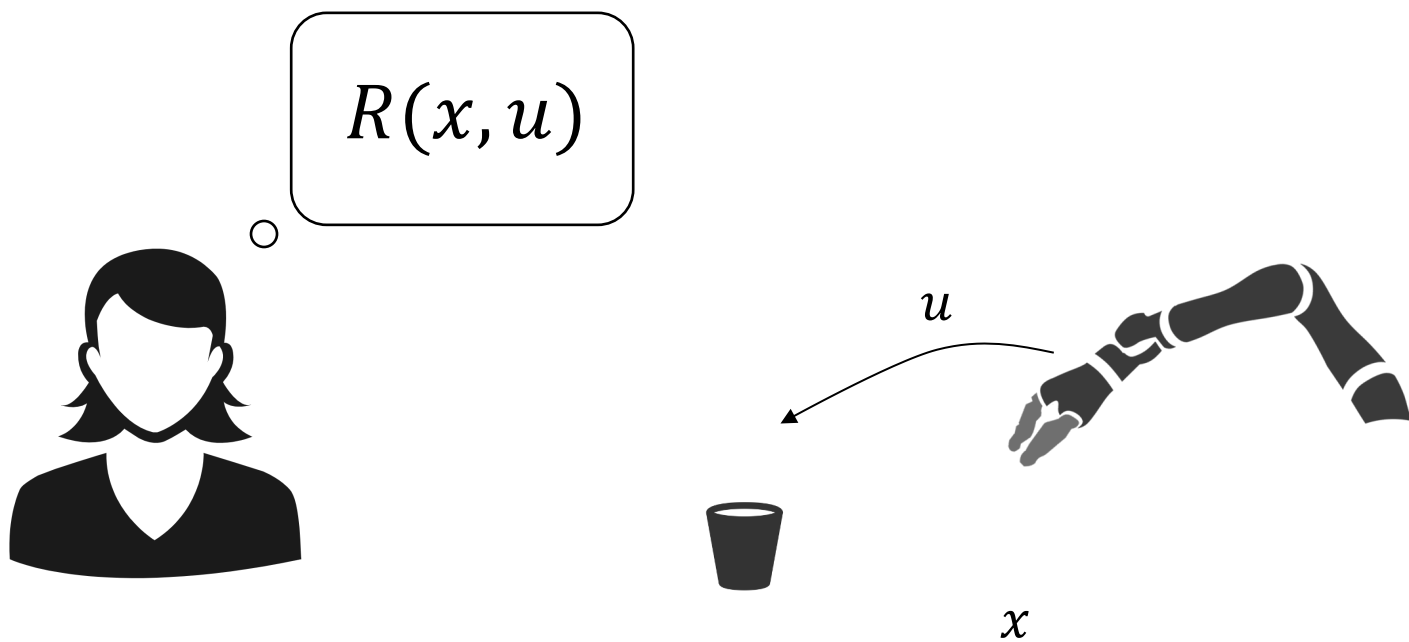
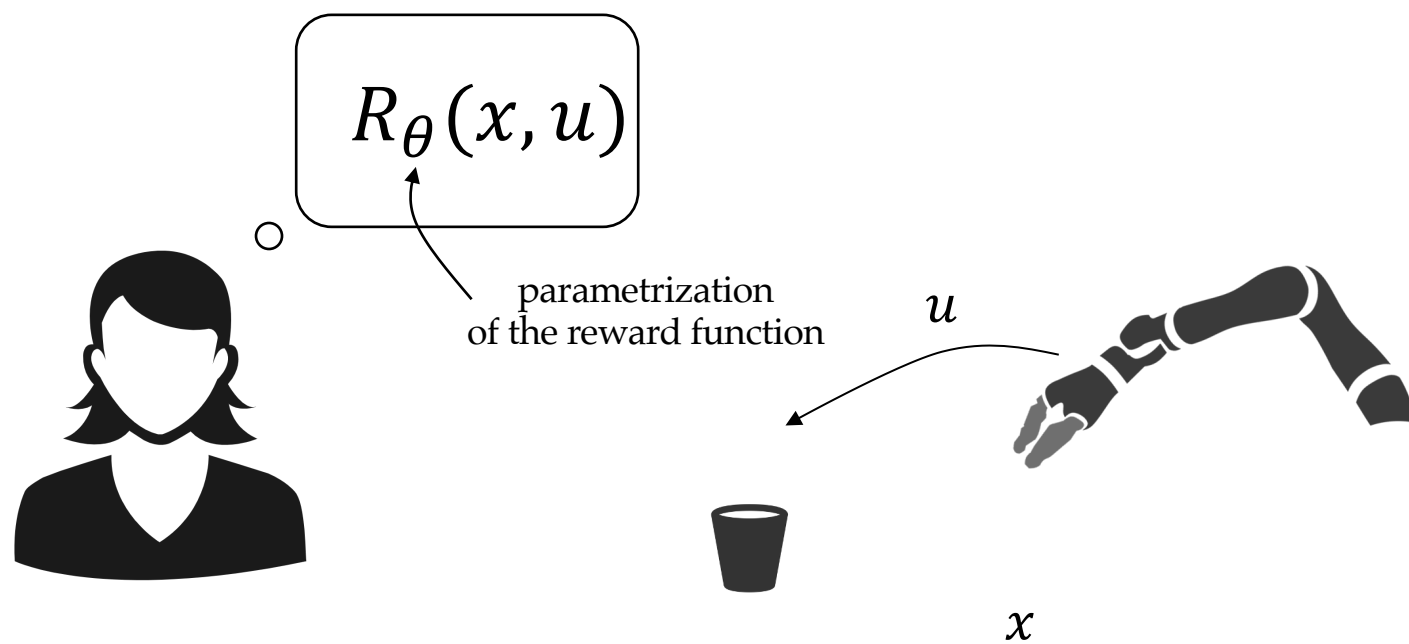
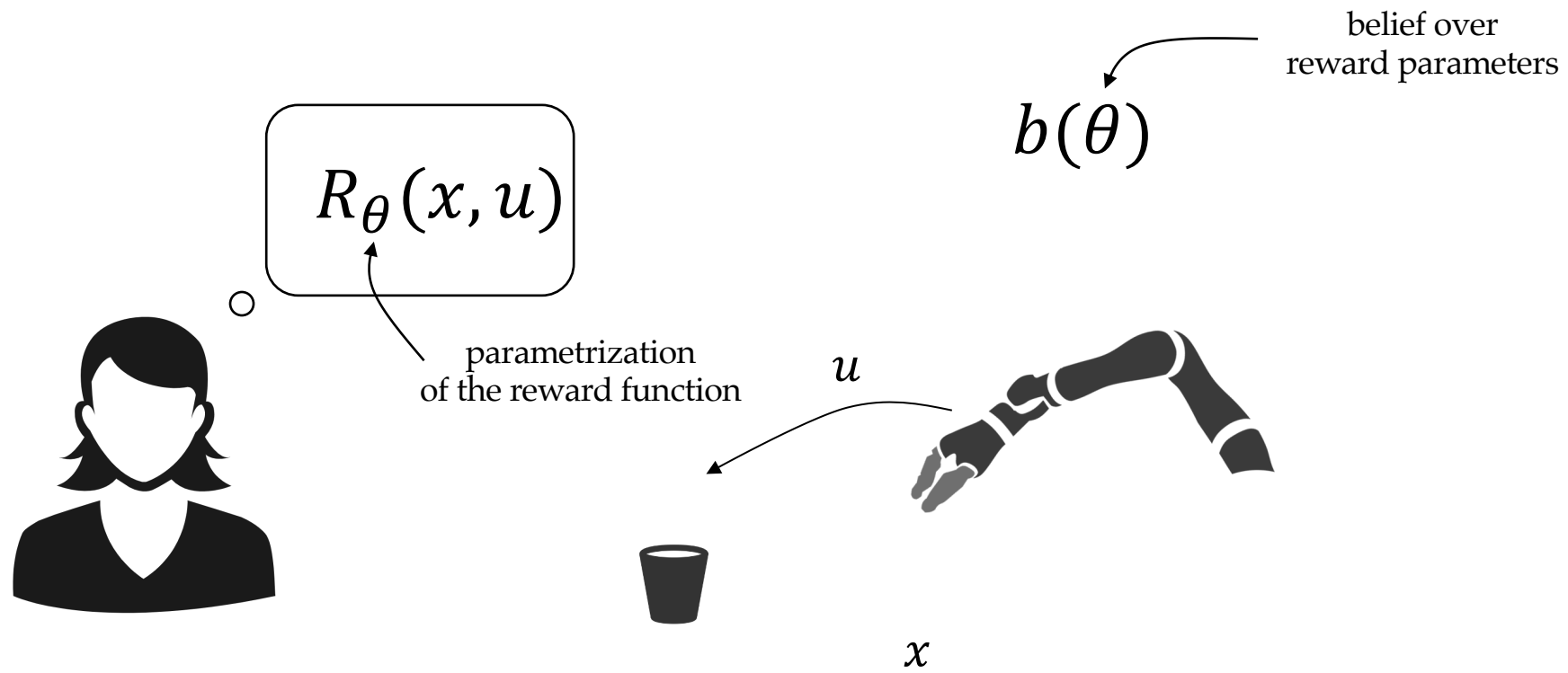


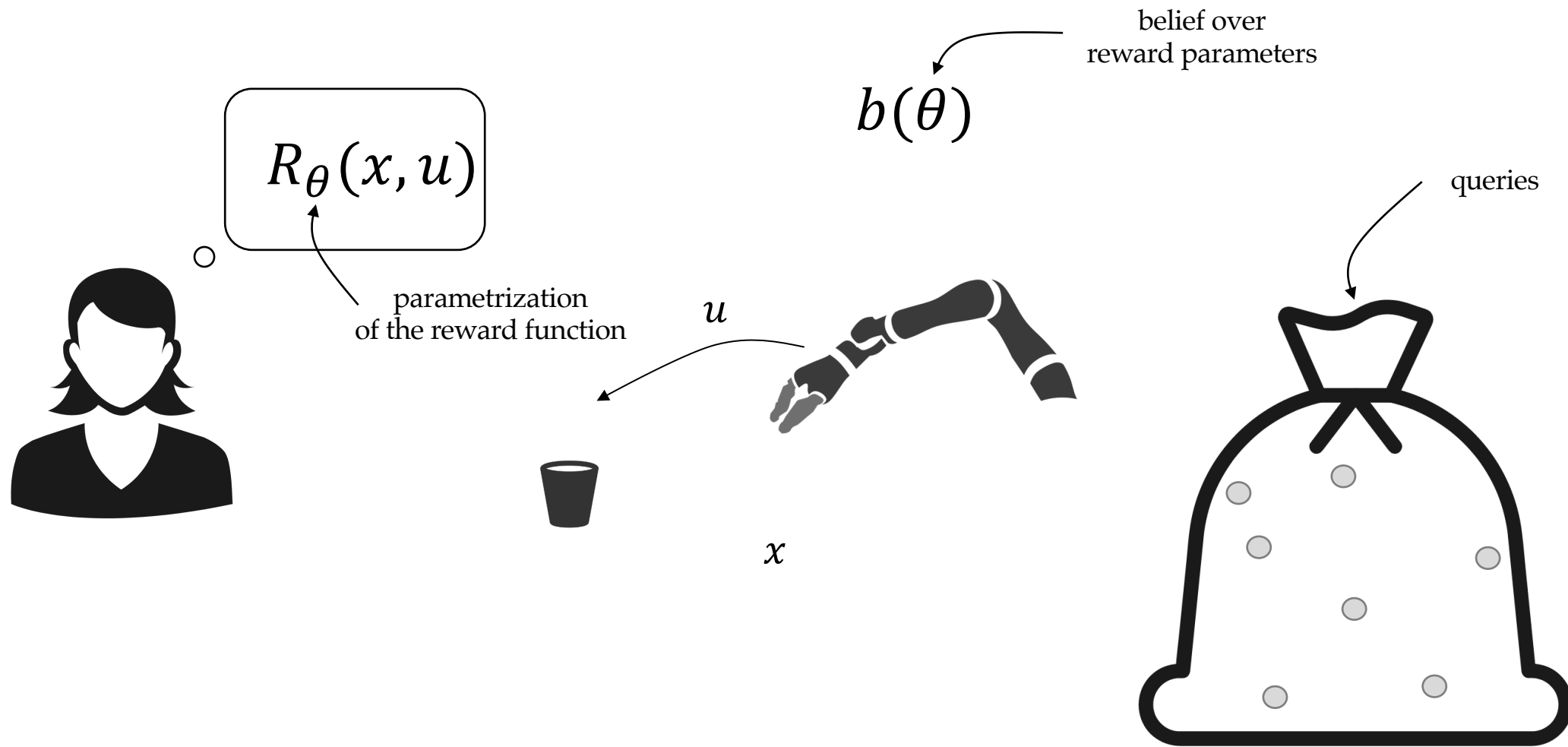
Learning Objectives and Preferences:

