# Alignment and Active Learning in HRI
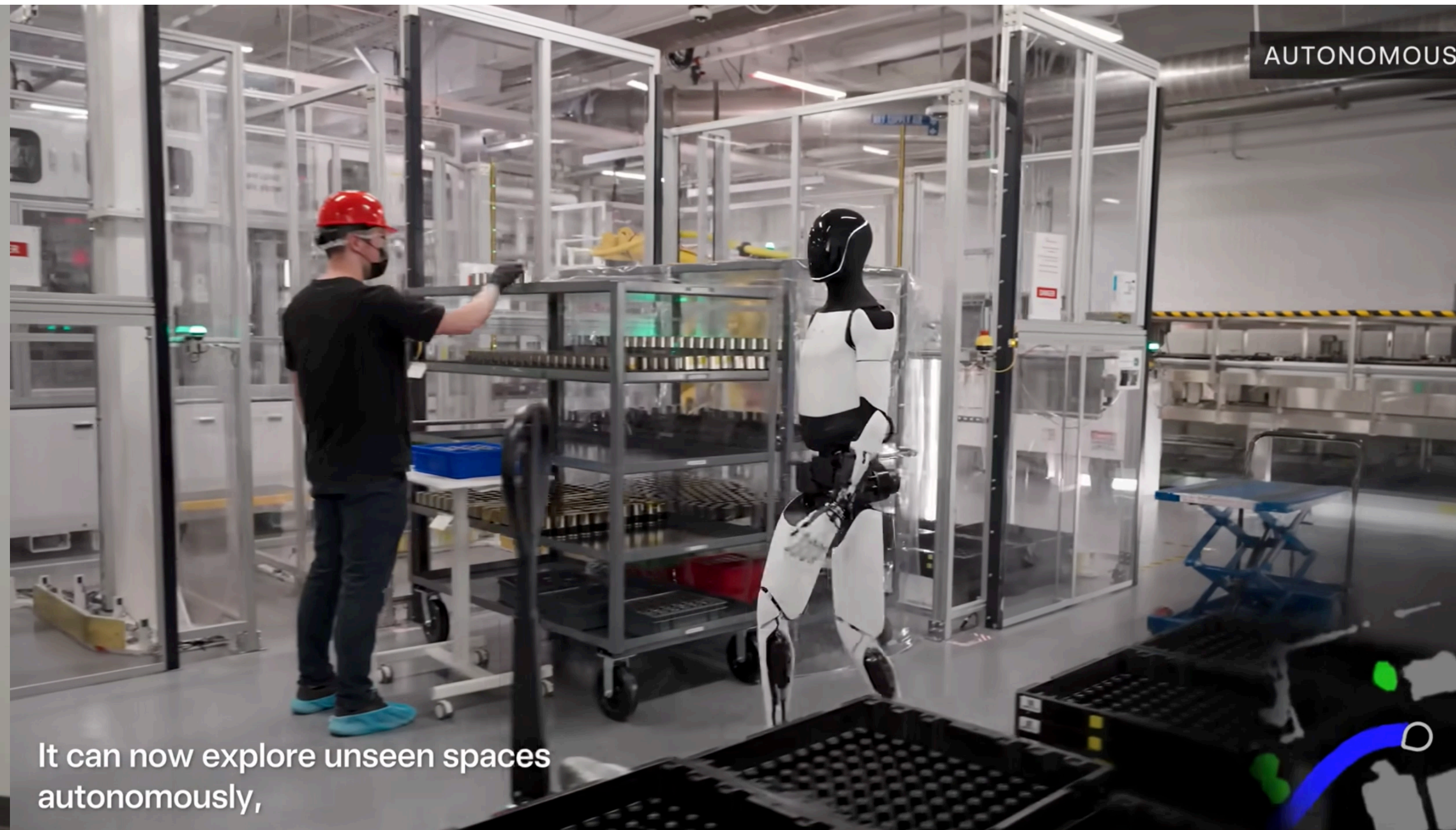
**Michelle Zhao,**
**October 29, 2024**

# Outline

- Alignment problem

- Alignment process: Learning from human feedback

- Case Study 1: Learning from preferences

- Active Learning: Why and How?

- Revisiting Case Study 1: Making learning from preference *active*

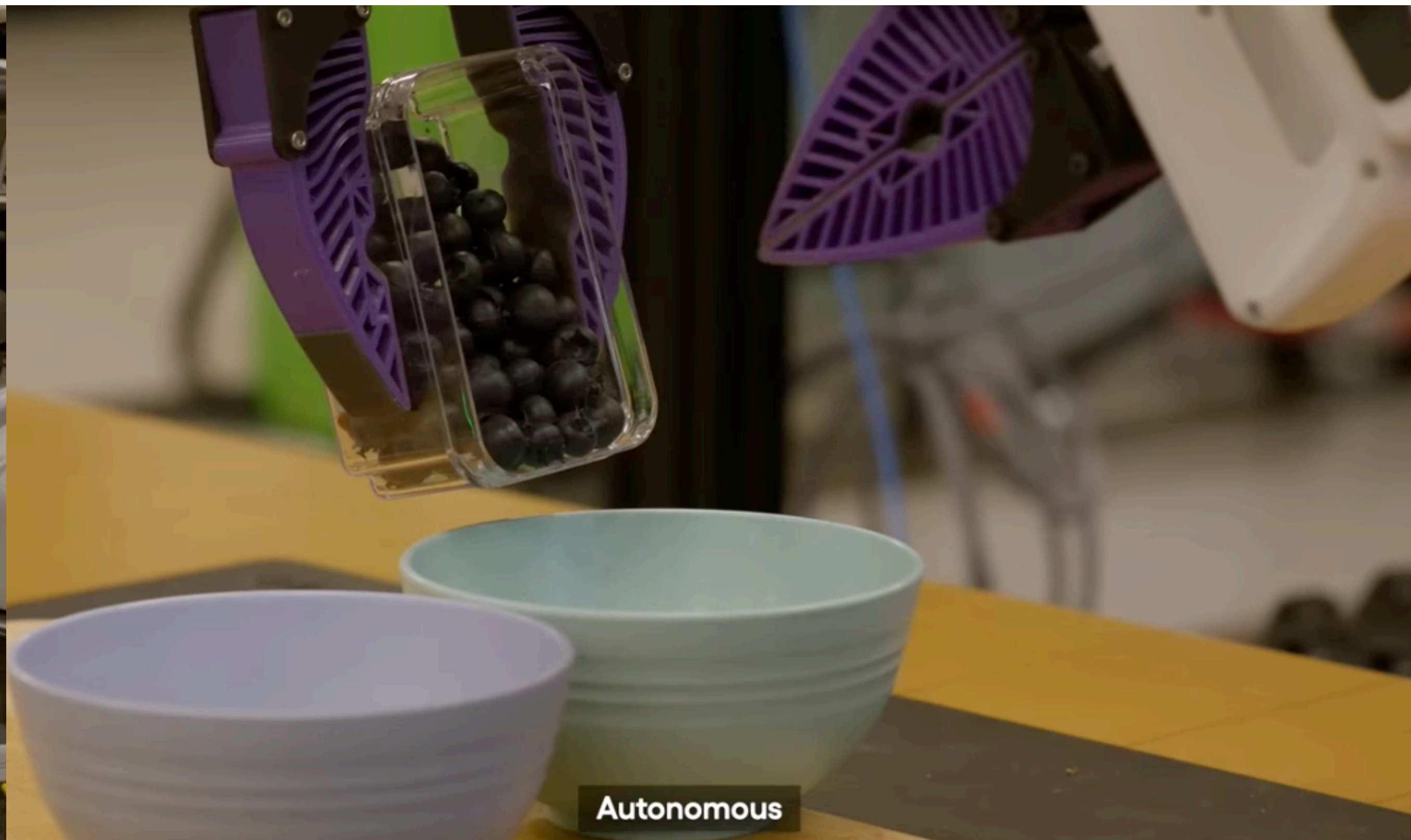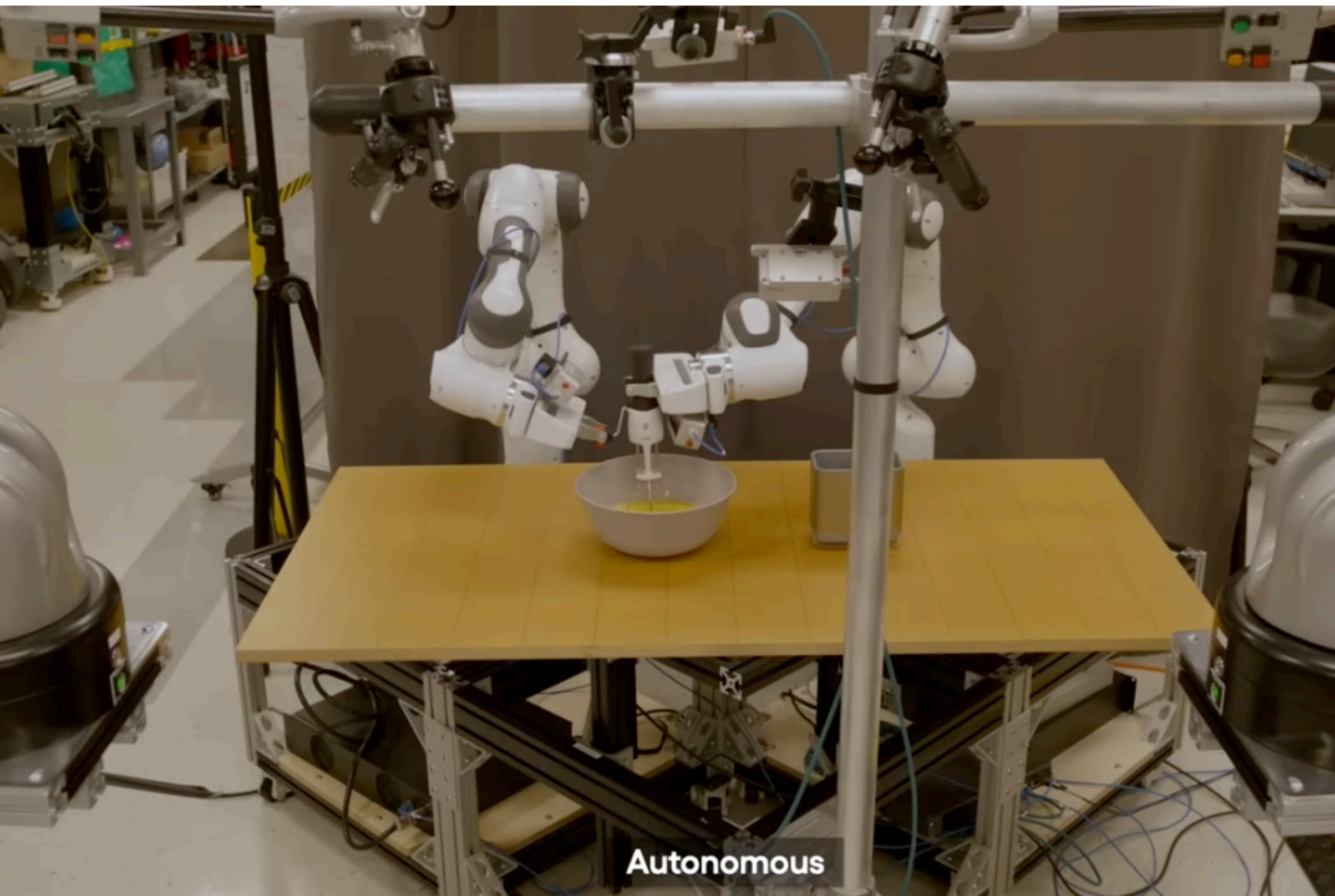- Case Study 2: Active learning for black-box policies

# Outline

- Alignment problem

- Alignment process: Learning from human feedback

- Case Study 1: Learning from preferences

- Active Learning: Why and How?

- Revisiting Case Study 1: Making learning from preference *active*

- Case Study 2: Active learning for black-box policies

# We're starting to see remarkable strides in learning for robotics

# We're starting to see remarkable strides in learning for robotics

# Underlying Aim: Robots that behave as we want them to!



Source

Source

# Alignment in Robotics

## How can we get robots to do what we want them to?

Leike, Jan, et al. "Scalable agent alignment via reward modeling: a research direction."
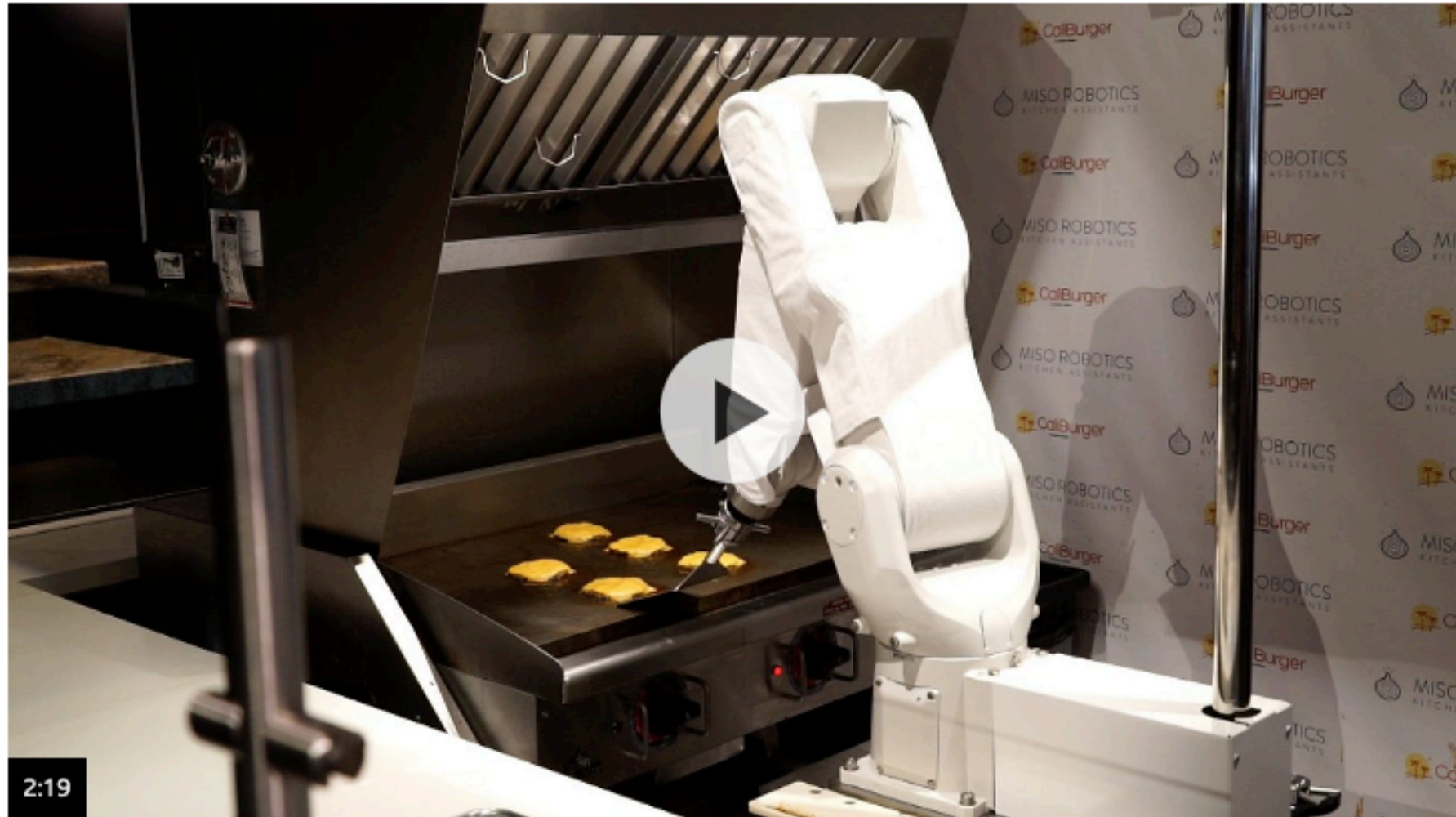
# Robots don't know what we want



**BBC**

Home   News   Sport   Business   Innovation   Culture   Travel   Earth   Video   Live

## Burger-flipping robot taken offline after one day

9 March 2018

Share

WATCH: Flippy the burger robot gets to work

**Flippy the burger-flipping robot that started work this week in a California restaurant has been forced to take a break because it was too slow.**

The robot was installed at a Cali Burger outlet in Pasadena and replaced human cooks.



