

Learning from Human Feedback: From Demonstrations to Comparative Language

Erdem Biyık

Assistant Professor

Thomas Lord Department of Computer Science
Ming Hsieh Department of Electrical and Computer Engineering

Machine Learning Solution to Robotics



Tons of text

(OpenAI 2023)

Language model that
produces human-like texts

14,000,000 images

(Deng et al. 2009)

Image recognition models
at human-level proficiency

44,000,000 chess games

(Silver et al. 2017)

Super-human
chess engines

We do not have large datasets in robotics



Comparing datasets

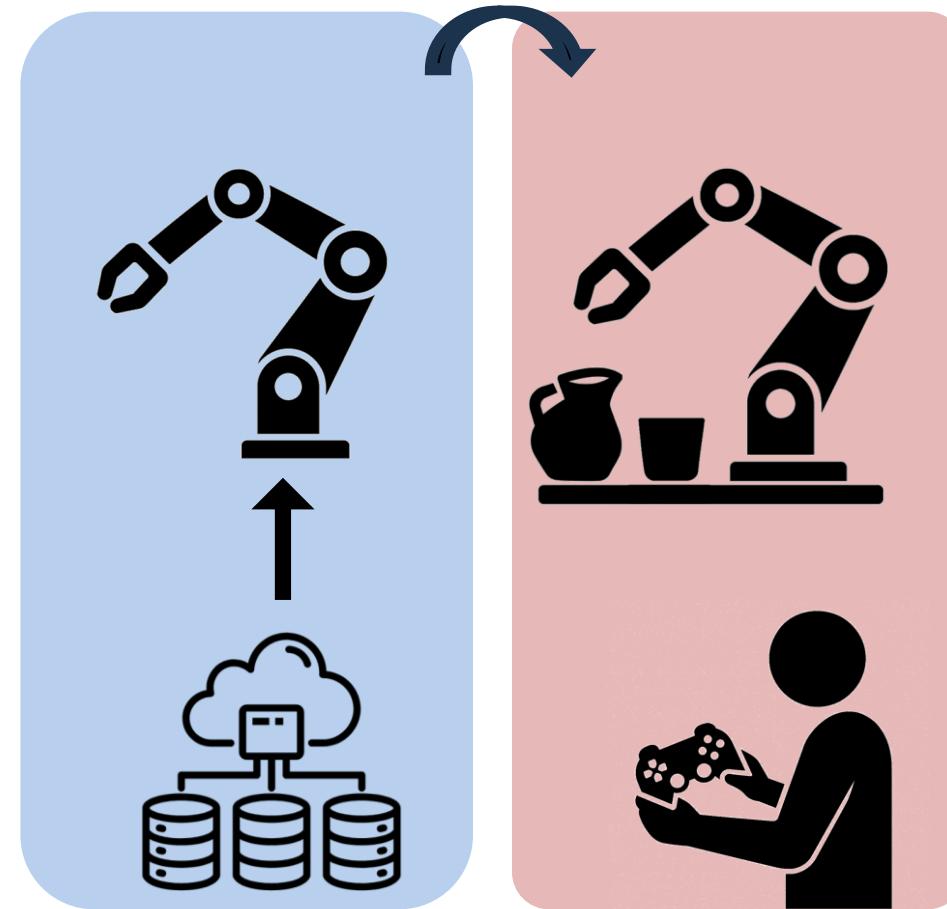
Assuming 238 words/minute, 1.33 tokens/word

OXE
4k hours

π data
10k hours

GPT-2
475k hours

Llama 3
790m hours



Learning from demonstrations (LfD)



Codevilla et al. ICRA'18

Cao et al. RSS'20

Chen et al. IROS'19

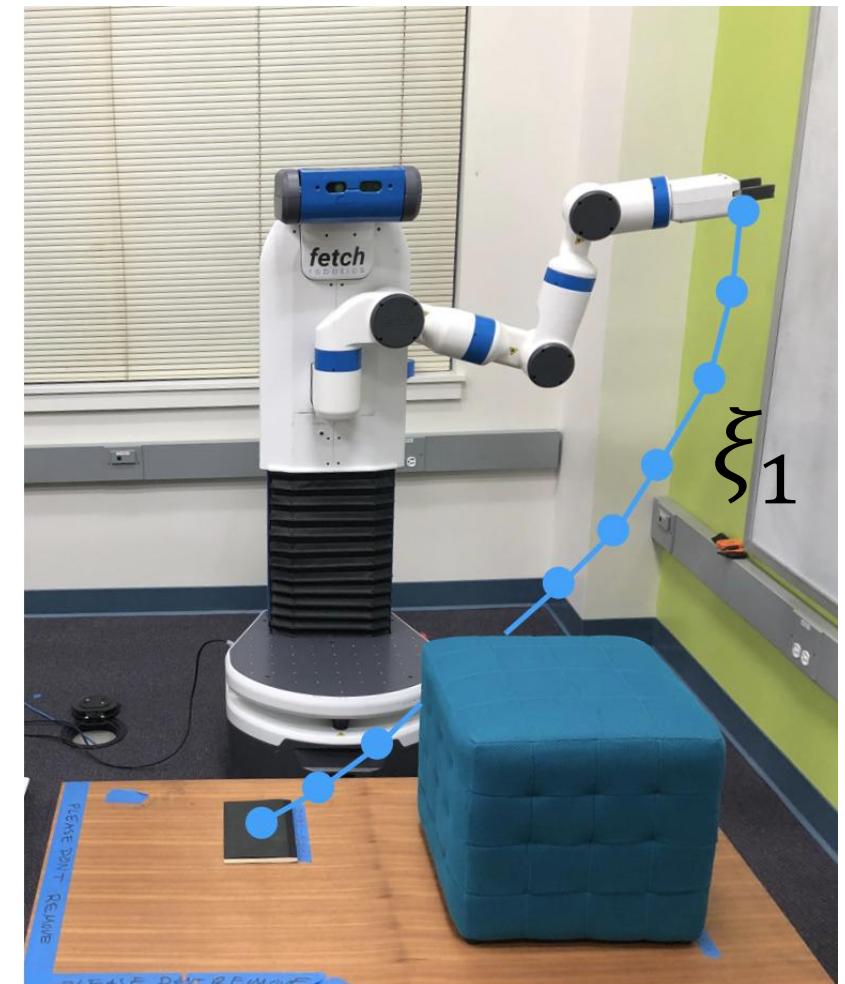
Why does LfD fail?

Demonstrations: $\mathcal{D} = \{\xi_1, \xi_2, \dots, \xi_L\}$

Trajectory features: $\phi(\xi_i) = \phi_i \in \mathbb{R}^d$

- Final distance to the notebook
- Minimum distance to the obstacle
- Average speed
- ...

Reward function : $R(\xi_i) = f_w(\underline{\phi}_i)$



Bayesian inverse reinforcement learning

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D})$$

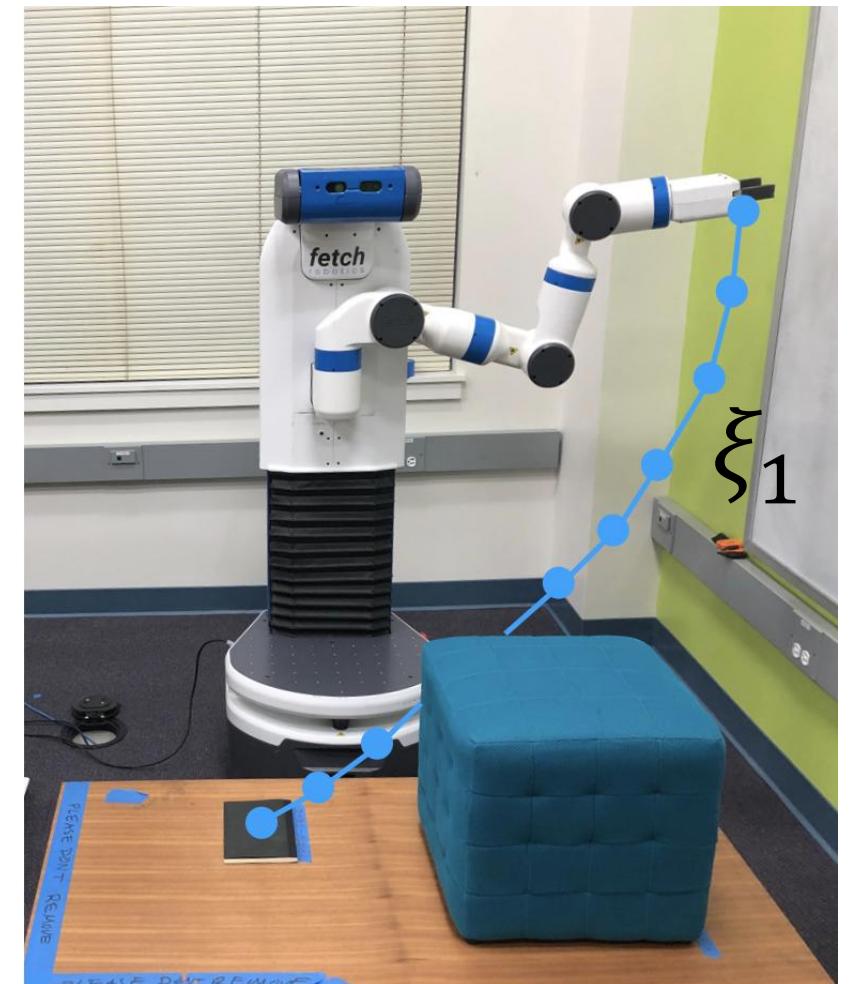
$$P(w \mid \mathcal{D}) \propto P(w) \underline{P(\mathcal{D} \mid w)}$$

$$= P(w) \prod_{i=1}^L P(\xi_i \mid w)$$



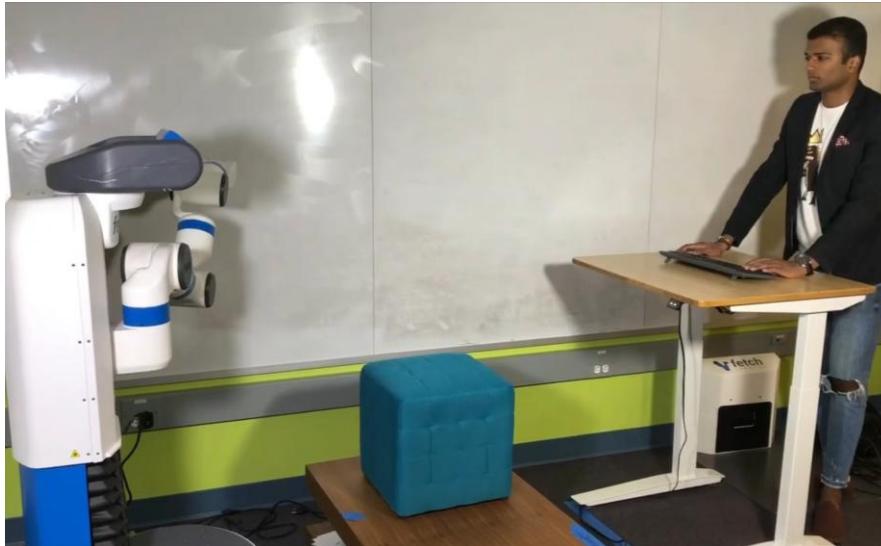
$$\propto P(w) \prod_{i=1}^L \exp f_w(\xi_i)$$

(Noisy humans)



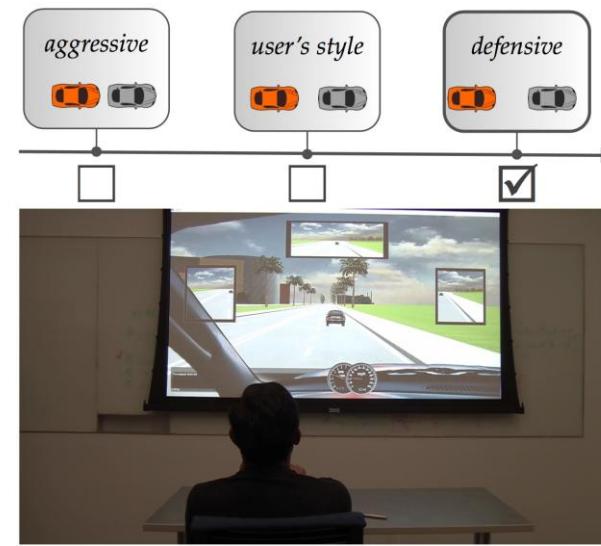
Humans are Suboptimal

Robots with high degrees of freedom are hard to teleoperate.



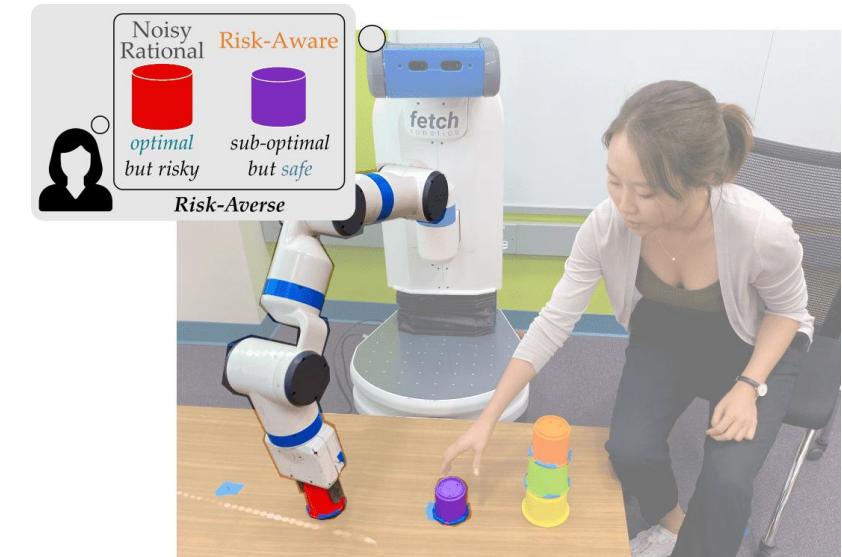
Palan et al. RSS'19

Humans do not like their own demonstrations.