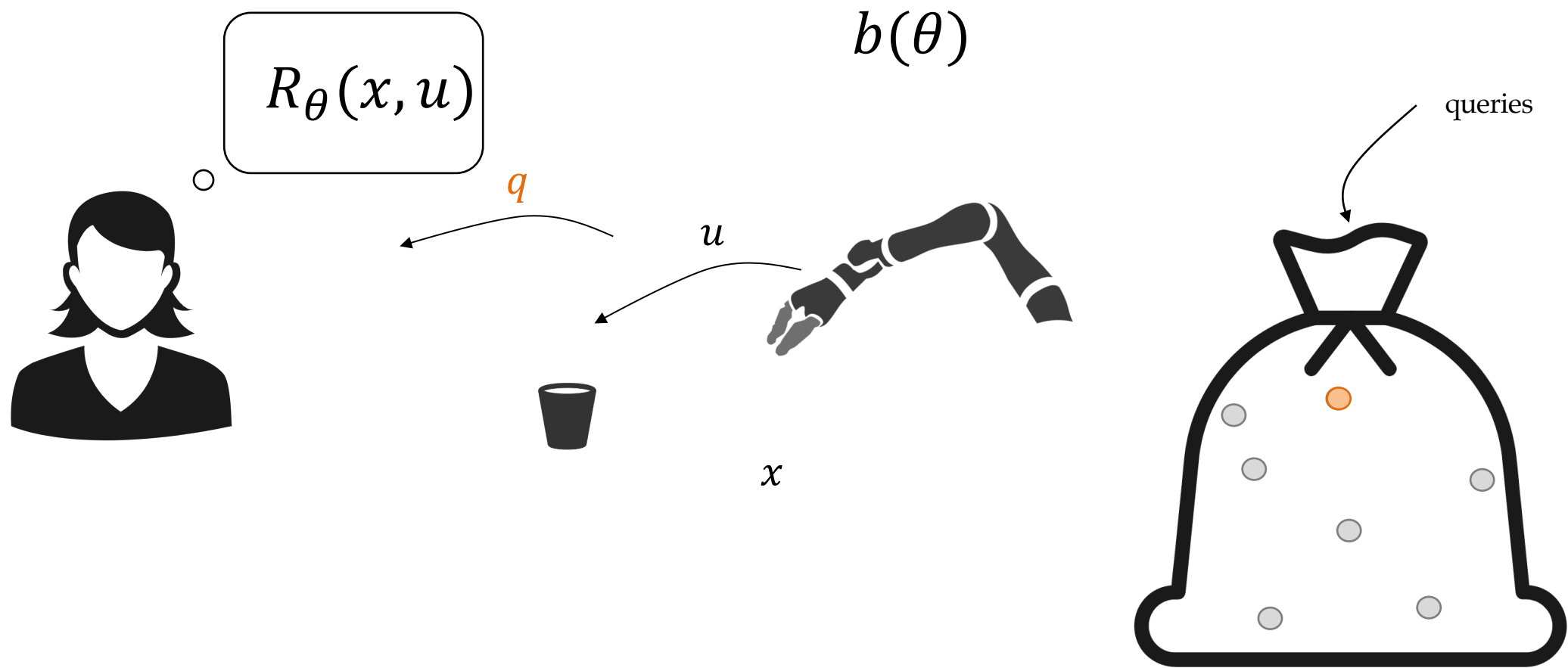
How? *Actively*

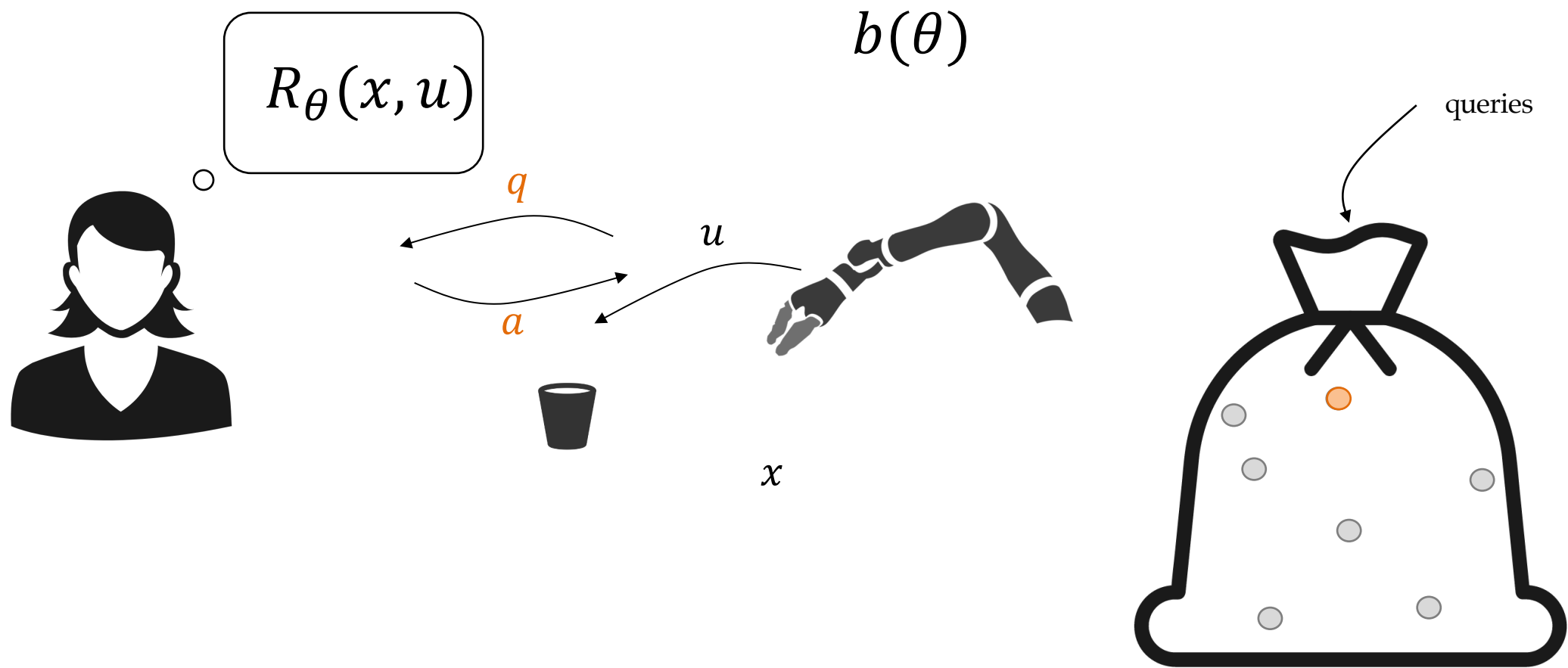


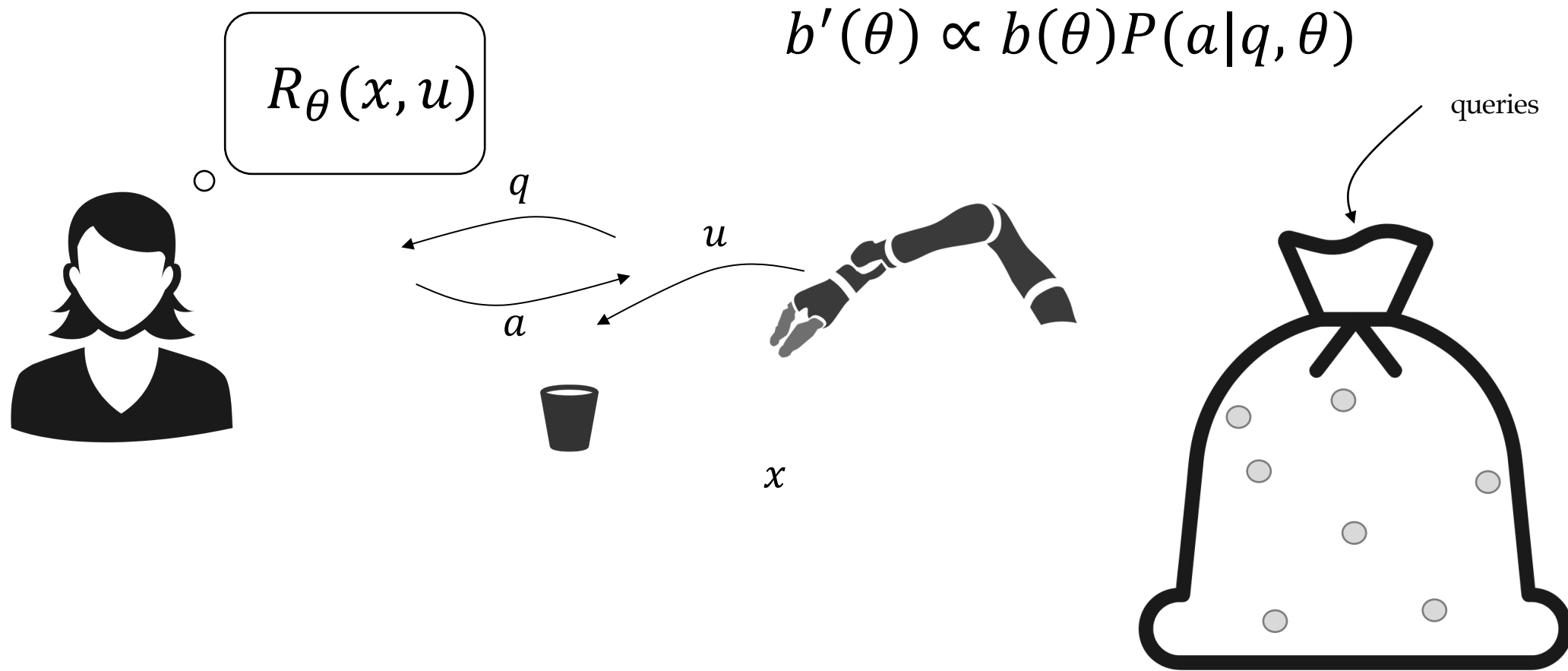




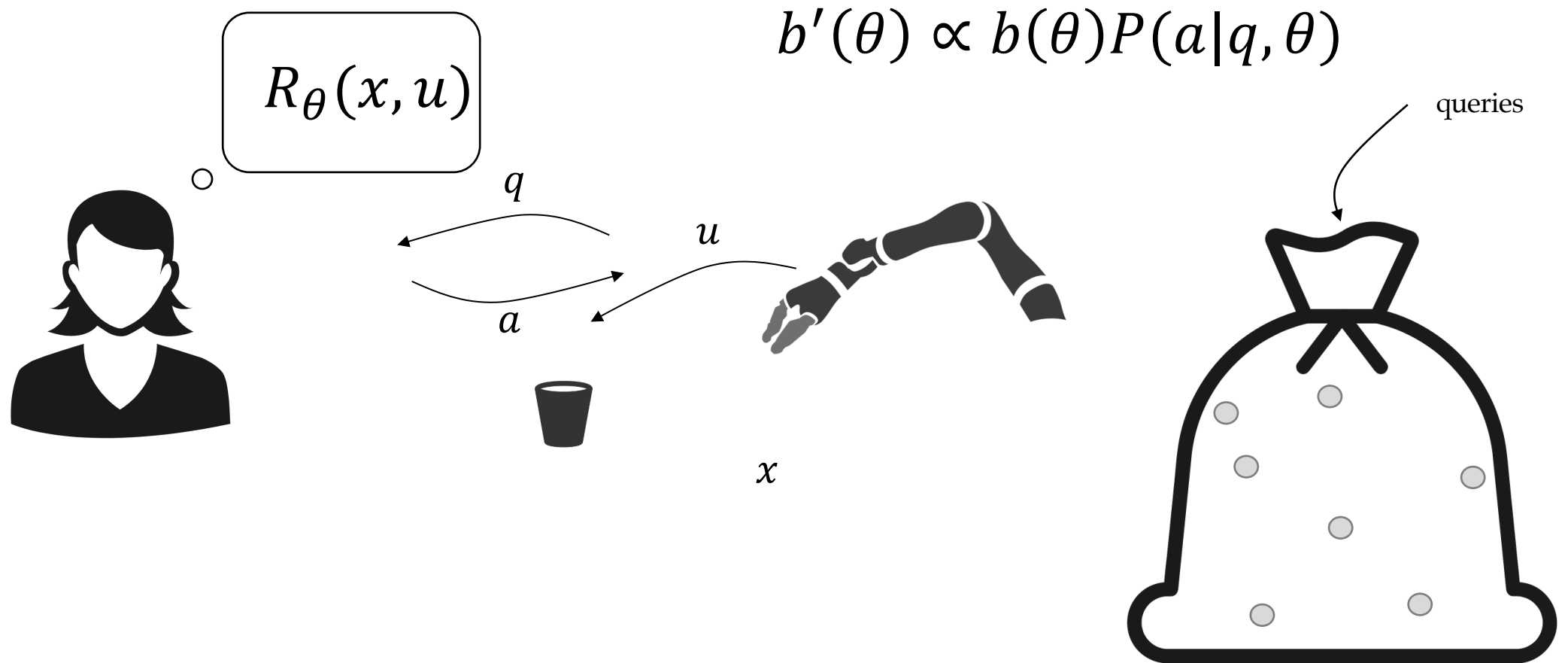


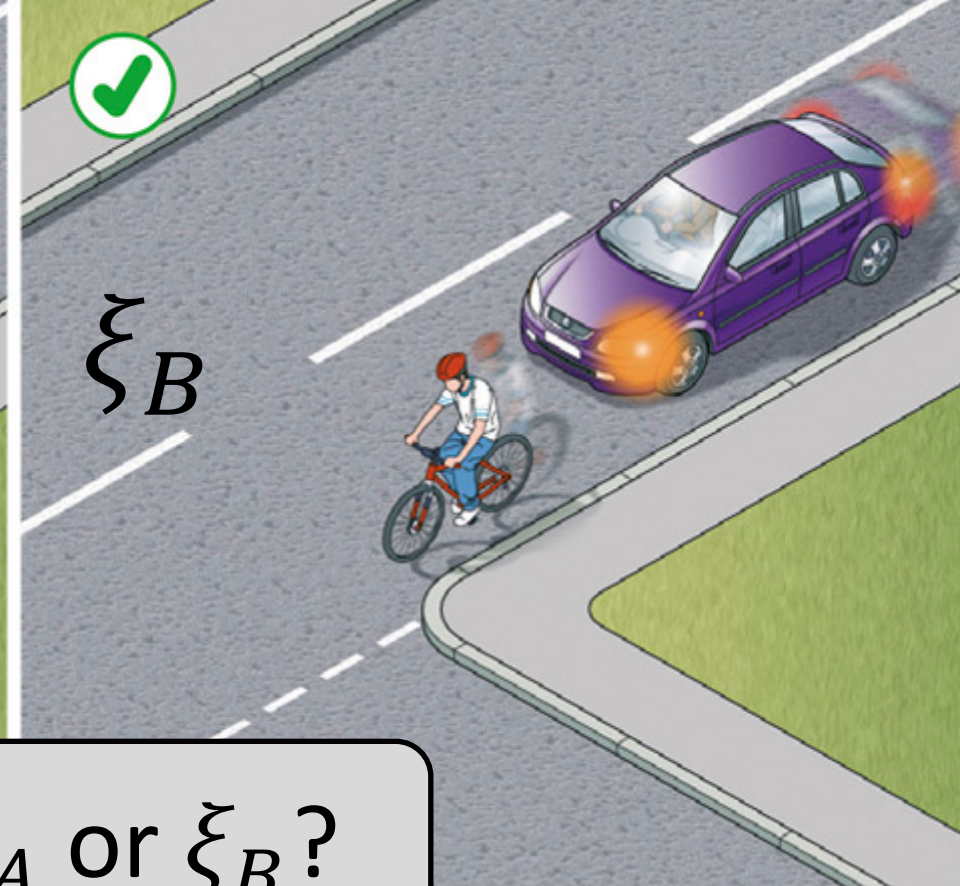
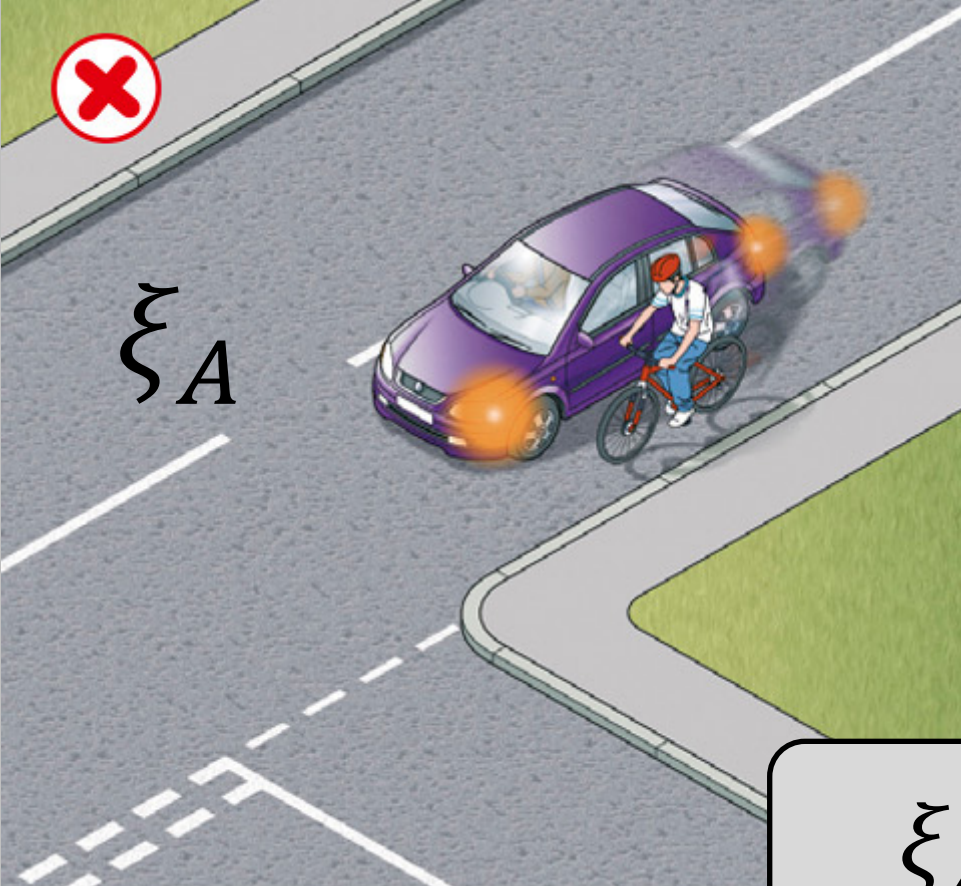






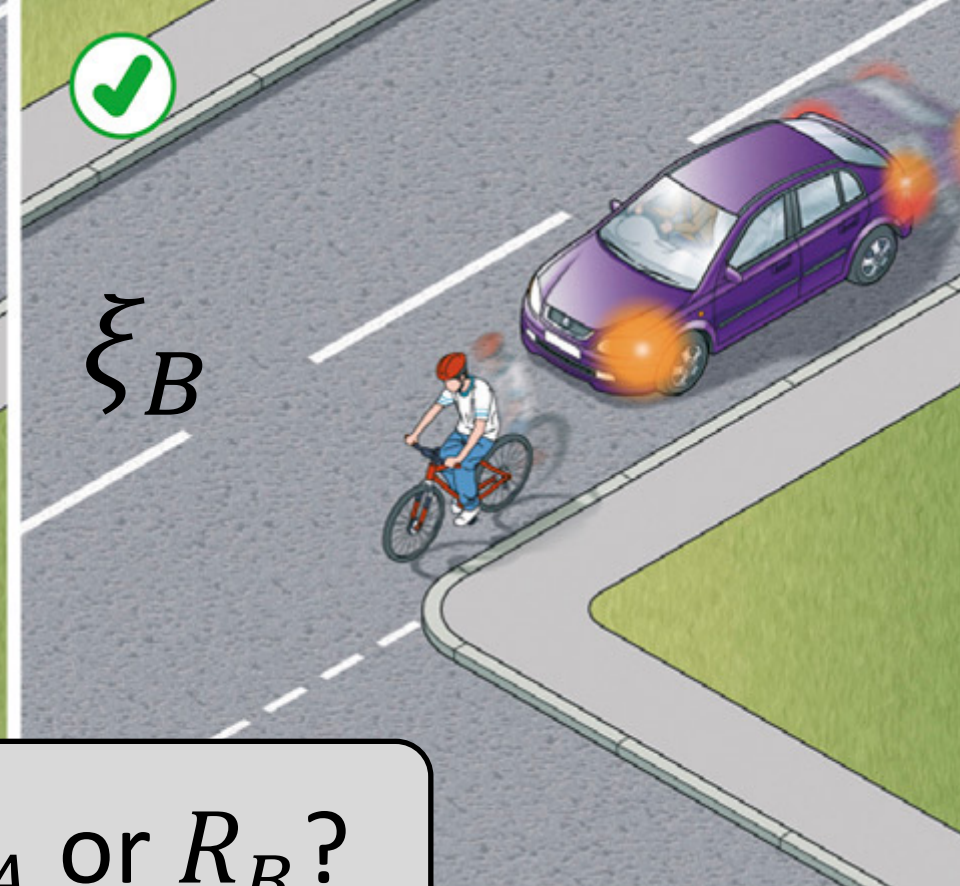
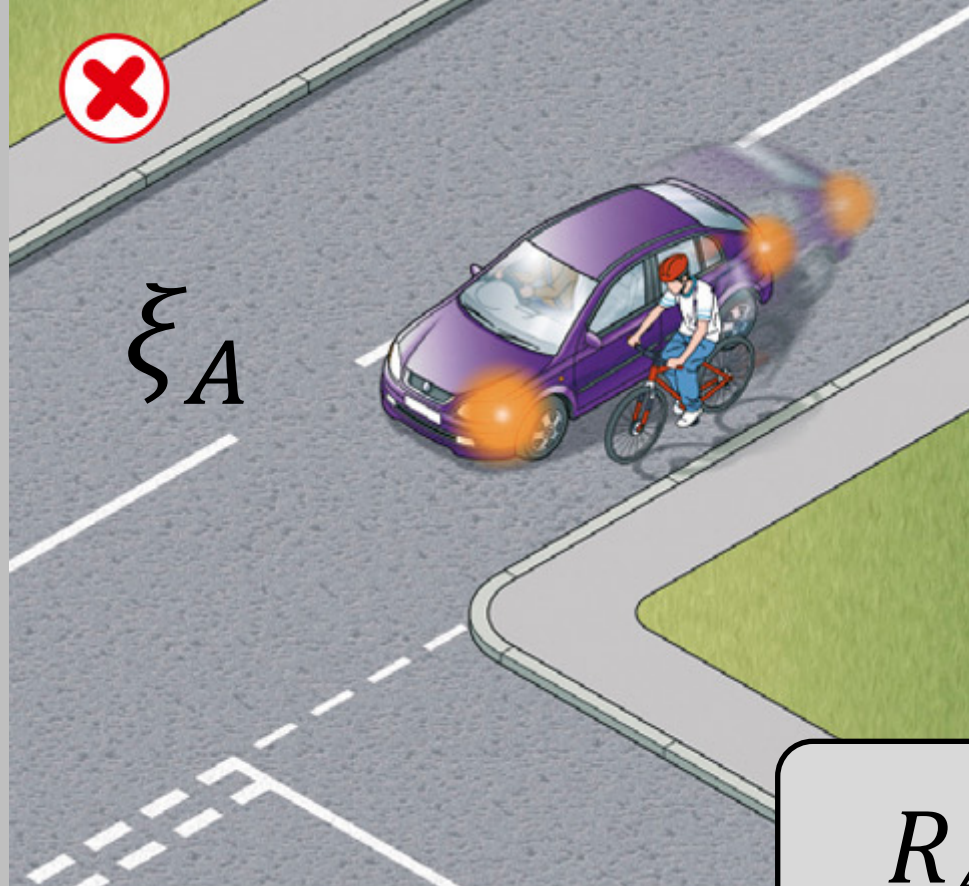
Where do queries come from?



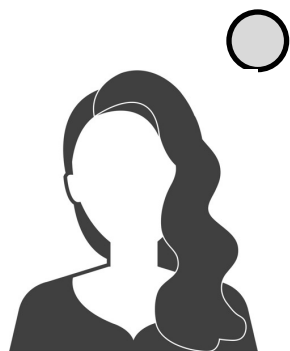


ξ_A or ξ_B ?

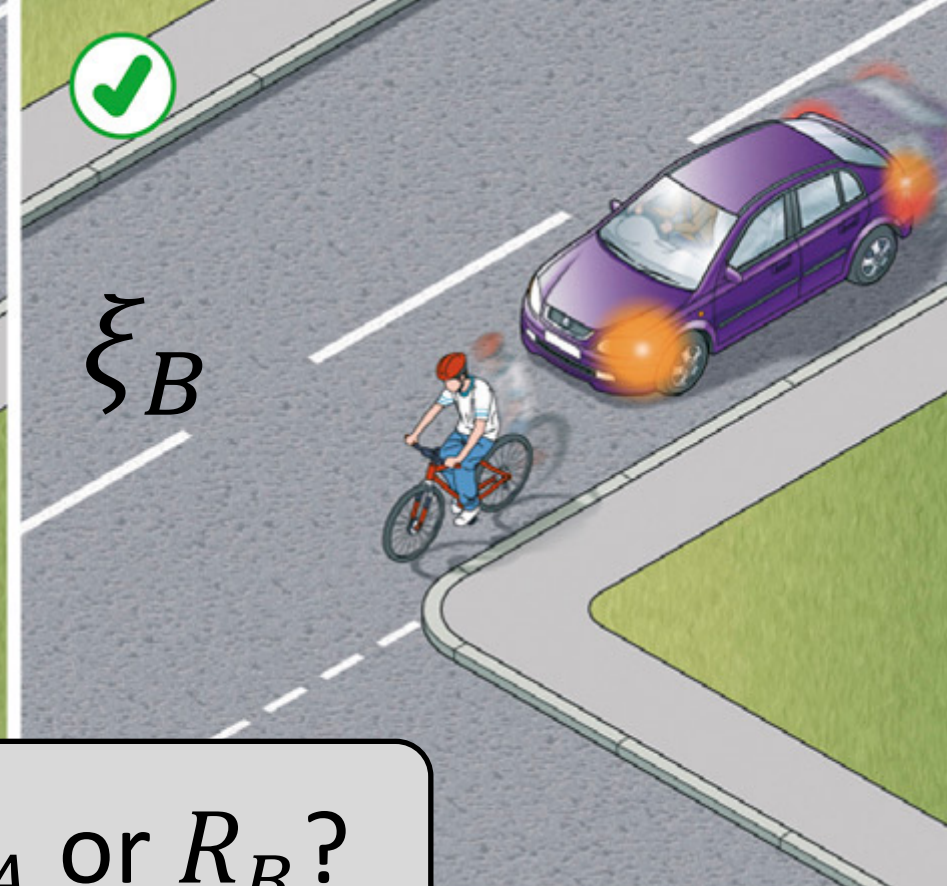
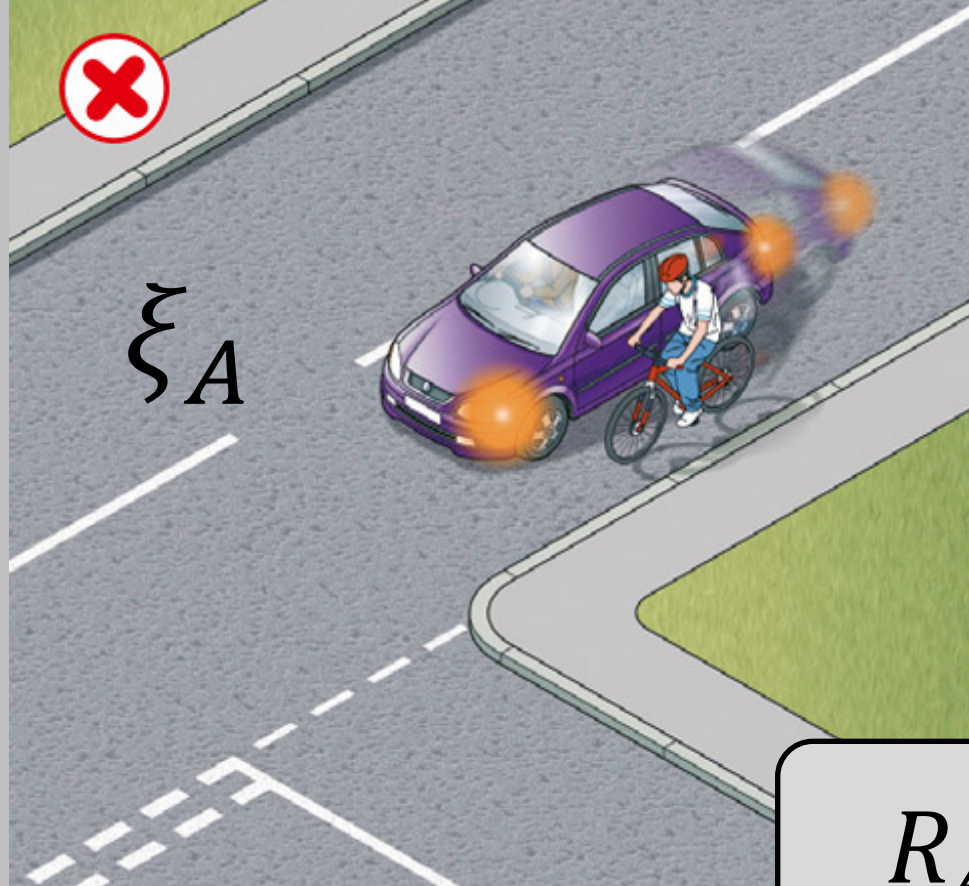




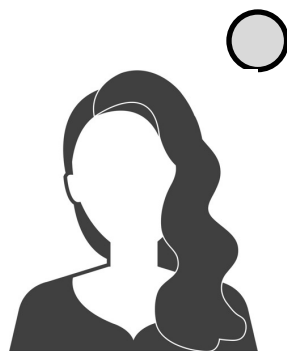
R_A or R_B ?



$$R = \theta \cdot \phi(\xi)$$



R_A or R_B ?

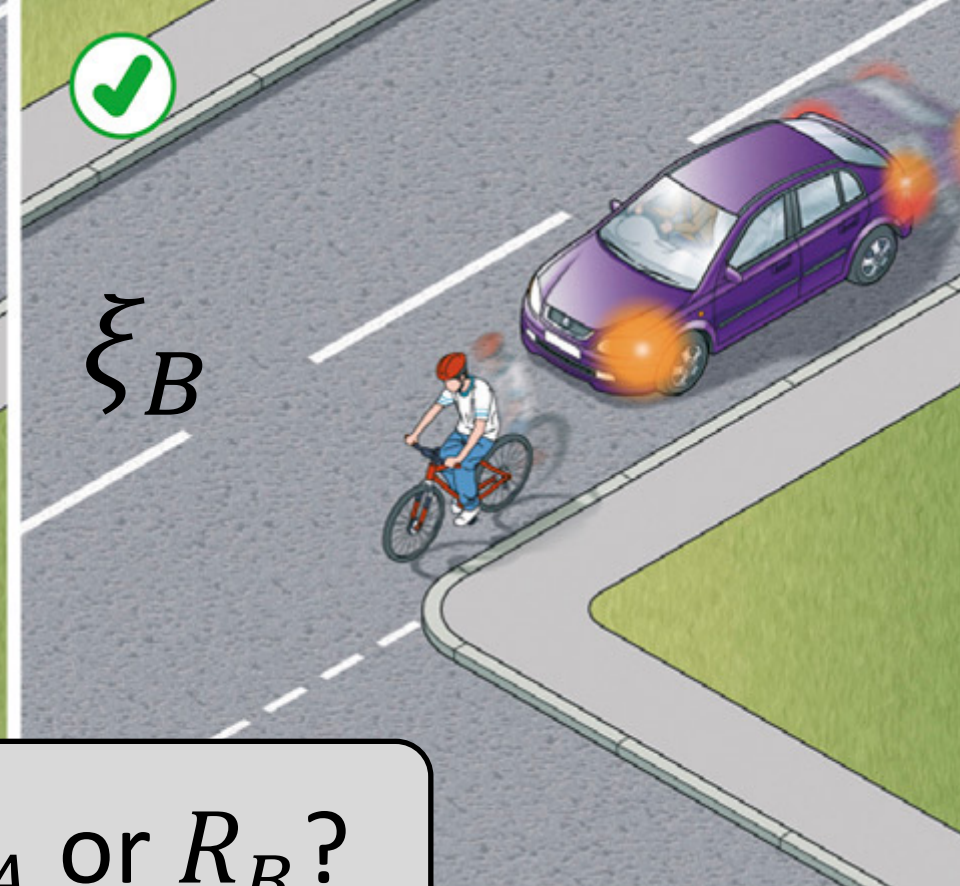
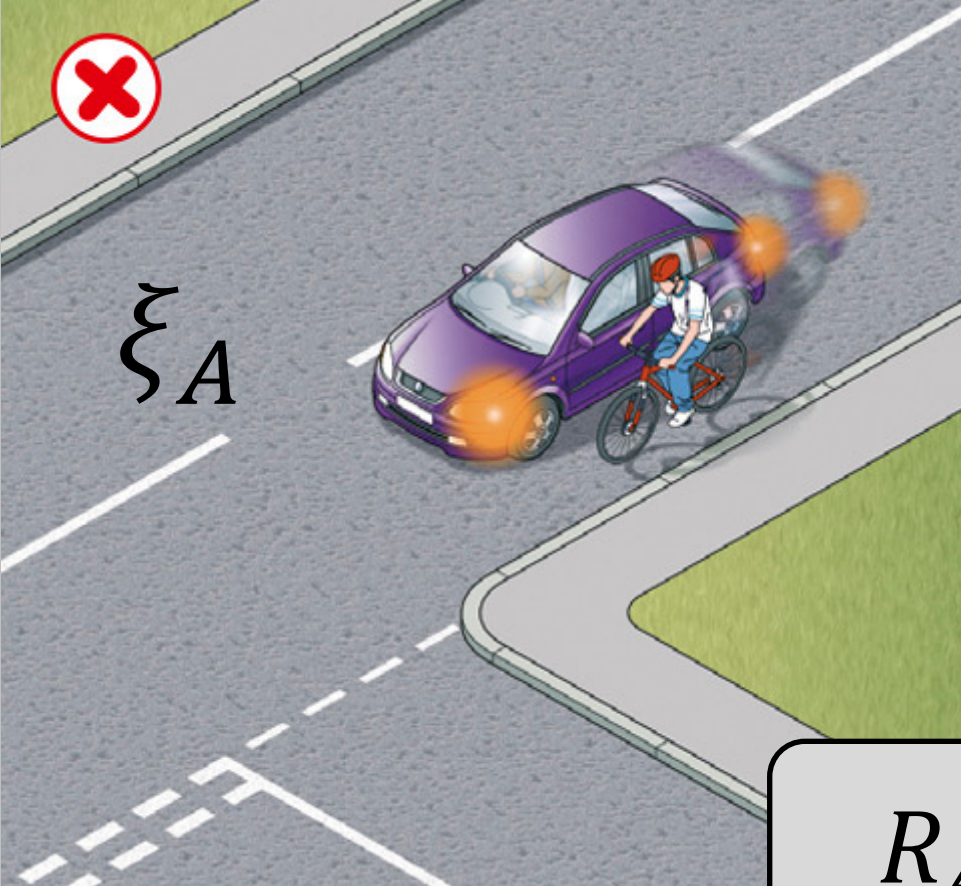


$$R = \theta \cdot \phi$$

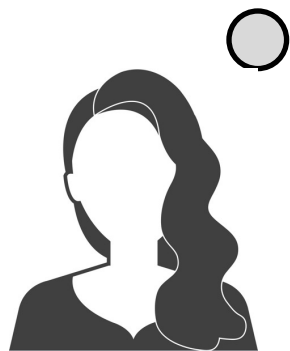
or

$$\theta \cdot \phi_A > \theta \cdot \phi_B$$

$$\theta \cdot \phi_A < \theta \cdot \phi_B$$



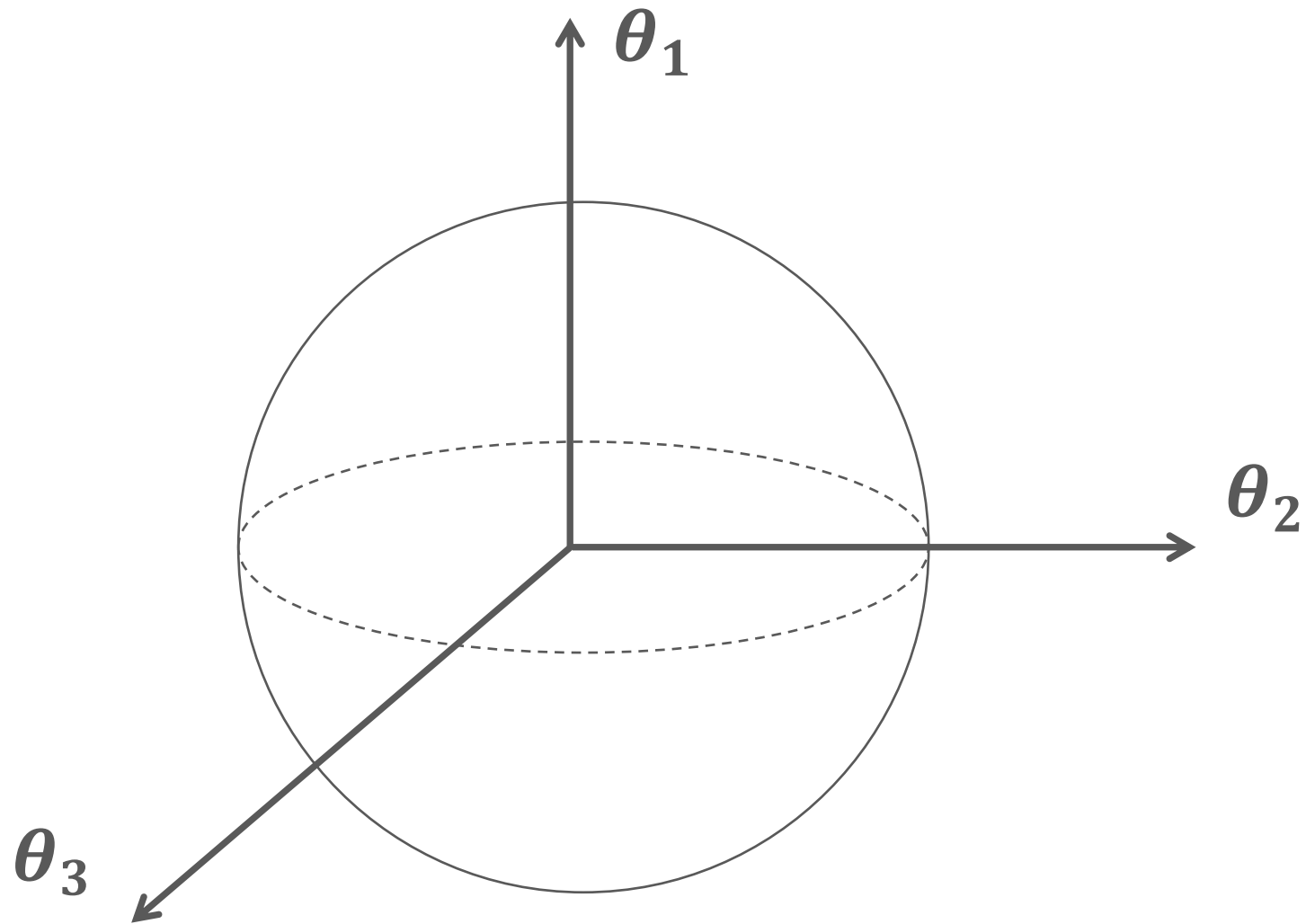
R_A or R_B ?