**BBC**

Home   News   Sport   Business   Innovation   Culture   Travel   Earth   Video   Live

Bacon ice cream and nugget overload sees misfiring McDonald's AI withdrawn

Asia / East Asia

## AI fail: Japan's Henn-na Hotel dumps 'annoying' robot staff, hires humans

- Dinosaur receptionists are a thing of the past as Japan's first robot hotel concludes there "are places where they are just not needed"

🎧 Listen to this article ▶

Julian Ryall   + FOLLOW
Published: 12:32pm, 16 Jan 2019

Why you can trust SCMP

Getty Images

# Outline

- **Alignment problem**

- Alignment process: Learning from human feedback

- Case Study 1: Learning from preferences

- Active Learning: Why and How?

- Revisiting Case Study 1: Making learning from preference *active*

- Case Study 2: Active learning for black-box policies
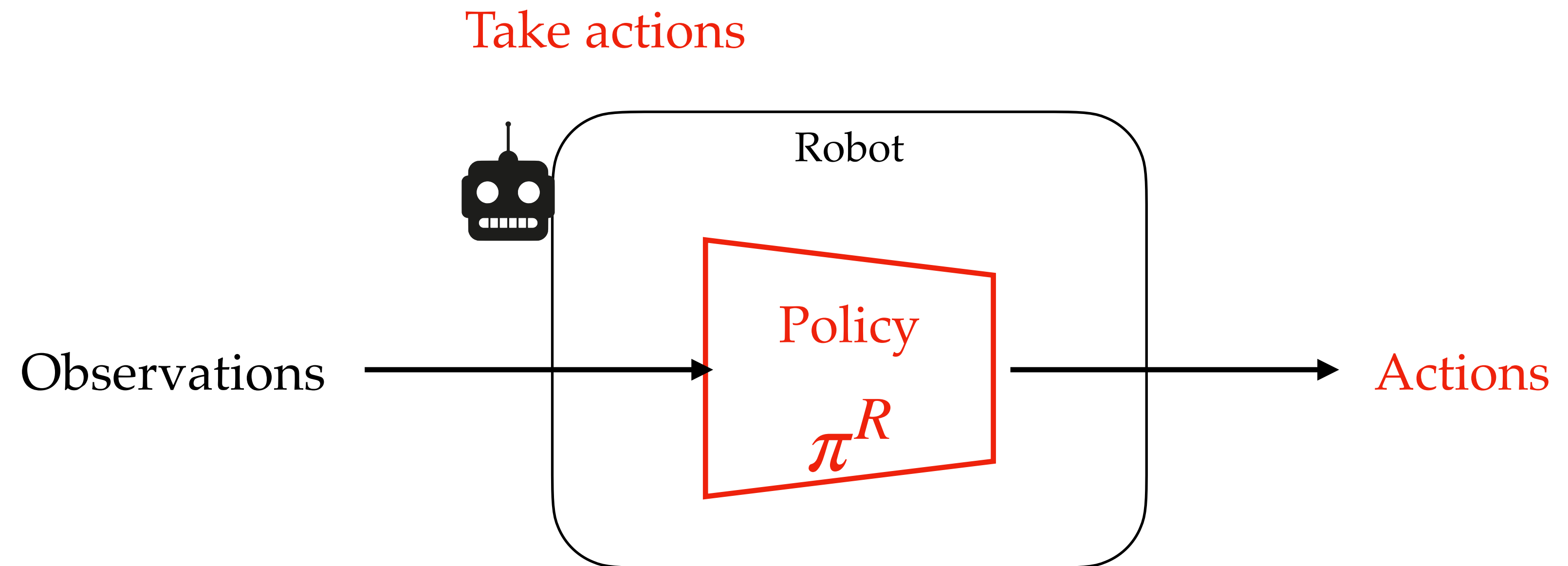
# Outline

• Alignment problem

• **Alignment process: Learning from human feedback**

• Case Study 1: Learning from preferences

• Active Learning: Why and How?

• Revisiting Case Study 1: Making learning from preference *active*

• Case Study 2: Active learning for black-box policies

# Alignment in Robotics

How can we get robots to do what we want them to?

Leike, Jan, et al. "Scalable agent alignment via reward modeling: a research direction."

# Alignment in Robotics

How can we get robots to <span style="color:red">do</span> what we want them to?

# Alignment in Robotics

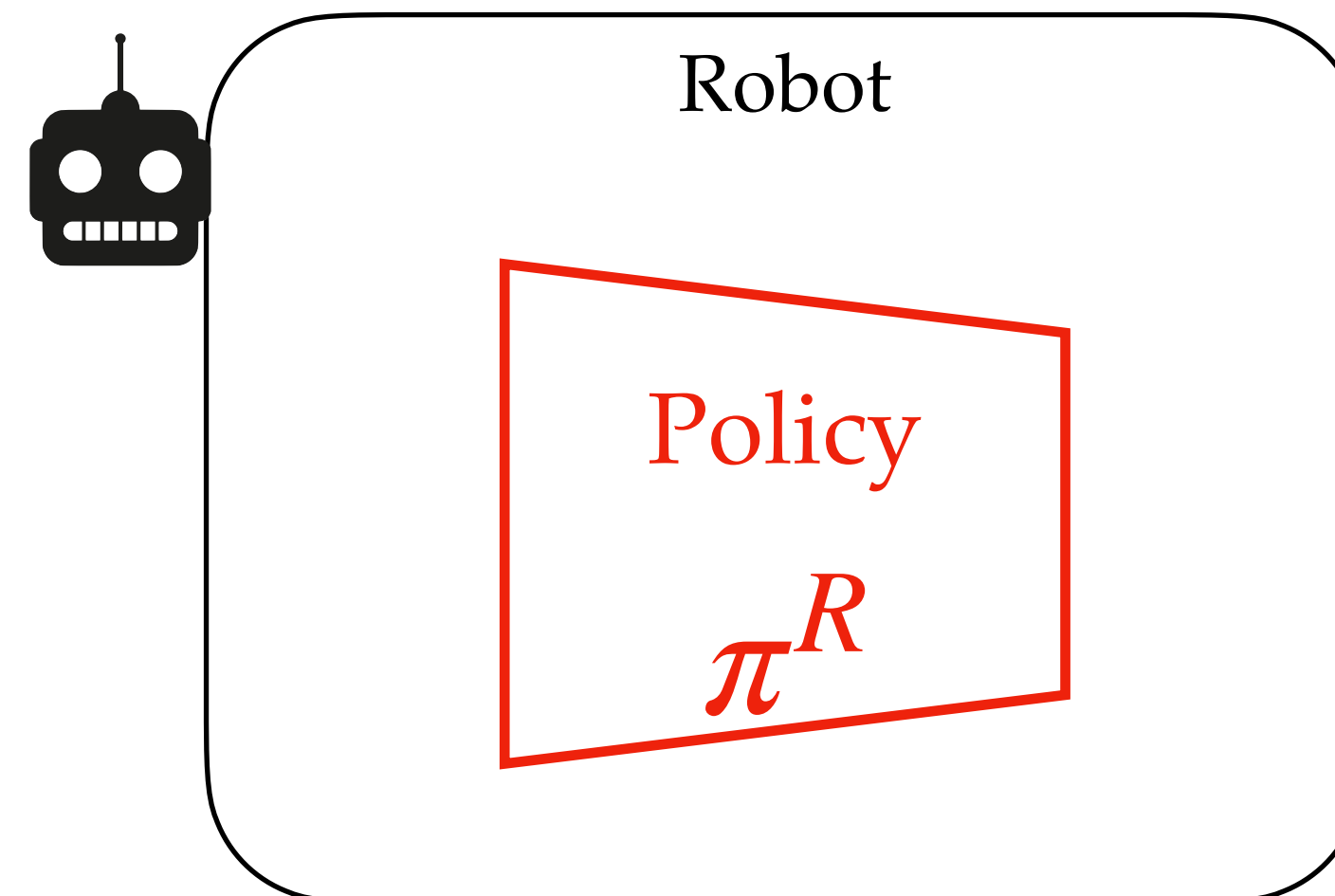How can we get robots to <span style="color:red">do</span> what <span style="color:green">we want them to</span>?

<span style="color:red">Take actions</span>

<span style="color:green">Achieve human objectives</span>

Robot

<span style="color:red">Policy</span>

$$\pi^R$$

# Alignment in Robotics

How can we get robots to <span style="color:red">do</span> what <span style="color:green">we want them to</span>?
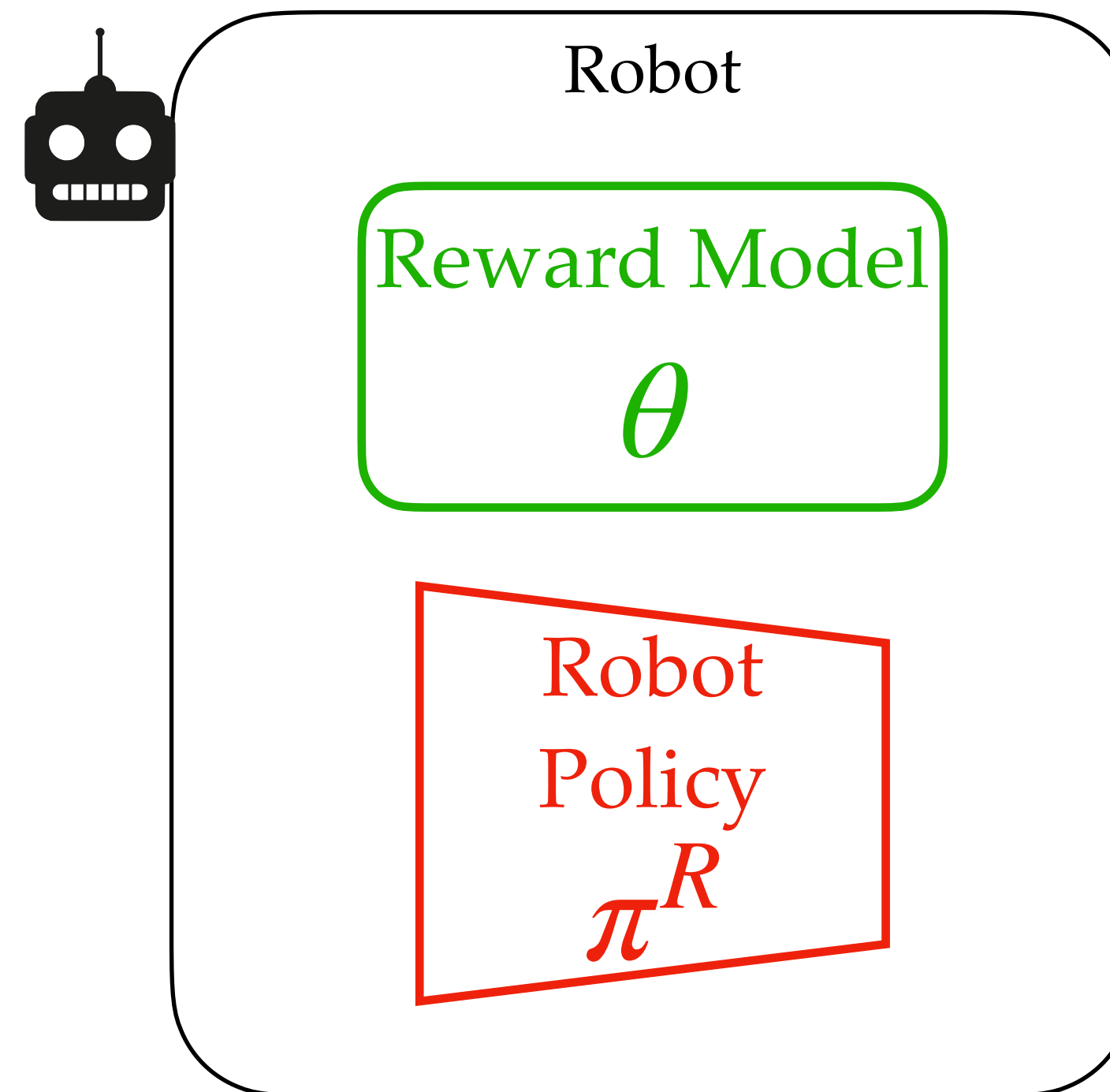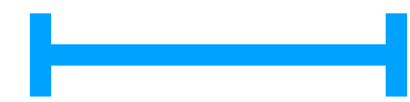
<span style="color:#1E90FF">Learning from
human feedback</span>

<span style="color:red">Take actions</span>

<span style="color:green">Achieve human
objectives</span>

Robot

<span style="color:green">Reward Model

$\theta$</span>

<span style="color:red">Robot
Policy

$\pi^R$</span>

# Alignment in Robotics

How can we get robots to do what we want them to?