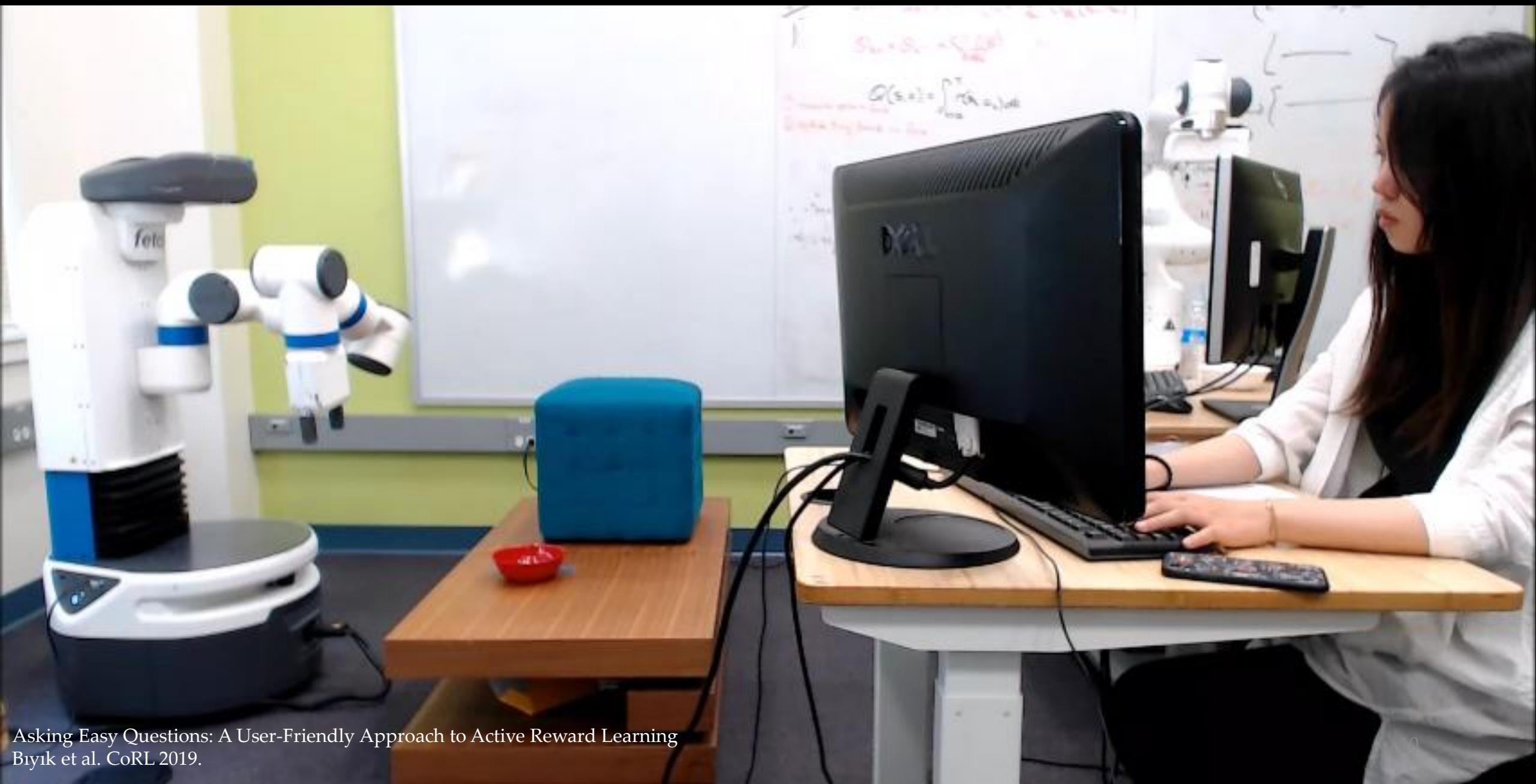
Basu et al. HRI'17

Humans take suboptimal actions in risky situations.



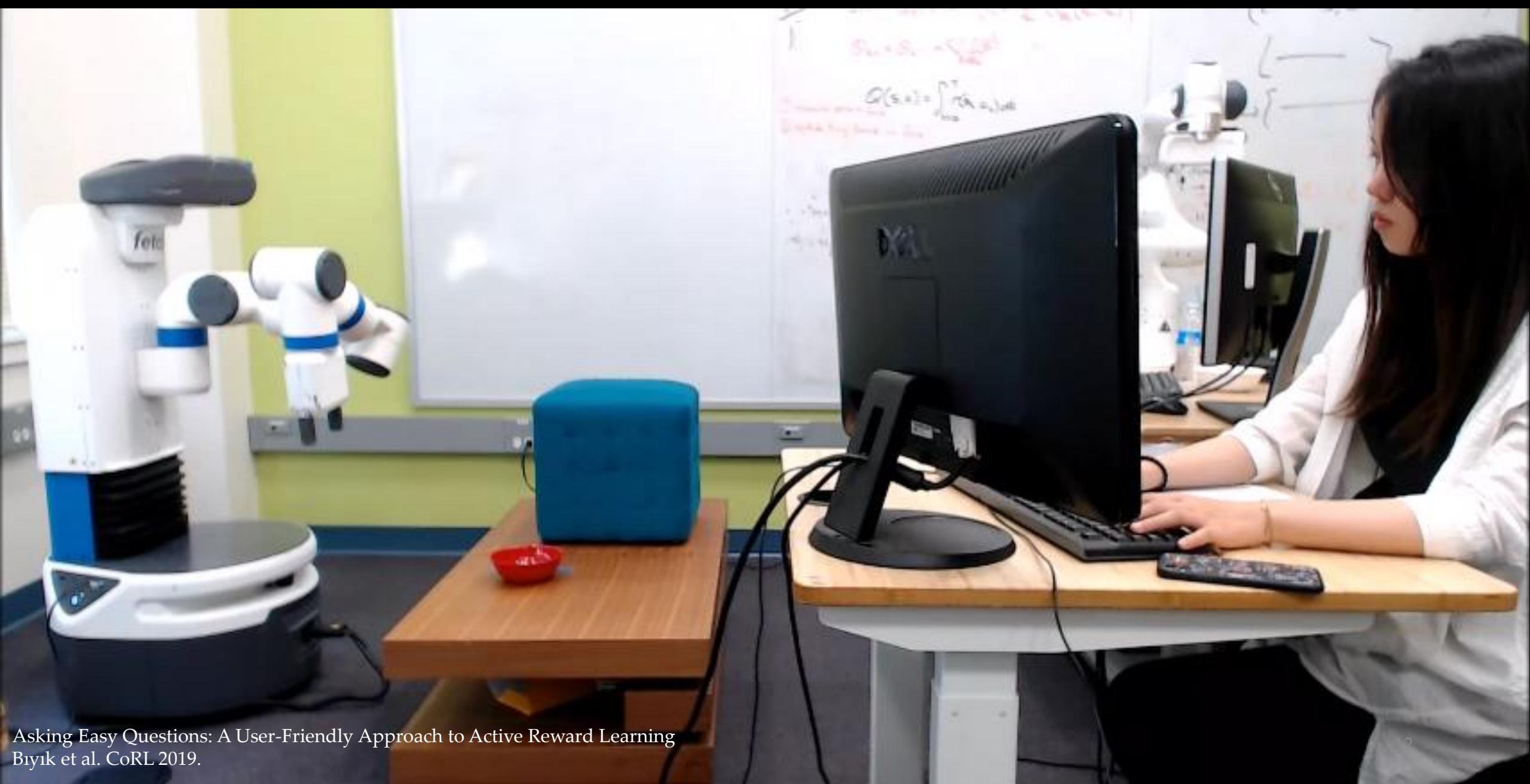
Kwon et al. HRI'20

We can let the human evaluate a robot demonstration



How dark is this blue?

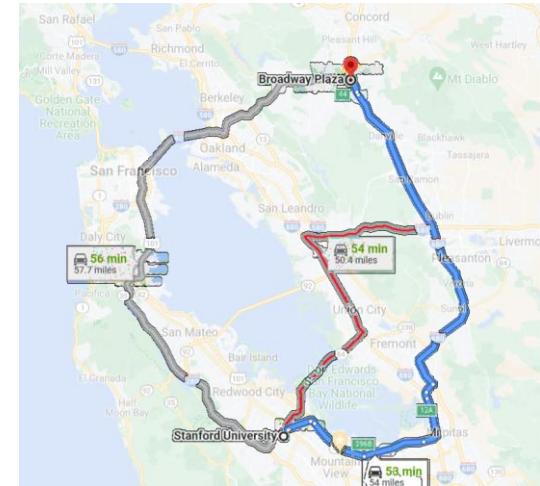
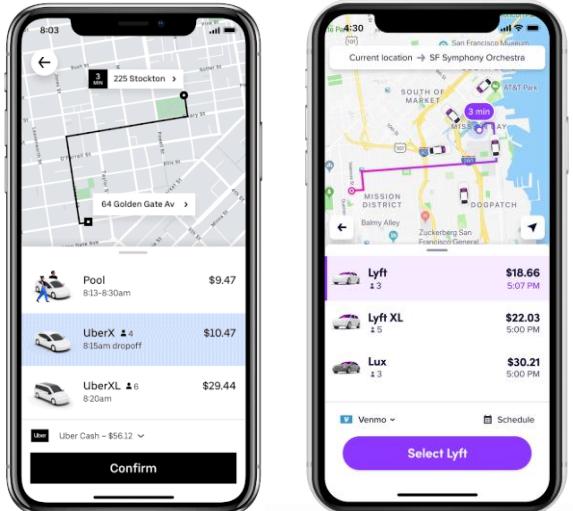
Human evaluations are often unreliable





Which blue is darker?

Comparison data



 **Mohsen Namjoo & Nederlands Blazers Ensemble - Nobahaari" @ Theaters...**

Café Nim
71K views • 5 years ago



 **how my deaf cat meows**

Thundy
7.6M views • 11 months ago



 **John Lowe 9-dart finish FIRST EVER ON TV**

Unicorn Darts
6.6M views • 4 years ago



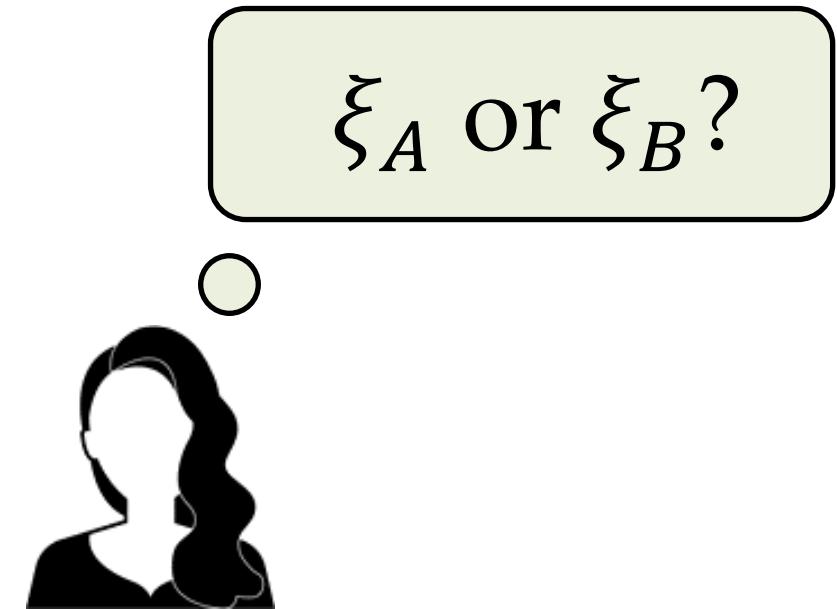
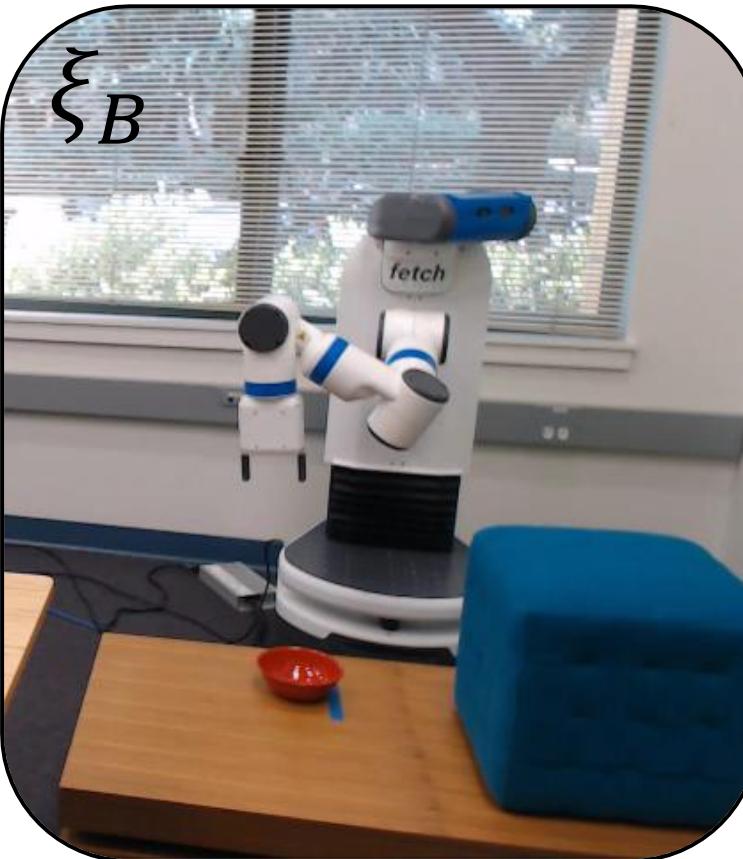
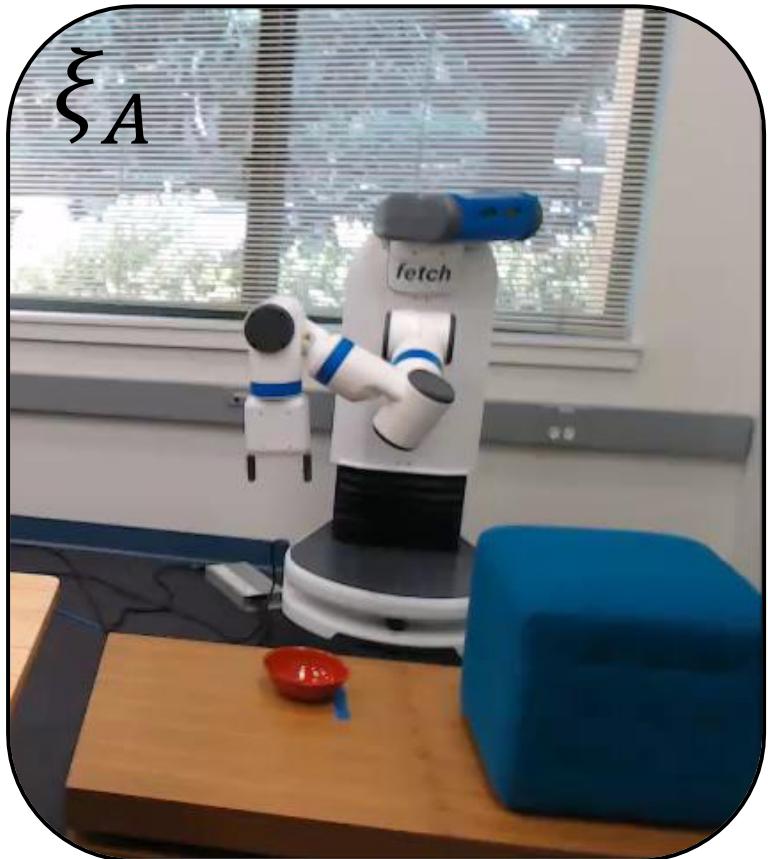
 **Carlsen - Nepomniachtchi | Game 8 | World Chess Championship | Howell,...**

chess24
24K watching

Today...