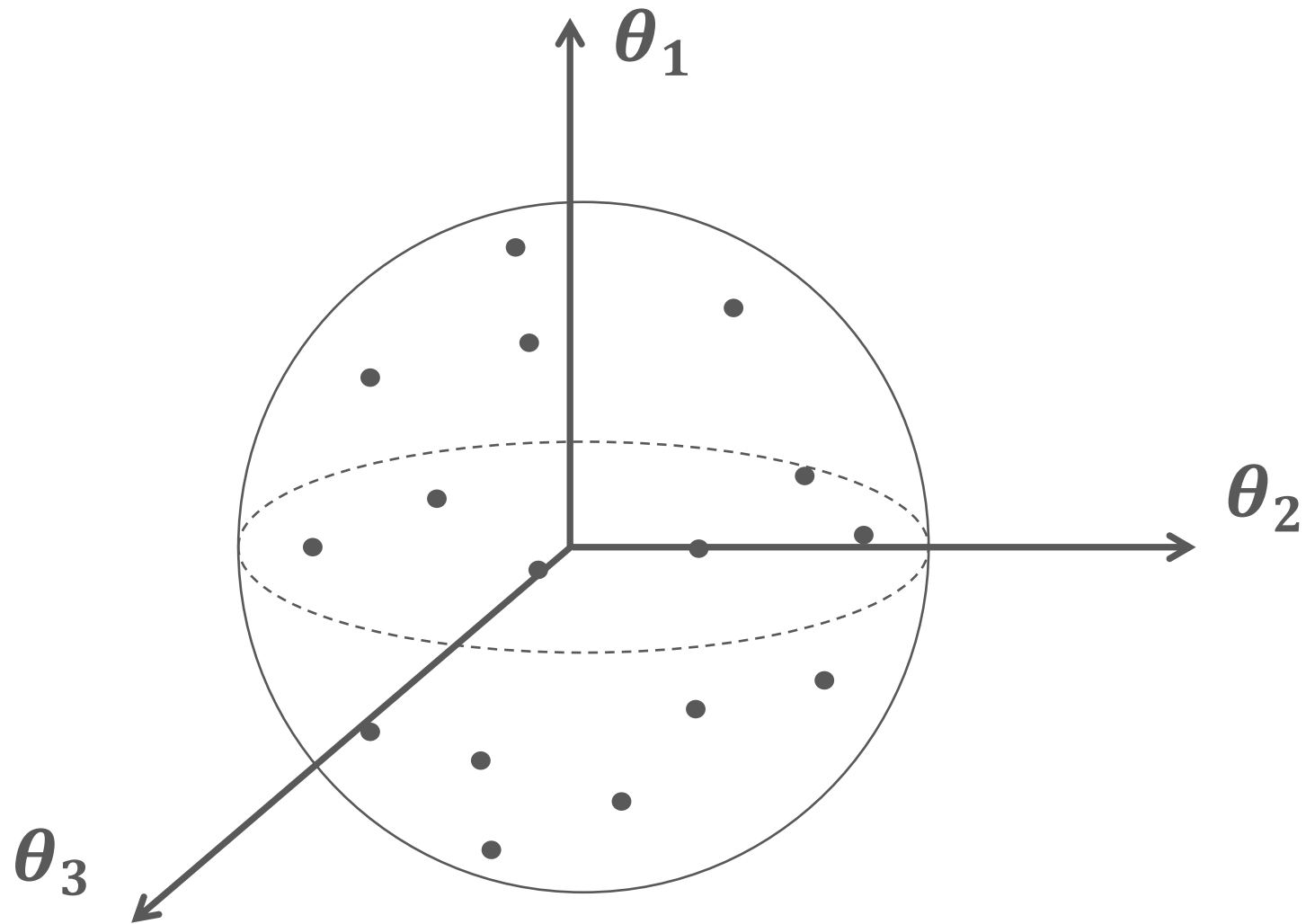
$$R = \theta \cdot \phi$$

$$\psi = (\phi_A - \phi_B)$$

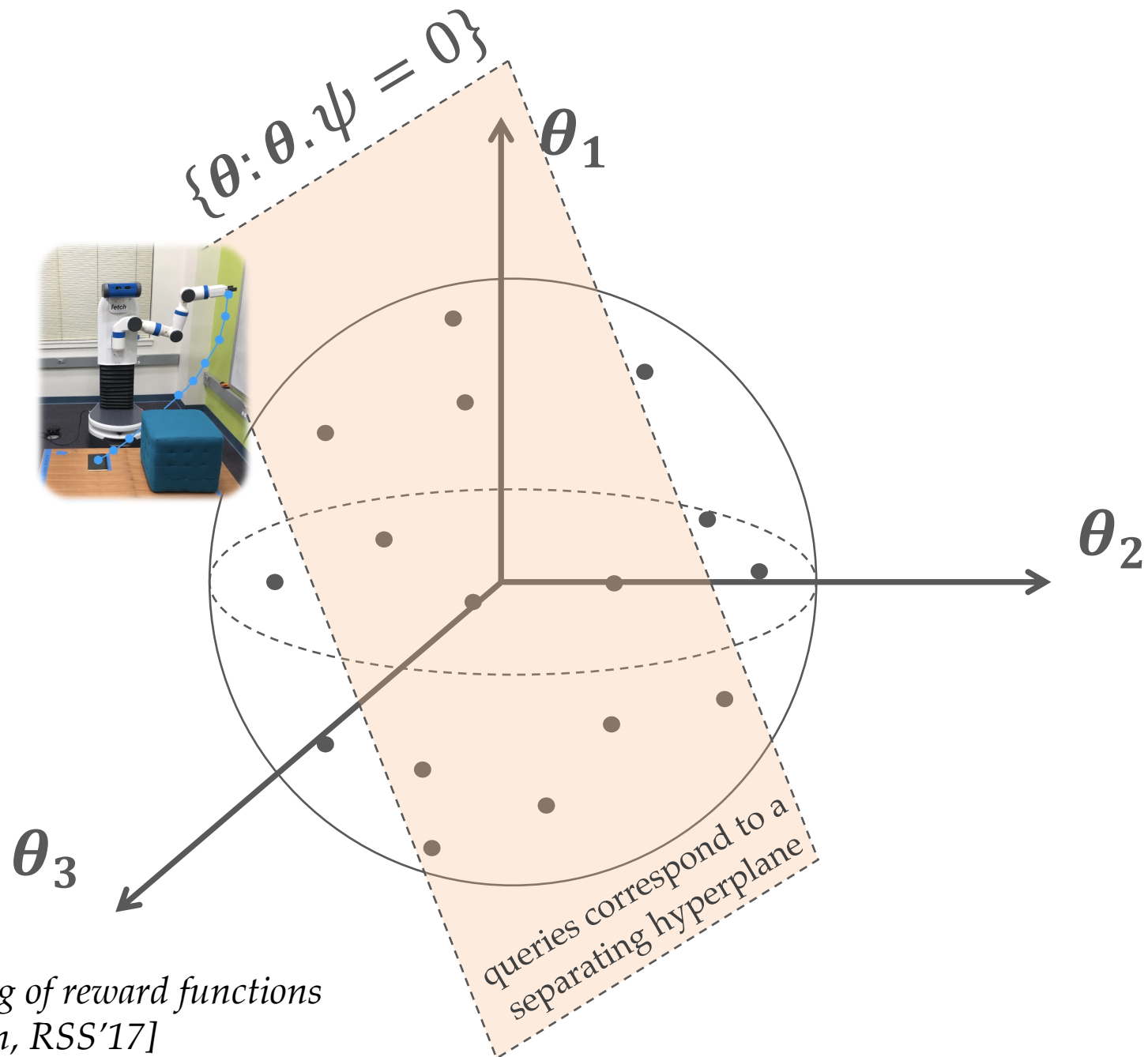
$$\theta \cdot \psi > 0 \text{ or } \theta \cdot \psi < 0$$



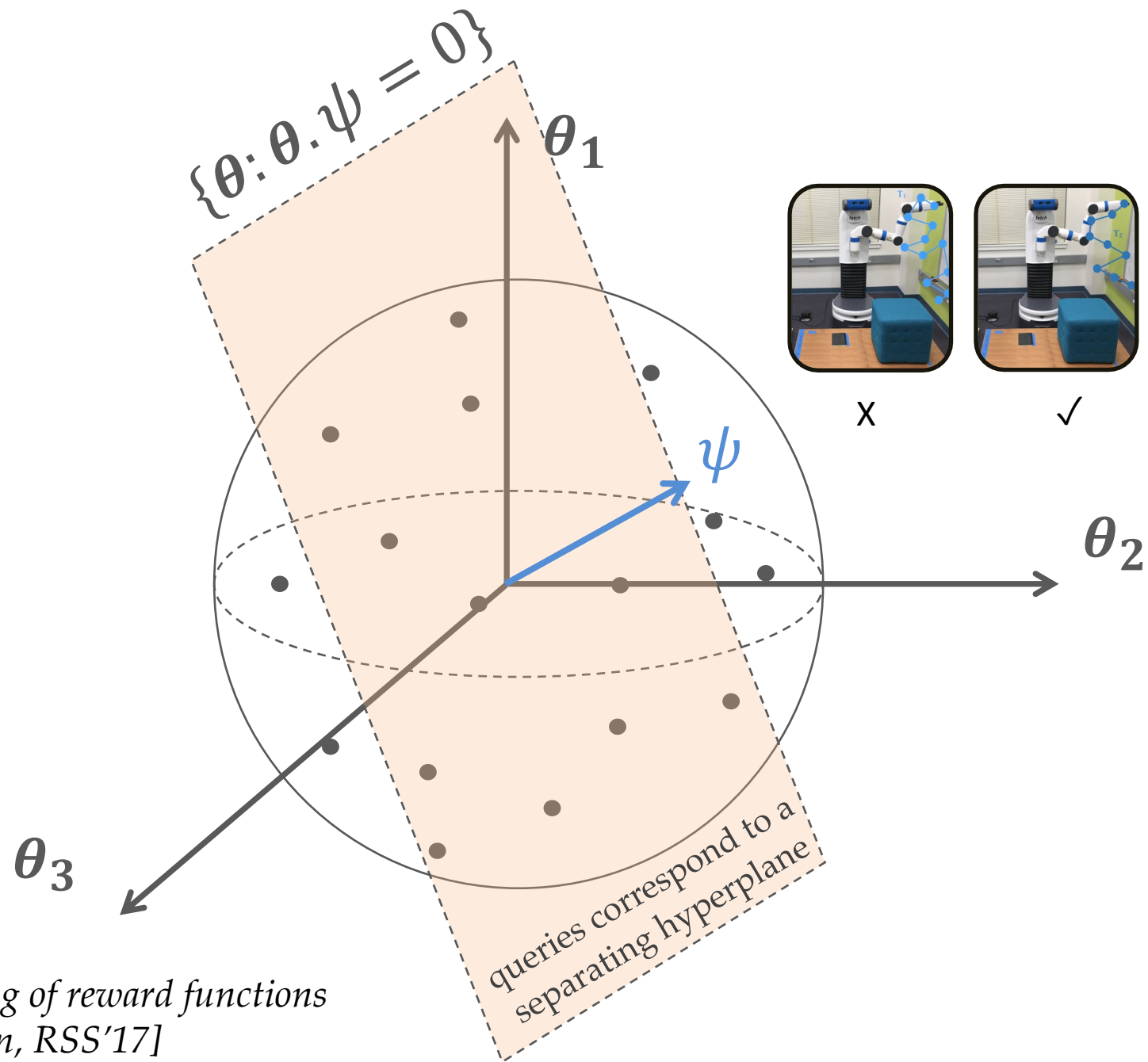
Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]



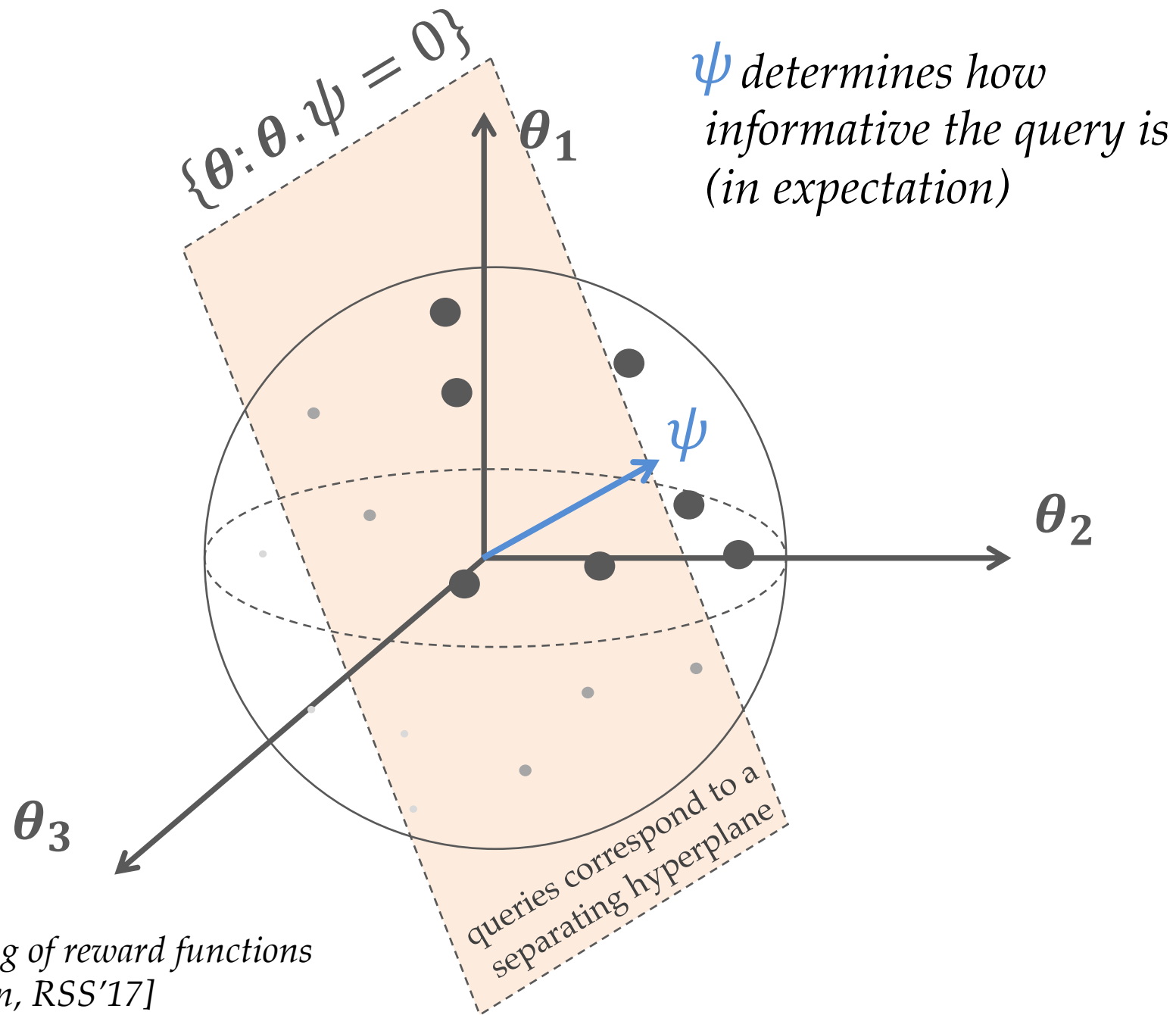
Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]



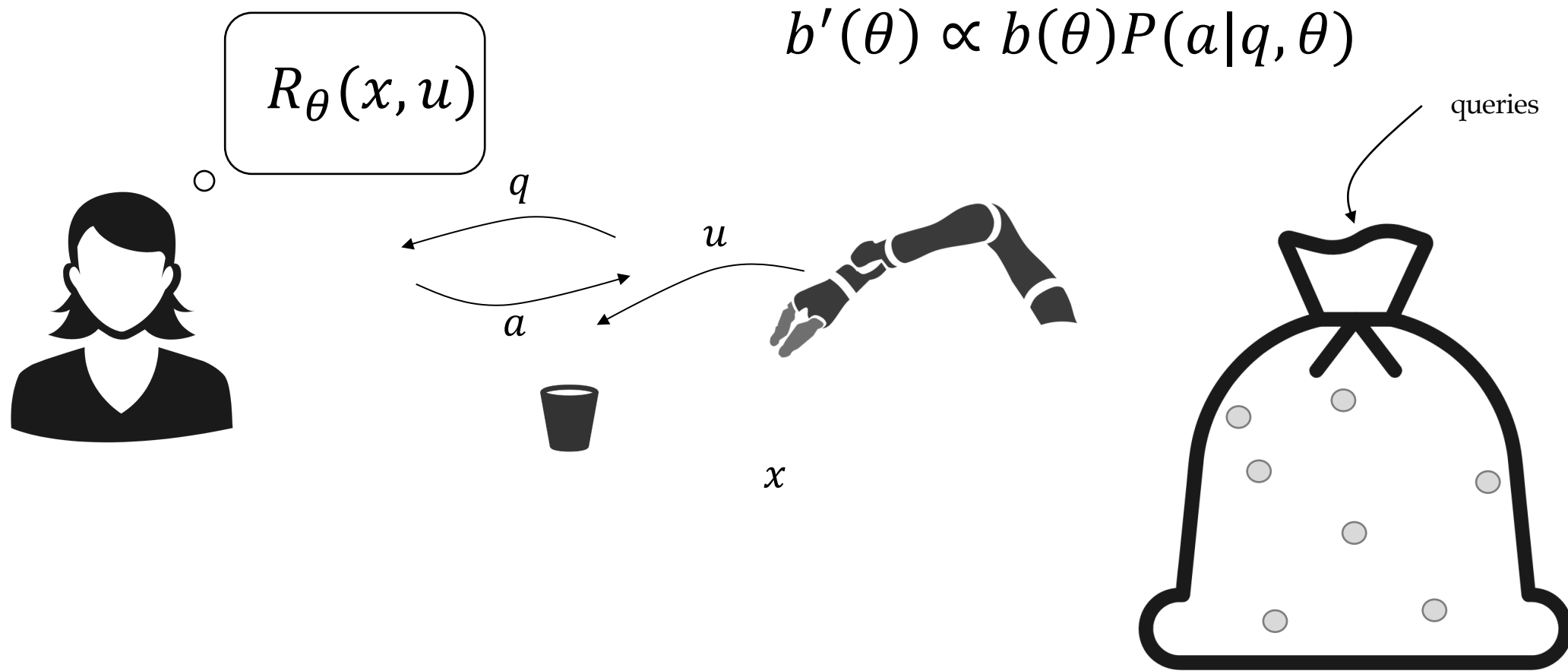
Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