Learning from
human feedback

*to have
robots*

Take actions

*that*
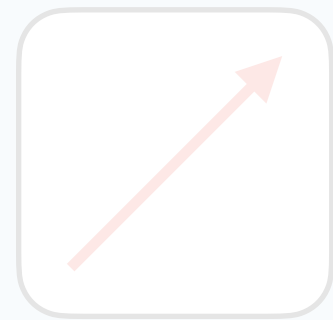
Achieve human
objectives

# Learning from human feedback

Human Data → Robot

# Learning from human feedback
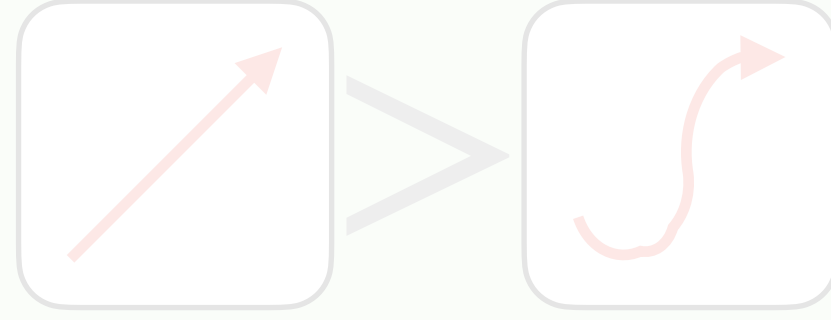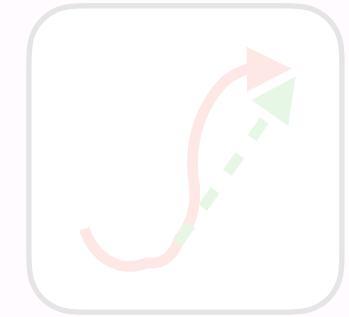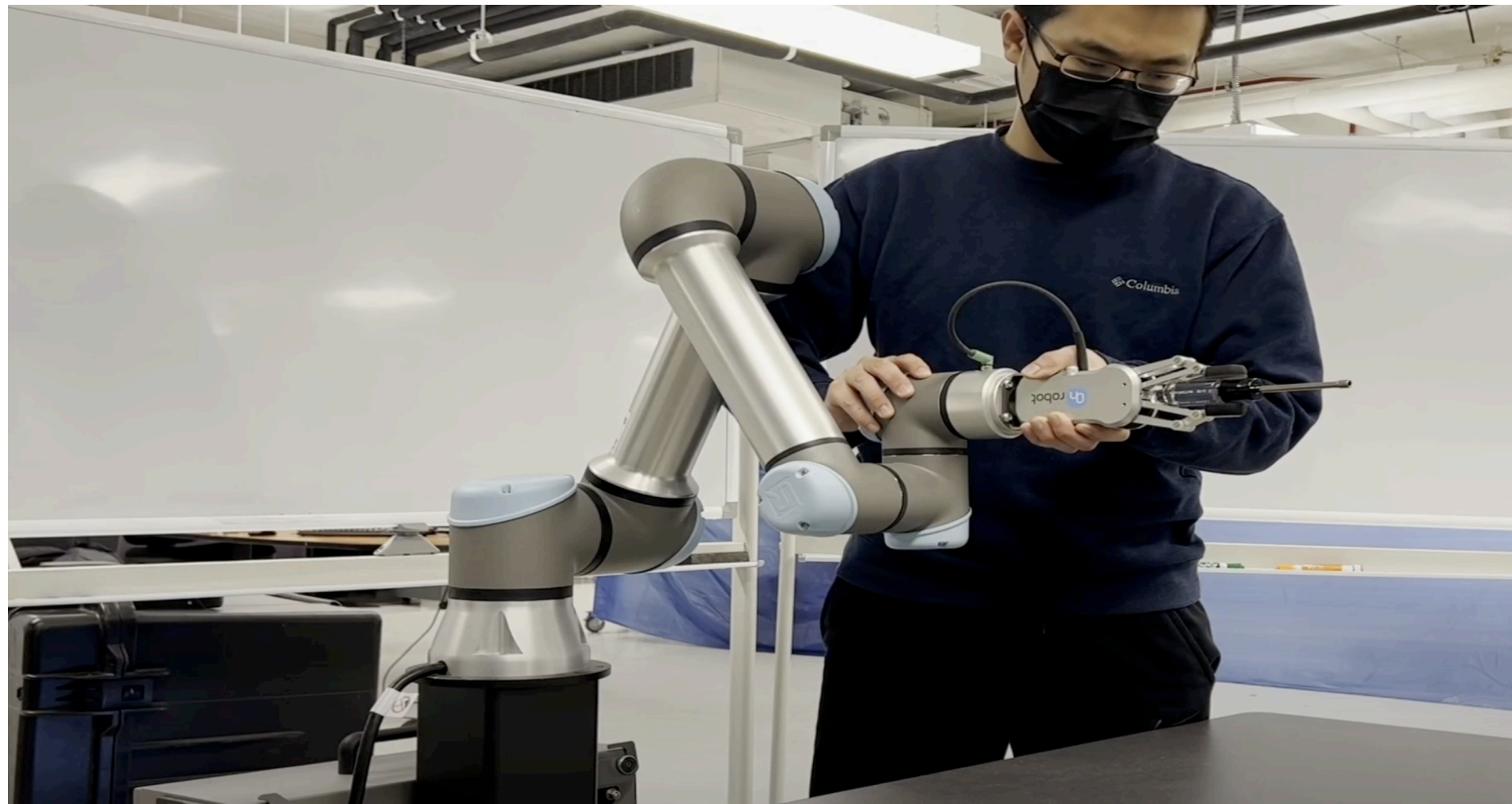
Human Data

Demonstrations

Preferences

Evaluations

★★★★★

Corrections

# Learning from human feedback

Human Data

Demonstrations | Preferences | Evaluations | Corrections

Robot Learning from Demonstration

# Learning from human feedback



Human Data

| Demonstrations | Preferences | Evaluations | Corrections |

# Learning from human feedback
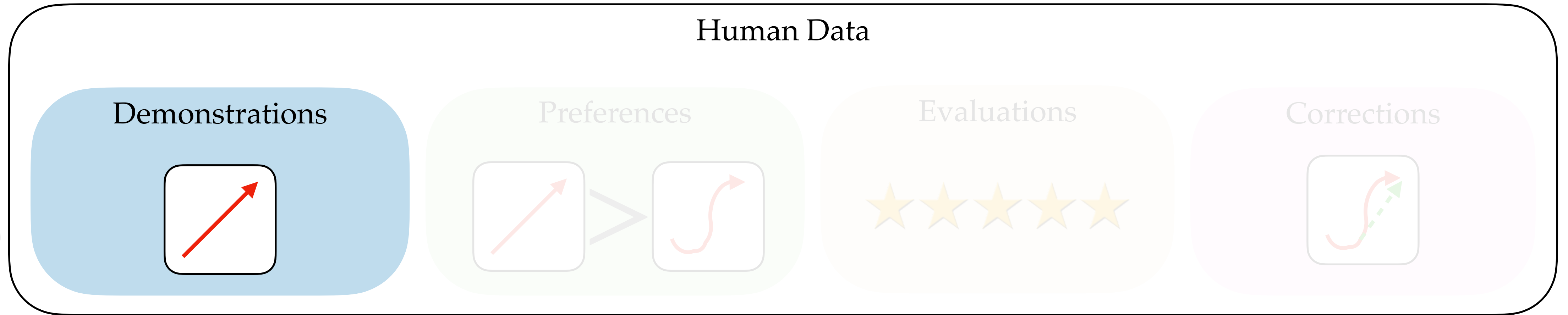
Human Data

Demonstrations

Preferences

Evaluations

★★★★★

Corrections

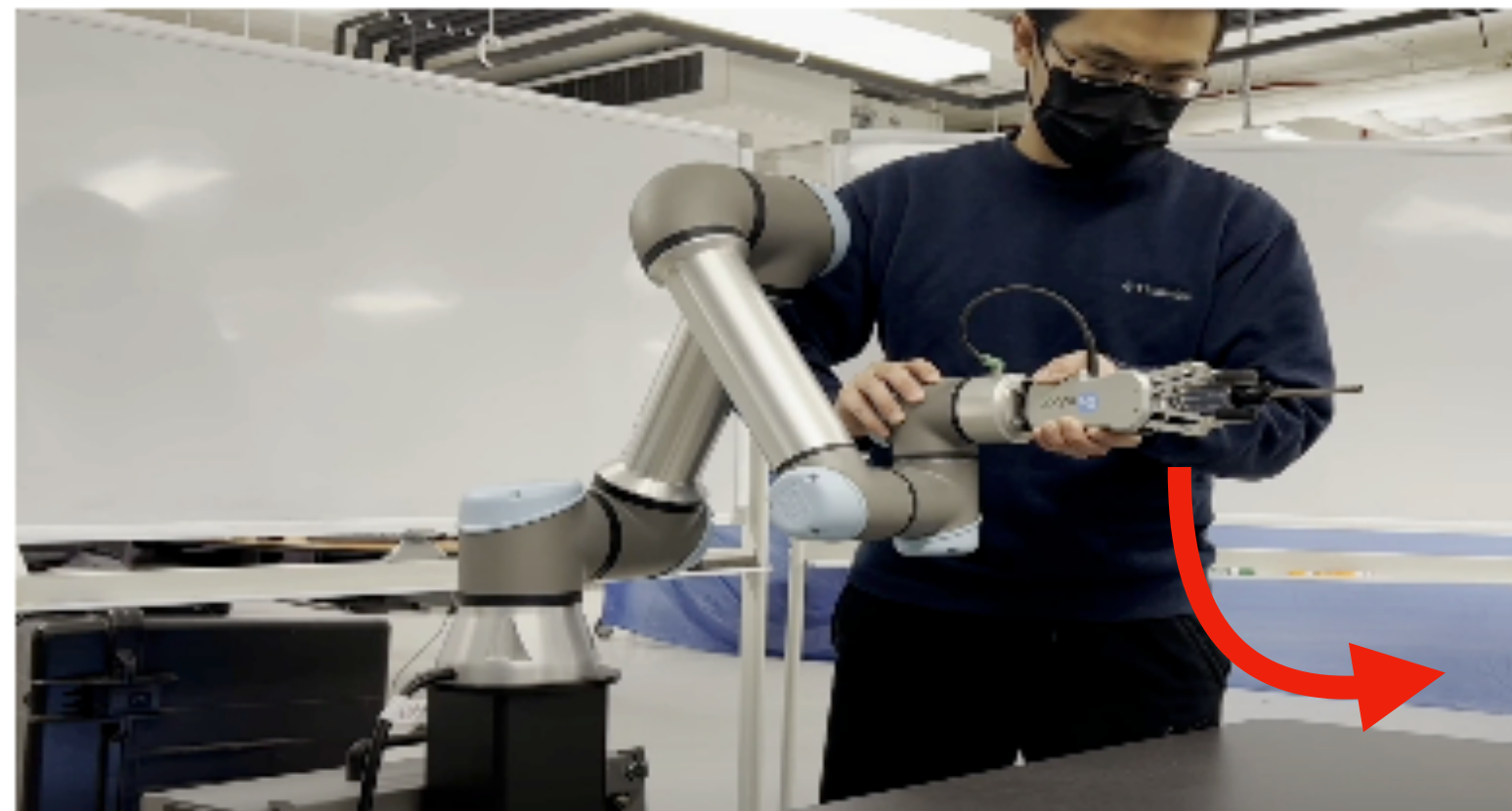# Learning from human feedback

Human Data

Demonstrations

Preferences

Evaluations

Corrections

Learning from Physical Human Corrections, One Feature at a Time

# Learning from human feedback

Human
Input

Robot

# Learning from human feedback



Human
Input

Robot

Reward Model
$\theta$

Robot
Policy
$\pi^R$

# Outline

● Alignment problem

● **Alignment process: Learning from human feedback**

● Case Study 1: Learning from preferences

● Active Learning: Why and How?

● Revisiting Case Study 1: Making learning from preference *active*

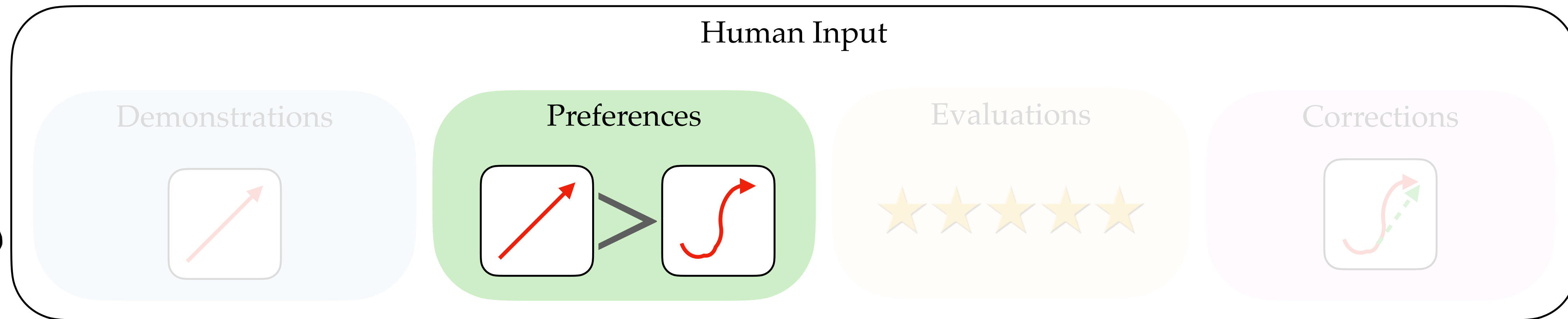● Case Study 2: Active learning for black-box policies

# Outline

• Alignment problem

• Alignment process: Learning from human feedback

• **Case Study 1: Learning from preferences**

• Active Learning: Why and How?

• Revisiting Case Study 1: Making learning from preference *active*

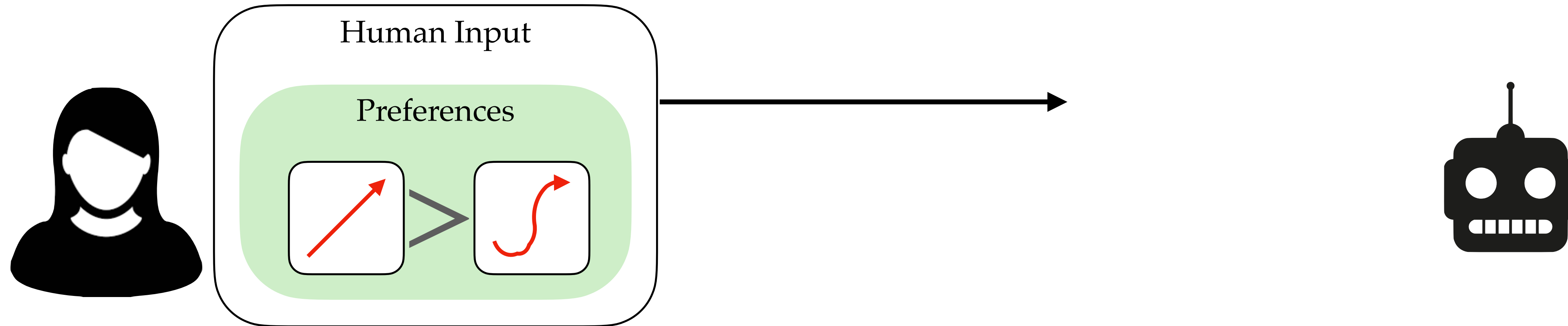• Case Study 2: Active learning for black-box policies

Let's take a closer look 👀 at
Active **Preference-Based Learning of Reward Functions**