- Learning from human feedback
 - Pairwise comparisons
 - Reinforcement learning from human feedback (RLHF)
 - Comparative language feedback

Incorporating Comparisons



Incorporating Comparisons

Demonstrations: $\mathcal{D} = \{\xi_1, \xi_2, \dots, \xi_L\}$

Comparisons: $\mathcal{C} = \left\{ \left(\xi_A^{(1)}, \xi_B^{(1)}, q^{(1)} \right), \dots, \left(\xi_A^{(N)}, \xi_B^{(N)}, q^{(N)} \right) \right\}$

Trajectory features: $\phi(\xi_i) = \phi_i \in \mathbb{R}^d$

- Final distance to the notebook
- Minimum distance to the obstacle
- Average speed
- ...

Reward function : $R(\xi_i) = f_w(\phi_i)$

Incorporating Comparisons

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D})$$

$$P(w \mid \mathcal{D}) \propto P(w)P(\mathcal{D} \mid w)$$

$$= P(w) \prod_{i=1}^L P(\xi_i \mid w)$$

Incorporating Comparisons

$$\operatorname{argmax}_w P(w \mid \mathcal{D}, \mathcal{C})$$

$$\begin{aligned} P(w \mid \mathcal{D}, \mathcal{C}) &\propto P(w)P(\mathcal{D} \mid w)P(\mathcal{C} \mid w) \\ &= P(w) \prod_{i=1}^L P(\xi_i \mid w) \prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)}) \end{aligned}$$

Incorporating Comparisons

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D}, \mathcal{C})$$

$$P(w \mid \mathcal{D}, \mathcal{C}) \propto P(w)P(\mathcal{D} \mid w)P(\mathcal{C} \mid w)$$

$$= P(w) \prod_{i=1}^L P(\xi_i \mid w) \boxed{\prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)})}$$

How do we compute this?

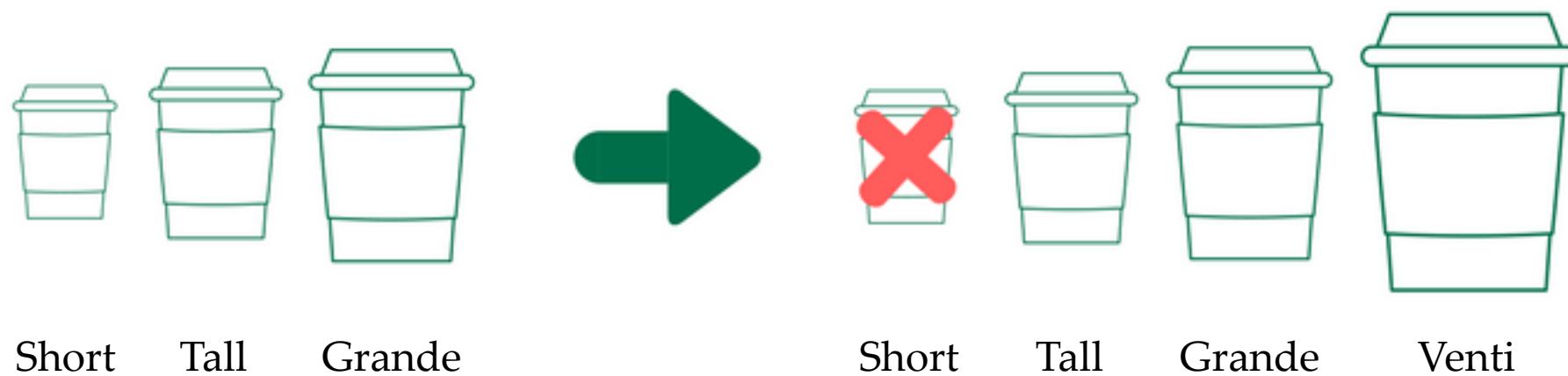
Luce's choice axiom

The probability of selecting one item over another from a pool of many items is not affected by the presence or absence of other items in the pool.

Selection of this kind is said to have *independence from irrelevant alternatives*.

Counterexamples for fun

- Starbucks: “*Compromise Effect*”



Counterexamples for fun

- Coca Cola vs. Pepsi



1985 (Spring)



1985 (Summer)

Regardless...

The probability of selecting one item over another from a pool of many items is not affected by the presence or absence of other items in the pool.

Selection of this kind is said to have *independence from irrelevant alternatives*.

Corollary

$$P(\xi_i \geq \xi_j \geq \xi_k) = P(\xi_i \geq \xi_j, \xi_k)P(\xi_j \geq \xi_k)$$

We only need to model the probability that the human chooses trajectory ξ over a pool of many trajectories.

Incorporating comparisons

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D}, \mathcal{C})$$

$$P(w \mid \mathcal{D}, \mathcal{C}) \propto P(w)P(\mathcal{D} \mid w)P(\mathcal{C} \mid w)$$

$$= P(w) \prod_{i=1}^L P(\xi_i \mid w) \boxed{\prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)})}$$

How do we compute this?

Models from discrete choice theory

$$P(q \mid w, \xi_A, \xi_B)$$

Thurstonian Model:

- Add Gaussian noise to the rewards:

$$\cdot u_A = f_w(\phi(\xi_A)) + z_A$$

$$\cdot u_B = f_w(\phi(\xi_B)) + z_B$$

where $z_A, z_B \sim \mathcal{N}(0, \sigma^2)$.

- The human choice is the noisy winner:

$$\cdot q = \begin{cases} A, & \text{if } u_A > u_B \\ B, & \text{otherwise} \end{cases}$$

$$\begin{aligned} P(q = A) &= P(u_A > u_B) \\ &= P(f_w(\phi(\xi_A)) + z_A > f_w(\phi(\xi_B)) + z_B) \\ &= P(z_A - z_B > f_w(\phi(\xi_B)) - f_w(\phi(\xi_A))) \end{aligned}$$

Models from discrete choice theory

$$P(q \mid w, \xi_A, \xi_B)$$

Bradley-Terry Model:

- The probability that the user chooses an option is proportional to the exponentials of the rewards:

$$P(q = A) = \frac{e^{\beta f_w(\phi(\xi_A))}}{e^{\beta f_w(\phi(\xi_A))} + e^{\beta f_w(\phi(\xi_B))}}$$

Incorporating comparisons

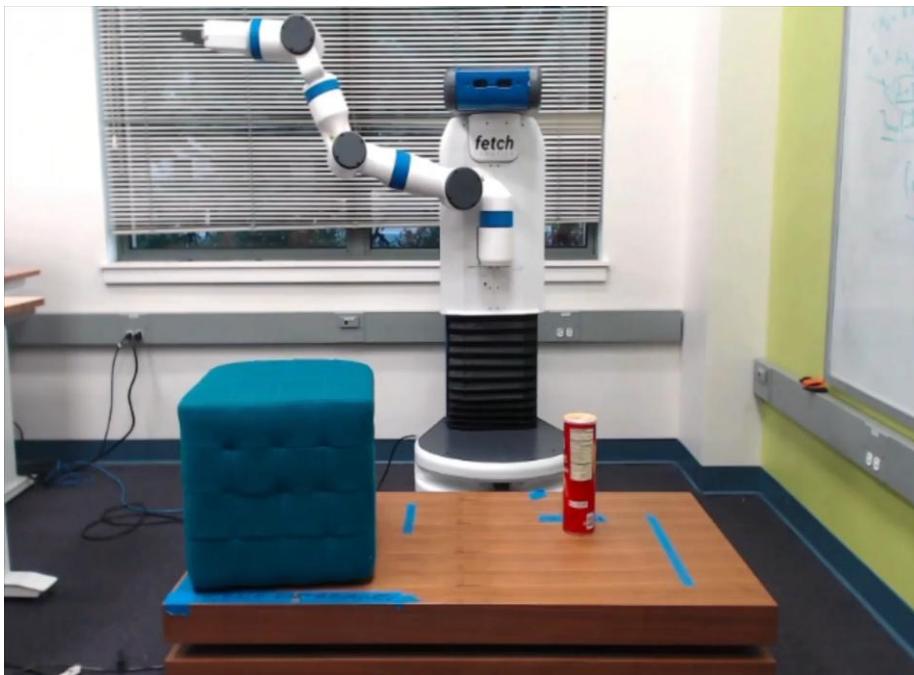
$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D}, \mathcal{C})$$

$$P(w \mid \mathcal{D}, \mathcal{C}) \propto P(w) P(\mathcal{D} \mid w) P(\mathcal{C} \mid w)$$

$$\begin{aligned} &= P(w) \prod_{i=1}^L P(\xi_i \mid w) \prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)}) \\ &\propto P(w) \prod_{i=1}^L \exp f_w(\xi_i) \prod_{i=1}^N \frac{\exp f_w(\xi_{q^{(i)}}^{(i)})}{\exp f_w(\xi_{q^{(i)}}^{(i)}) + \exp f_w(\xi_{\neg q^{(i)}}^{(i)})} \end{aligned}$$