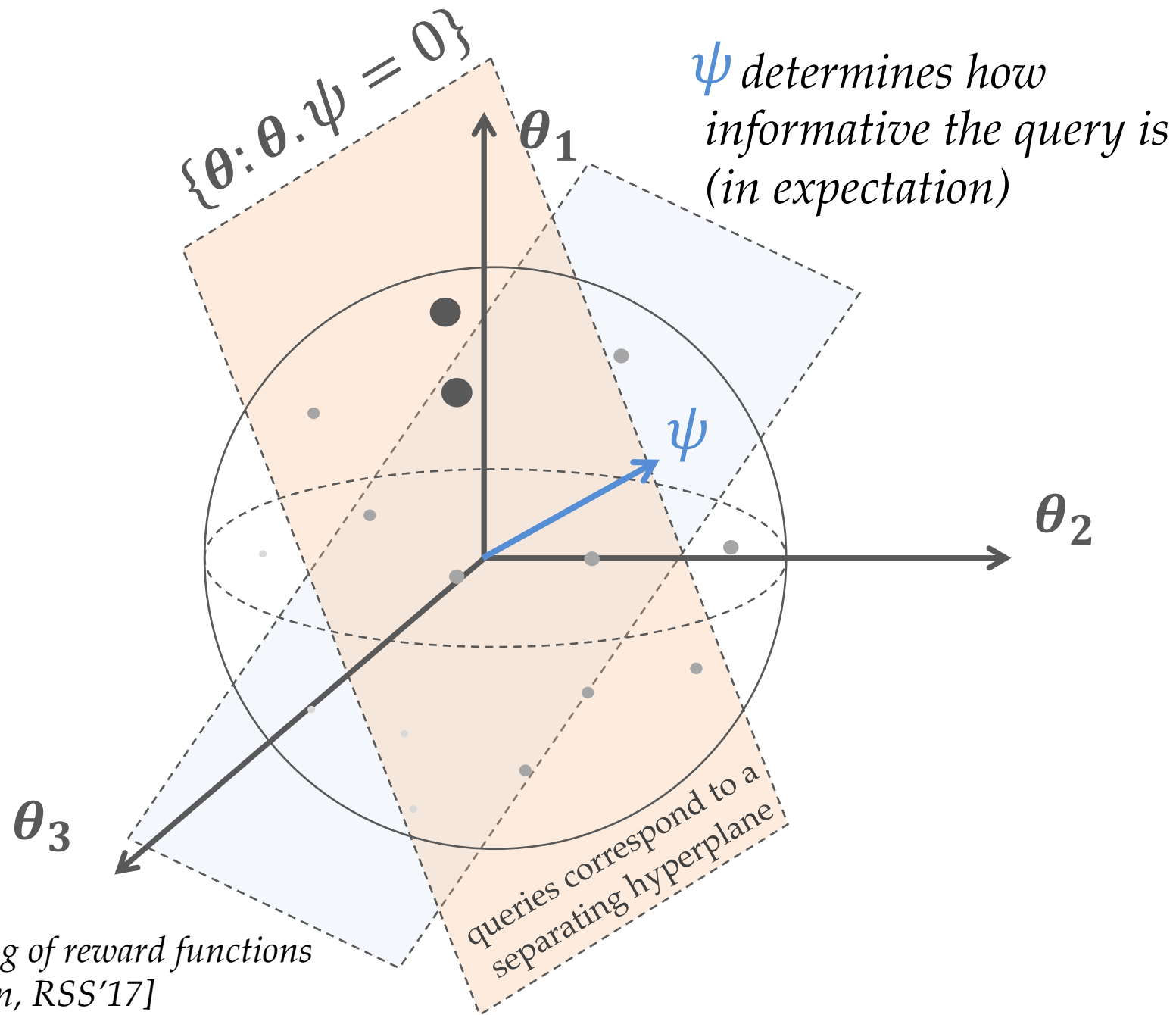


Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

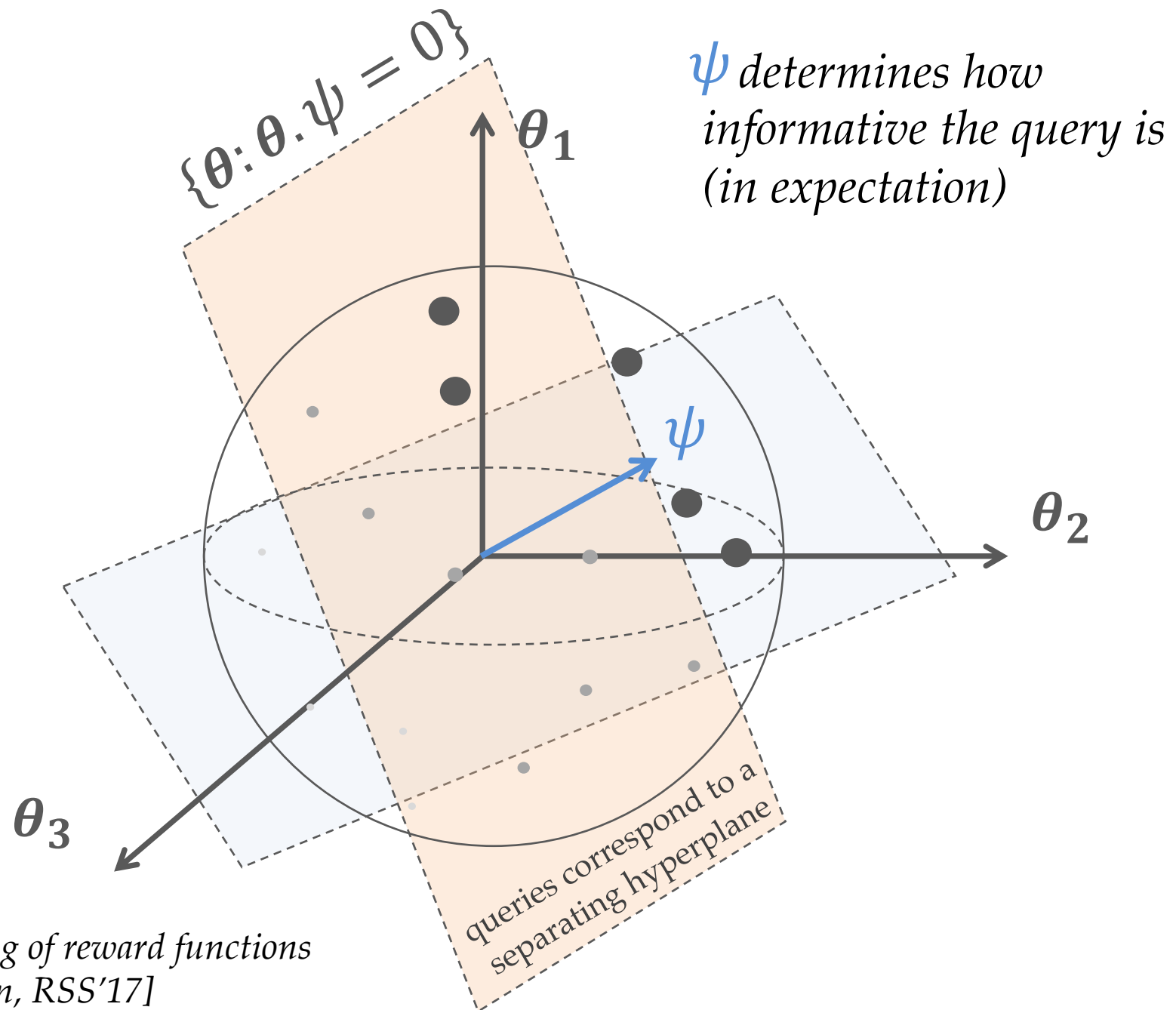


Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

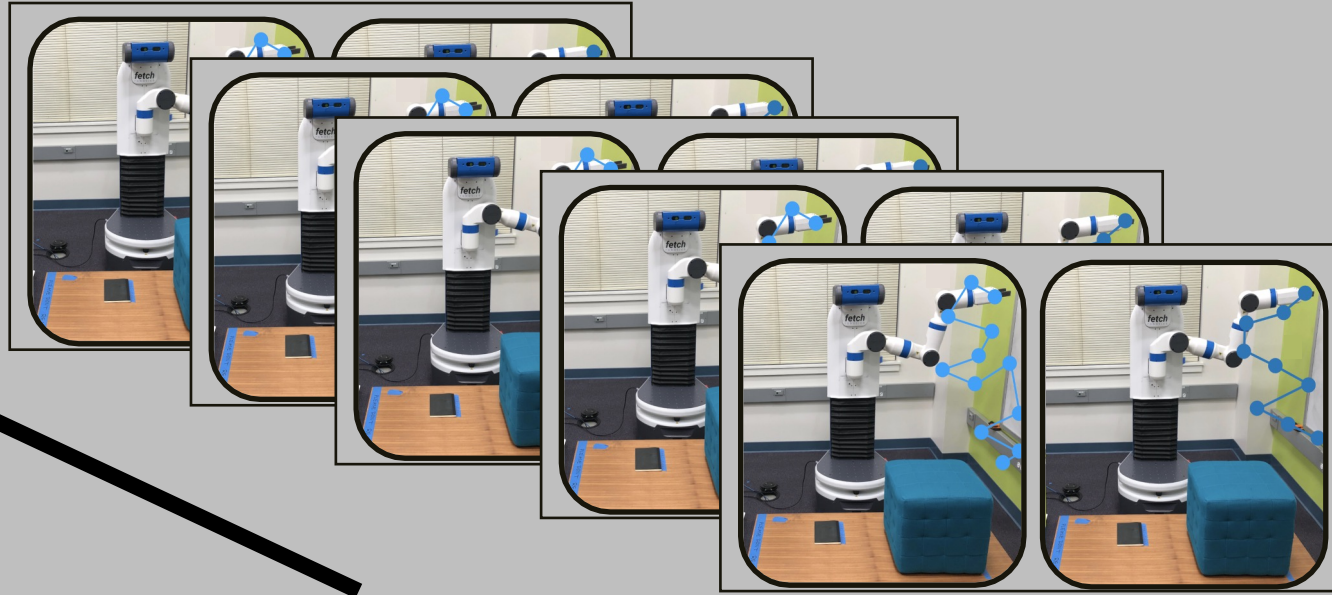




Active preference based learning of reward functions
 [Sadigh, Seshia, Sastry, Dragan, RSS'17]



Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

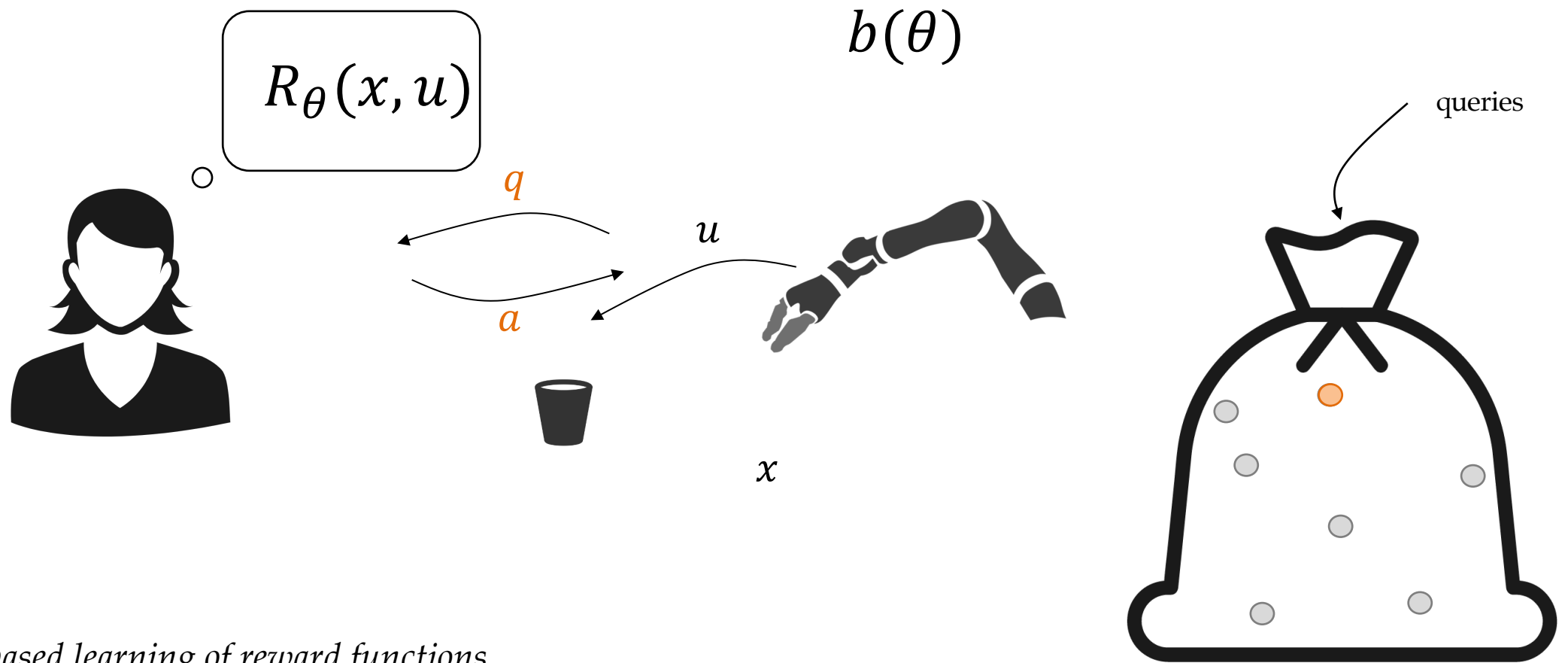


*Most informative,
diverse sequence of queries*

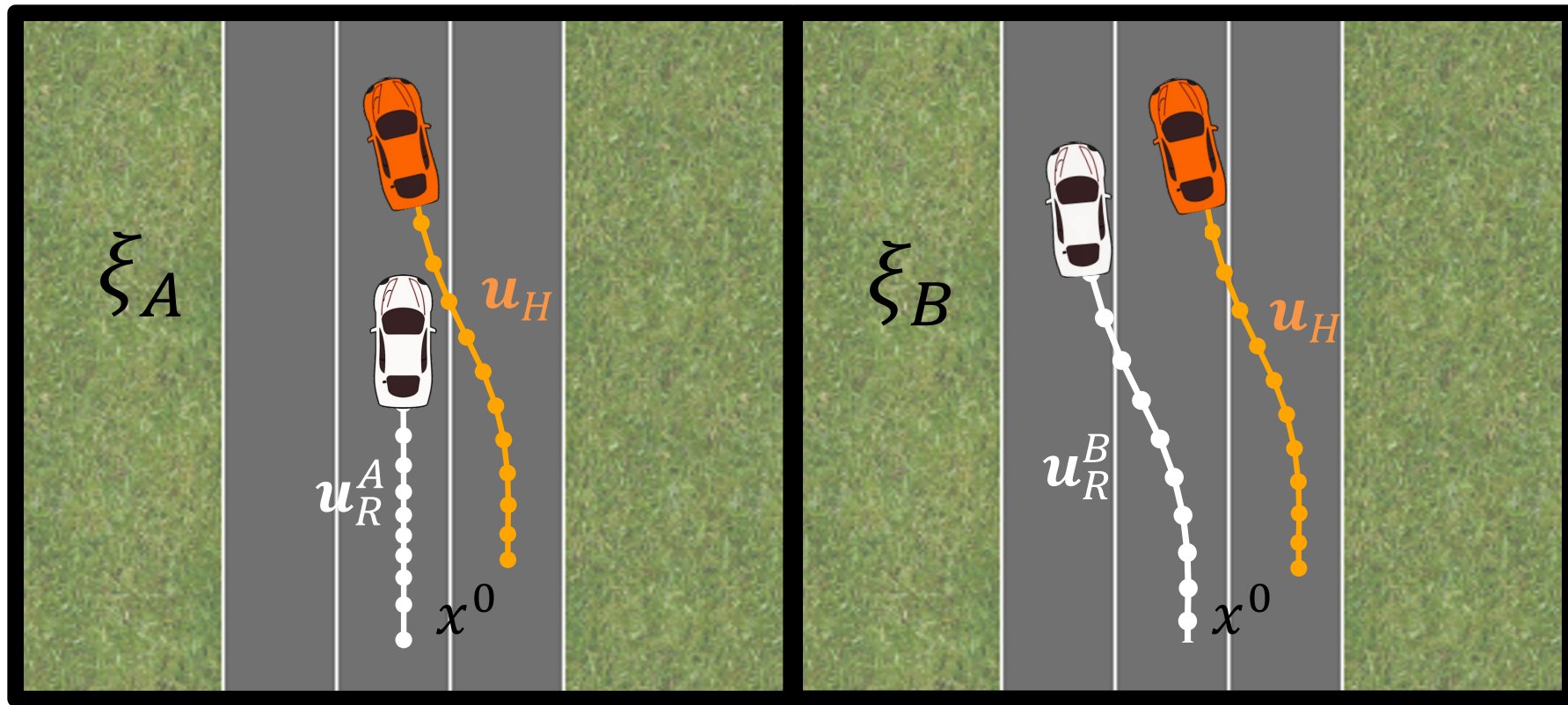


ξ_A or ξ_B ?

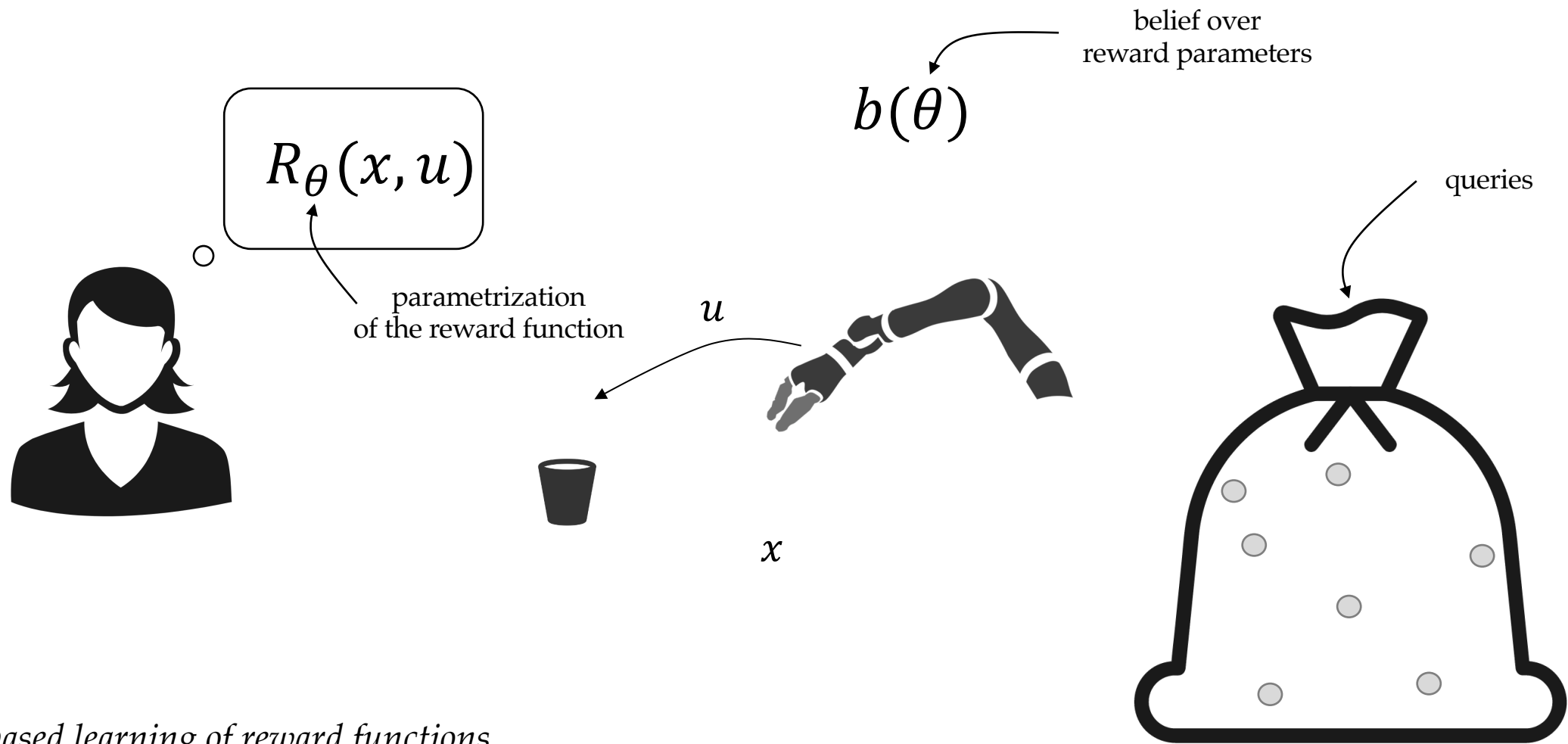
Queries should be actively selected.



Challenge:
Queries lie in a continuous and high-dimensional space.

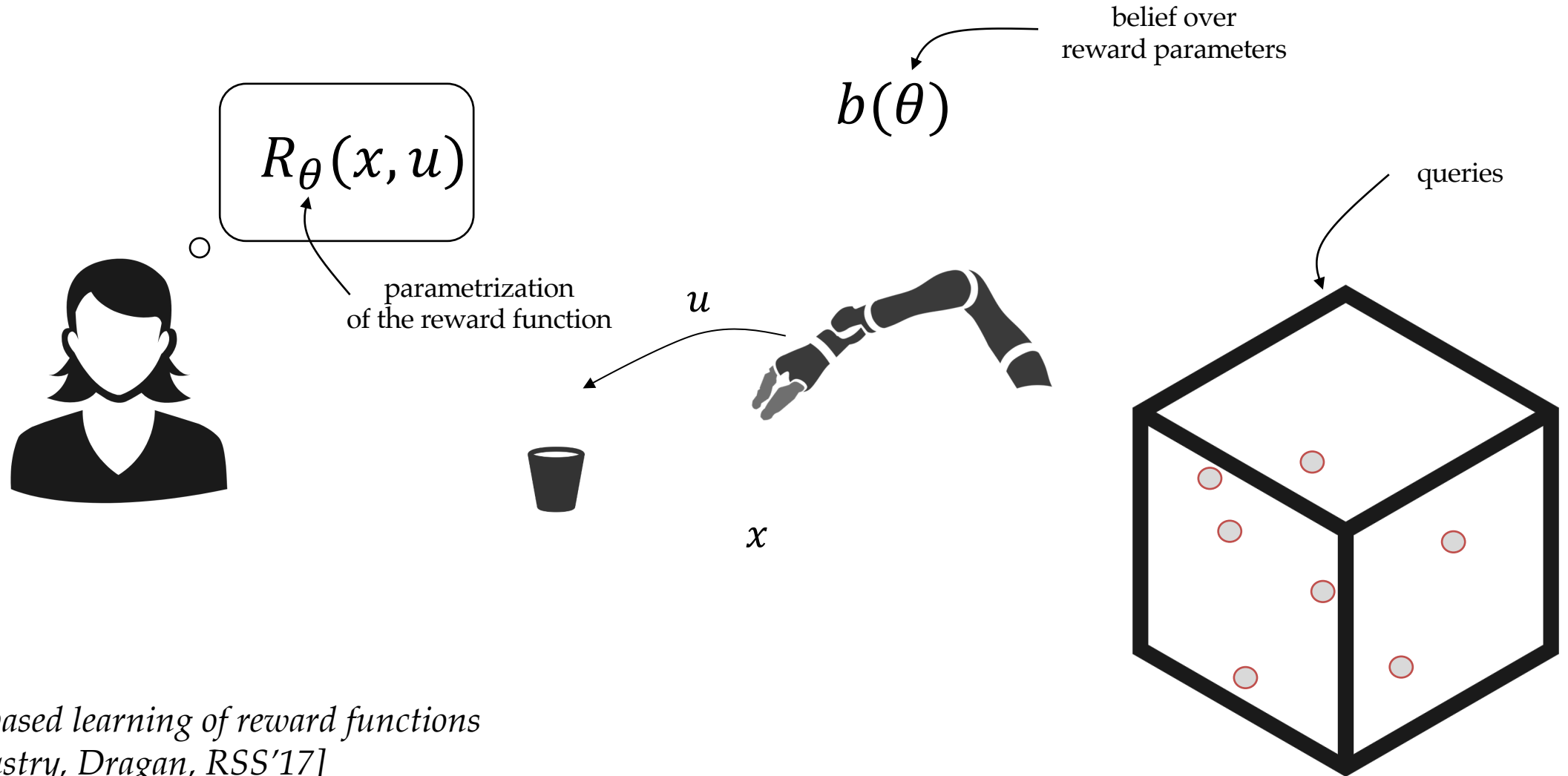


Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

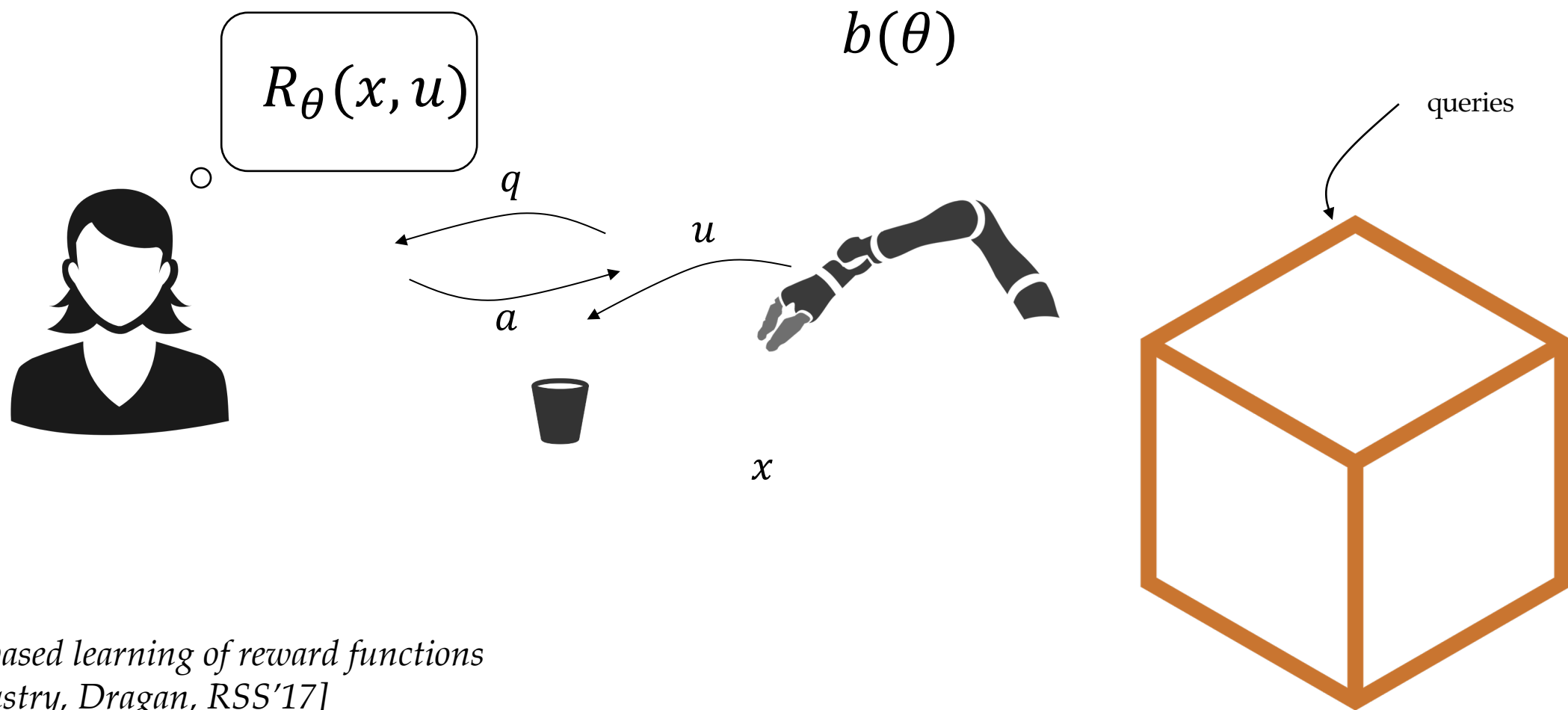


Active preference based learning of reward functions
[Sadigh, Seshia, Sastry, Dragan, RSS'17]

Real robots don't get handed a neat query set



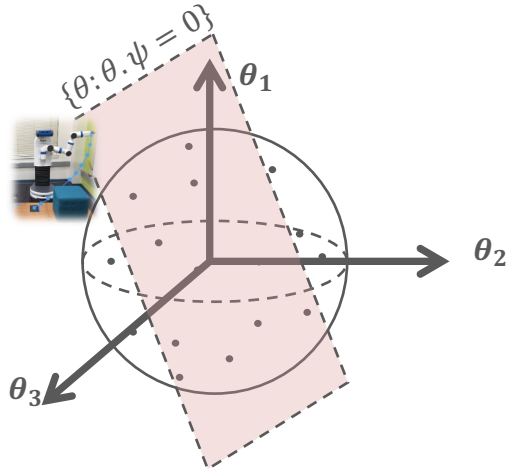
They have to synthesize their queries from scratch.