Sadigh, Dorsa, et al. *Active preference-based learning of reward functions*. 2017.
Biyik, Erdem, and Dorsa Sadigh. "Batch active preference-based learning of reward functions." *Conference on robot learning*. PMLR, 2018.
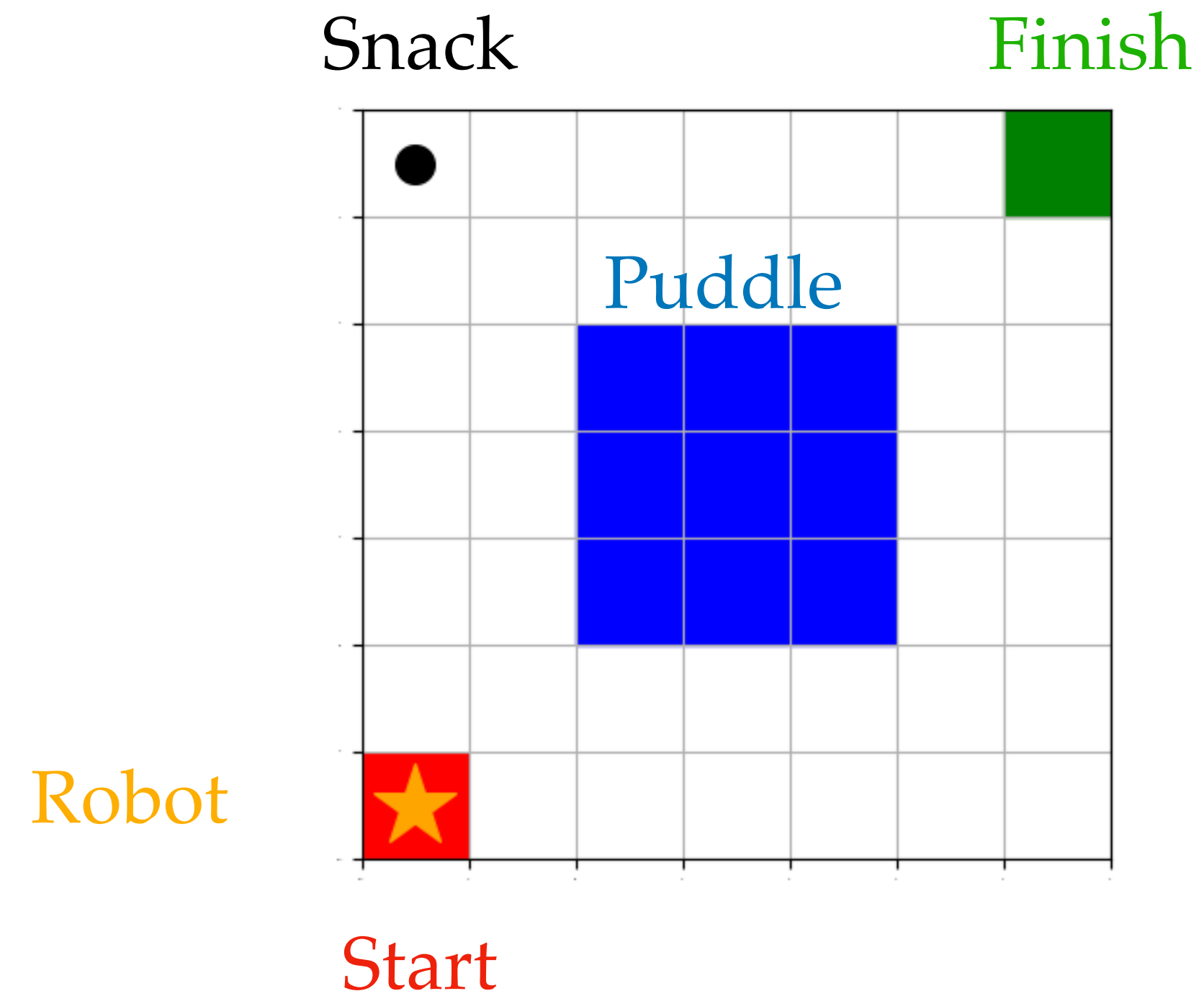Bıyık, Erdem, et al. "Asking easy questions: A user-friendly approach to active reward learning." *arXiv preprint arXiv:1910.04365* (2019).

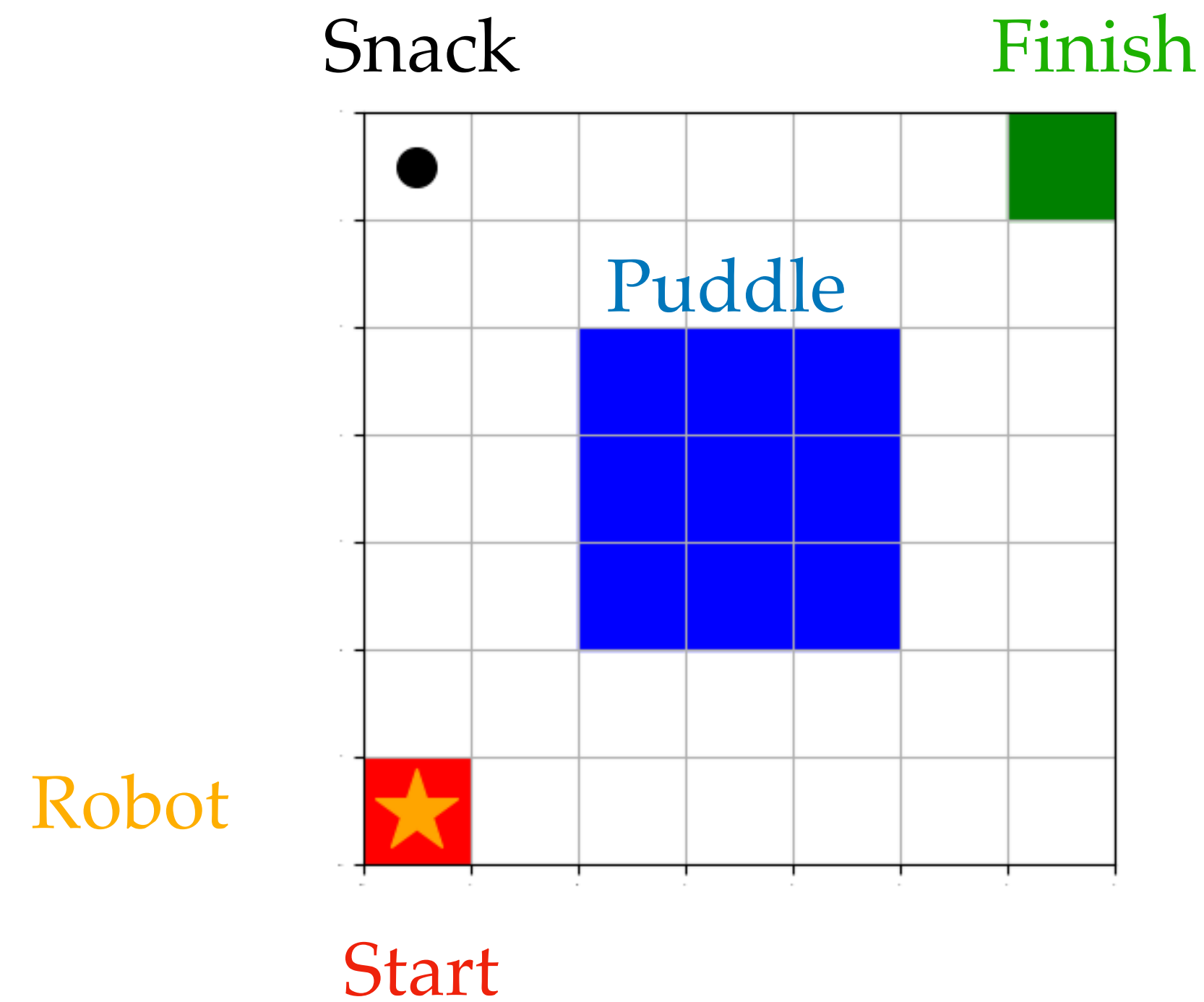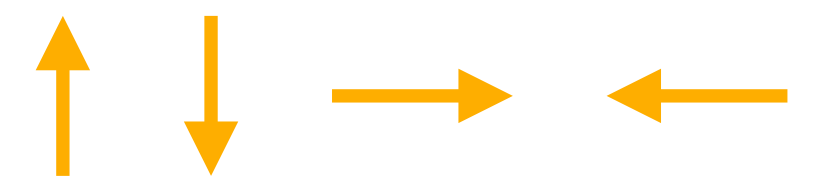# Preference-based learning: Interaction Setup

Human Input

Preferences

# Step 1: Formalizing the Objective

**Let's decide what we want**

# Step 1: Formalizing the Objective

**Let's decide what we want**

Snack        Finish

Puddle

Robot

Start

Actions $a \in A$

# Step 1: Formalizing the Objective
## Let's decide what we want

Task Objectives I want to teach the robot:

1. Snack: Good *(want to eat a snack)*

# Step 1: Formalizing the Objective
## Let's decide what we want

Task Objectives I want to teach the robot:

1. **Snack**: Good *(want to eat a snack)*
2. **Puddle**: Bad *(want to avoid puddles)*
3. **Finish**: Good *(want to get to the finish)*
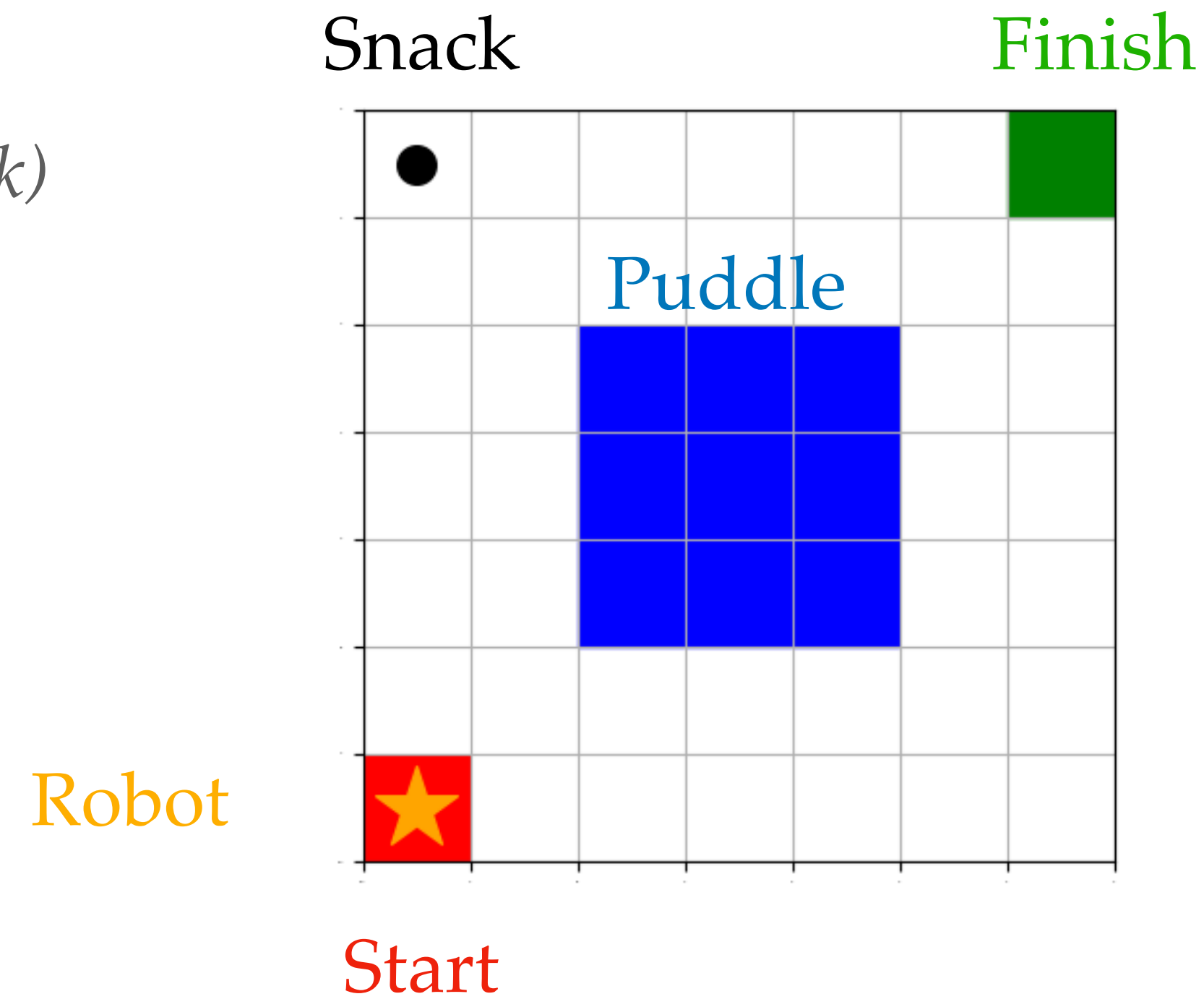4. **Steps**: Bad *(want to take as few steps as possible)*

# Step 1: Formalizing the Objective
## Let's decide what we want

Task Objectives I want to teach the robot:

1. Snack: Good *(want to eat a snack)*
2. Puddle: Bad *(want to avoid puddles)*
3. Finish: Good *(want to get to the finish)*
4. Steps: Bad *(want to take as few steps as possible)*

What matters?

How does it matter?

Reward function
$$R(s) = \theta^T \phi(s)$$

Weights    Set of selected features

# Step 1: Formalizing the Objective
## Let's decide what we want

Task Objectives I want to teach the robot:

$$R(s) = \theta^T \phi(s)$$

Weights    Set of selected features

1. Snack: Good *(want to eat a snack)*
2. Puddle: Bad *(want to avoid puddles)*
3. Finish: Good *(want to get to the finish)*
4. Steps: Bad *(want to take as few steps as possible)*

$\phi(s)$: [# snacks,
distance from puddle,
distance from finish,
# timesteps occurred]



What matters?

# Step 1: Formalizing the Objective

## Let's decide what we want

Task Objectives I want to teach the robot:

$$R(s) = \theta^T \phi(s)$$

Weights   Set of selected features

1. Snack: Good *(want to eat a snack)*
2. Puddle: Bad *(want to avoid puddles)*
3. Finish: Good *(want to get to the finish)*
4. Steps: Bad *(want to take as few steps as possible)*

$\phi(s)$: [# snacks,
distance from puddle,
distance from finish,
# timesteps occurred]

$\phi(s) = [\ 0\ ,\ 3.6\ ,\ 7.8,\ \ 1\ \ ]$

# Step 1: Formalizing the Objective
## Let's decide what we want

Task Objectives I want to teach the robot:

1. Snack: Good *(want to eat a snack)*
2. Puddle: Bad *(want to avoid puddles)*
3. Finish: Good *(want to get to the finish)*
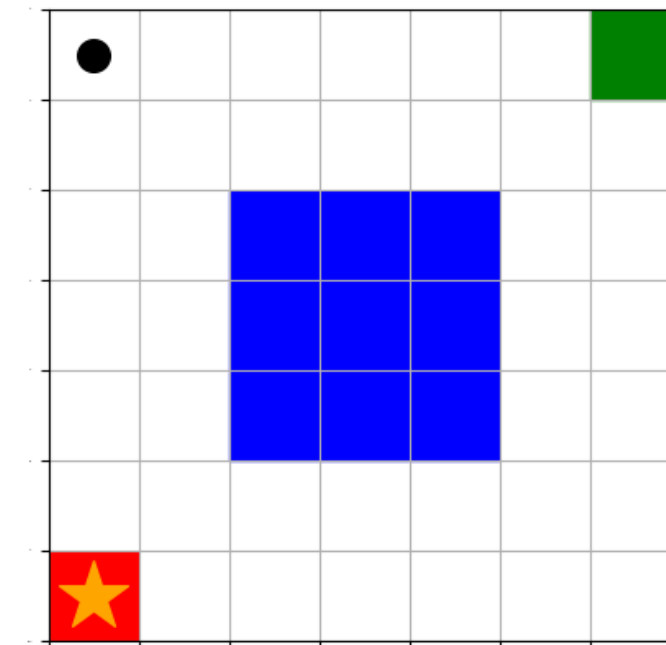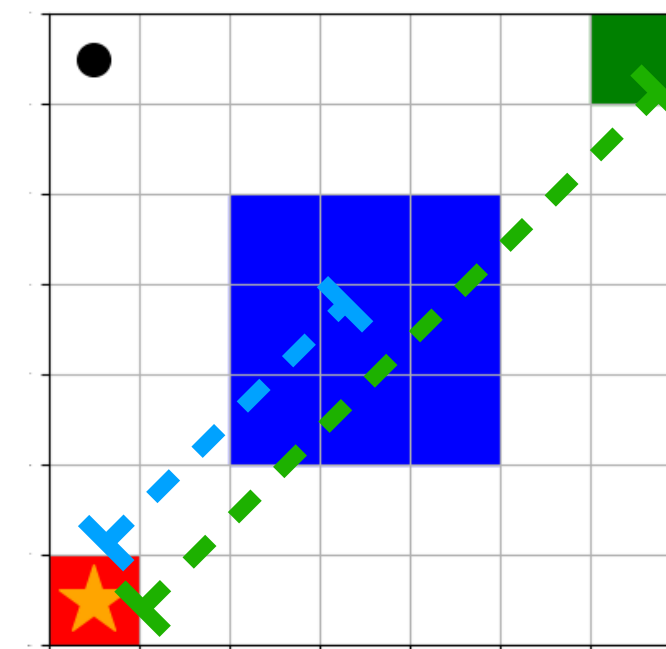4. Steps: Bad *(want to take as few steps as possible)*
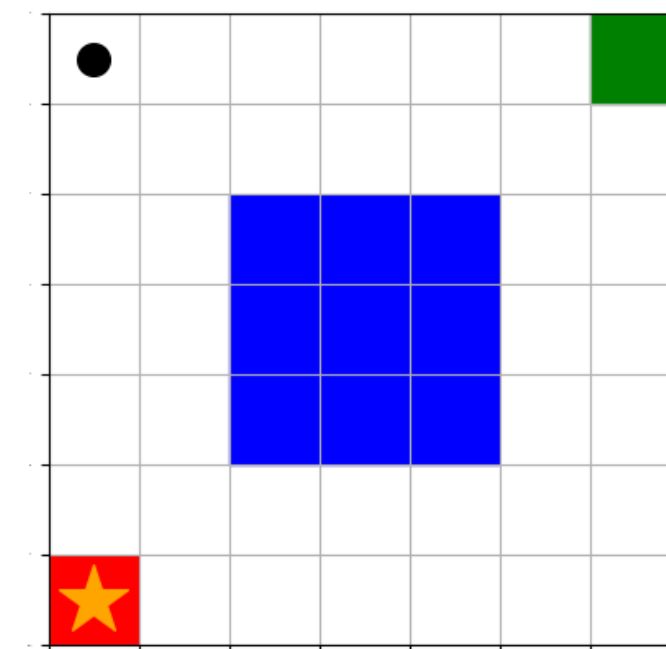
How does it matter?