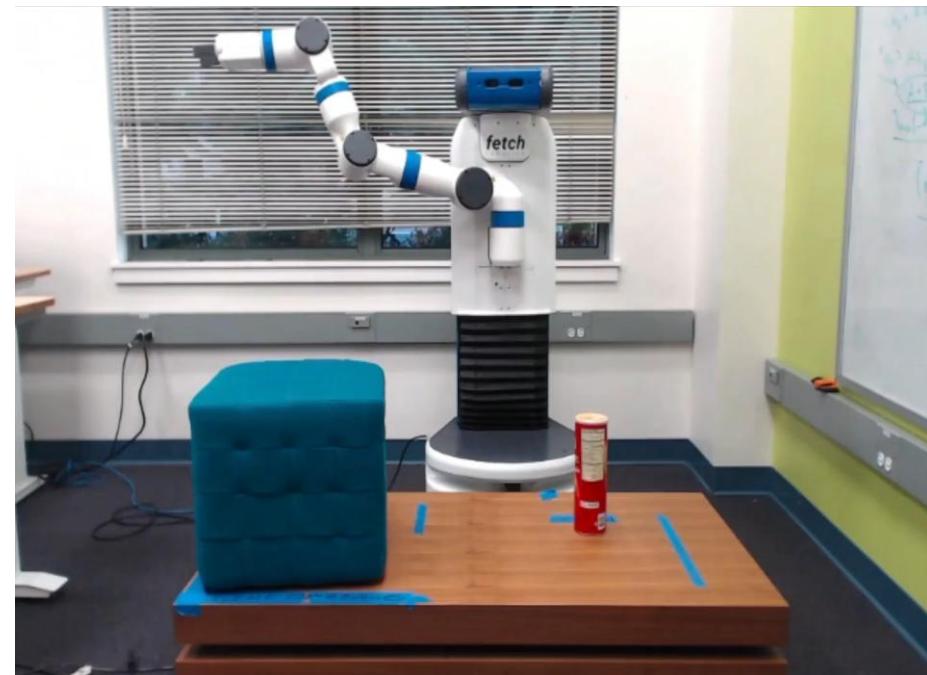
Benefit of comparisons

Bayesian IRL

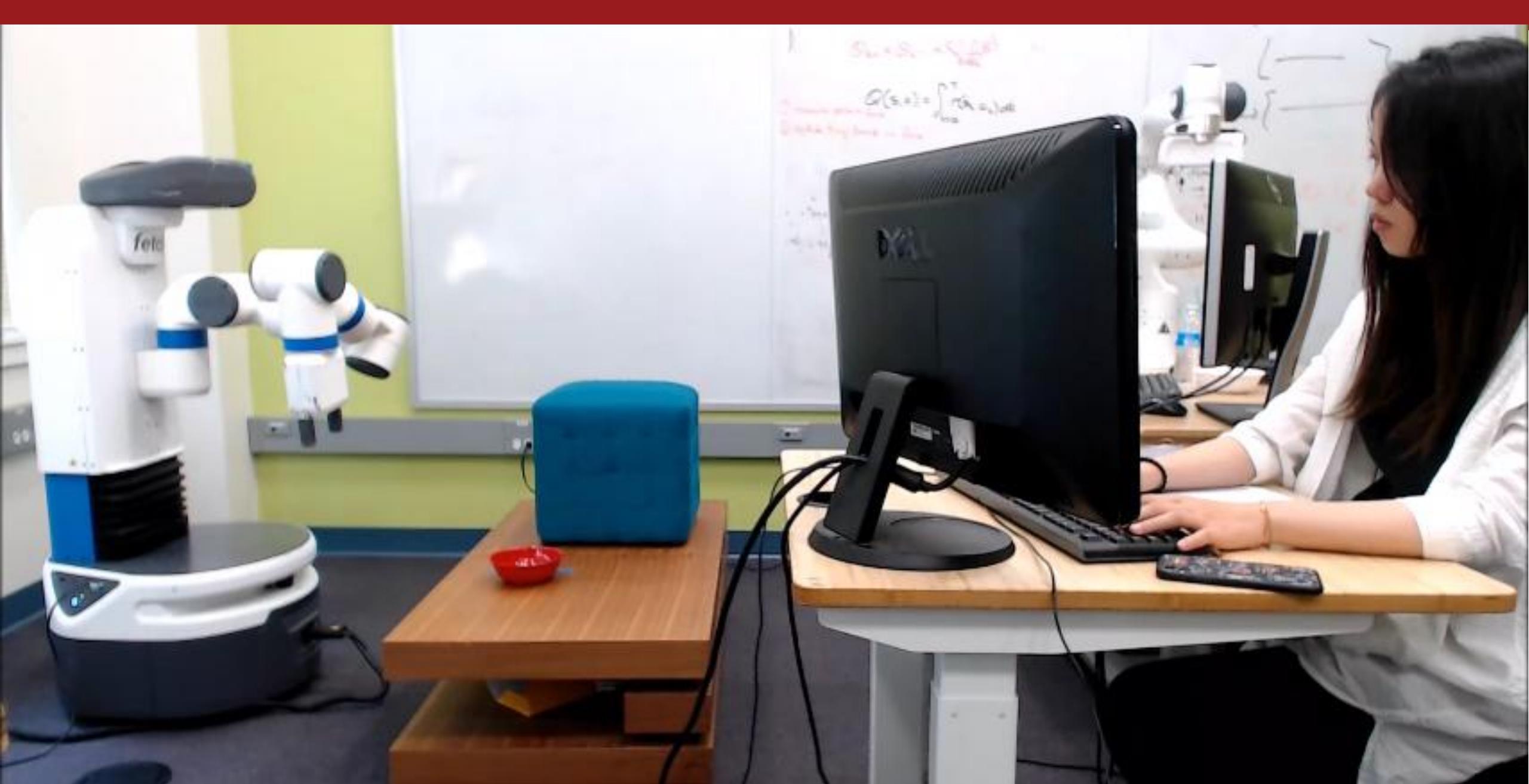


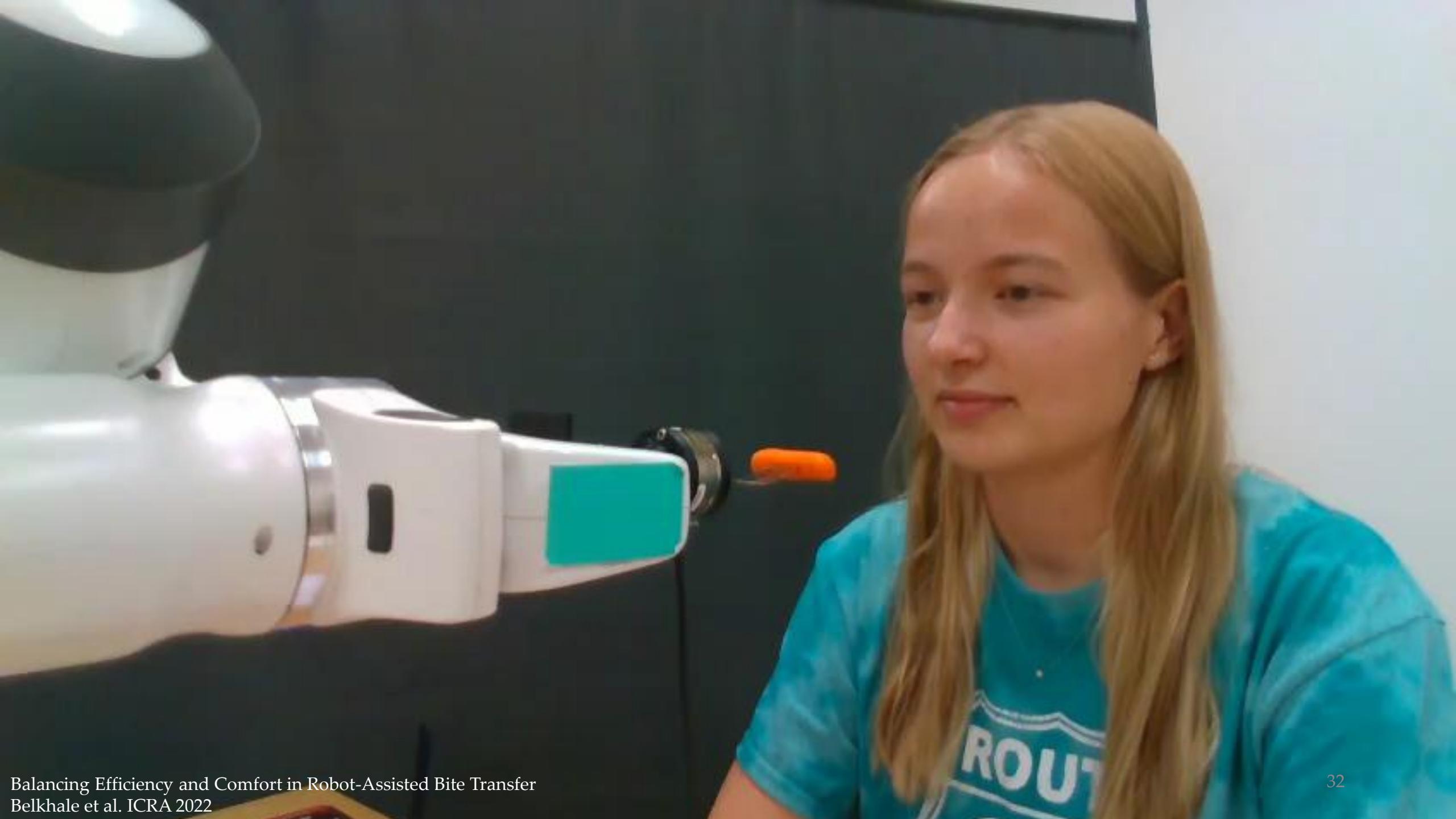
5 demonstrations

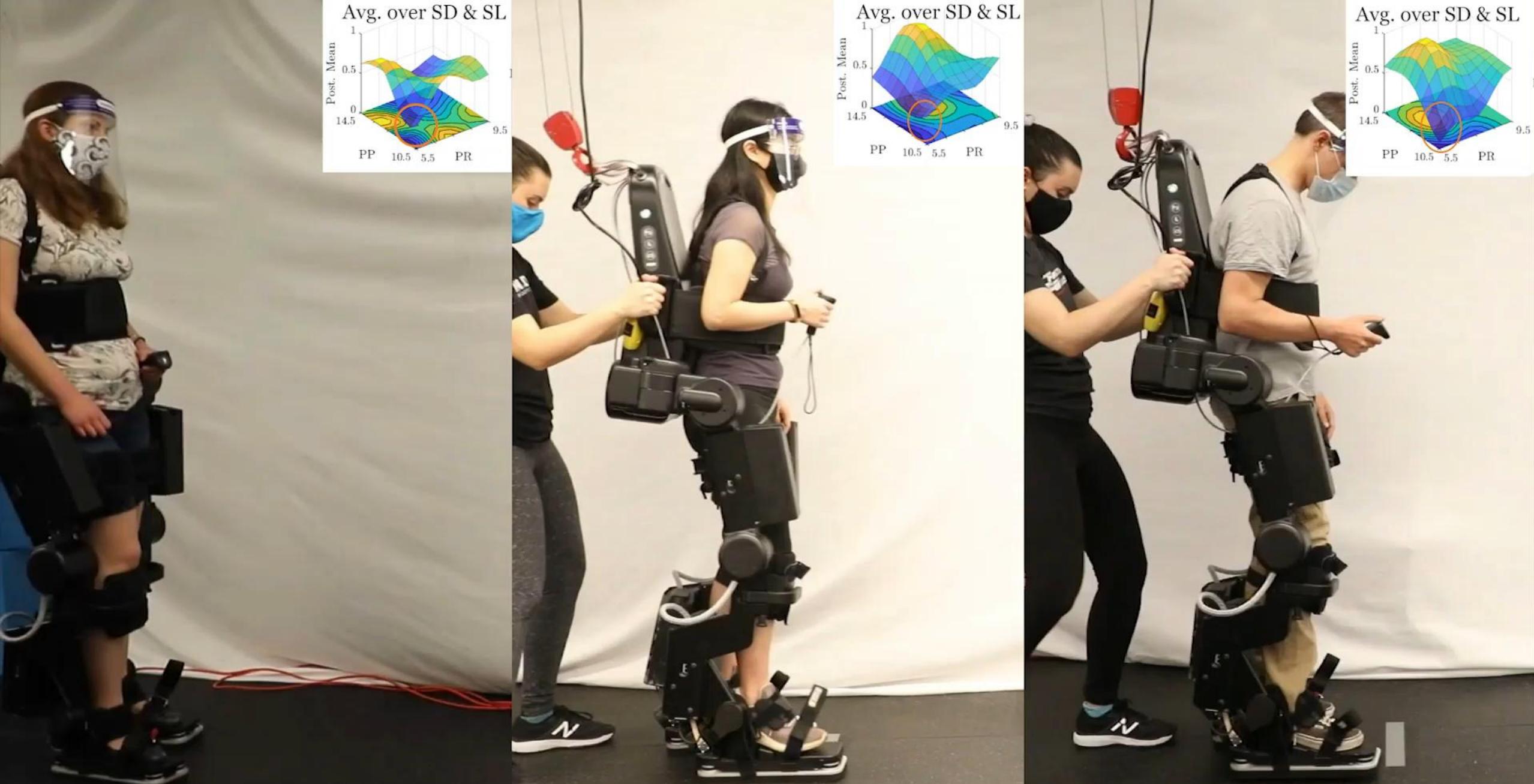
Ours



1 demonstration + 15 comparisons



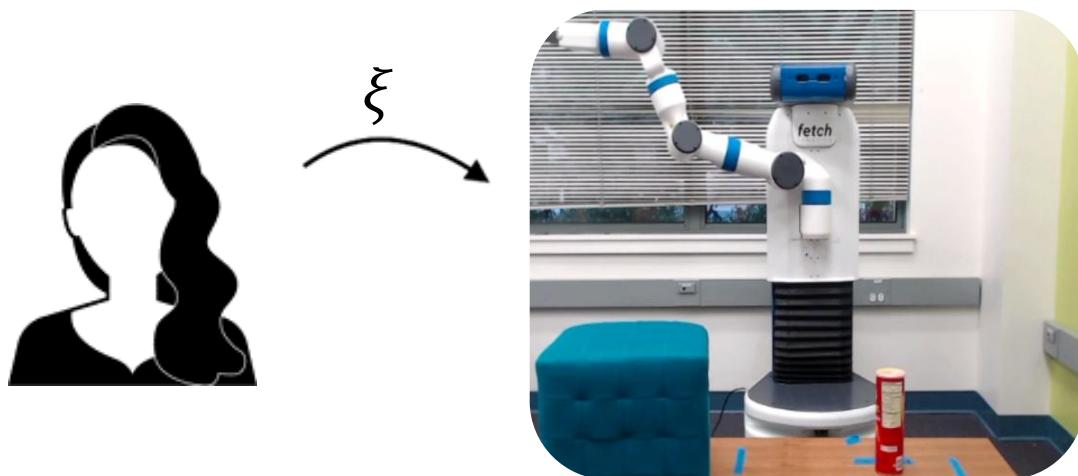




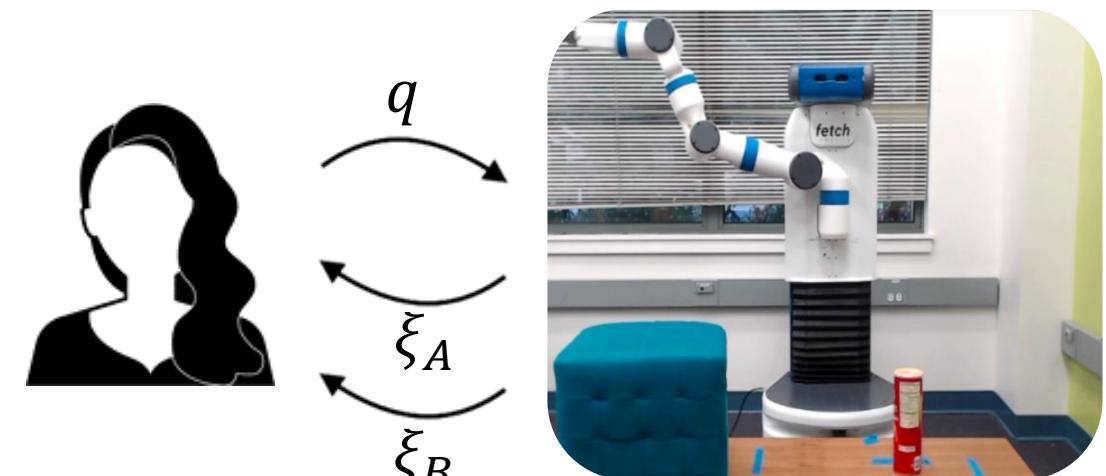
ROIAL: Region of Interest Active Learning for Characterizing
Exoskeleton Gait Preference Landscapes
Li et al., ICRA 2021

Choosing Queries

Demonstrations



Comparisons



How do we quantify information?