Queries should be actively synthesized.

Actively synthesizing queries

minimum volume removed

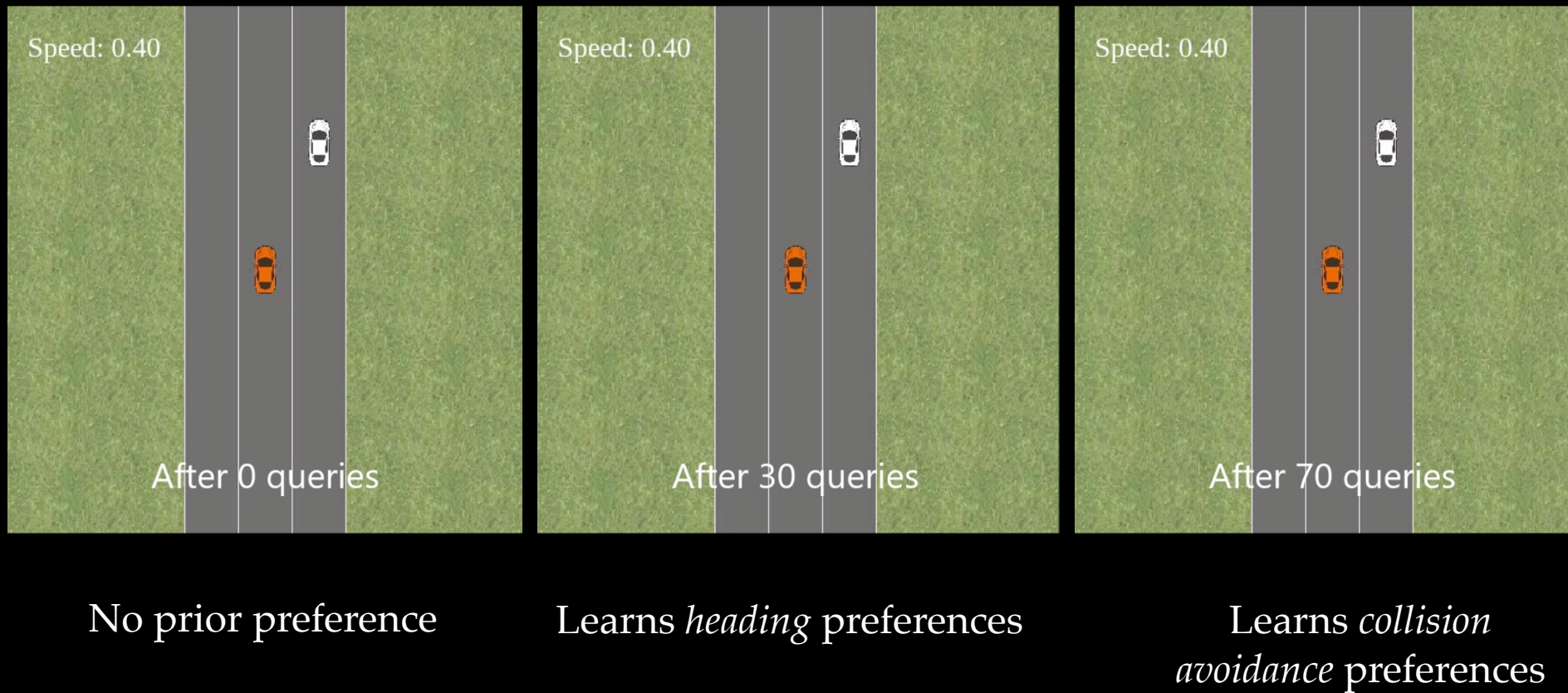


$$\max_{\psi} \min\{\mathbb{E}[1 - f_{\psi}(\theta)], \mathbb{E}[1 - f_{-\psi}(\theta)]\}$$

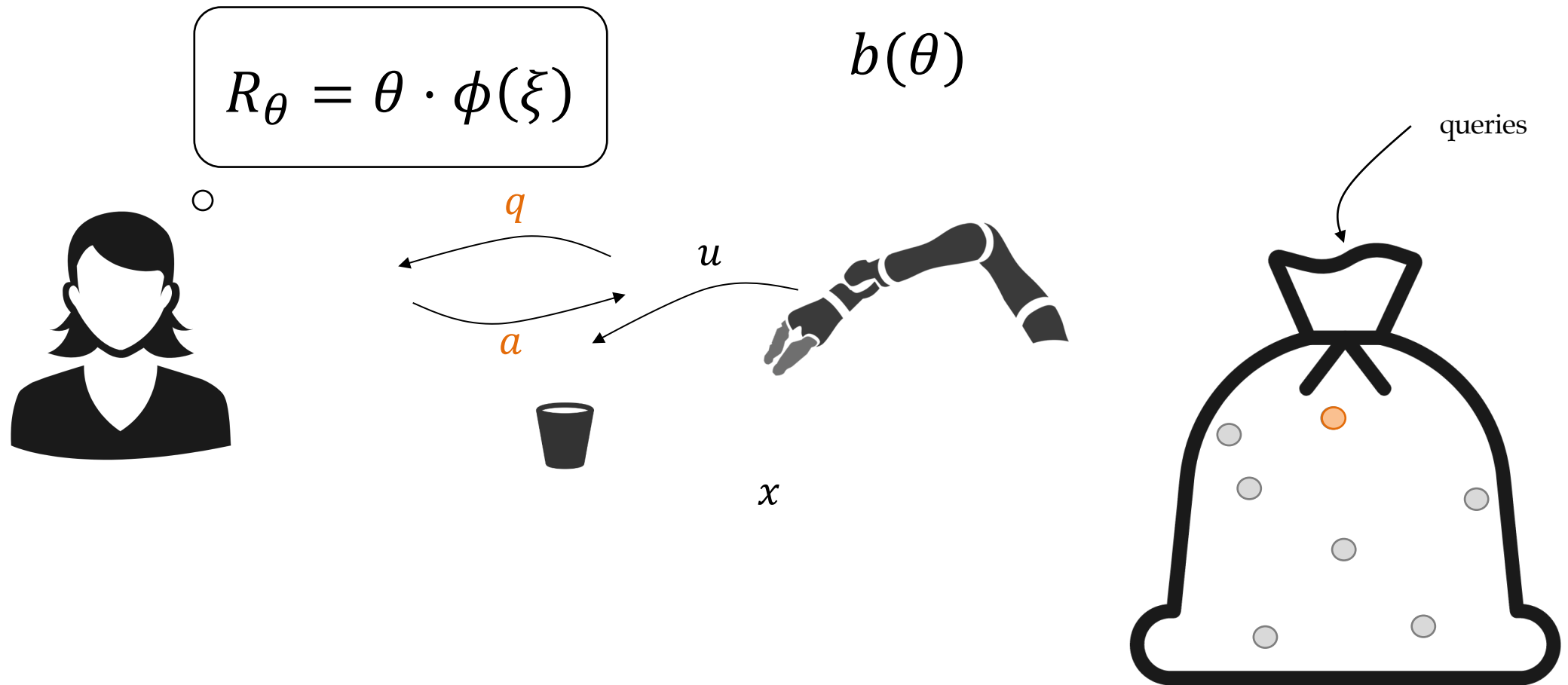
Subject to $\psi \in \mathbb{F}$

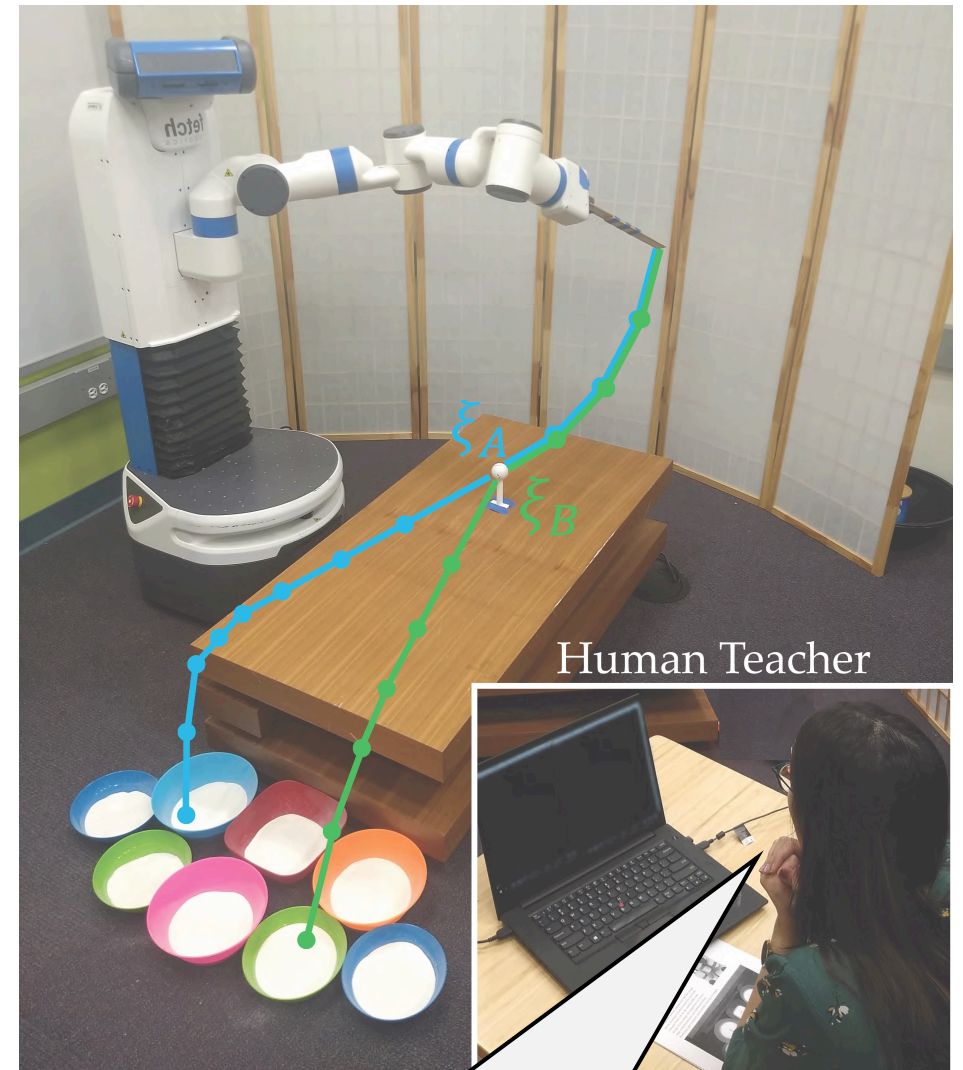
$$\mathbb{F} = \{\psi: \psi = \Phi(\xi_A) - \Phi(\xi_B), \xi_A, \xi_B \in \Xi\}$$

Human update function $f_{\psi}(\boldsymbol{\theta}) = \min(1, \exp(I_t \boldsymbol{\theta}^T \psi))$



Queries should be actively synthesized.





$$R(\xi_A) > R(\xi_B)$$

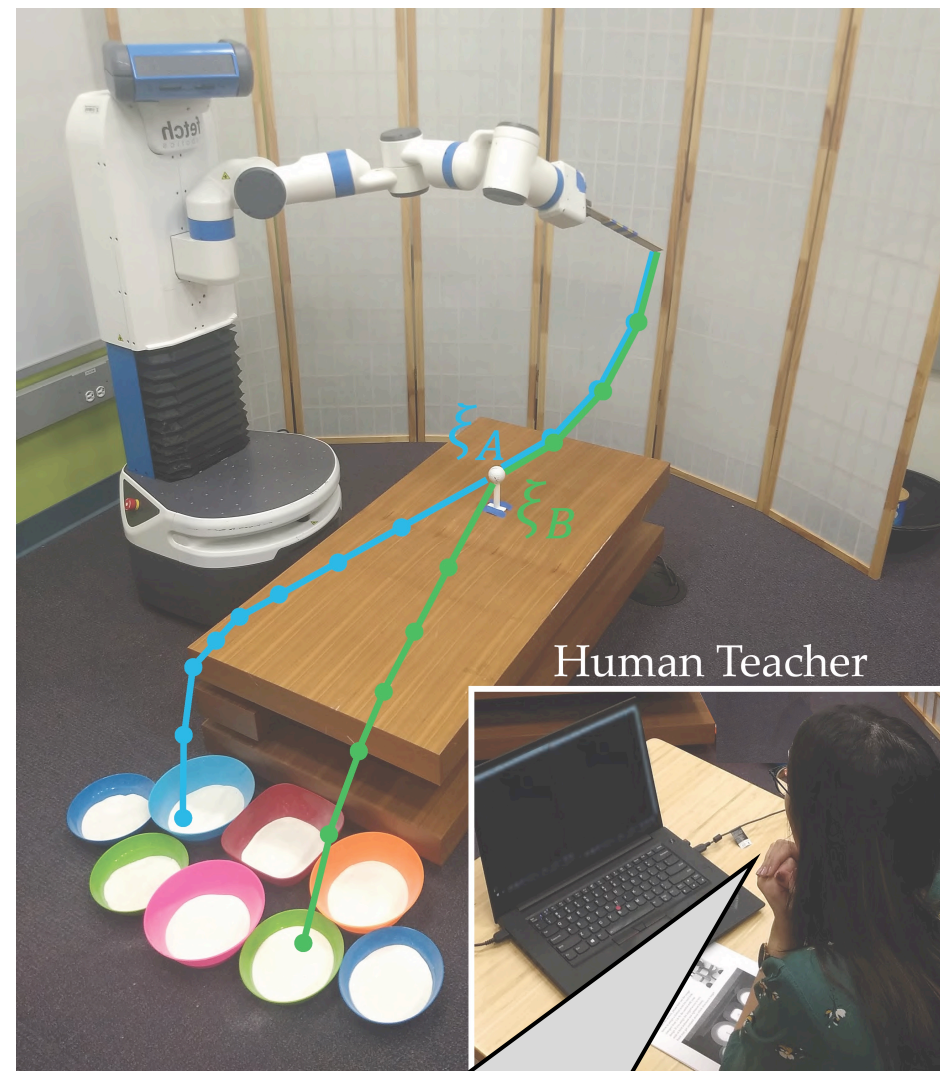
$$R(\xi_A) = \theta \cdot \phi(\xi_A)$$



Designing features is hard.

$e^{-c_4 d_4}$ where d_4 is the final horizontal distance between the object and the center of the closest basket, and $c_4 = 3$.

The average of $e^{-c_2 d_2^2 - c_3 d_3^2}$ over the trajectory, where d_2 and d_3 are the horizontal and vertical distances between the ego car and the other car, respectively; and $c_2 = 7, c_3 = 3$

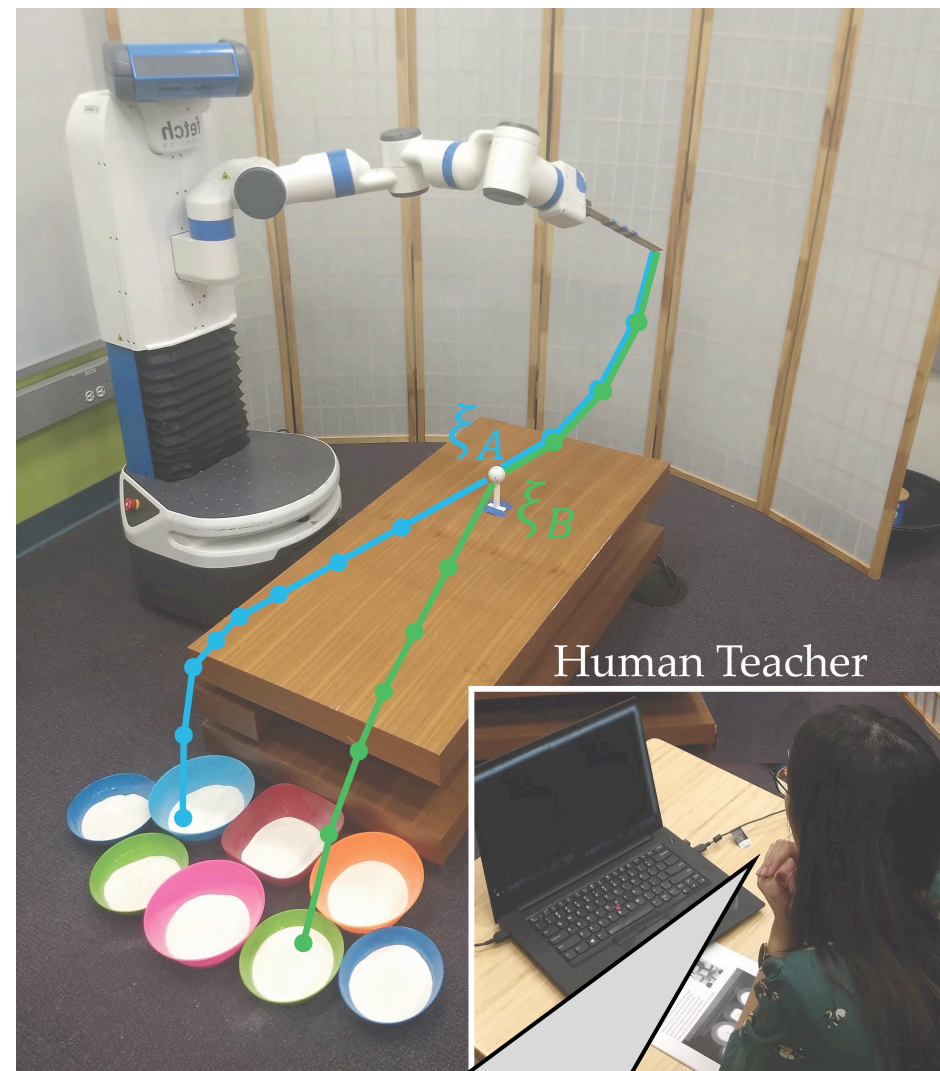


$$R(\xi_A) > R(\xi_B)$$

[Sadigh'17]
[Basu'18]
[Biyik'18,'19]
[Katz'19]
[Palan'19]
[Wilde'20]

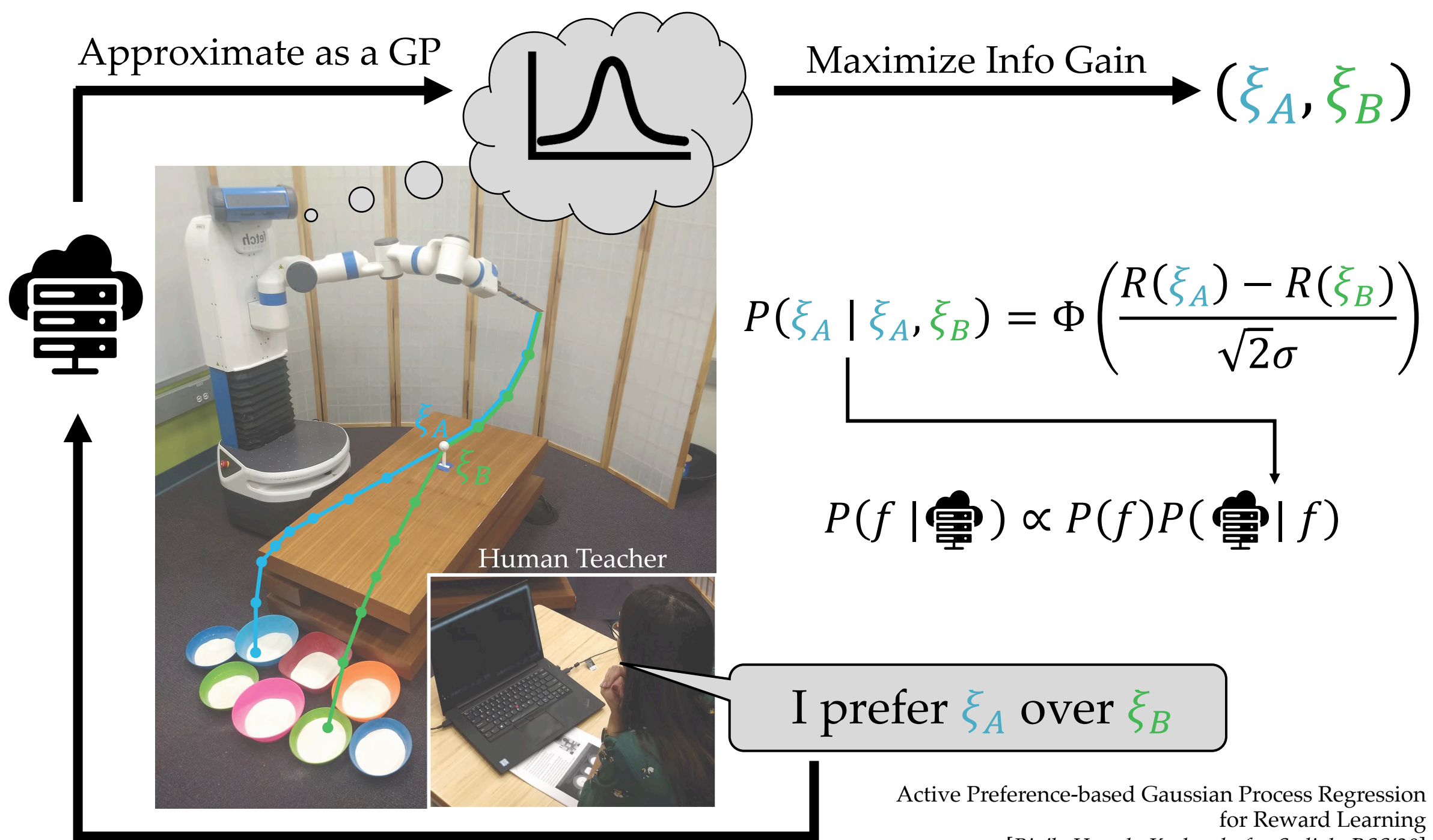
Trajectory Features: Shot Speed, Shot Angle