$$R(s) = \theta^T \phi(s)$$

Weights   Set of selected features

$\phi(s)$: [# snacks, distance from puddle, distance from finish, # timesteps occurred]

$\phi(s) = [ \ 0 \ , \ 3.6 \ , \ 7.8, \quad 1 \quad ]$

Weights   $\theta \quad = [ \ 5 \ , \quad 2 \ , \quad -1, \quad -1 \quad ]$

# Preference-based learning: Interaction Setup

Human Input

Preferences

Robot

Reward Model
$\hat{\theta}$

Robot
Policy
$\pi^R$

Reward $\theta$

$\theta \quad = [\ 5\ ,\quad 2\ ,\quad \text{-}1,\quad \text{-}1\ ]$

# Let's try giving the robot a preference together!

# Step 2: What does the robot do with this information?



Human Input

Reward $\theta$

Robot

Reward Model
$\hat{\theta}$

Robot
Policy
$\pi^R$

# Step 2: Bayes update to learn from preference

I have no idea what $\theta$ might be!
It could be anything in $\mathbb{R}^4$

# Step 2: Bayes update to learn from preference

We first initialize a distribution over $\Theta$

$$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$$

# Step 2: Bayes update to learn from preference

We first initialize a distribution over $\Theta$

$P(\theta_i)$

1/6

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

# Step 2: Bayes update to learn from preference

$P(\theta_i)$

$1/6$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

We just received data from the user in the form of a preference

Query $Q = \{\xi_A, \xi_B\}$

Choice c = $\xi_A$
(*the one the user preferred*)

$>$

We will use Bayes Rule to obtain a posterior probability

# Building up: Bayes Update

$$P(Y)\, P(X|Y) \quad = \quad P(X, Y)$$

Chain Rule: $\quad P(c, Q)\, P(\theta | c, Q) = P(c, Q, \theta)$

# Building up: Bayes Update

$$P(Y)\ P(X|Y) \quad = \quad P(X, Y)$$

Chain Rule: $\quad P(c, Q)\ P(\theta | c, Q) = P(c, Q, \theta)$

$$P(\theta | c, Q) = \frac{P(c, Q, \theta)}{P(c, Q)}$$

# Building up: Bayes Update

$$P(Y)\ P(X\,|\,Y) \quad = \quad P(X, Y)$$

Chain Rule:  $P(c, Q)\ P(\theta\,|\,c, Q) = P(c, Q, \theta)$

$$P(\theta\,|\,c, Q) = \frac{\boxed{P(c, Q, \theta)}}{P(c, Q)}$$

$$= \frac{\boxed{P(c\,|\,Q, \theta)\ P(Q, \theta)}}{P(c, Q)}$$

# Building up: Bayes Update

$$P(Y)\ P(X\,|\,Y) \quad = \quad P(X, Y)$$

Chain Rule: $\quad P(c, Q)\ P(\theta\,|\,c, Q) = P(c, Q, \theta)$

$$P(\theta\,|\,c, Q) = \frac{P(c, Q, \theta)}{P(c, Q)}$$

$$= \frac{P(c\,|\,Q, \theta)\ P(Q, \theta)}{P(c, Q)}$$

$$= \frac{P(c\,|\,Q, \theta)\ P(Q)\ P(\theta)}{P(c\,|\,Q)\ P(Q)}$$

# Building up: Bayes Update

$$P(Y)\,P(X\,|\,Y) \quad = \quad P(X, Y)$$

Chain Rule: $\quad P(c, Q)\,P(\theta\,|\,c, Q) = P(c, Q, \theta)$

$$P(\theta\,|\,c, Q) = \frac{P(c, Q, \theta)}{P(c, Q)}$$

$$= \frac{P(c\,|\,Q, \theta)\,P(Q, \theta)}{P(c, Q)}$$

$$= \frac{P(c\,|\,Q, \theta)\,P(Q)\,P(\theta)}{P(c\,|\,Q)\,P(Q)}$$

$$= \frac{P(c\,|\,Q, \theta)\,\,P(\theta)}{P(c\,|\,Q)}$$

# Building up: Bayes Update

Bayes Rule:

$$P(\theta \mid c, Q) = \frac{P(c \mid Q, \theta) \; P(\theta)}{P(c \mid Q)}$$

P(rew | query choice)

P(choice | query, rew)

P(rew)

Uniform prior

Normalization

P(choice | query)

# Boltzmann: Likelihood of Human Decision | Model

$$P(c \,|\, Q, \theta) = \frac{e^{R(c)}}{\sum_{q \in Q} e^{R(q)}}$$

P(choice | query, rew)

Boltzmann Rational Model
(Might also see this as Bradley-Terry model of preferences)

# Step 2: Bayes update to learn from preference

$Q, c$  $>$ 

Use Bayes to compute prob.
model given data

$$P(\theta \mid c, Q) = \frac{P(c \mid Q, \theta) P(\theta)}{P(c \mid Q)}$$



$P(\theta_i)$

$1/6$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

$P(\theta_i \mid c, Q)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

# Step 2: Bayes update to learn from preference



$Q, c$

$P(\theta_i \,|\, c, Q)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

MAP (maximum a posteriori) estimate serves as $\hat{\theta}$

Reward Model
$\hat{\theta}$

# Preference-Based Learning of Reward Functions

# Outline

• Alignment problem

• Alignment process: Learning from human feedback

• **Case Study 1: Learning from preferences**

• Active Learning: Why and How?

• Revisiting Case Study 1: Making learning from preference *active*

• Case Study 2: Active learning for black-box policies

# Challenges in the [Passive] Learning from Feedback Paradigm

- The agent's ability to learn relies on good training data.

# Let's consider another pair of trajectories

# Challenges in the [Passive] Learning from Feedback Paradigm

- The agent's ability to learn relies on good training data.

- The onus to provide the good training data falls completely on the user to know what the robot needs.

- What else?

# Challenges in the [Passive] Learning from Feedback Paradigm

- The agent's ability to learn relies on good training data.

- The onus to provide the good training data falls completely on the user to know what the robot needs.

- What else?

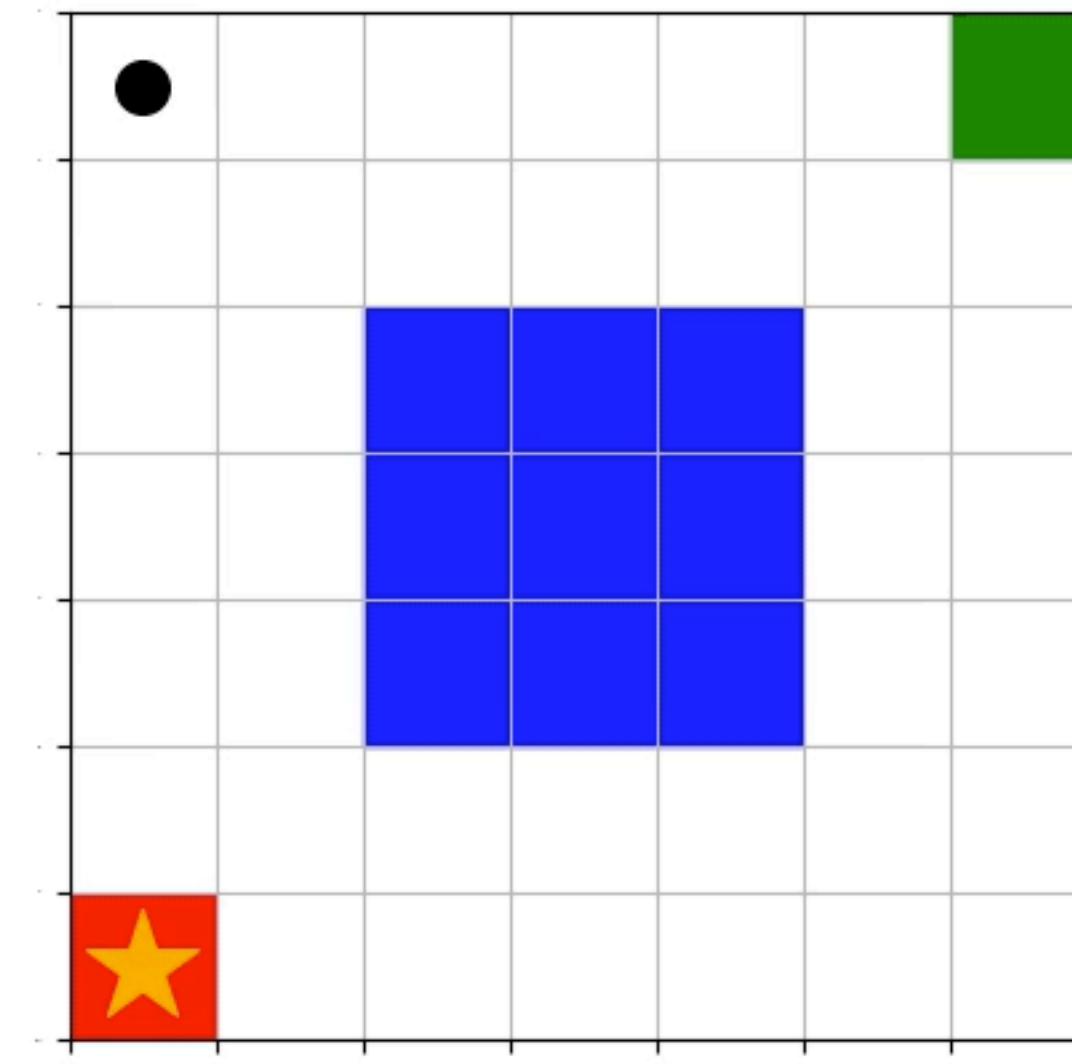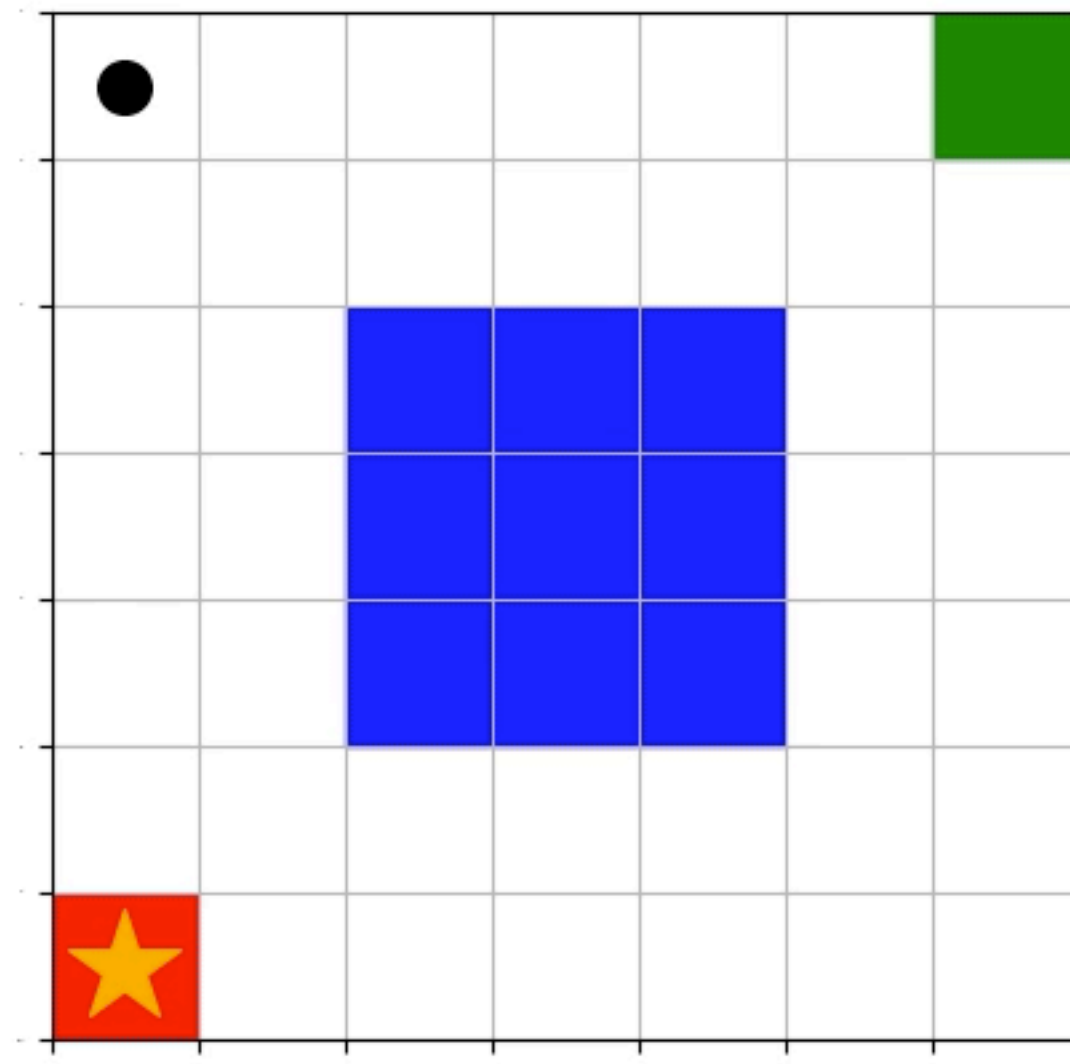- At scale, it can require fleets of highly trained users.

