
hopper-med-replay	14.9	102.8	48.6	73.1	31.4	94.7	11.8
halfcheetah-med-exp	64.3	95.5	62.4	90.0	97.9	86.7	35.8
walker2d-med-exp	82.6	116.3	98.7	96.1	101.1	109.6	6.4
hopper-med-exp	64.7	112.6	111.0	111.1	112.2	91.5	111.9
pen-human	207.8	79.3	37.5	-	-	71.5	34.4
hammer-human	25.4	6.7	4.4	-	-	1.4	1.5
door-human	28.6	8.7	9.9	-	-	4.3	0.5
relocate-human	86.1	0.3	0.2	-	-	0.1	0.0
pen-cloned	107.7	73.9	39.2	-	-	37.3	56.9
hammer-cloned	8.1	2.3	2.1	-	-	2.1	0.8
door-cloned	12.1	8.2	0.4	-	-	1.6	-0.1
relocate-cloned	28.7	0.8	-0.1	-	-	-0.2	-0.1
pen-exp	105.7	159.5	107.0	-	-	-	85.1
hammer-exp	96.3	128.4	86.7	-	-	-	125.6
door-exp	100.5	105.5	101.5	-	-	-	34.9
relocate-exp	101.6	106.5	95.0	-	-	-	101.3

Datasets where ATAC is the best performing algorithm, with 9% improvement (median) compared with the best baseline algorithm.