Surprise

$95\% \rightarrow X = \text{Heads}$



Surprise: $\log_2 \frac{1}{0.95} \cong 0.074$

$5\% \rightarrow X = \text{Tails}$



Surprise: $\log_2 \frac{1}{0.05} \cong 4.322$

Entropy (a measure of uncertainty)

$$\begin{aligned} 95\% \rightarrow X = \text{Heads} &\longrightarrow \text{Surprise: } \log_2 \frac{1}{0.95} \cong 0.074 \\ 5\% \rightarrow X = \text{Tails} &\longrightarrow \text{Surprise: } \log_2 \frac{1}{0.05} \cong 4.322 \end{aligned}$$

Entropy is the expected surprise.

$$\text{Entropy: } H(X) = 0.95 \times \log_2 \frac{1}{0.95} + 0.05 \times \log_2 \frac{1}{0.05} \cong 0.286$$

Another example

$50\% \rightarrow X = \text{Heads}$

$50\% \rightarrow X = \text{Tails}$

Another example

$$50\% \rightarrow X = \text{Heads} \quad \longrightarrow \quad \text{Surprise: } \log_2 \frac{1}{0.50} = 1$$
$$50\% \rightarrow X = \text{Tails} \quad \longrightarrow \quad \text{Surprise: } \log_2 \frac{1}{0.50} = 1$$

Another example

$$50\% \rightarrow X = \text{Heads} \quad \longrightarrow \quad \text{Surprise: } \log_2 \frac{1}{0.50} = 1$$
$$50\% \rightarrow X = \text{Tails} \quad \longrightarrow \quad \text{Surprise: } \log_2 \frac{1}{0.50} = 1$$

$$\text{Entropy: } H(X) = 0.50 \times \log_2 \frac{1}{0.50} + 0.50 \times \log_2 \frac{1}{0.50} \cong 1$$

Mutual information

Uncertainty

$$H(X) = 1$$



Alice



Bob

$50\% \rightarrow X = \text{Heads}$

$50\% \rightarrow X = \text{Tails}$

Mutual information

Uncertainty

$$H(X | X) = 0$$

$$0 \times \log \frac{1}{0} + 1 \times \log \frac{1}{1} = 0$$

This is 0 in
information theory.



Alice

What is X ?



Bob

Tails!

50% $\rightarrow X = \text{Heads}$

50% $\rightarrow X = \text{Tails}$

Mutual information

Uncertainty

$$H(X | X) = 0$$



Alice

What is X ?



Bob

Tails!

50% $\rightarrow X = \text{Heads}$

50% $\rightarrow X = \text{Tails}$

$$\begin{aligned}\text{Mutual Information} &= \text{Reduction in Entropy}: I(X; X) &= H(X) - H(X | X) \\ &= 1 - 0 = 1\end{aligned}$$

Mutual information

Uncertainty

$$H(X) = 1$$



Alice



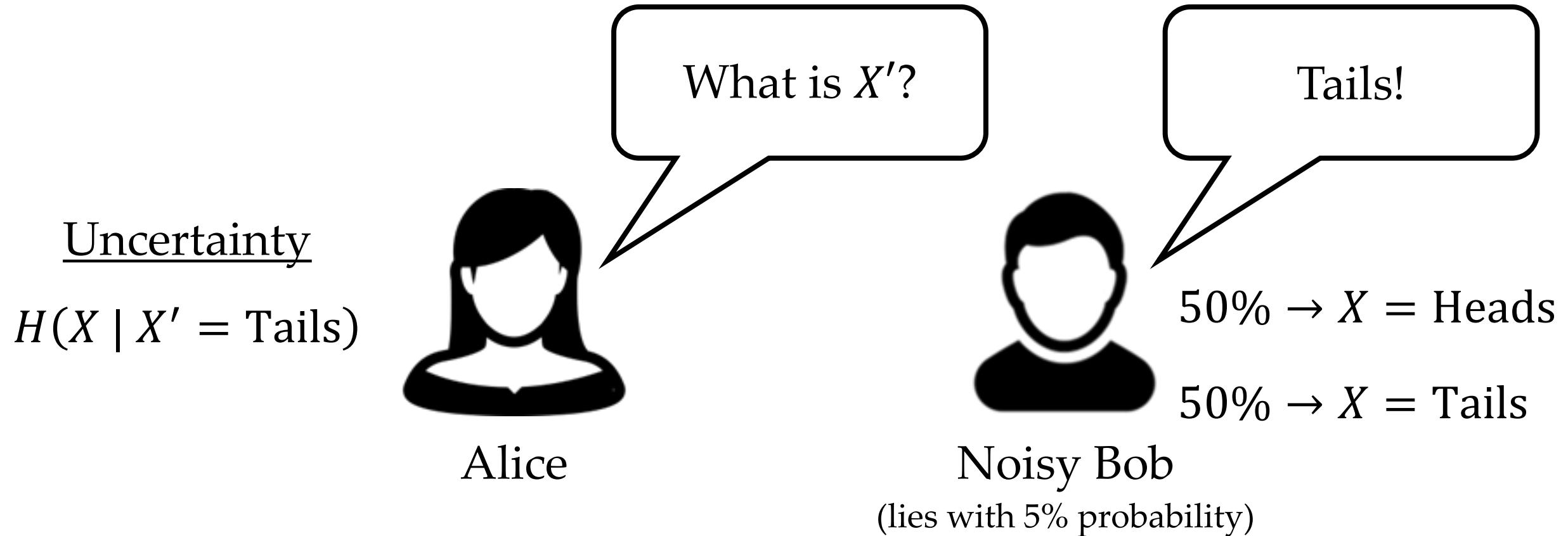
Noisy Bob

(lies with 5% probability)

$50\% \rightarrow X = \text{Heads}$

$50\% \rightarrow X = \text{Tails}$

Mutual information



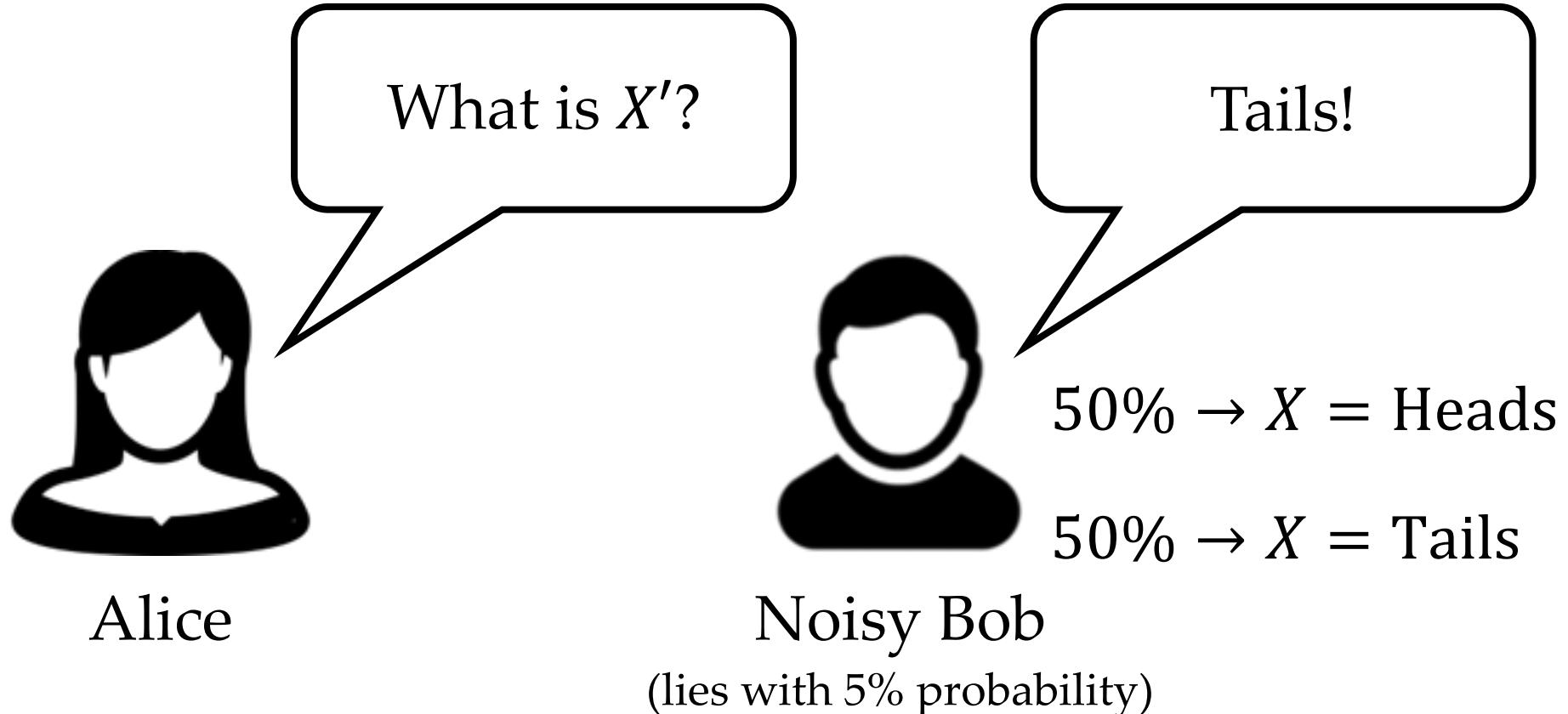
$$P(X = \text{Tails} | X' = \text{Tails}) \propto P(X' = \text{Tails} | X = \text{Tails})P(X = \text{Tails})$$

$$P(X = \text{Heads} | X' = \text{Tails}) \propto P(X' = \text{Tails} | X = \text{Heads})P(X = \text{Heads})$$

Mutual information

Uncertainty

$$H(X | X' = \text{Tails})$$



$$P(X = \text{Tails} | X' = \text{Tails}) = 0.95$$

$$P(X = \text{Heads} | X' = \text{Tails}) = 0.05$$

Mutual information

Uncertainty

$$H(X | X' = \text{Tails}) \\ \approx 0.286$$



Alice

What is X' ?



Noisy Bob
(lies with 5% probability)

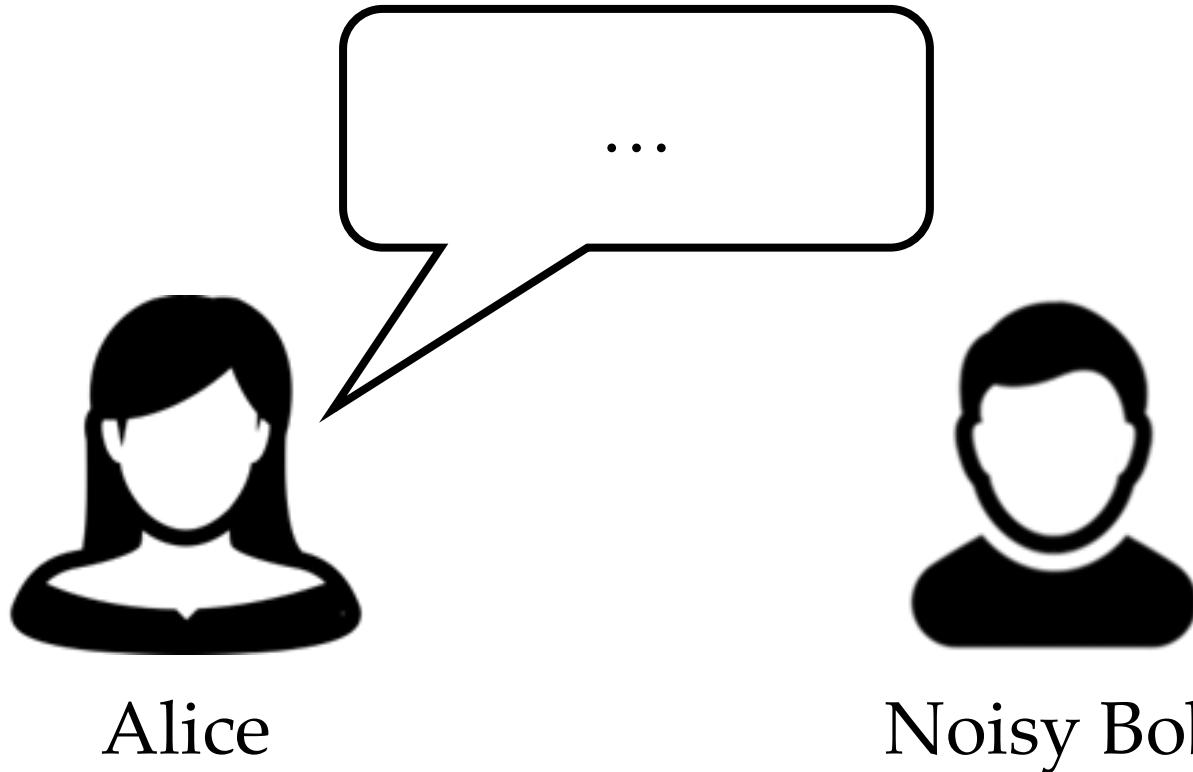
Tails!

$$50\% \rightarrow X = \text{Heads} \\ 50\% \rightarrow X = \text{Tails}$$

Mutual Information = Reduction in Entropy: $I(X; X') = H(X) - H(X | X')$

$$\approx 1 - 0.286 = 0.714$$

Mutual information: what do you ask?