# Outline

• Alignment problem

• Alignment process: Learning from human feedback

• Case Study 1: Learning from preferences

• **Active Learning: Why and How?**

• Revisiting Case Study 1: Making learning from preference *active*

• Case Study 2: Active learning for black-box policies

# Active Learning



Teacher ← → Learner

Provide labels for requested queries

Select queries to request labels for

Learn a Model

The learner (robot) remains in control and requests annotated data from the human teacher.

The learner can be curious and request information from the teacher based on different query strategies.

# The key decision in active learning: Query Strategy

Provide labels for requested queries

Learn a Model

Select queries to request labels for

Teacher

Learner

# To design: an active robot learner who asks for help

Provide labels for requested queries
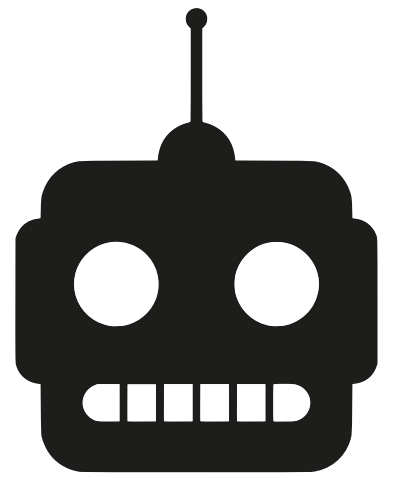
Teacher

Select queries to request labels for

Learner

I know that I need help

*How do I know that I need help?*

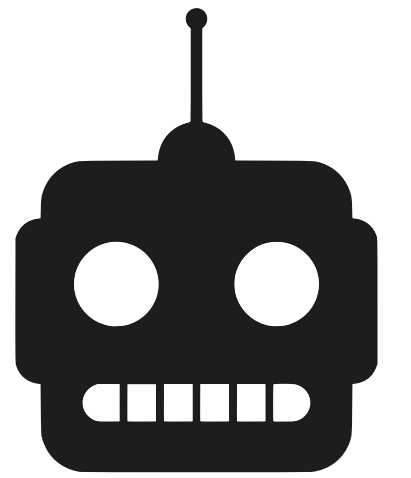# To design: an active robot learner who asks for help

Provide labels for requested queries

Teacher

Learner

Select queries to request labels for

I know what help to ask for

I know that I need help

*How do I ask for help?*

*How do I know that I need help?*

# Query Strategy: how do I ask for help?

## Uncertainty Minimization (Gaining Information)

Selects unlabeled items whose labels (once received) will reduce the robot's uncertainty over the model.

✓ Volume Removal

✓ Information Gain

## Diversity Sampling (Exploration)

Selects unlabeled items that differ from or are unseen in the data the robot has already seen.

✓ Variety of diversity metrics

✓ Different exploration objectives

## Random

# Outline

• Alignment problem

• Alignment process: Learning from human feedback

• Case Study 1: Learning from preferences

• **Active Learning: Why and How?**

• Revisiting Case Study 1: Making learning from preference *active*
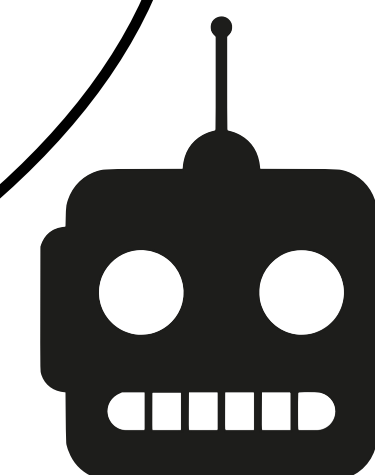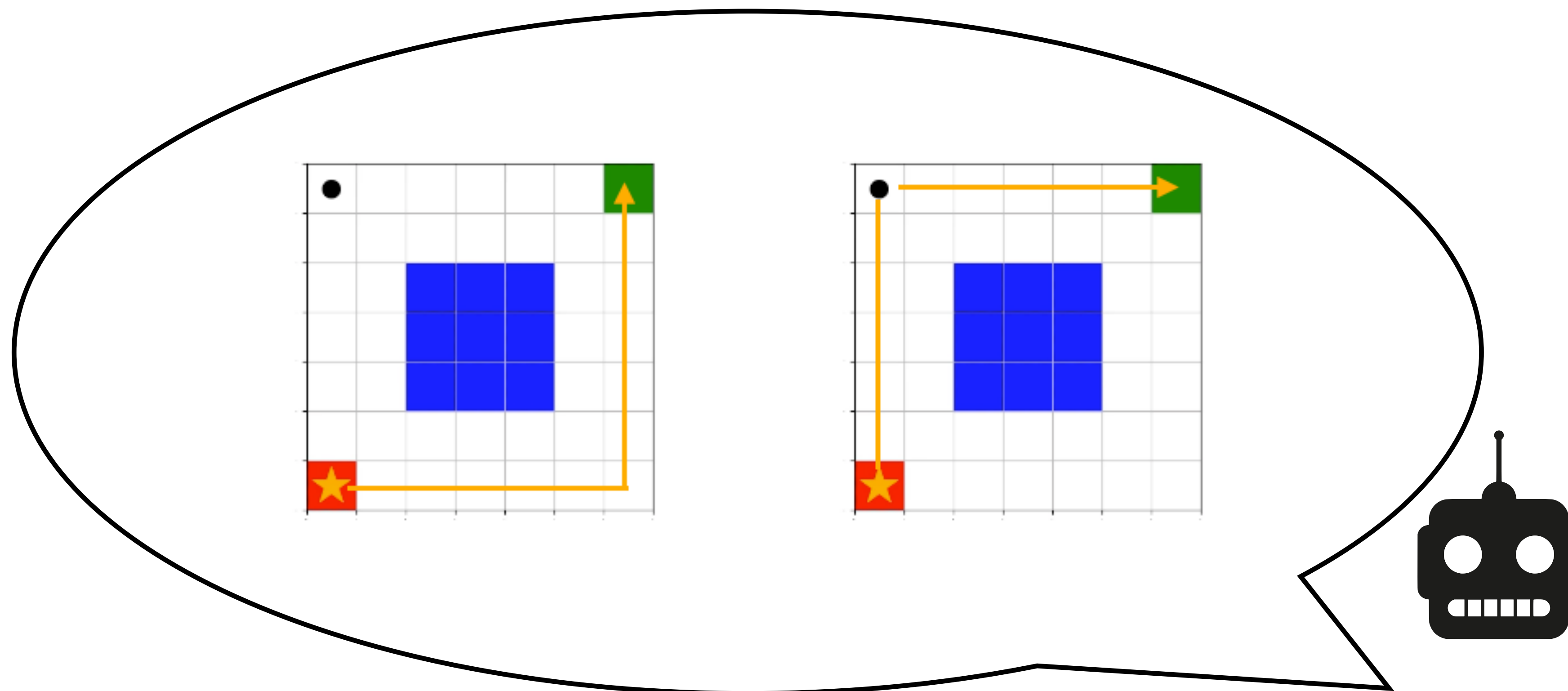
• Case Study 2: Active learning for black-box policies

# Outline

- Alignment problem

- Alignment process: Learning from human feedback

- Case Study 1: Learning from preferences

- Active Learning: Why and How?

- **Revisiting Case Study 1: Making learning from preference *active***

- Case Study 2: Active learning for black-box policies

Let's take a closer look 👀 at

**Active** Preference-Based Learning of Reward Functions

# Let's take a closer look 👀 at
# **Active** Preference-Based Learning of Reward Functions

I know what help to ask for

I know that I need help

*How do I ask for help?*

*How do I know that I need help?*

# Let's take a closer look 👀 at
# **Active** Preference-Based Learning of Reward Functions

| I know what help to ask for | I know that I need help |

*How do I ask for help?*   *How do I know that I need help?*

*Pick query that maximally reduces uncertainty*

*High uncertainty*

I know that I need help

$P(\theta_i)$

1/6

$\theta_1$  $\theta_2$  $\theta_3$  $\theta_4$  $\theta_5$  $\theta_6$

**Uncertainty = Entropy**

$P(\theta_i | c, Q)$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \theta_4 \quad \theta_5 \quad \theta_6$

I know that I need help

**Low Entropy** → **High Entropy**

The robot wants to find the query that will reduce the its uncertainty the most

$$\underset{Q \in \substack{\text{Possible} \\ \text{Queries}}}{\arg\max} \quad \text{Uncertainty Reduction } (Q)$$

The robot wants to find the query that will reduce the its uncertainty the most

$$\underset{Q \in \text{Possible Queries}}{\arg\max} \quad \text{Uncertainty Reduction} (Q)$$

The robot wants to find the query that will reduce the its uncertainty the most

$$\underset{Q \in \substack{\text{Possible} \\ \text{Queries}}}{\arg\max} \quad \overbrace{H(\theta)}^{\substack{\text{Uncertainty prior to} \\ \text{query}}} - \overbrace{\mathbb{E}_c[H(\theta \,|\, c, Q)]}^{\substack{\text{Uncertainty after} \\ \text{human response}}}$$

Uncertainty Reduction $(Q)$

The robot wants to find the query that will reduce the its uncertainty the most

$$\underset{Q \in \substack{\text{Possible} \\ \text{Queries}}}{\arg\max} \quad \overset{\text{Uncertainty prior to query}}{H(\theta)} - \overset{\text{Uncertainty after human response}}{\mathbb{E}_c[H(\theta \,|\, c, Q)]}$$

Uncertainty Reduction $(Q)$

Sadigh, Dorsa, et al. *Active preference-based learning of reward functions*. 2017.

Bıyık, Erdem, and Dorsa Sadigh. "Batch active preference-based learning of reward functions." *Conference on robot learning*. PMLR, 2018.
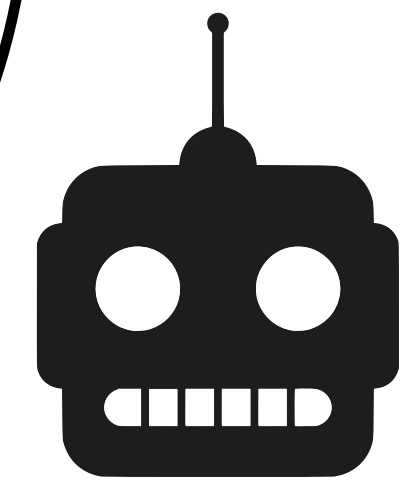
Bıyık, Erdem, et al. "Asking easy questions: A user-friendly approach to active reward learning." *arXiv preprint arXiv:1910.04365* (2019).

Choosing the query that will reduce the its uncertainty the most maximizes information gain

$$\underset{Q \in \substack{\text{Possible} \\ \text{Queries}}}{\arg \max} \quad \text{Information Gain}(Q)$$

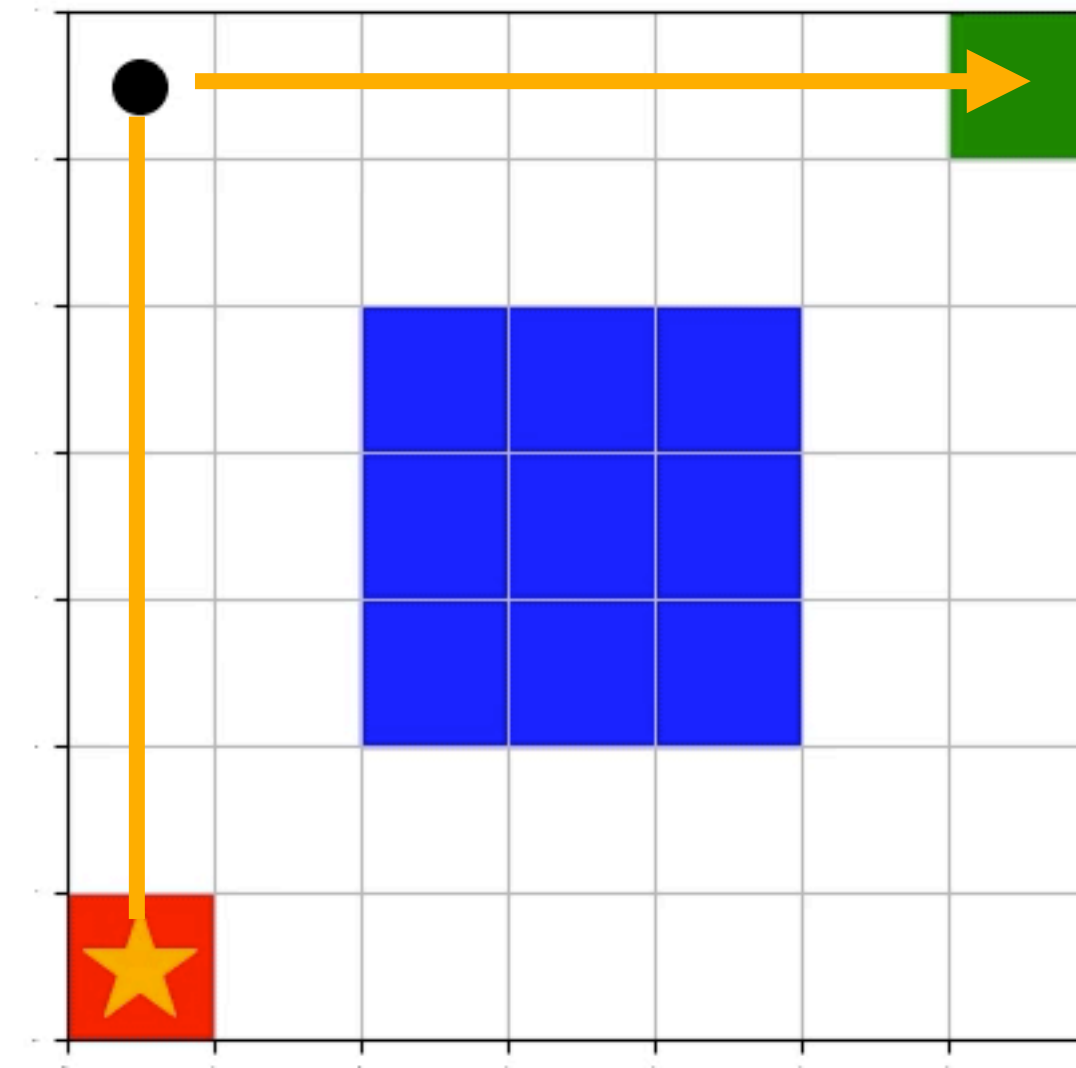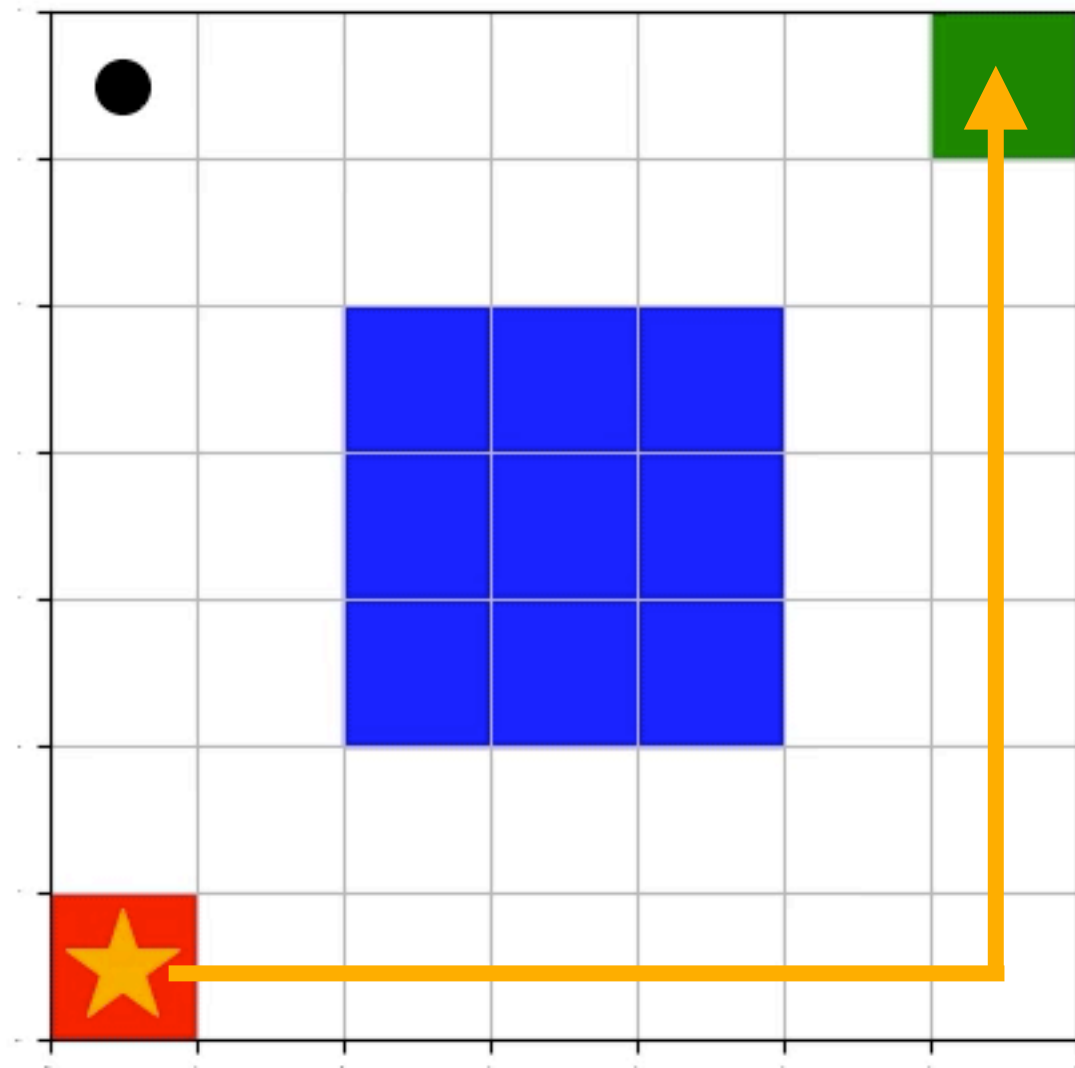When we optimize for information gain, we simultaneously produce queries that seek to be easy for the human to answer.