$$R(\xi_A) = \theta(\phi(\xi_A))$$



$$R(\xi_A) > R(\xi_B)$$

1. What if our reward function is **nonlinear**?



Learned Policies



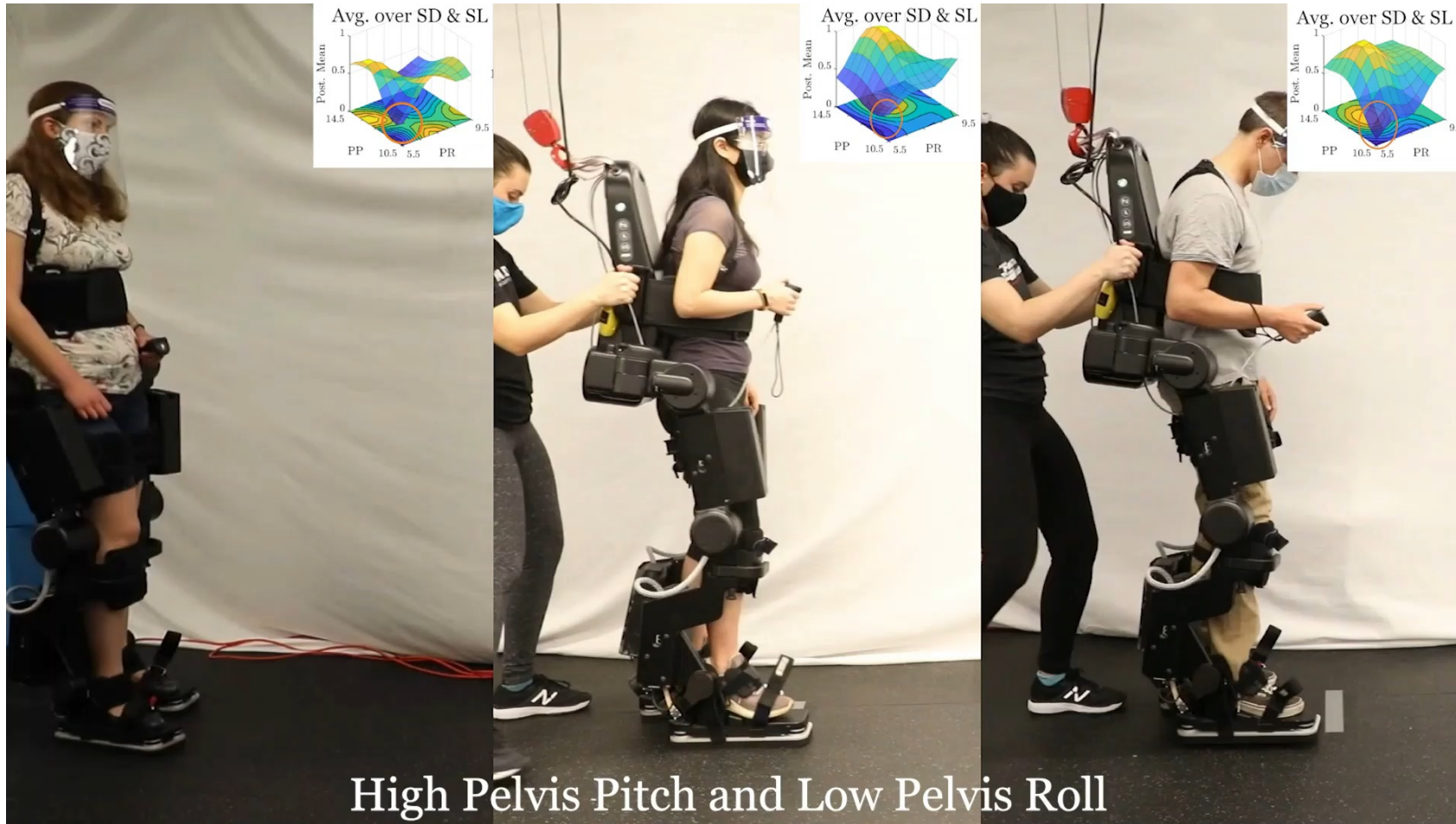
Linear Reward



GP Reward

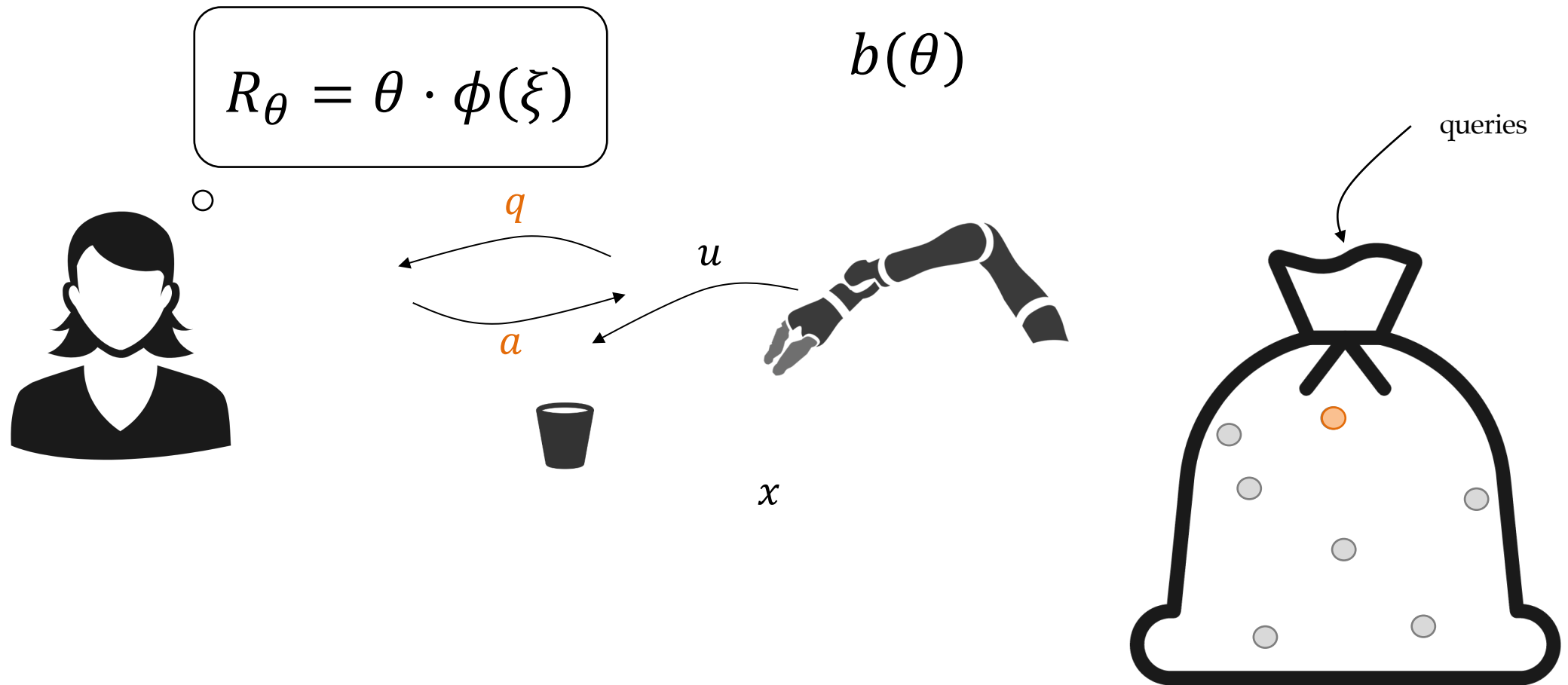
Active Preference-based Gaussian Process Regression for Reward Learning
[Biyik, Huynh, Kochenderfer, Sadigh. RSS'20]

Nonlinear Rewards for Exoskeletons



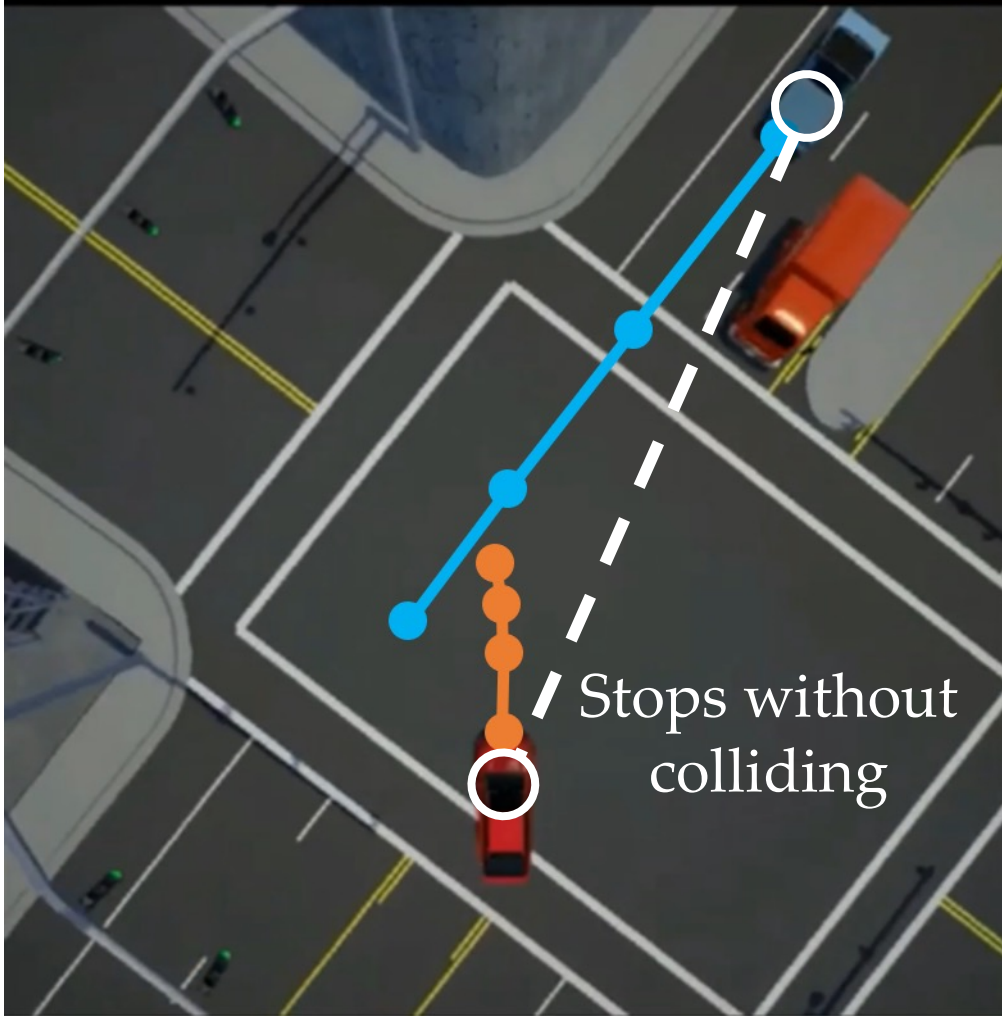
ROIAL: Region of Interest Active Learning for Characterizing Exoskeleton Gait Preference Landscapes
[Li, Tucker, Biyik, Novoseller, Burdick, Sui, Sadigh, Yue, Ames, ICRA'21]

Queries should be actively synthesized.



Learning from

The red car will make an unprotected left turn...

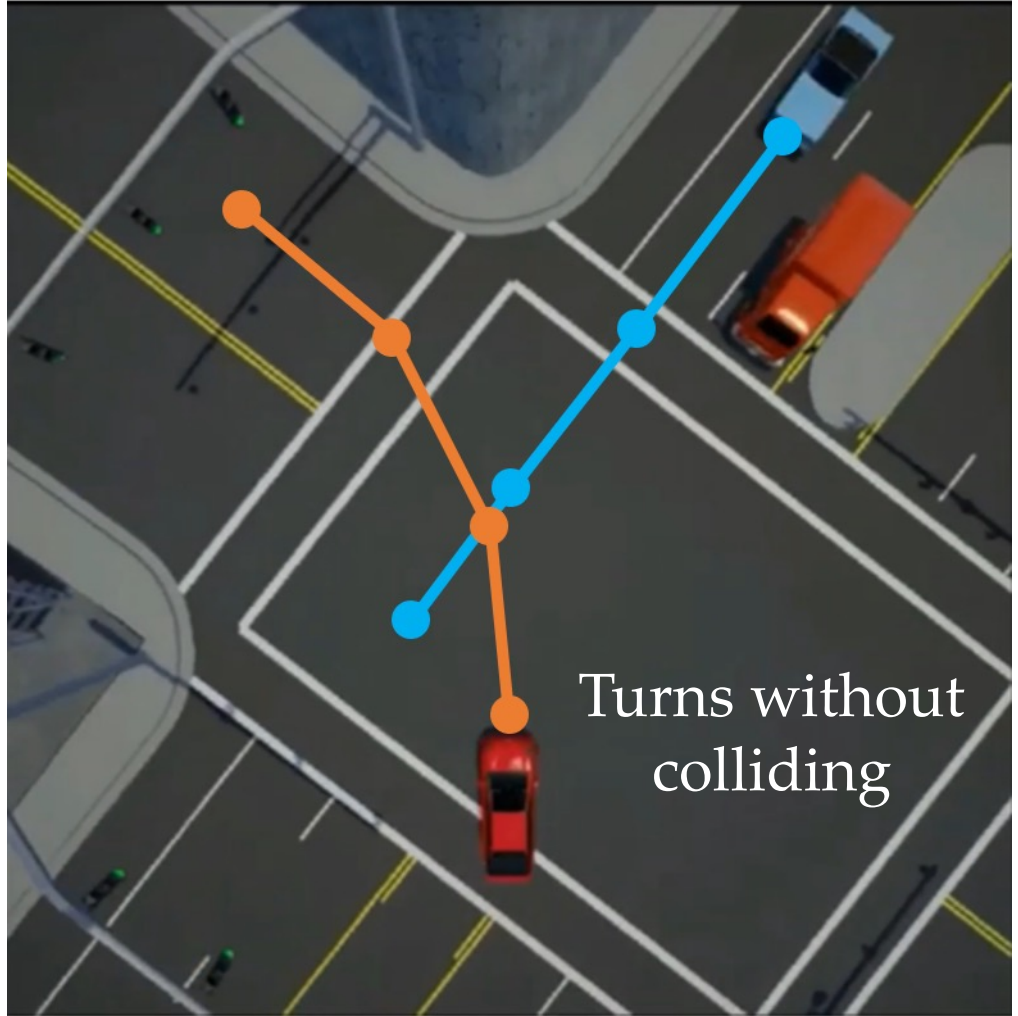


...but realizes the blue car is coming.



A timid driver

The red car will make an unprotected left turn...



...but realizes the blue car is coming.

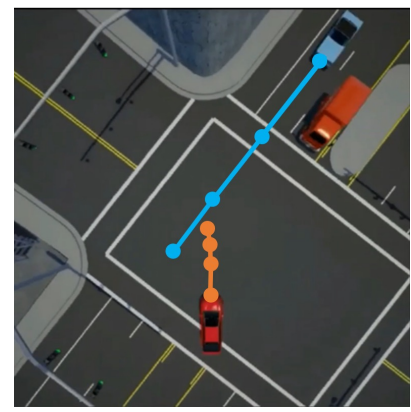
Learning from



A timid driver

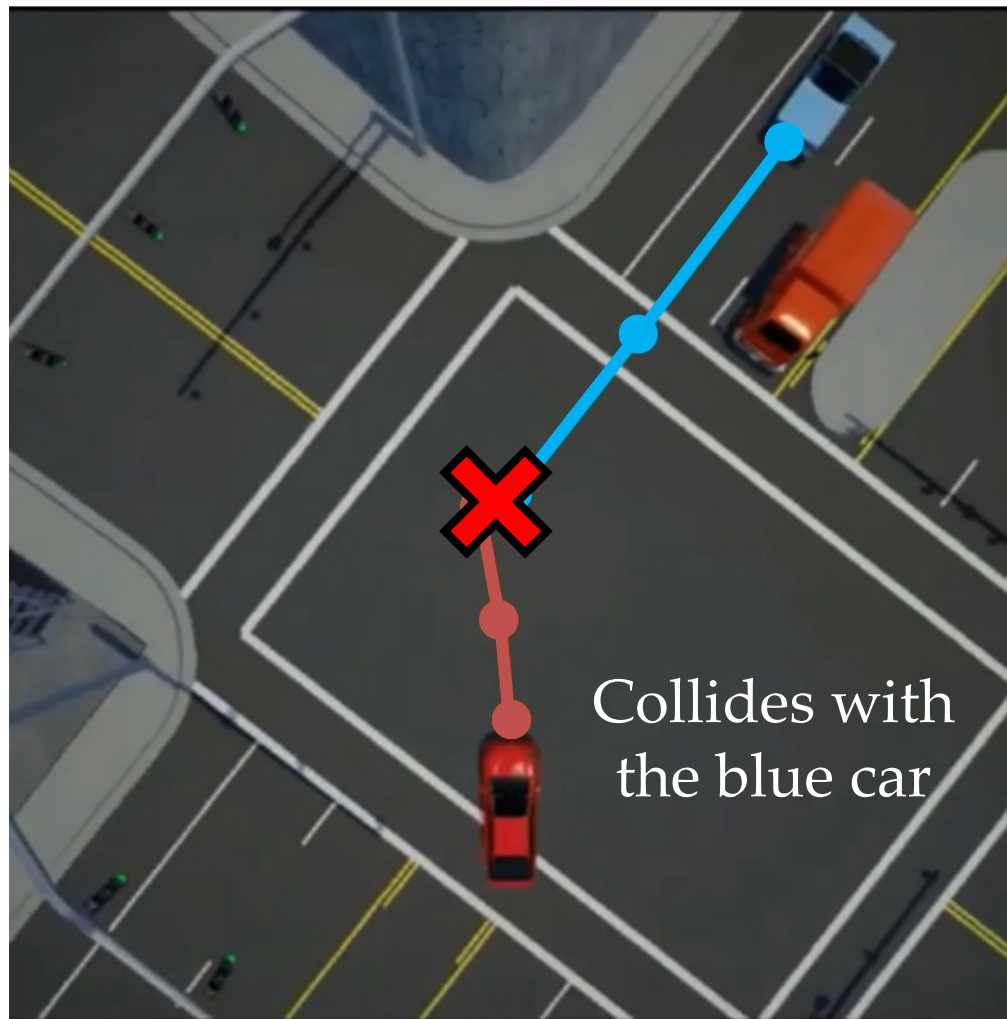


An aggressive driver



Learning from

The red car will make an unprotected left turn...



...but realizes the blue car is coming.



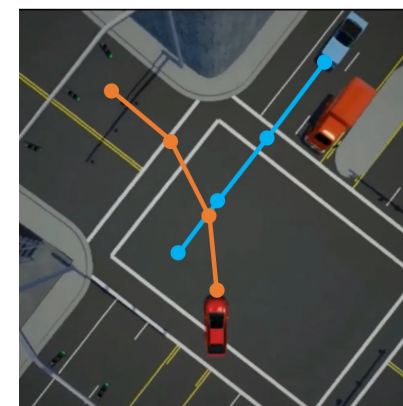
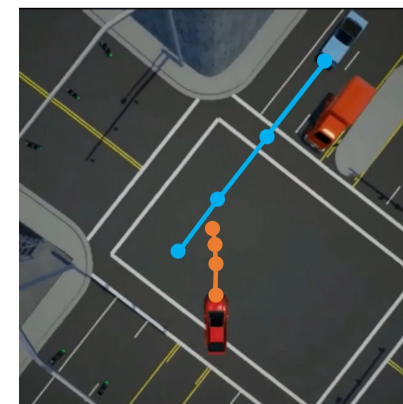
A timid driver



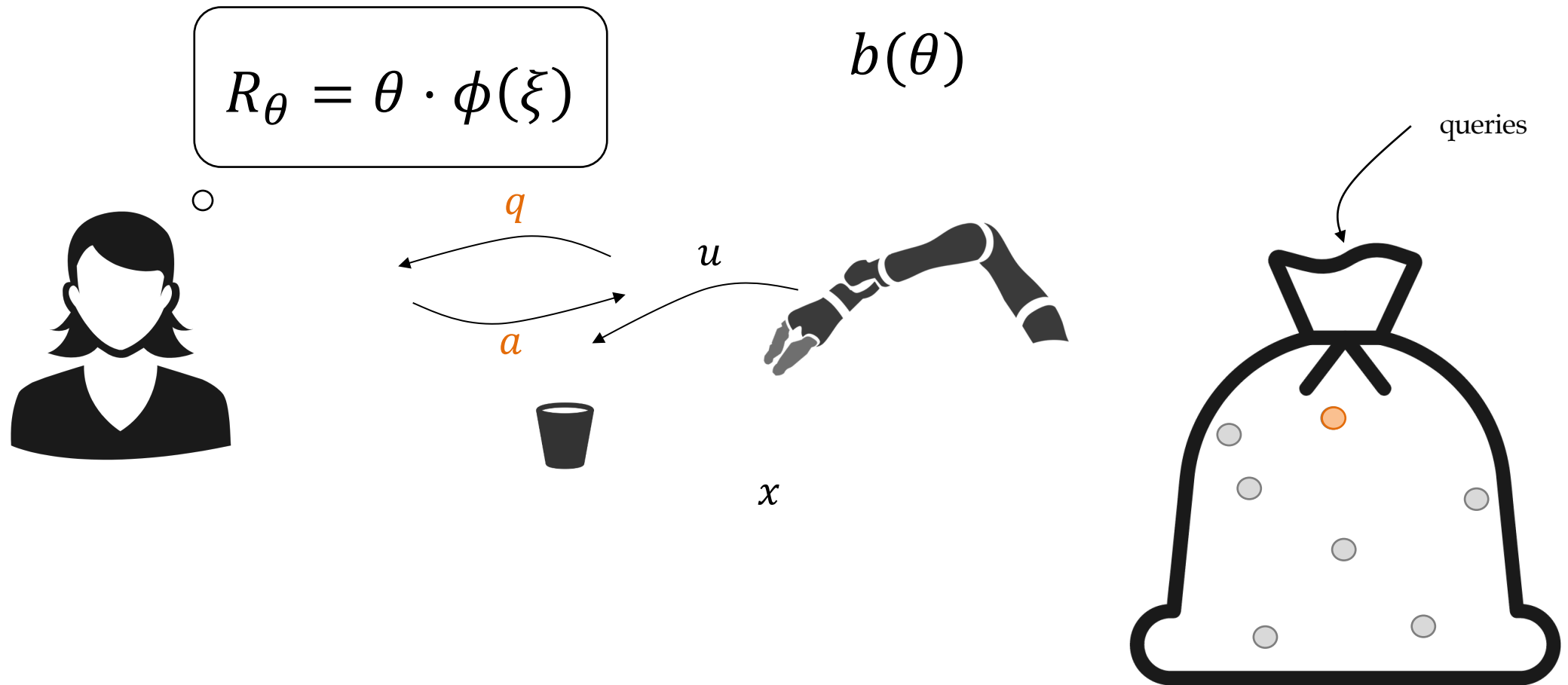
An aggressive driver



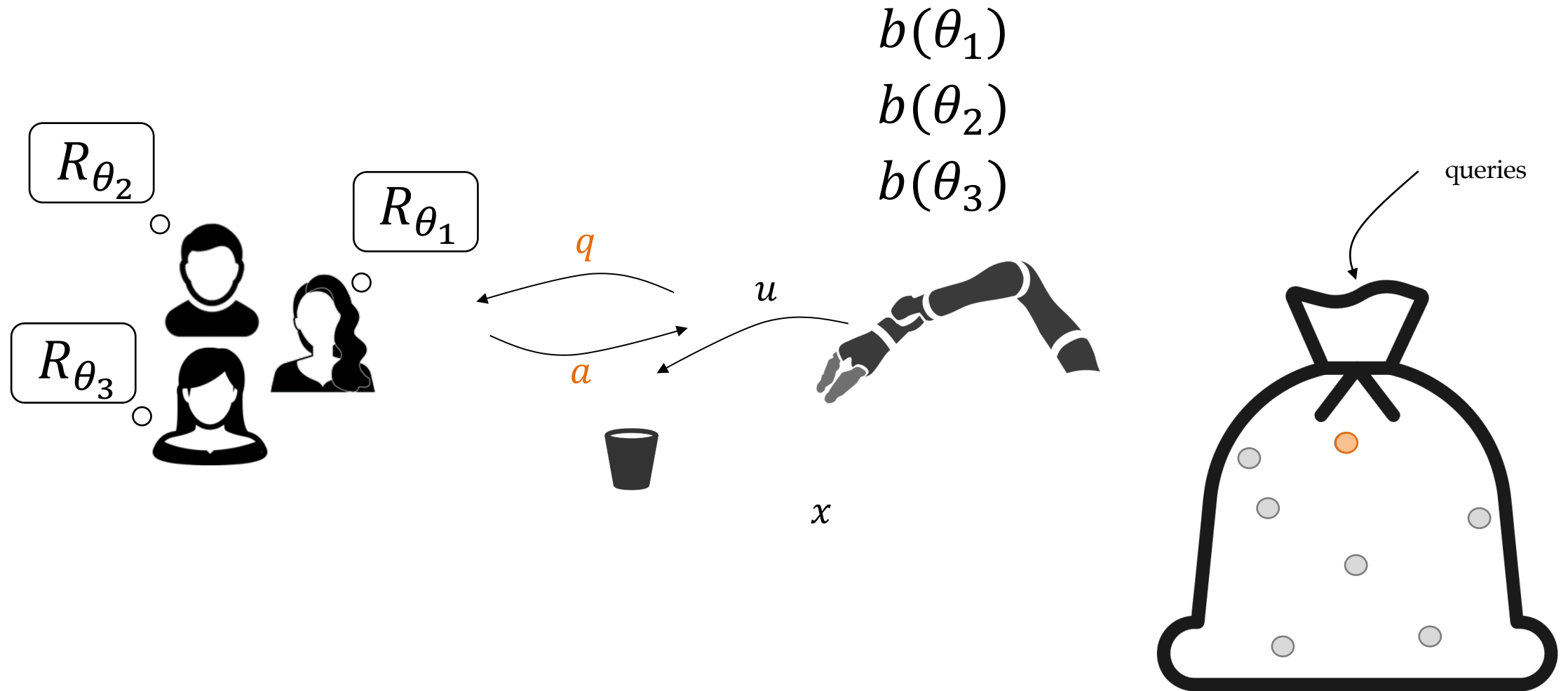
Both types of drivers



Queries should be actively synthesized.



Queries should be actively synthesized.



1. What if our reward function is **nonlinear**?
2. What if our reward function is **multimodal**?

Learning from Comparisons



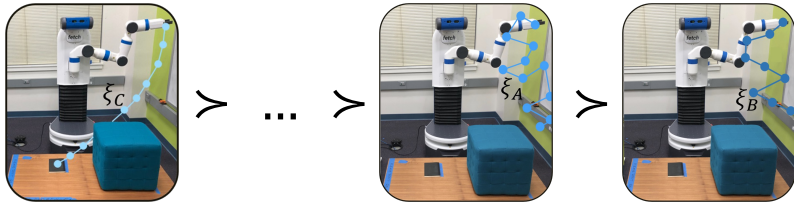
Unimodal

Gervasio et al., 1999
Akroun et al., 2012
Sadigh et al., 2017
Christiano et al., 2017
Bıyık et al., 2020
Tucker et al., 2020
Wilde et al., 2021

Multimodal

Impossible
(Zhao et al., 2016)

Learning from Rankings



Unimodal

Gervasio et al., 1999
Akroun et al., 2012
Sadigh et al., 2017
Christiano et al., 2017
Bıyık et al., 2020
Tucker et al., 2020
Wilde et al., 2021

Multimodal

This work



Please rank these trajectories



Users have their own reward functions.