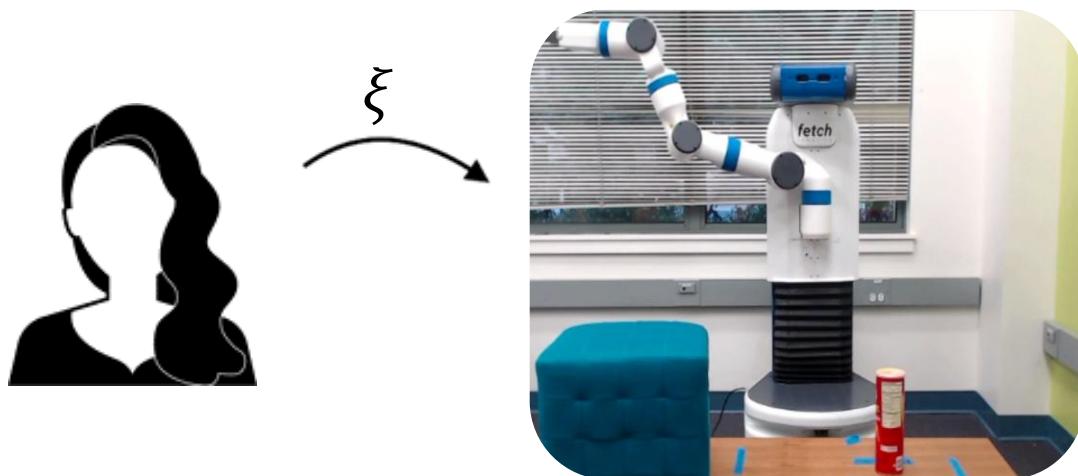


(tells the truth for X_1 and X_2 ,
lies with 5% probability for X_3)

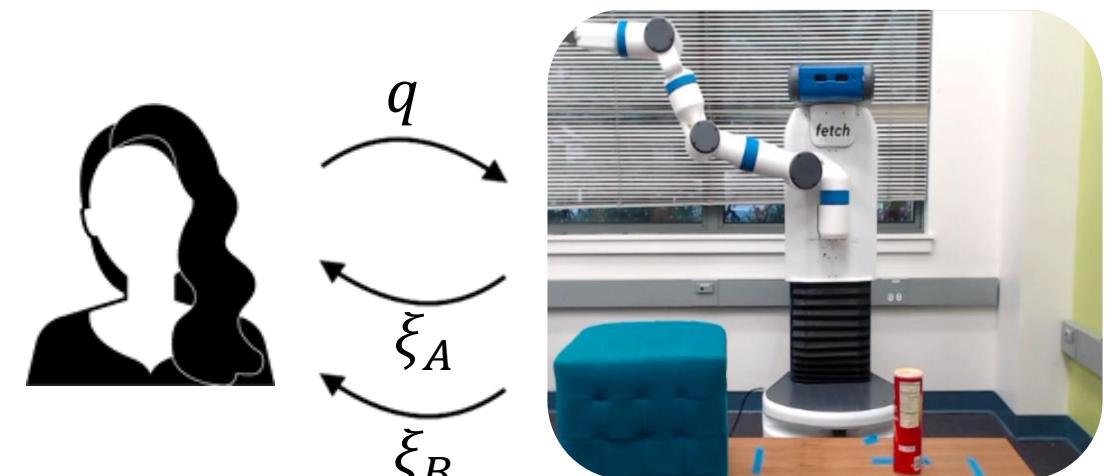
- 20% $\rightarrow X = (\text{H}, \text{H}, \text{H})$
- 20% $\rightarrow X = (\text{H}, \text{H}, \text{T})$
- 20% $\rightarrow X = (\text{H}, \text{T}, \text{H})$
- 20% $\rightarrow X = (\text{H}, \text{T}, \text{T})$
- 20% $\rightarrow X = (\text{T}, \text{T}, \text{T})$

Choosing queries

Demonstrations



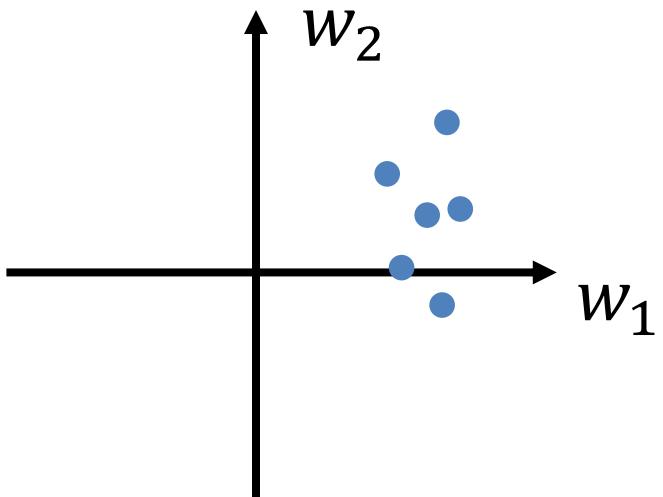
Comparisons



The robot can query the user with the query that will give the **most information**.

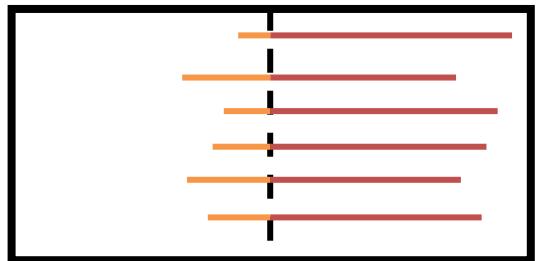
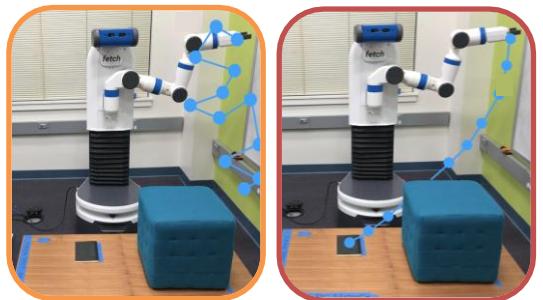
Maximum volume removal

Posterior $P(w | \mathcal{C})$

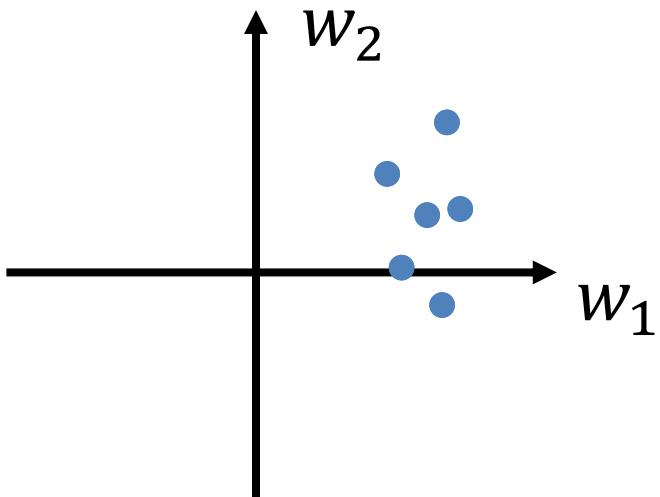


Maximum volume removal

Posterior $P(w | \mathcal{C})$

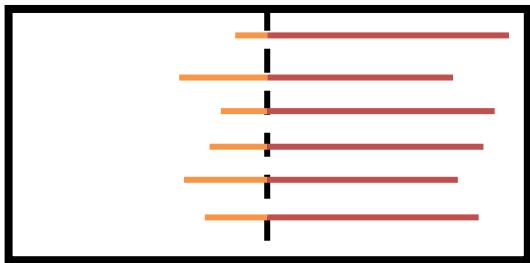
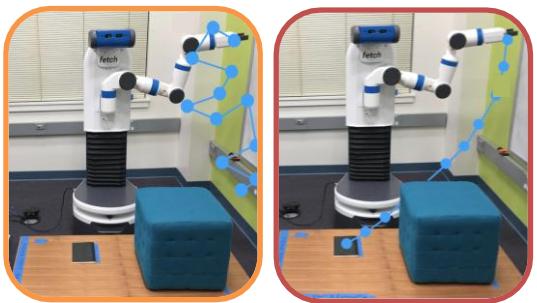


User Choice

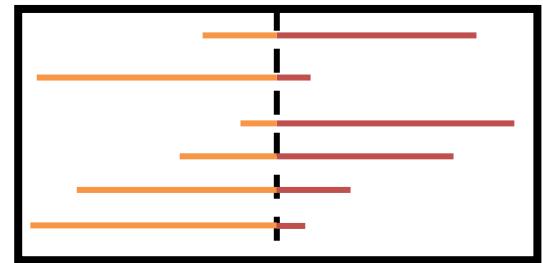
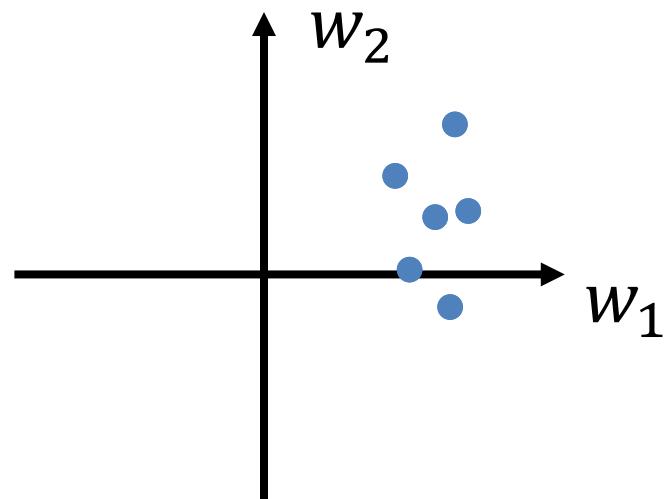


Maximum volume removal

Posterior $P(w | \mathcal{C})$



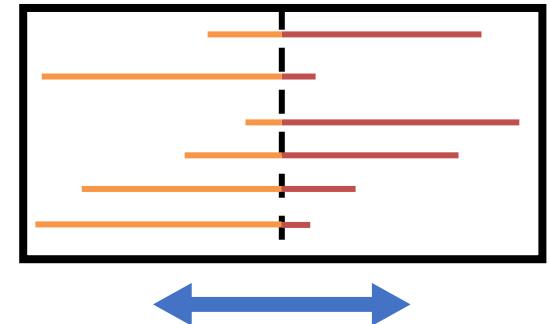
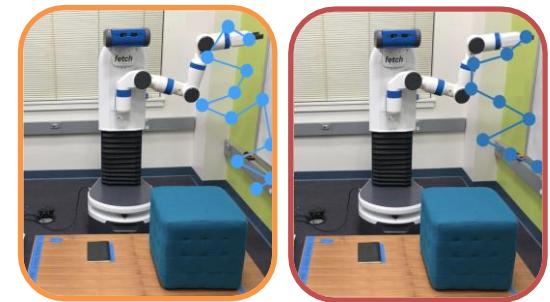
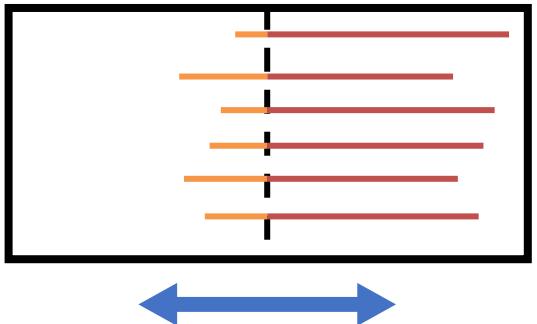
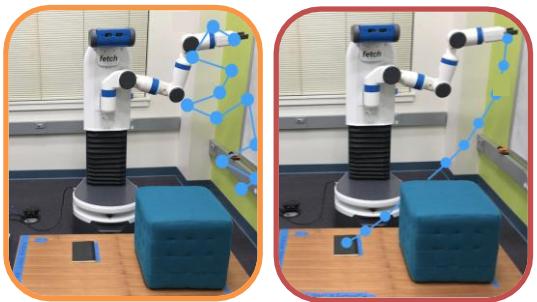
User Choice



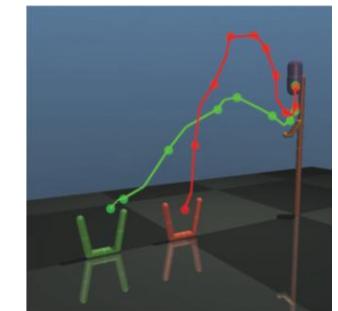
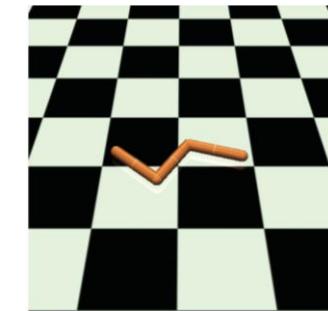
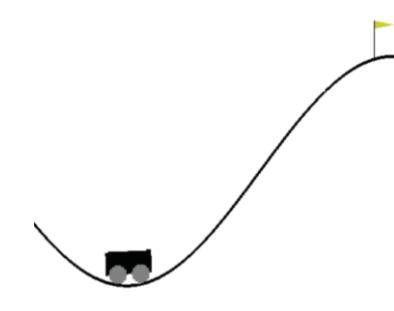
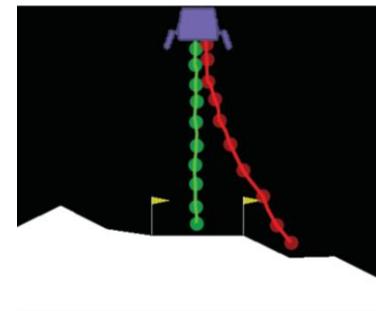
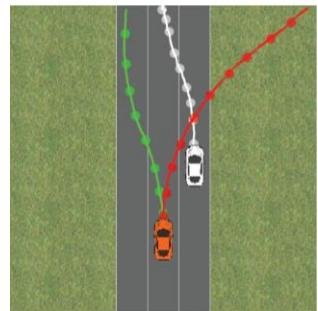
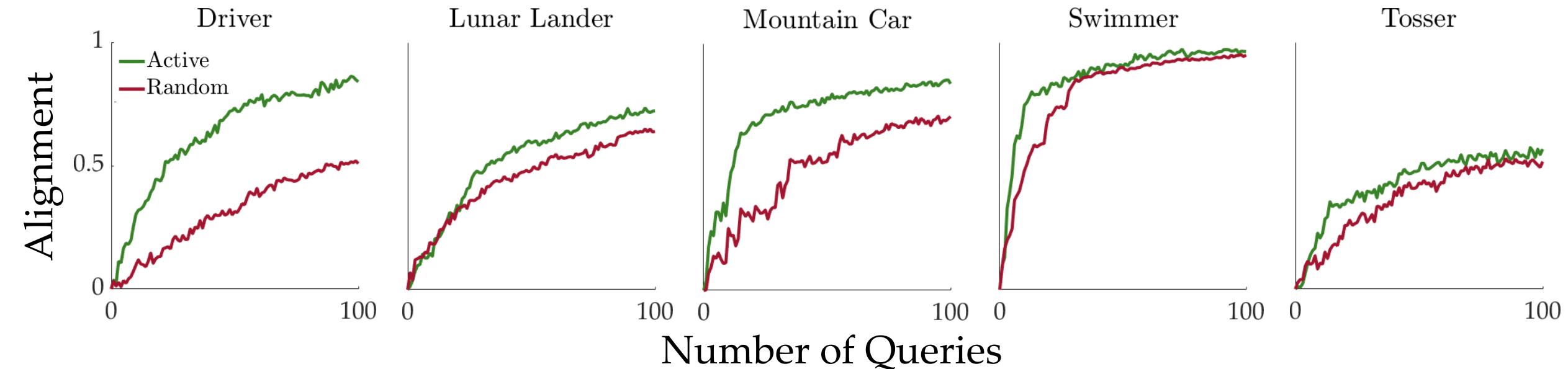
User Choice