Bıyık, Erdem, et al. "Asking easy questions: A user-friendly approach to active reward learning." *arXiv preprint arXiv:1910.04365* (2019).

Applying information gain, We have for our most informative query:

# These preference learning techniques are key in training models like ChatGPT

# More robots that ask for help



Hey robot, could you put the bowl on the small counter in the microwave?

Large Language Model Planner

Next-Step Prediction with Scores

Put plastic bowl in recycling bin - 0.08
Put metal bowl in microwave - 0.41
Put plastic bowl in microwave - 0.44
Put metal bowl in landfill bin - 0.03

Ren, Allen Z., et al. "Robots that ask for help: Uncertainty alignment for large language model planners." *CoRL* (2023).

- Robots capable of self-assessments ability and a priori competency predictions can help improve overall team performance and trust.



Bridgwater, Tom, et al. "Examining profiles for robotic risk assessment: Does a robot's approach to risk affect user trust?." Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction. 2020.

# Outline

• Alignment problem

• Alignment process: Learning from human feedback

• Case Study 1: Learning from preferences

• Active Learning: Why and How?

• **Revisiting Case Study 1: Making learning from preference *active***

• Case Study 2: Active learning for black-box policies

# Summary thus far:

- Learning from human feedback seeks to align robot behaviors to human intentions

  - Human feedback can occur through multiple types (i.e. preferences)

# Summary thus far:

- Learning from human feedback seeks to align robot behaviors to human intentions

  - Human feedback can occur through multiple types (i.e. preferences)

- Active learning shifts the interaction paradigm by allowing the learner to *actively* request feedback (labels) from the human.

  - Effective learning occurs by coming up with a good strategy for querying the human (*strategically asking for help*)

# Summary thus far and Questions?

- Learning from human feedback seeks to align robot behaviors to human intentions

  - Human feedback can occur through multiple types (i.e. preferences)

- Active learning shifts the interaction paradigm by allowing the learner to *actively* request feedback (labels) from the human.

  - Effective learning occurs by coming up with a good strategy for querying the human (*strategically asking for help*)

- To enable an effective active robot learner, 2 necessary design decisions:
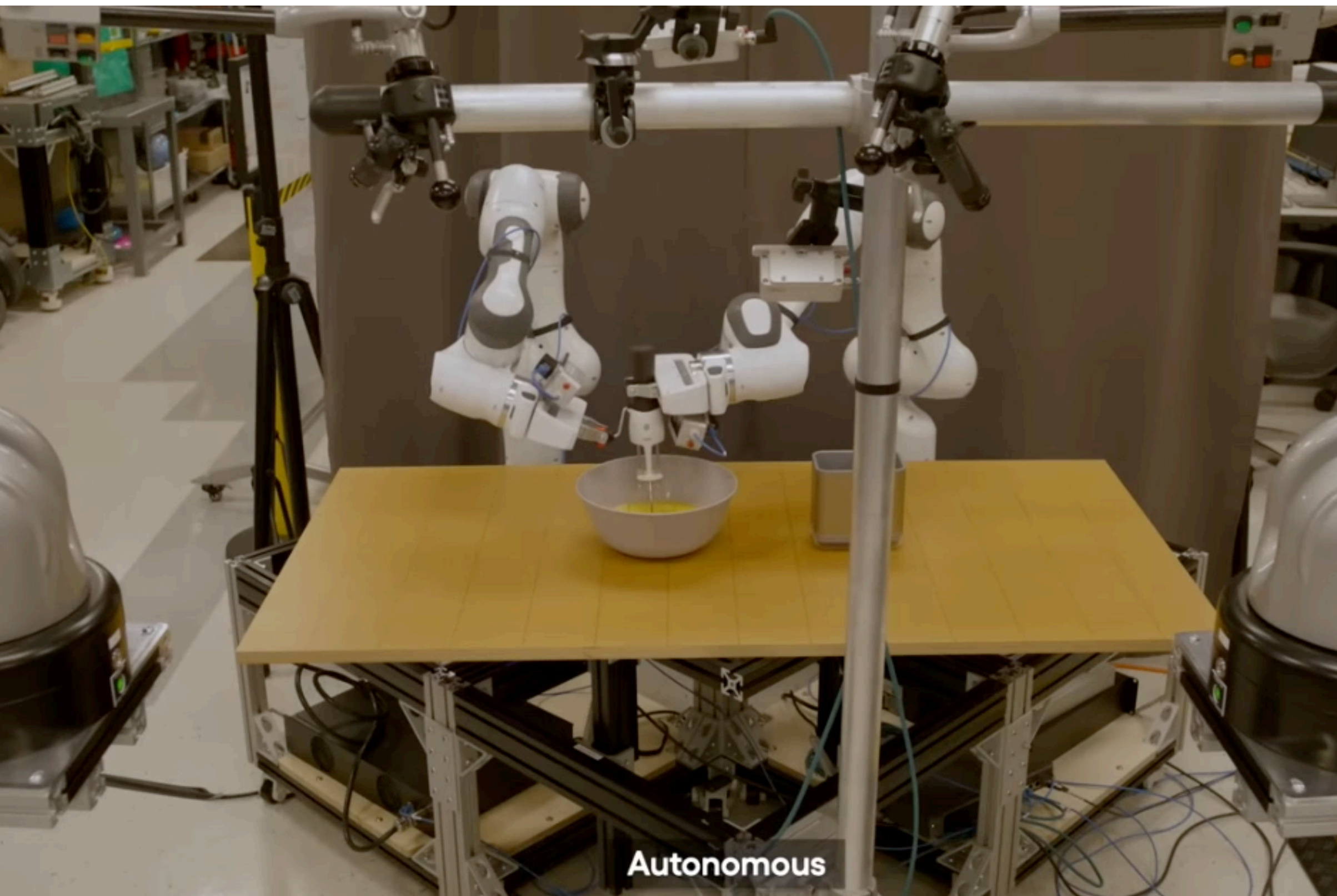
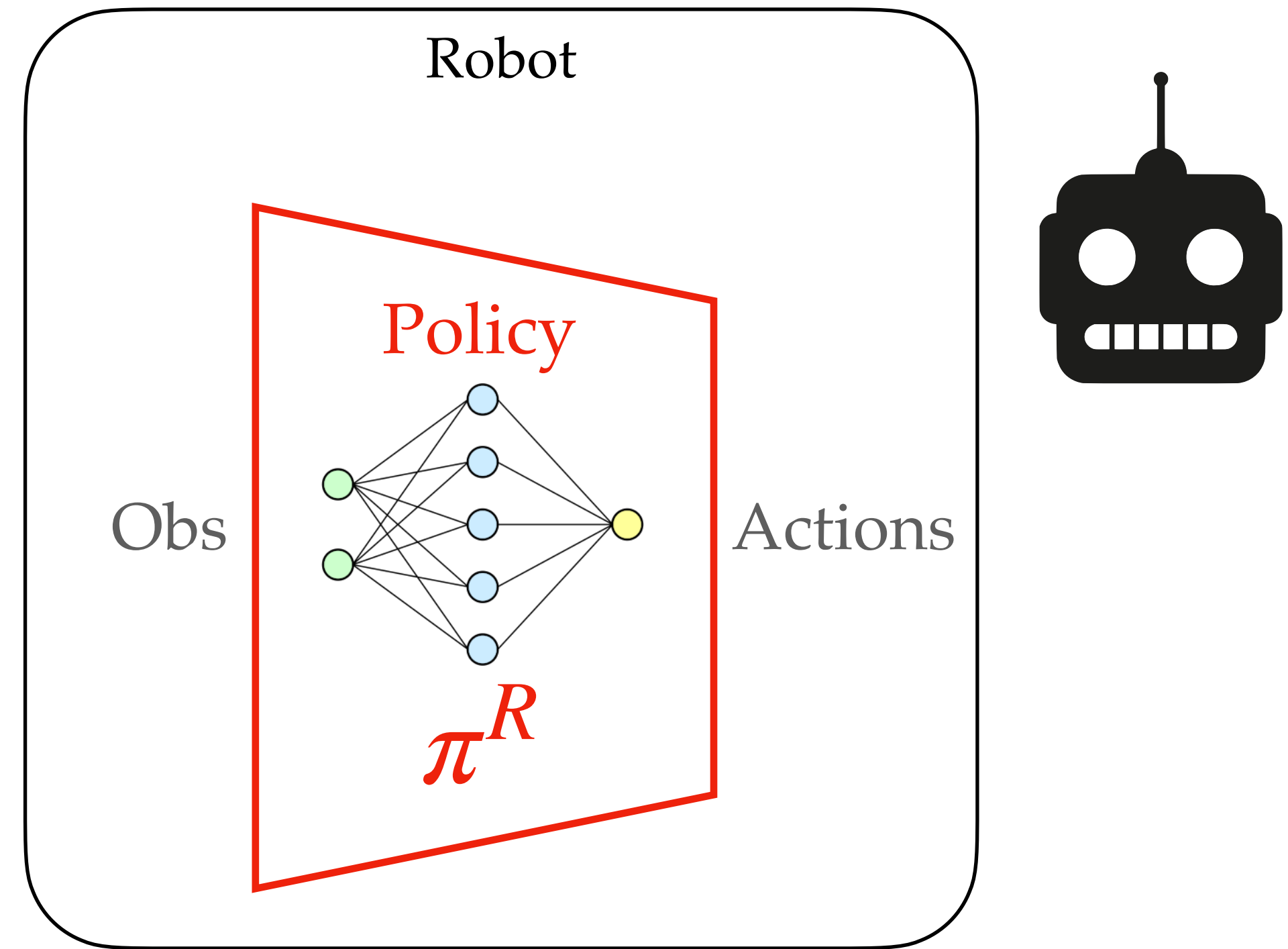  | I know what help to ask for | I know that I need help |
  |---|---|
  | *How do I ask for help?* | *How do I know that I need help?* |

# Outline

● Alignment problem

● Alignment process: Learning from human feedback

● Case Study 1: Learning from preferences

● Active Learning: Why and How?

● Revisiting Case Study 1: Making learning from preference *active*

● **Case Study 2: Active learning for black-box policies**

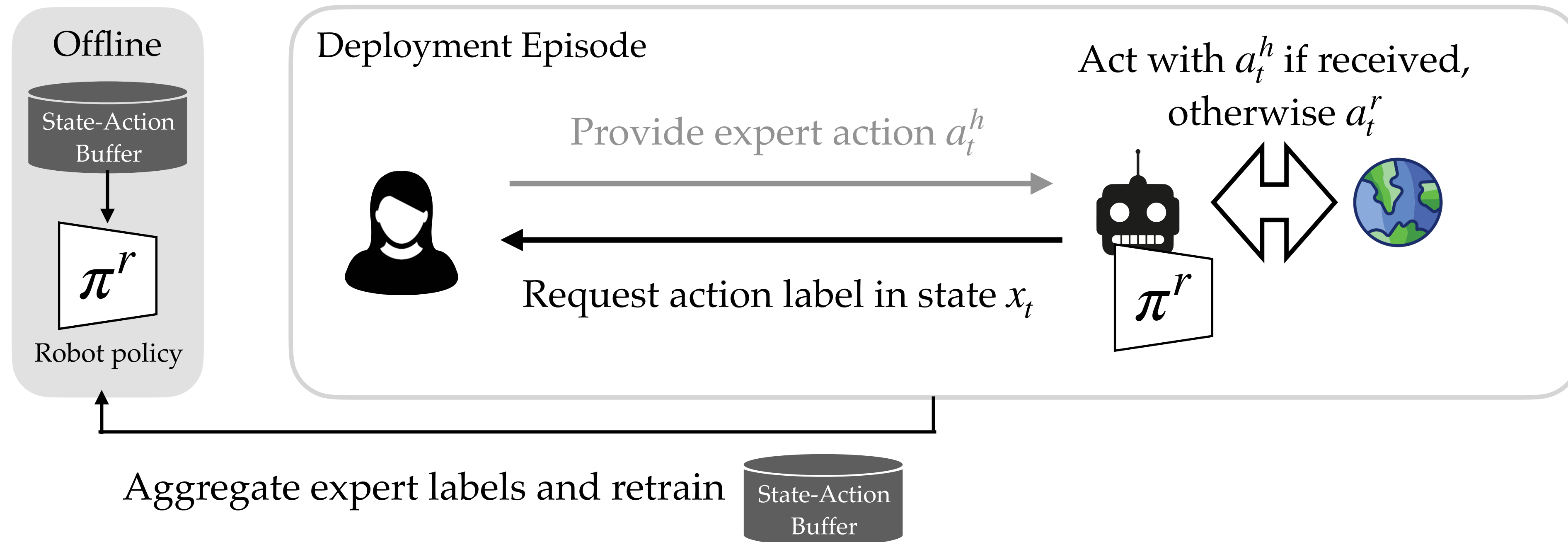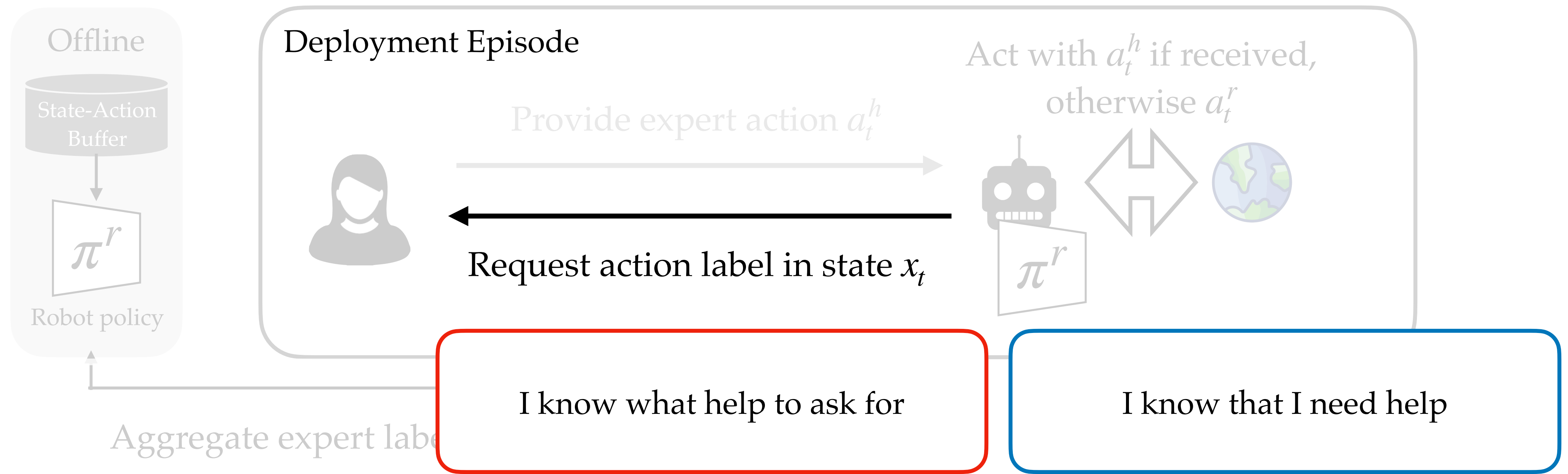# What about black-box policies?