With probability $\frac{\alpha_1}{\alpha_1 + \alpha_2}$



$$R_1(\xi) = \theta_1^\top \phi(\xi)$$

With probability $\frac{\alpha_2}{\alpha_1 + \alpha_2}$

$$R_2(\xi) = \theta_2^\top \phi(\xi)$$





Please rank these trajectories



With probability $\frac{\alpha_1}{\alpha_1 + \alpha_2}$



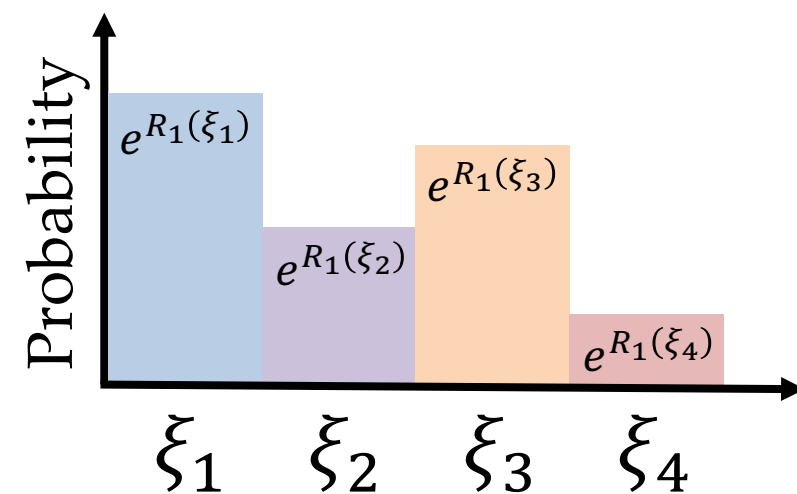
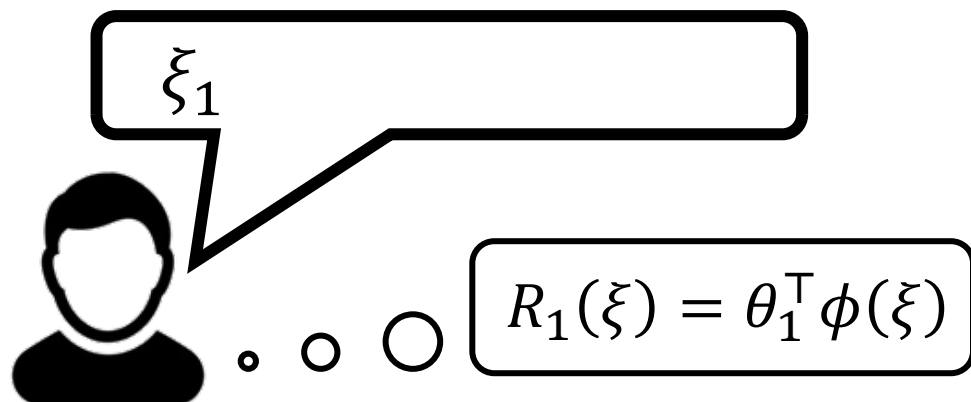
$$R_1(\xi) = \theta_1^\top \phi(\xi)$$



Please rank these trajectories



The user noisily chooses his best option: He chooses ξ_1





Please rank these trajectories

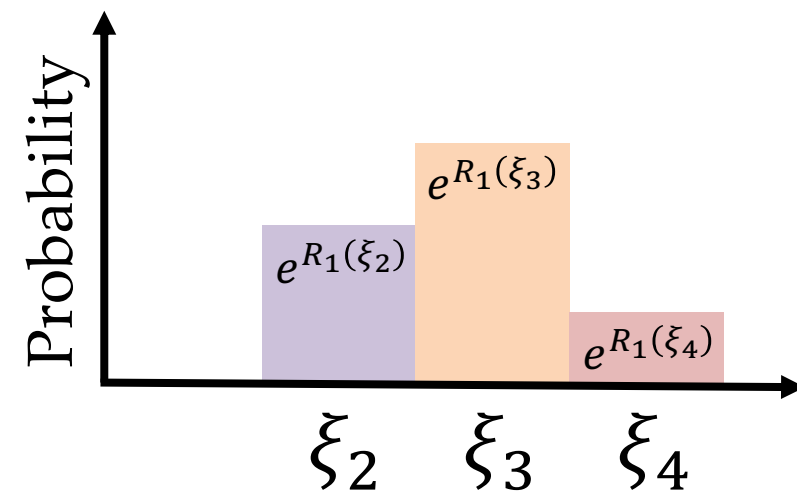


The user noisily chooses his second best option: He chooses ξ_3

$$\xi_1 \succ \xi_3$$



$$R_1(\xi) = \theta_1^\top \phi(\xi)$$





Please rank these trajectories

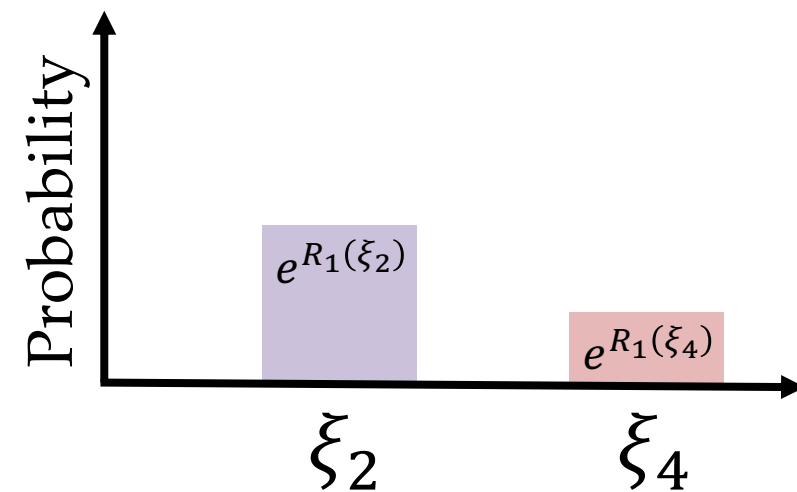


The user noisily chooses his third best option: He chooses ξ_4

$$\xi_1 > \xi_3 > \xi_4$$



$$R_1(\xi) = \theta_1^\top \phi(\xi)$$





Please rank these trajectories

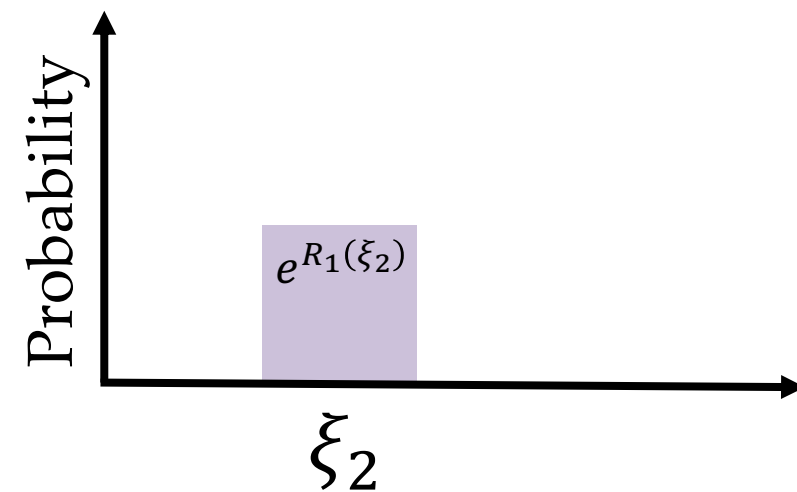


There is only one option left.

$$\xi_1 > \xi_3 > \xi_4 > \xi_2$$



$$R_1(\xi) = \theta_1^\top \phi(\xi)$$





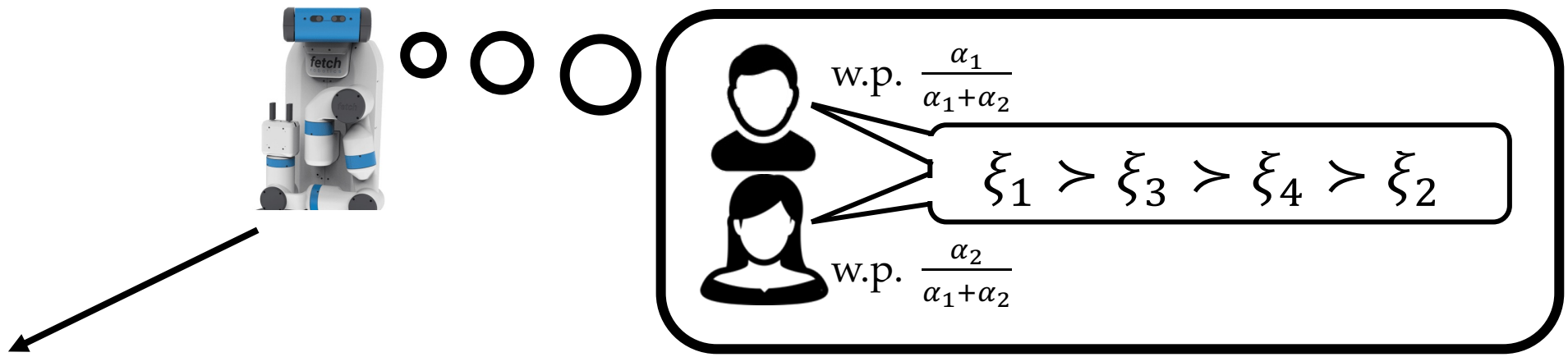
Please rank these trajectories



$\xi_1 \succ \xi_3 \succ \xi_4 \succ \xi_2$



$$R_1(\xi) = \theta_1^\top \phi(\xi)$$



Starts with a prior $P(\alpha, \theta_1, \theta_2)$

Posterior: $P(\alpha, \theta_1, \theta_2 \mid \xi_1 > \xi_3 > \xi_4 > \xi_2) \propto P(\alpha, \theta_1, \theta_2) \underbrace{P(\xi_1 > \xi_3 > \xi_4 > \xi_2 \mid \alpha, \theta_1, \theta_2)}_{\text{Probability of observing the ranking}}$

$$P(\xi_1 > \xi_3 > \xi_4 > \xi_2 \mid \alpha, \theta_1, \theta_2) = \\ P(\text{man} \mid \alpha)P(\xi_1 > \xi_3 > \xi_4 > \xi_2 \mid \text{man}, \theta_1) + P(\text{woman} \mid \alpha)P(\xi_1 > \xi_3 > \xi_4 > \xi_2 \mid \text{woman}, \theta_2)$$



Please rank these trajectories



How does the robot choose which trajectories to show to the users?

We **actively query** the users by maximizing **information gain**.

Maximize Information Gain

Unimodal: $\max_{\text{query}} I(\theta ; \text{response} \mid \text{query})$

Multimodal: $\max_{\text{query}} I(\alpha, (\theta_i)_{i=1}^M ; \text{response} \mid \text{query})$

$$\min_{\text{query}} \mathbb{E}_{\text{response}, \alpha, (\theta_i)_{i=1}^M \mid \text{query}} \log \frac{\mathbb{E}_{\alpha', (\theta'_i)_{i=1}^M} P(\text{response} \mid \text{query}, \alpha', (\theta'_i)_{i=1}^M)}{P(\text{response} \mid \text{query}, \alpha, (\theta_i)_{i=1}^M)}$$

Maximize Information Gain

Unimodal: $\max_{\text{query}} I(\theta ; \text{response} \mid \text{query})$

Multimodal: $\max_{\text{query}} I(\alpha, (\theta_i)_{i=1}^M ; \text{response} \mid \text{query})$

$$\min_{\text{query}} \mathbb{E}_{\text{response}, \alpha, (\theta_i)_{i=1}^M \mid \text{query}} \log \frac{\mathbb{E}_{\alpha', (\theta'_i)_{i=1}^M} P(\text{response} \mid \text{query}, \alpha', (\theta'_i)_{i=1}^M)}{P(\text{response} \mid \text{query}, \alpha, (\theta_i)_{i=1}^M)}$$

Maximize Information Gain

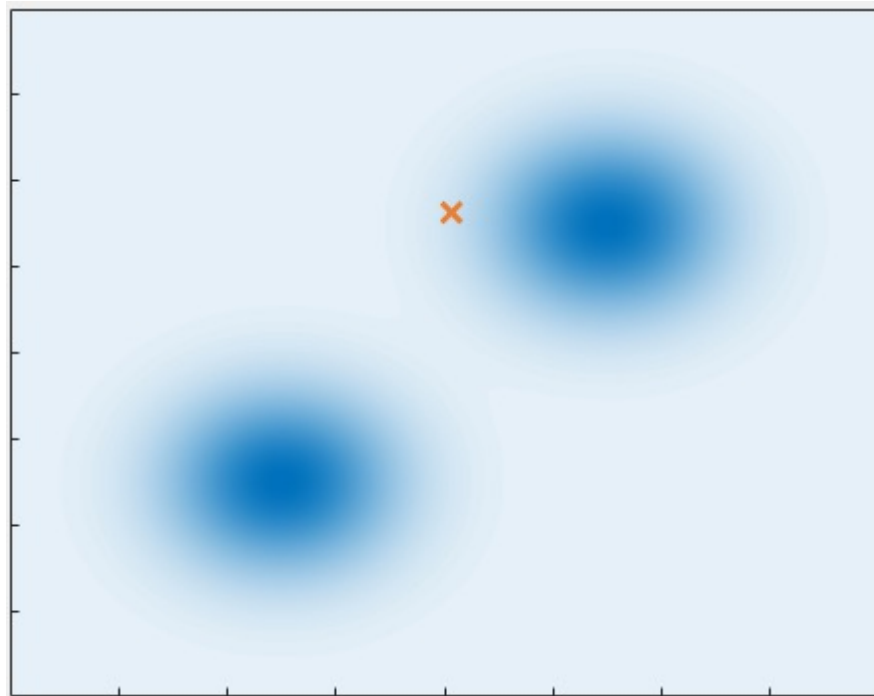
Unimodal: $\max_{\text{query}} I(\theta ; \text{response} \mid \text{query})$

Multimodal: $\max_{\text{query}} I(\alpha, (\theta_i)_{i=1}^M ; \text{response} \mid \text{query})$

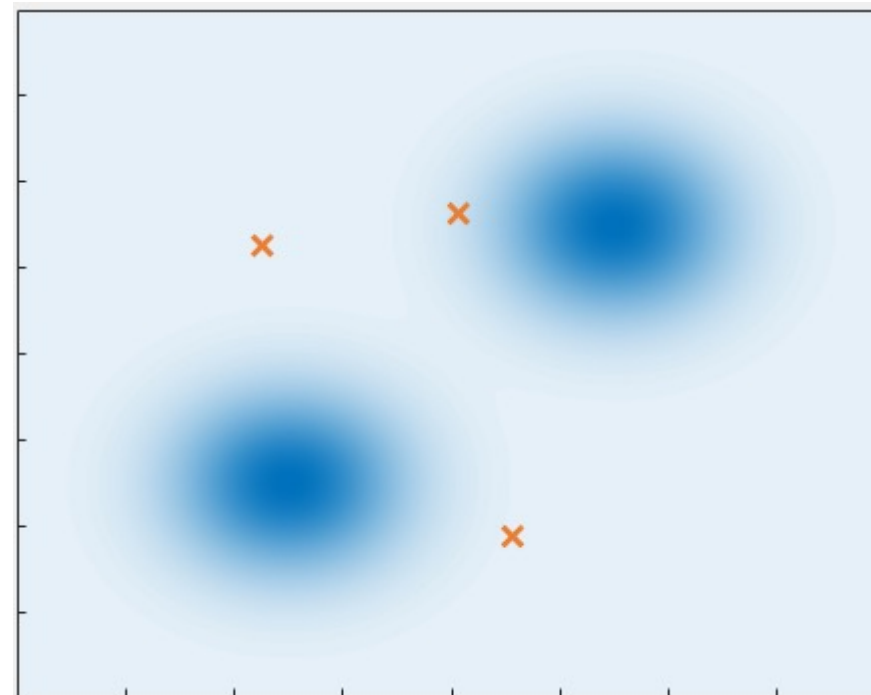
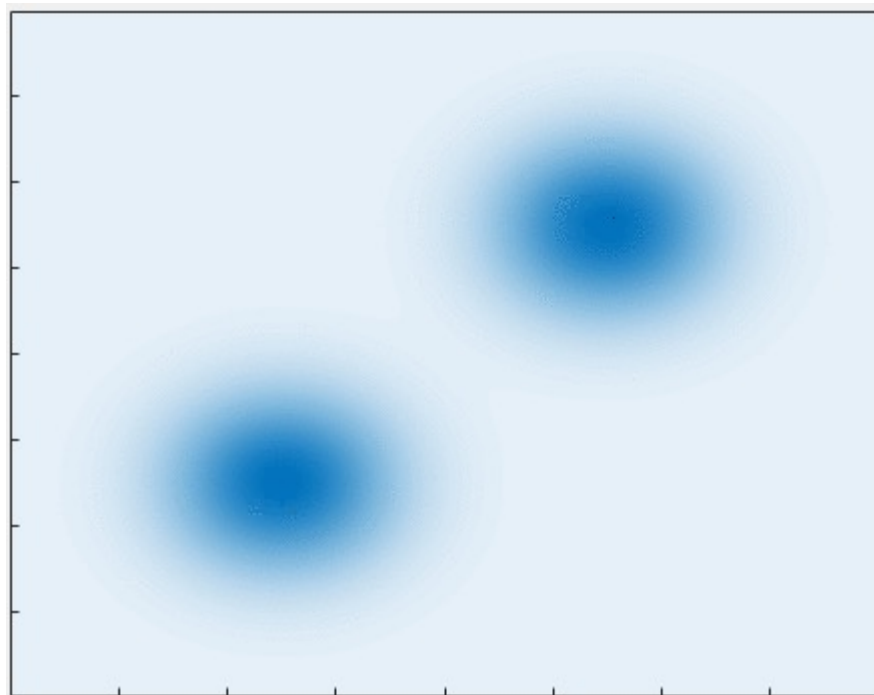
$$\min_{\text{query}} \mathbb{E}_{\text{response}, \alpha, (\theta_i)_{i=1}^M \mid \text{query}} \log \frac{\mathbb{E}_{\alpha', (\theta'_i)_{i=1}^M} P(\text{response} \mid \text{query}, \alpha', (\theta'_i)_{i=1}^M)}{P(\text{response} \mid \text{query}, \alpha, (\theta_i)_{i=1}^M)}$$

Sample from: $P(\alpha, (\theta_i)_{i=1}^M)$ Multimodal!

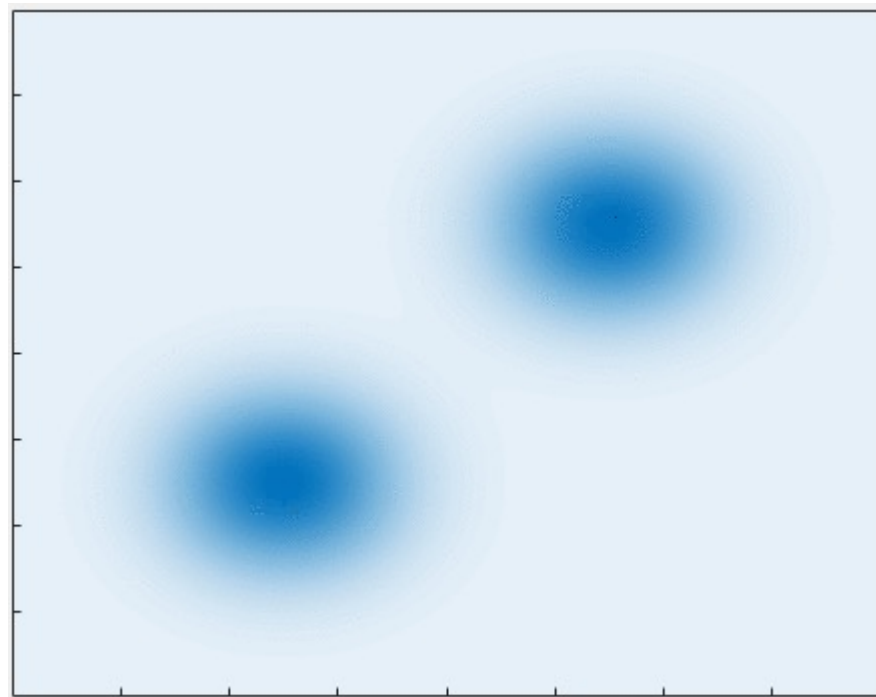
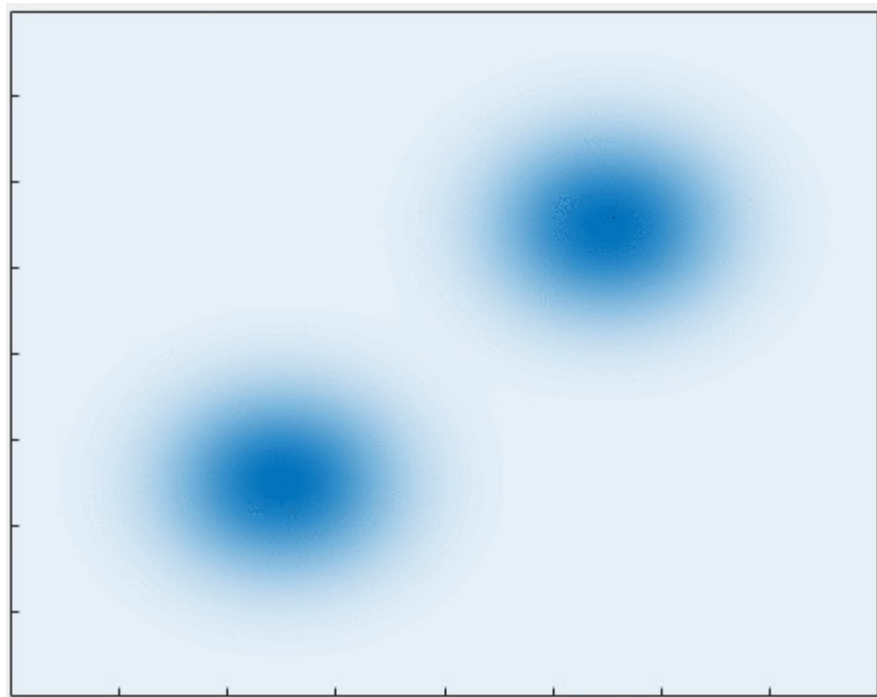
Sampling from a Multimodal Distribution



Sampling from a Multimodal Distribution



Sampling from a Multimodal Distribution



Maximize Information Gain

Unimodal: $\max_{\text{query}} I(\theta ; \text{response} \mid \text{query})$

Multimodal: $\max_{\text{query}} I(\alpha, (\theta_i)_{i=1}^M ; \text{response} \mid \text{query})$

$$\min_{\text{query}} \mathbb{E}_{\text{response}, \alpha, (\theta_i)_{i=1}^M \mid \text{query}} \log \frac{\mathbb{E}_{\alpha', (\theta'_i)_{i=1}^M} P(\text{response} \mid \text{query}, \alpha', (\theta'_i)_{i=1}^M)}{P(\text{response} \mid \text{query}, \alpha, (\theta_i)_{i=1}^M)}$$

Maximize Information Gain

Unimodal: $\max_{\text{query}} I(\theta ; \text{response} \mid \text{query})$

Multimodal: $\max_{\text{query}} I(\alpha, (\theta_i)_{i=1}^M ; \text{response} \mid \text{query})$

$$\min_{\text{query}} \mathbb{E}_{\text{response}, \alpha, (\theta_i)_{i=1}^M \mid \text{query}} \log \frac{\mathbb{E}_{\alpha', (\theta'_i)_{i=1}^M} P(\text{response} \mid \text{query}, \alpha', (\theta'_i)_{i=1}^M)}{P(\text{response} \mid \text{query}, \alpha, (\theta_i)_{i=1}^M)}$$

Active Method Better Learns

After 10 random queries:

Middle Shelf



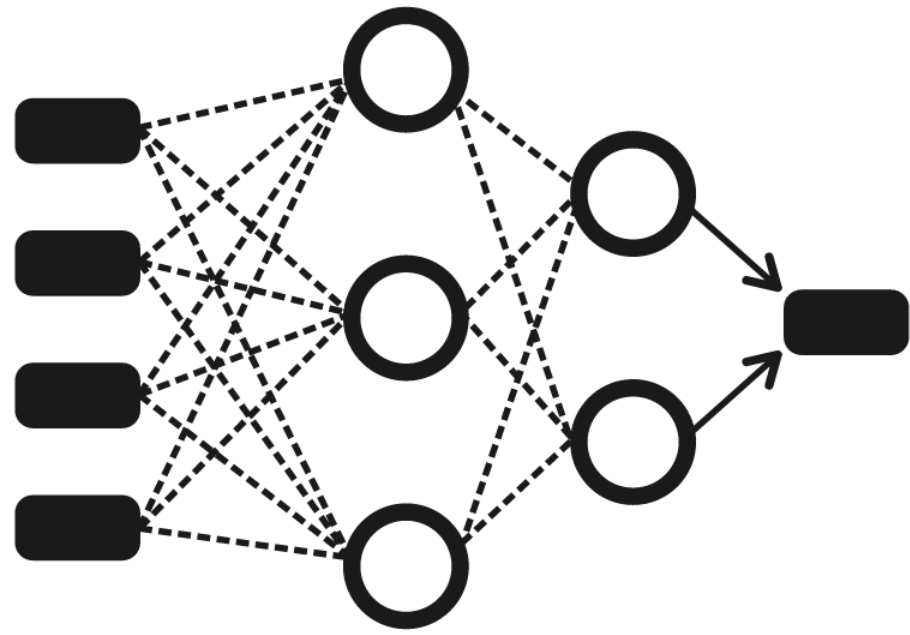
Random querying
did not properly
learn not dropping
items yet