Maximum volume removal

Posterior $P(w | \mathcal{C})$

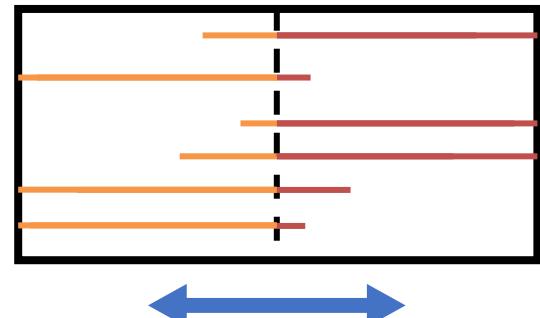
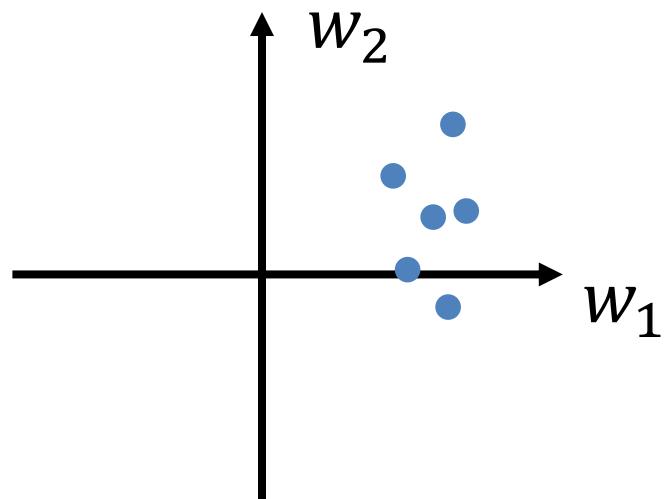
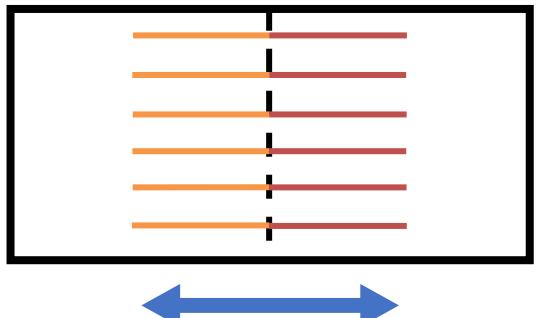


Active vs. random querying



Maximum volume removal

Posterior $P(w | \mathcal{C})$



Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q; w | \mathcal{C}, \xi_A, \xi_B)$$

The diagram illustrates the components of the mutual information expression. The top part shows the formula $\max_{\xi_A, \xi_B} I(q; w | \mathcal{C}, \xi_A, \xi_B)$. Below the formula, three labels are positioned: "User response" under q , "Dataset" under w , and "Query" under \mathcal{C} . Arrows point from "User response" to q and from "Dataset" to w . A bracket under ξ_A and ξ_B points to \mathcal{C} .

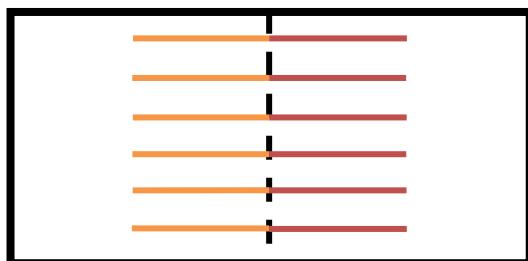
Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w | \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q | \mathcal{C}, \xi_A, \xi_B) - H(q | \mathcal{C}, \xi_A, \xi_B, w)$$



Model
Uncertainty



User Choice



User
Uncertainty



User Choice

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} -\mathbb{E}_{q \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \mathcal{C}, \xi_A, \xi_B)] + \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \mathcal{C}, \xi_A, \xi_B, w)]$$

No w here!

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} -\mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \mathcal{C}, \xi_A, \xi_B)] + \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \mathcal{C}, \xi_A, \xi_B, w)]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \mathcal{C}, \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\begin{aligned} \max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} & \left[\log P(q \mid \xi_A, \xi_B, w) - \log \int \underline{P(q, w' \mid \mathcal{C}, \xi_A, \xi_B) dw'} \right] \\ & P(w' \mid \mathcal{C}, \xi_A, \xi_B) P(q \mid \mathcal{C}, \xi_A, \xi_B, w') \\ & = P(w' \mid \mathcal{C}) P(q \mid \xi_A, \xi_B, w') \end{aligned}$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log \int P(w' \mid \mathcal{C}) P(q \mid \xi_A, \xi_B, w') dw']$$

This is an expectation over $w' \mid \mathcal{C}$

Take samples from $w' \mid \mathcal{C}$ to compute.

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} \left[\log P(q \mid \xi_A, \xi_B, w) - \log \frac{1}{|\Omega|} \sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w') \right]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} \left[\log P(q \mid \xi_A, \xi_B, w) - \log \sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w') \right]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

$$\begin{aligned} P(q, w \mid \mathcal{C}, \xi_A, \xi_B) &= P(w \mid \mathcal{C}, \xi_A, \xi_B) P(q \mid \mathcal{C}, \xi_A, \xi_B, w) \\ &= P(w \mid \mathcal{C}) P(q \mid \xi_A, \xi_B, w) \end{aligned}$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \xi_A, \xi_B} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

$$\max_{\xi_A, \xi_B} \frac{1}{|\Omega|} \sum_{w \in \Omega} \mathbb{E}_{q \mid \xi_A, \xi_B, w} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \xi_A, \xi_B} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

$$\max_{\xi_A, \xi_B} \sum_{w \in \Omega} \mathbb{E}_{q \mid \xi_A, \xi_B, w} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

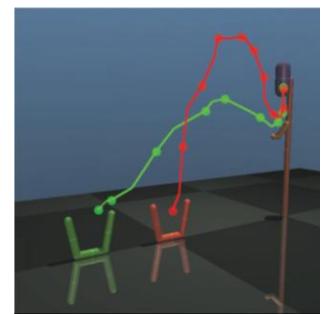
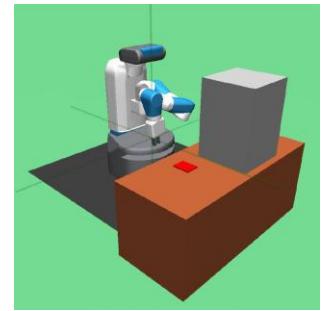
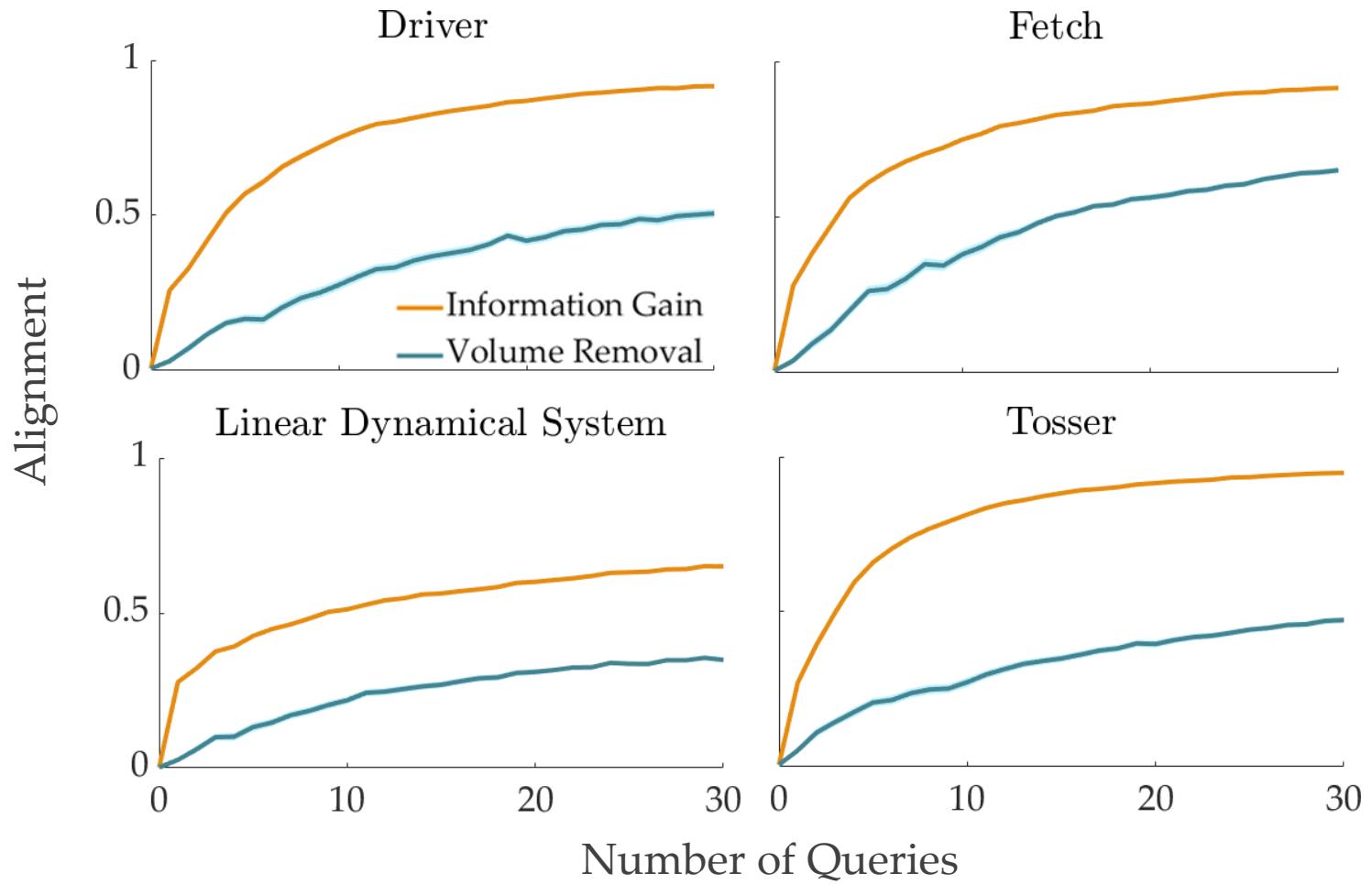
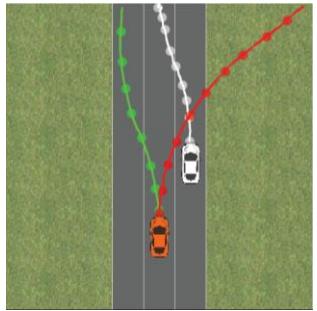
$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \mathcal{C}, \xi_A, \xi_B} [\log P(q \mid \xi_A, \xi_B, w) - \log P(q \mid \mathcal{C}, \xi_A, \xi_B)]$$

$$\max_{\xi_A, \xi_B} \mathbb{E}_{q, w \mid \xi_A, \xi_B} \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

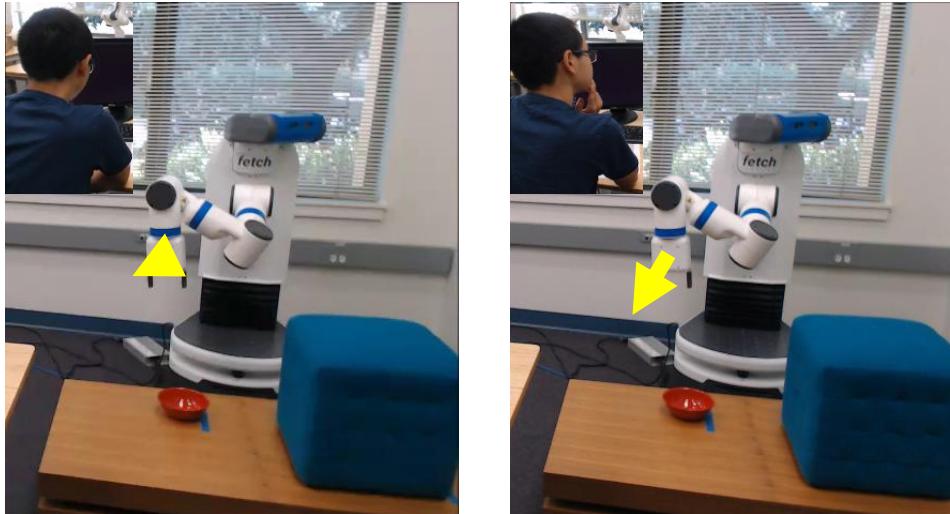
$$\max_{\xi_A, \xi_B} \sum_{w \in \Omega} \sum_q P(q \mid \xi_A, \xi_B, w) \left[\log \frac{P(q \mid \xi_A, \xi_B, w)}{\sum_{w' \in \Omega} P(q \mid \xi_A, \xi_B, w')} \right]$$