# We benefited from our robot representing the reward distribution

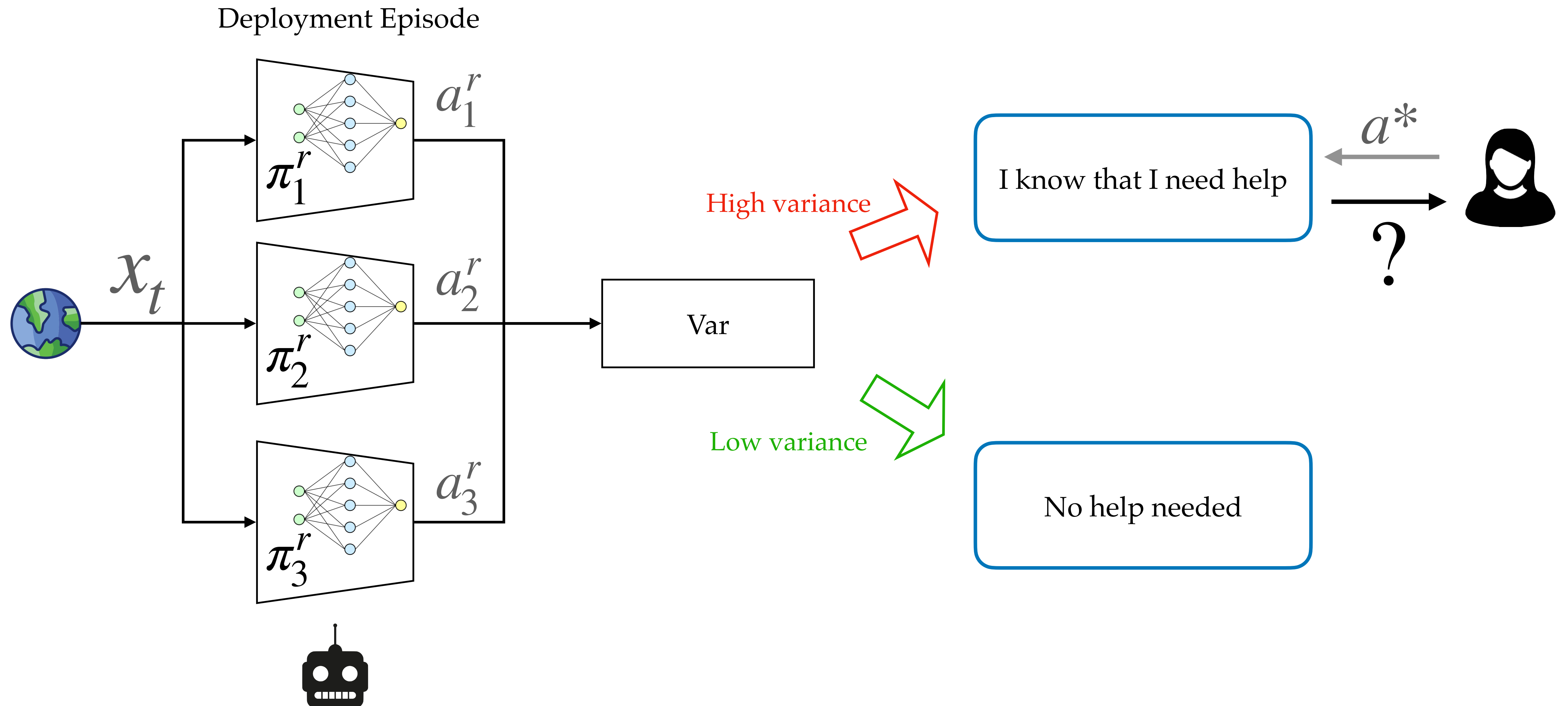# Online Interactive Imitation Learning

Offline

State-Action
Buffer

$\pi^r$

Robot policy

Deployment Episode

Act with $a_t^h$ if received,
otherwise $a_t^r$

Provide expert action $a_t^h$

Request action label in state $x_t$

$\pi^r$

Aggregate expert lab

I know what help to ask for

I know that I need help

*How do I ask for help?*

*How do I know that I need help?*

Expert Teleoperation

✓ Ensembles

✓ Conformal prediction

# Uncertainty Quantification: Ensemble Disagreement

Deployment Episode



$a_1^r$

$\pi_1^r$

$x_t$

$a_2^r$

$\pi_2^r$

$a_3^r$

$\pi_3^r$

Var

High variance

Low variance

I know that I need help

No help needed

$a*$

?

# Uncertainty Quantification: Conformal Prediction



Deployment Episode

$x_t$

$\pi^r$

$a^r$

Conformal Prediction Uncertainty

Big intervals

Small intervals

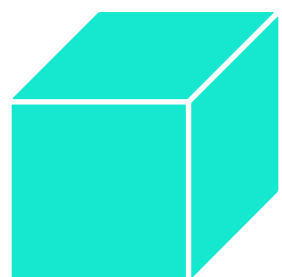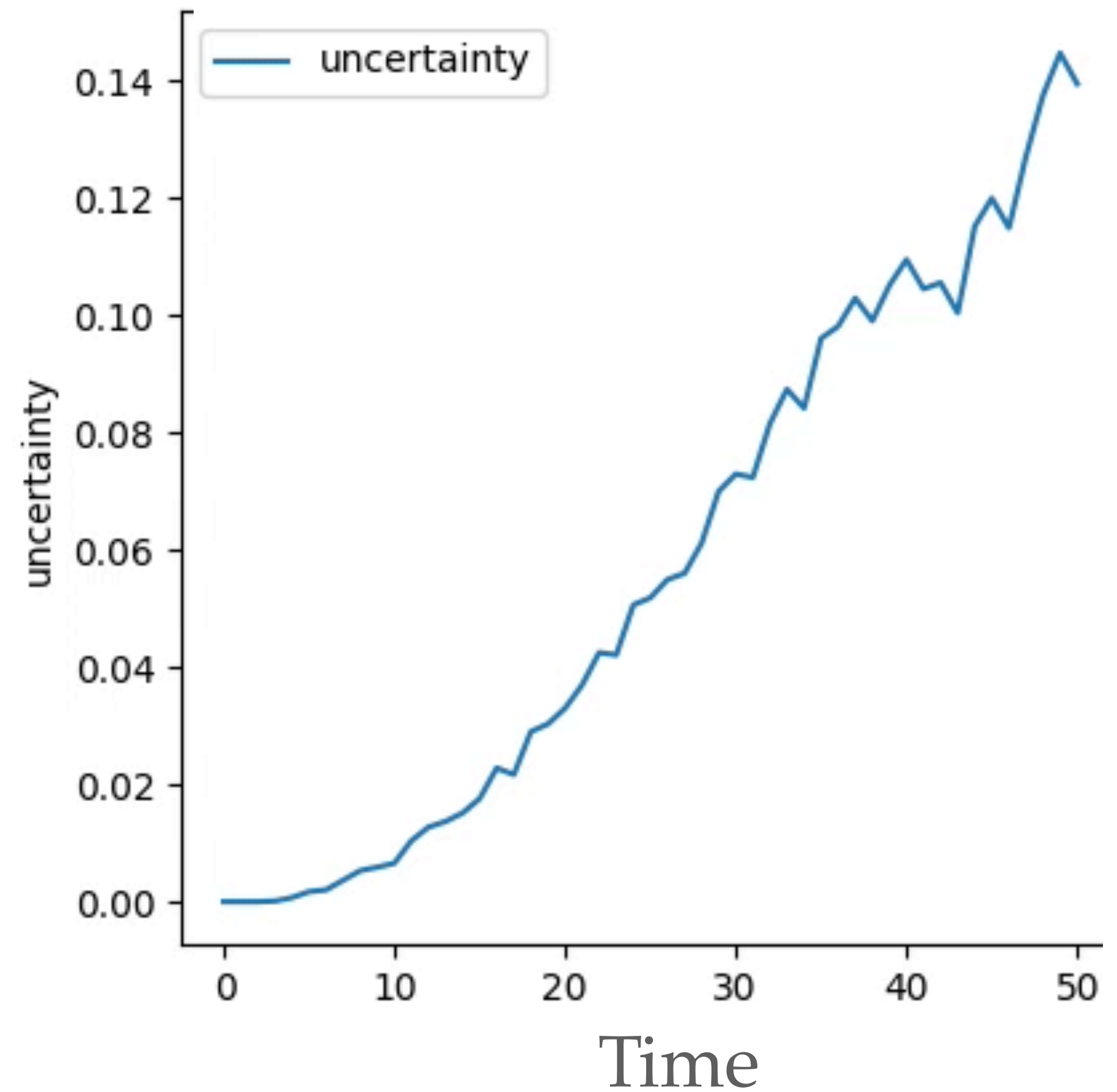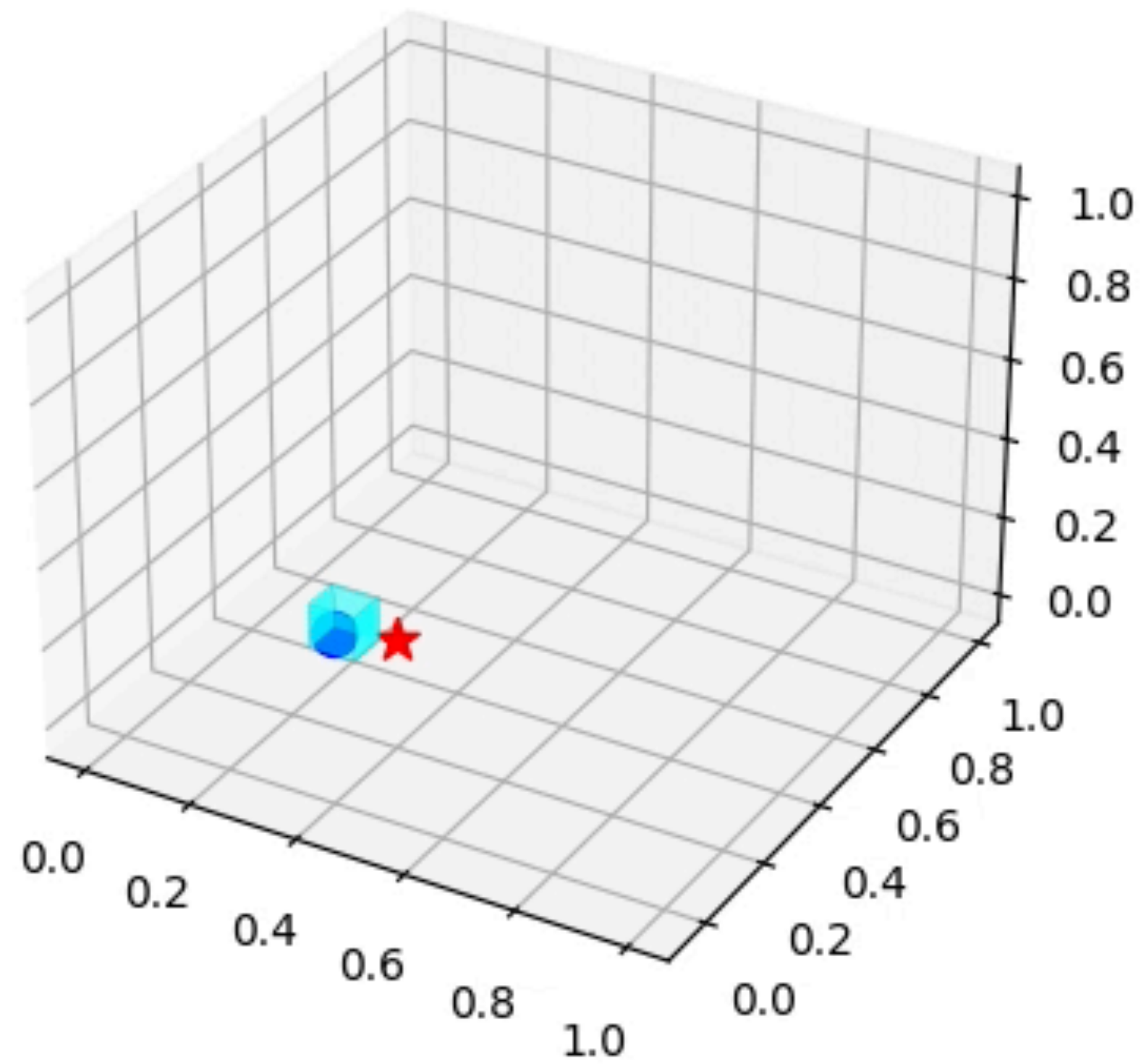I know that I need help

No help needed

$a*$

?

**Online Conformal:** Whenever we observe the human's ground truth action, use the prediction error to adaptively adjust uncertainty estimate.

# Uncertainty Quantification: Online Conformal Prediction

# Summary and Questions?

- Learning from human feedback seeks to align robot behaviors to human intentions

  - Human feedback can occur through multiple types (i.e. preferences)

- Active learning shifts the interaction paradigm by allowing the learner to *actively* request feedback (labels) from the human.

  - Effective learning occurs by coming up with a good strategy for querying the human (*strategically asking for help*)

- To enable an effective active robot learner, 2 necessary design decisions:

✔ Preferences

✔ Expert control

I know what help to ask for

*How do I ask for help?*

I know that I need help

✔ Information gain

✔ Ensembles

✔ Conformal prediction

*How do I know that I need help?*