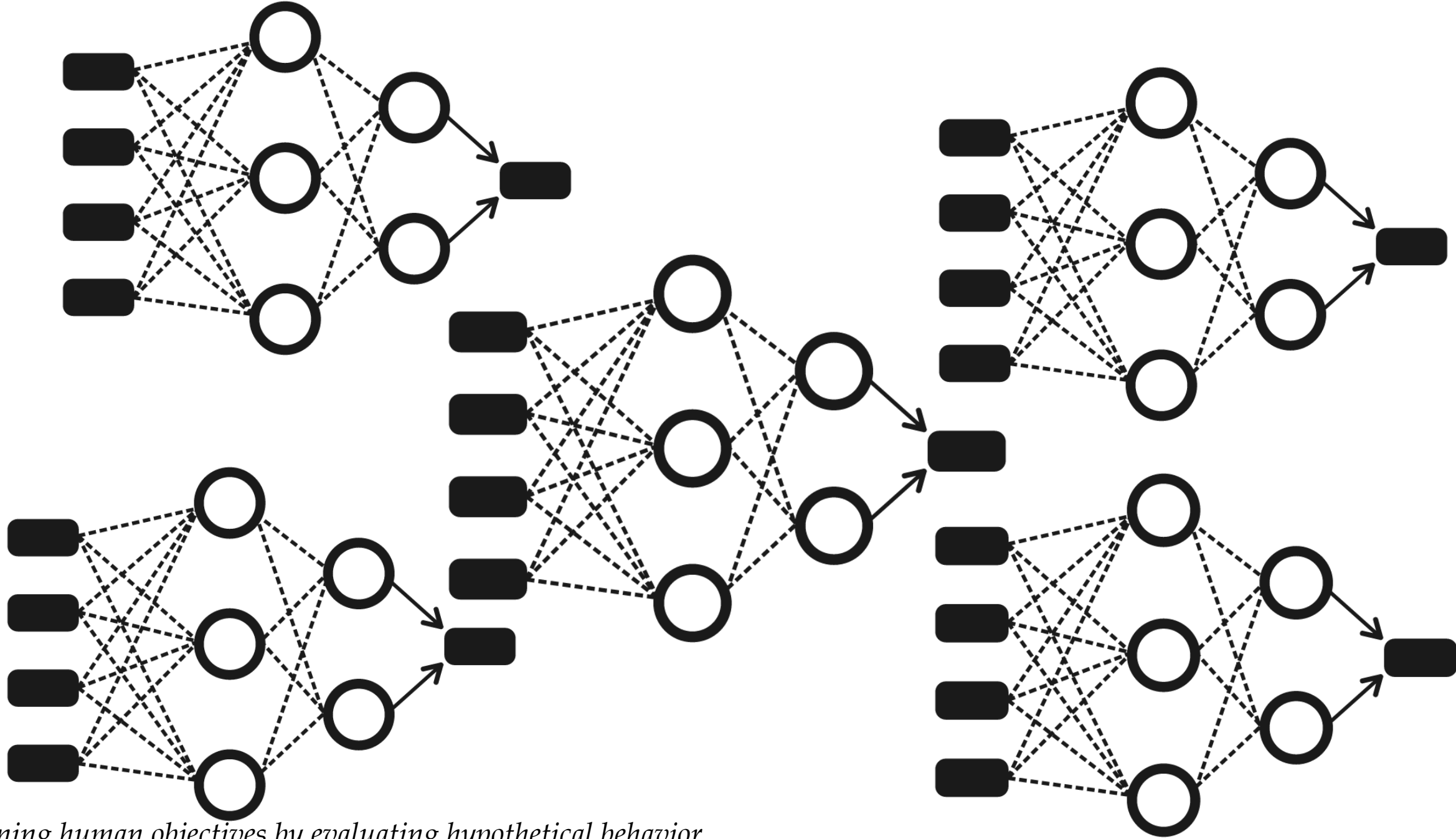
After 10 information-optimal queries:



1. What if our reward function is **nonlinear**?
2. What if our reward function is **multimodal**?

1. What if our reward function is **nonlinear**?
2. What if our reward function is **multimodal**?
3. What if we have a **neural** reward function?






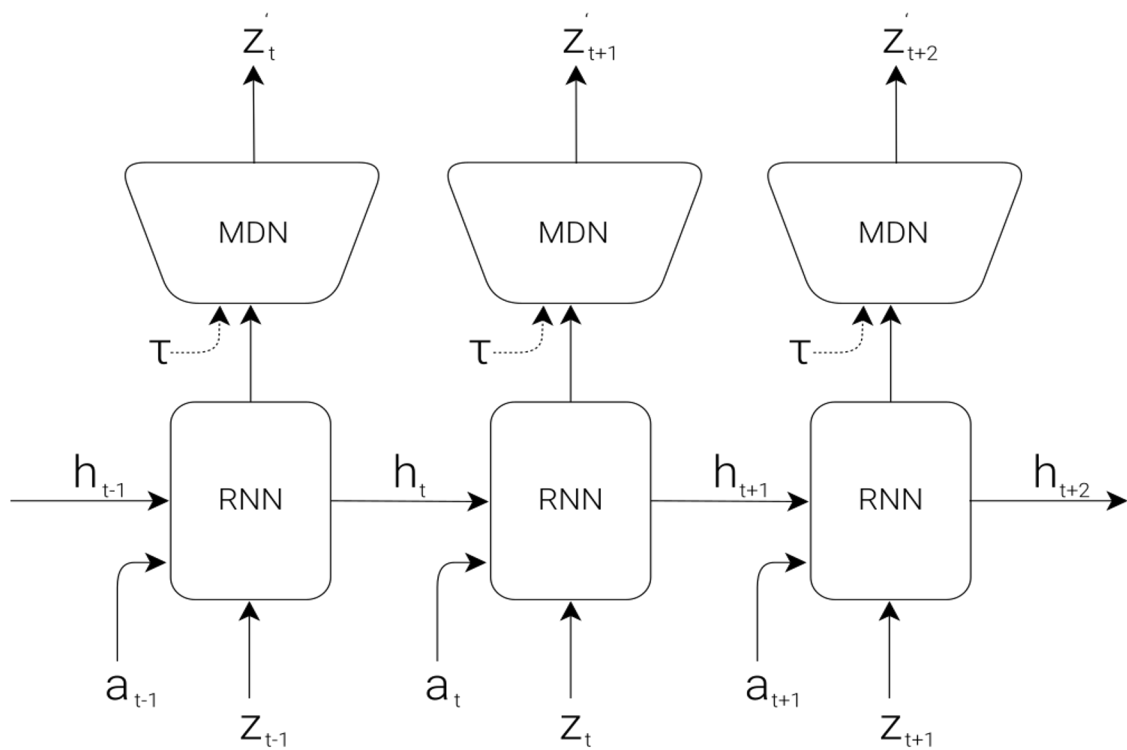
Learning human objectives by evaluating hypothetical behavior
[Reddy, Dragan, Levine, Legg, Leike, ICML'20]

Actively synthesizing queries

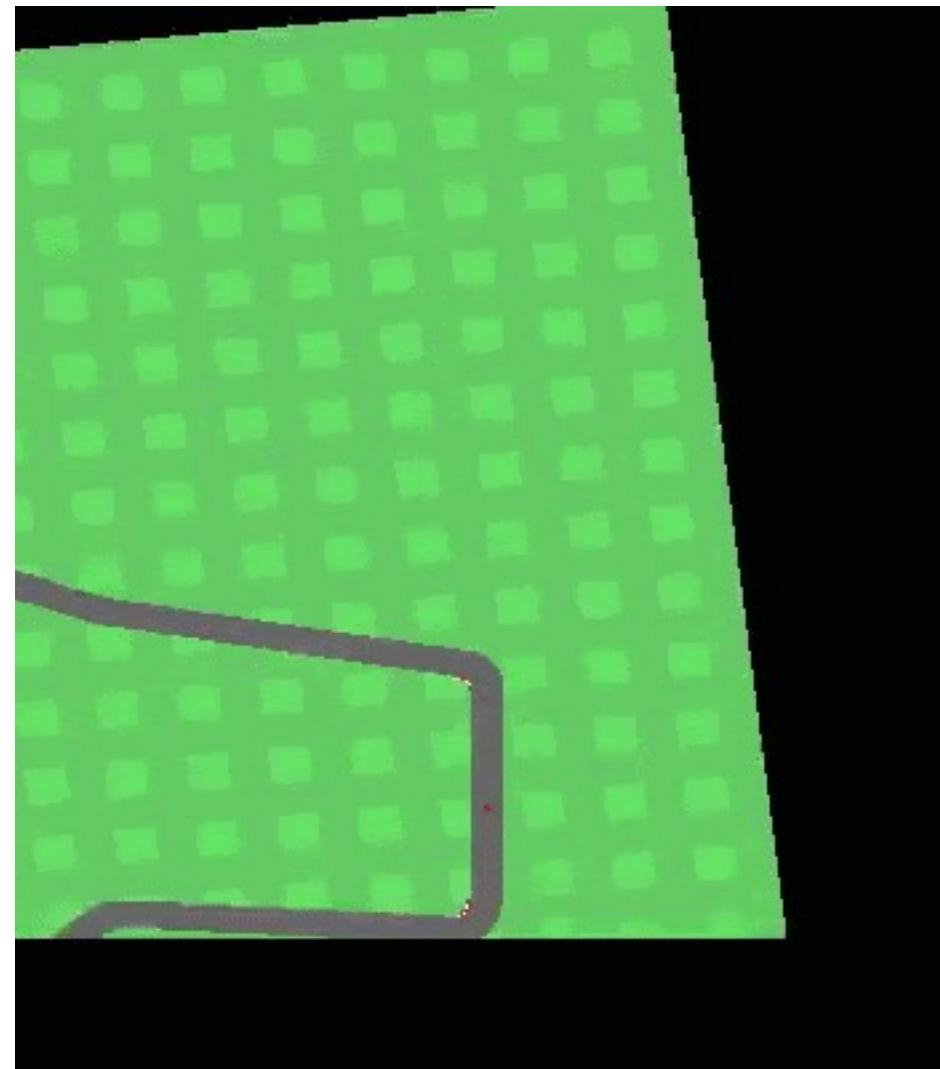
$$\max_{\xi} \text{EnsembleDisagreement}(\mathbf{R}_{\theta}(\xi))$$

Subject to $p_{\phi}(\xi) > \tau$

 learned dynamics model

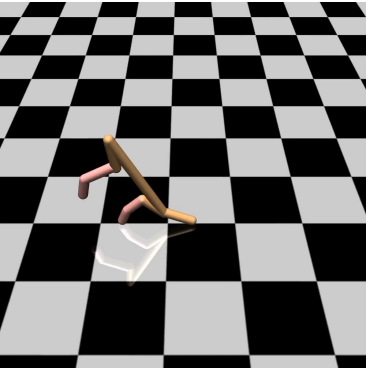


VAE + RNN + Mixture Density Network

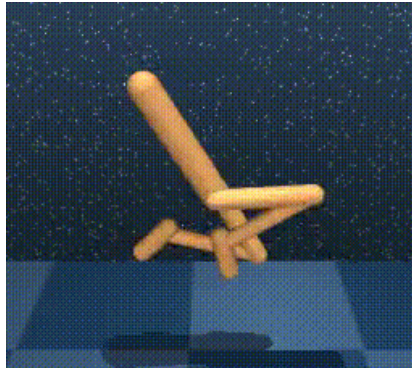
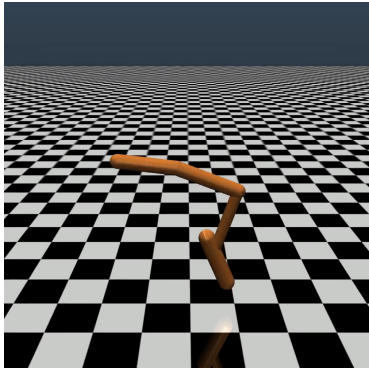


Learning human objectives by evaluating hypothetical behavior
 [Reddy, Dragan, Levine, Legg, Leike, ICML'20]

Active Learning of Neural Rewards



Christiano et al. NeurIPS'17

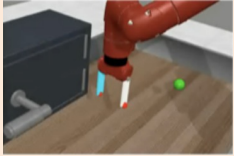


Lee et al. ICML'21

- Generating trajectories takes too much time.
- Human data are expensive.

Pre-training

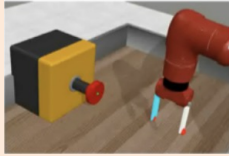
Prior Tasks



Door Open



Window Open



Button Press

Pre-training

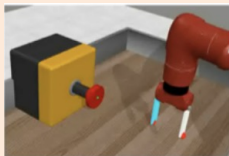
Prior Tasks



Door Open



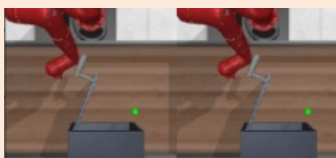
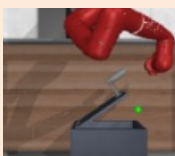
Window Open



Button Press



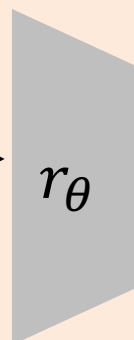
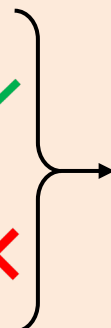
...



...



Segments σ_1, σ_2

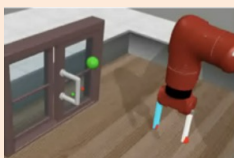


Pre-training

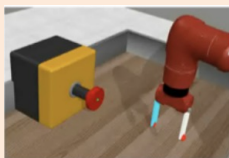
Prior Tasks



Door Open



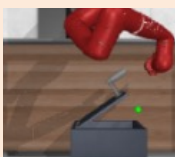
Window Open



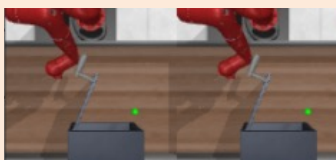
Button Press



...



✓

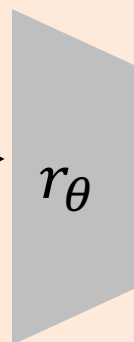
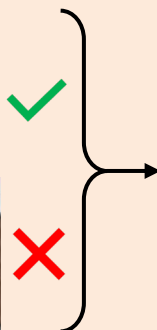


...



✗

Segments σ_1, σ_2



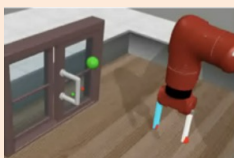
Online Adaptation

Pre-training

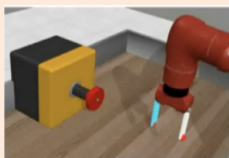
Prior Tasks



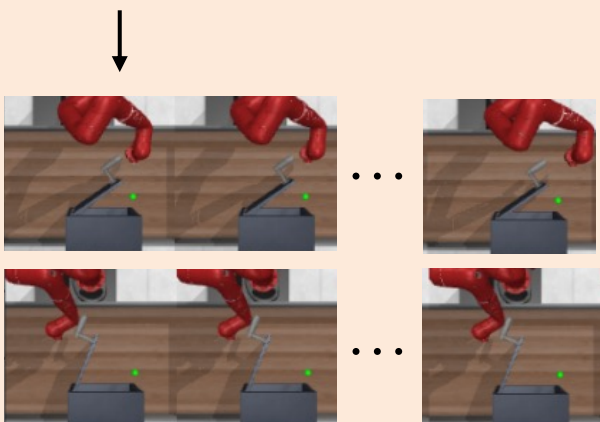
Door Open



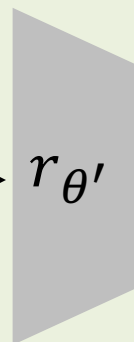
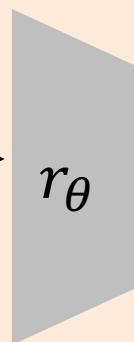
Window Open



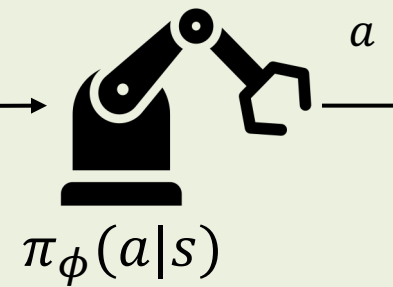
Button Press



Segments σ_1, σ_2



$$r_{\theta'}(s, a)$$



$$\pi_{\phi}(a|s)$$



Drawer Open
New Task

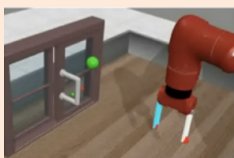
Online Adaptation

Pre-training

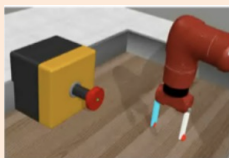
Prior Tasks



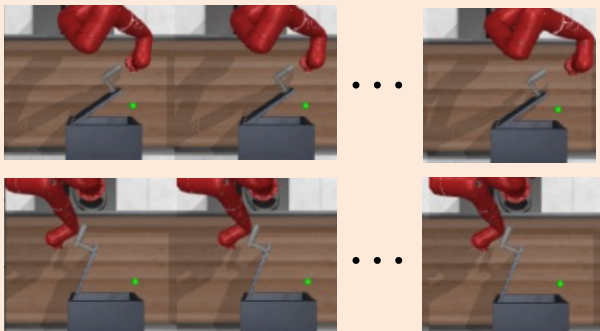
Door Open



Window Open



Button Press



Segments σ_1, σ_2

r_θ

$r_{\theta'}$

$r_{\theta'}(s, a)$

Online Adaptation

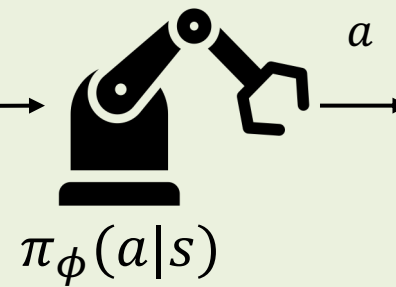
Segments σ_1, σ_2



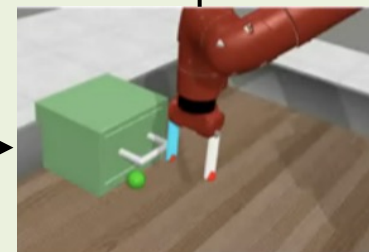
(s, a, s')

Replay Buffer

(s, a, s')



$\pi_\phi(a|s)$



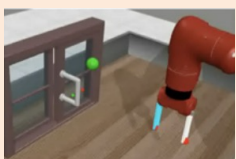
Drawer Open
New Task

Pre-training

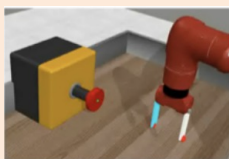
Prior Tasks



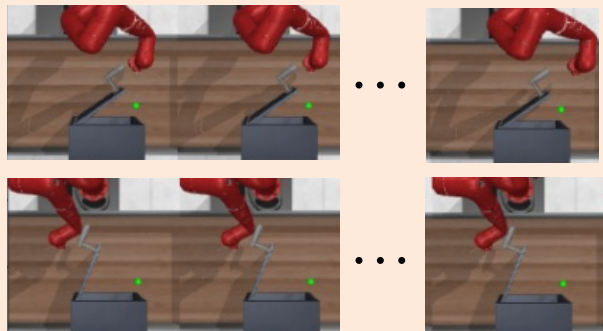
Door Open



Window Open



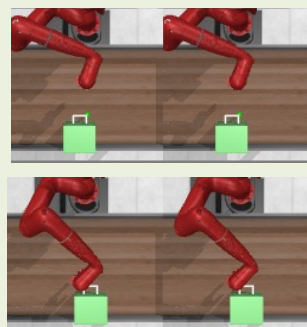
Button Press



Segments σ_1, σ_2

r_θ

Online Adaptation



...



...



Segments σ_1, σ_2

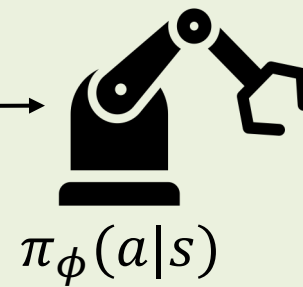
(s, a, s')



(s, a, s')

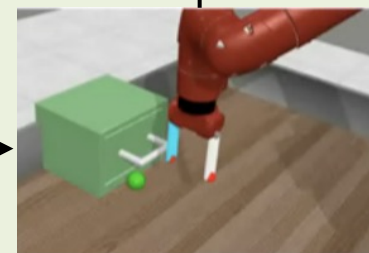
$r_{\theta'}(s, a)$

$r_{\theta'}$



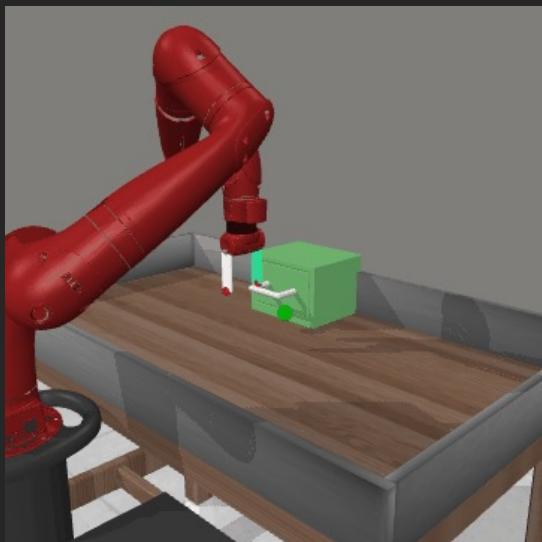
$\pi_\phi(a|s)$

a

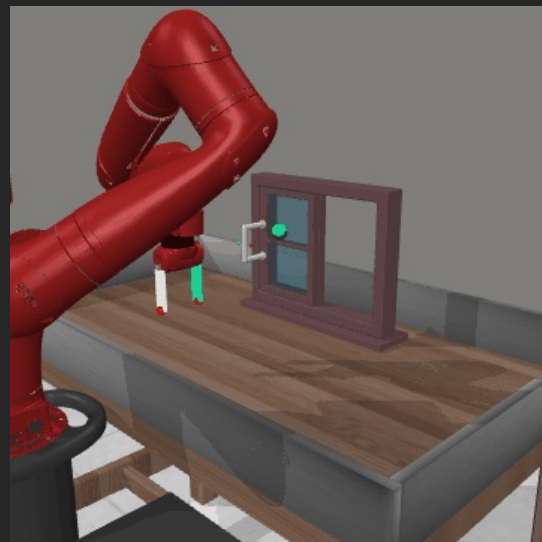


Drawer Open
New Task

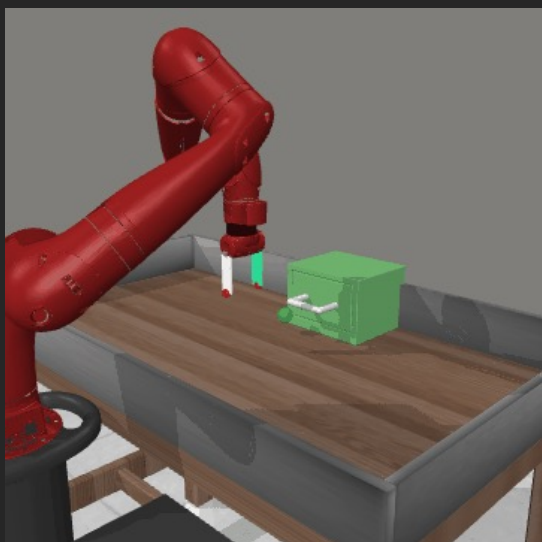
OURS



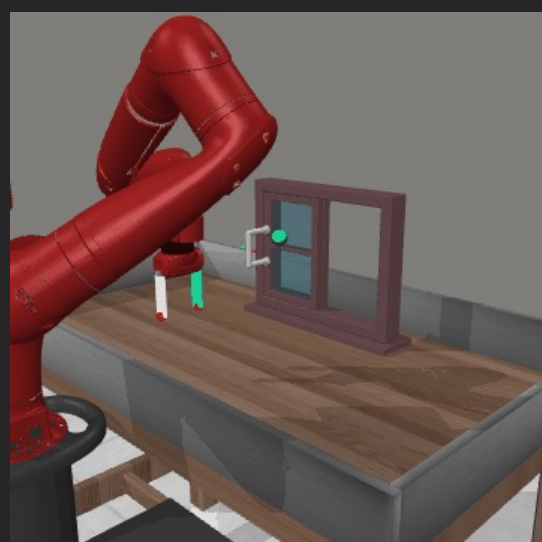
OURS



PEBBLE



PEBBLE



1. What if our reward function is **nonlinear**?
2. What if our reward function is **multimodal**?
3. What if we have a **neural** reward function?