Mutual information maximization



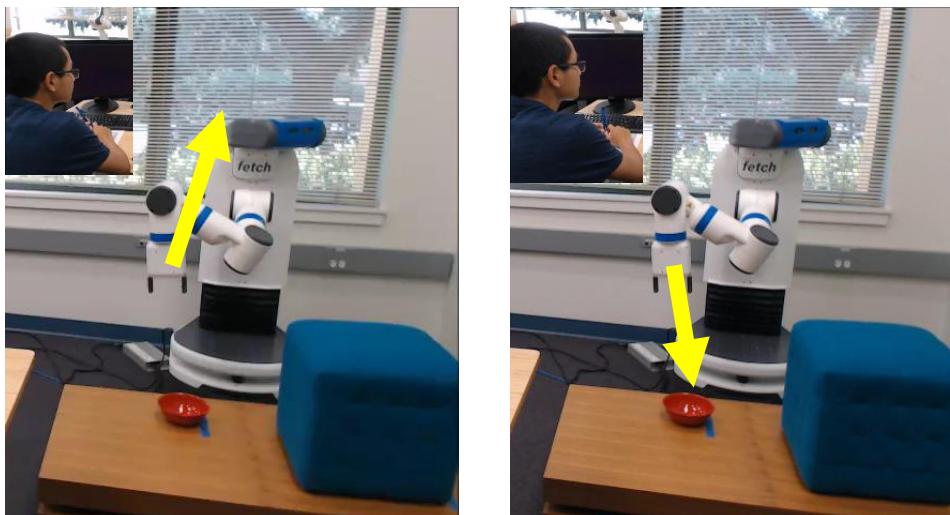
Similar
Trajectories

Volume Removal

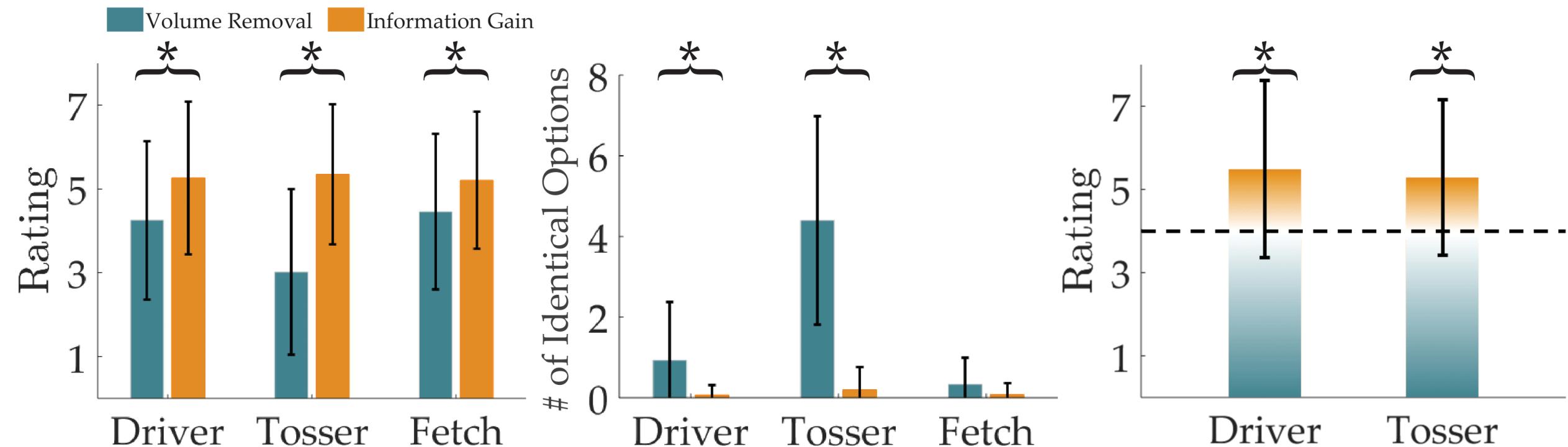


More
Distinguishable
Query

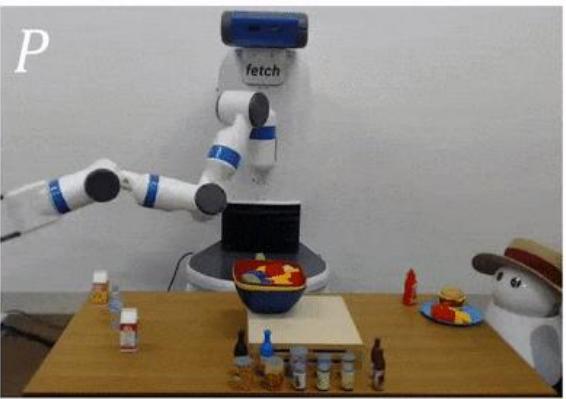
Information Gain



Mutual information maximization



P



Richer forms of comparative feedback enables learning more general reward functions.

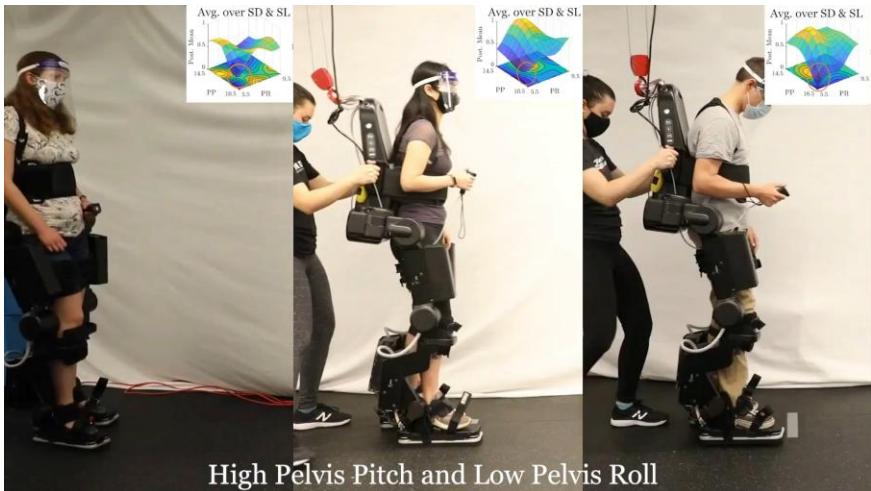
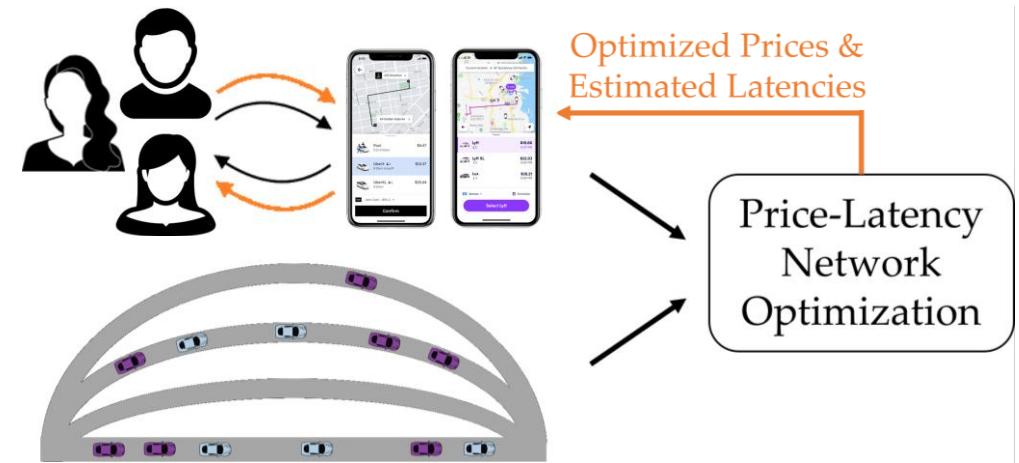
*Basu, Biyik, He, Singhal, Sadigh. IROS'19.
Wilde*, Biyik*, Sadigh, Smith. CoRL'21.*

Data from Uber/Lyft are comparative feedback.
We can use them to optimize traffic routing.

Biyik, Lazar, Sadigh, Pedarsani. CDC'19.

Biyik, Lazar*, Pedarsani, Sadigh. TCNS'21.*

Beliaev, Biyik, Lazar, Wang, Sadigh, Pedarsani. ICCPS'21.



Active preference-based Gaussian process regression enables learning complex rewards with small amounts of data.

Biyik, Huynh*, Kochenderfer, Sadigh. RSS'20.*

Li, Tucker, Biyik, Novoseller, Burdick, Sui, Sadigh, Yue, Ames. ICRA'21.

Biyik, Huynh, Kochenderfer, Sadigh. IJRR'23.

Incorporating comparisons

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D}, \mathcal{C})$$

$$P(w \mid \mathcal{D}, \mathcal{C}) \propto P(w) P(\mathcal{D} \mid w) P(\mathcal{C} \mid w)$$

$$\begin{aligned} &= P(w) \prod_{i=1}^L P(\xi_i \mid w) \prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)}) \\ &\propto P(w) \prod_{i=1}^L \exp f_w(\xi_i) \prod_{i=1}^N \frac{\exp f_w(\xi_{q^{(i)}}^{(i)})}{\exp f_w(\xi_{q^{(i)}}^{(i)}) + \exp f_w(\xi_{\neg q^{(i)}}^{(i)})} \end{aligned}$$

Other types of human feedback

Feedback	Constraint	Probabilistic
Comparisons	$r(\xi_1) \geq r(\xi_2)$	$\mathbb{P}(\xi_1 r, \mathcal{C}) = \frac{\exp(\beta \cdot r(\xi_1))}{\exp(\beta \cdot r(\xi_1)) + \exp(\beta \cdot r(\xi_2))}$
Demonstrations	$r(\xi_D) \geq r(\xi) \quad \forall \xi \in \Xi$	$\mathbb{P}(\xi_D r, \Xi) = \frac{\exp(\beta \cdot r(\xi_D))}{\sum_{\xi \in \Xi} \exp(\beta \cdot r(\xi))}$
Corrections	$r(\xi_R + A^{-1}\Delta q) \geq r(\xi_R + A^{-1}\Delta q') \quad \forall \Delta q' \in Q - Q$	$\mathbb{P}(\Delta q' r, Q - Q) = \frac{\exp(\beta \cdot r(\xi_R + A^{-1}\Delta q))}{\sum_{\Delta q \in Q - Q} \exp(\beta \cdot r(\xi_R + A^{-1}\Delta q))}$
Improvement	$r(\xi_{\text{improved}}) \geq r(\xi_R)$	$\mathbb{P}(\xi_{\text{improved}} r, \mathcal{C}) = \frac{\exp(\beta \cdot r(\xi_{\text{improved}}))}{\exp(\beta \cdot r(\xi_{\text{improved}})) + \exp(\beta \cdot r(\xi_R))}$
Off	$r(\xi_R^{0:t} \xi^t \dots \xi^t) \geq r(\xi_R)$	$\mathbb{P}(\text{off} r, \mathcal{C}) = \frac{\exp(\beta \cdot r(\xi_R^{0:t} \xi^t \dots \xi^t))}{\exp(\beta \cdot r(\xi_R^{0:t} \xi^t \dots \xi^t)) + \exp(\beta \cdot r(\xi_R))}$
Language	$\mathbb{E}_{\xi \sim \text{Unif}(G(\lambda^*))}[r(\xi)] \geq \mathbb{E}_{\xi \sim \text{Unif}(G(\lambda))}[r(\xi)] \quad \forall \lambda \in \Lambda$	$\mathbb{P}(\lambda^* r, \Lambda) = \frac{\exp(\beta \cdot \mathbb{E}_{\xi \sim \text{Unif}(G(\lambda^*))}[r(\xi)])}{\sum_{\lambda \in \Lambda} \exp(\beta \cdot \mathbb{E}_{\xi \sim \text{Unif}(G(\lambda))}[r(\xi)])}$
Proxy Rewards	$\mathbb{E}_{\tilde{\xi} \sim \pi(\tilde{\xi} \tilde{r})}[r(\tilde{\xi})] \geq \mathbb{E}_{\tilde{\xi} \sim \pi(\tilde{\xi} c)}[r(\tilde{\xi})] \quad \forall c \in \tilde{\mathcal{R}}$	$\mathbb{P}(\tilde{r} r, \tilde{\mathcal{R}}) = \frac{\exp(\beta \cdot \mathbb{E}_{\tilde{\xi} \sim \pi(\tilde{\xi} \tilde{r})}[r(\tilde{\xi})])}{\sum_{c \in \tilde{\mathcal{R}}} \exp(\beta \cdot \mathbb{E}_{\tilde{\xi} \sim \pi(\tilde{\xi} c)}[r(\tilde{\xi})])}$
Reward/Punish	$r(\xi_R) \geq r(\xi_{\text{expected}})$	$\mathbb{P}(+1 r, \mathcal{C}) = \frac{\exp(\beta \cdot r(\xi_R))}{\exp(\beta \cdot r(\xi_R)) + \exp(\beta \cdot r(\xi_{\text{expected}}))}$
Initial state	$\mathbb{E}_{\xi \sim \psi(s^*)}[r(s^*)] \geq \mathbb{E}_{\xi \sim \psi(s)}[r(s)] \quad \forall s \in \mathcal{S}$	$\mathbb{P}(s^* r, \mathcal{S}) = \frac{\exp(\beta \cdot \mathbb{E}_{\xi \sim \psi(s^*)}[r(\xi)])}{\sum_{s \in \mathcal{S}} \exp(\beta \cdot \mathbb{E}_{\xi \sim \psi(s)}[r(\xi)])}$
Meta-choice	$\mathbb{E}_{\xi \sim \psi(c_i)}[r(\xi)] \geq \mathbb{E}_{\xi \sim \psi(c_j)}[r(\xi)] \quad \forall j \in [n]$	$\mathbb{P}(c_i r, \mathcal{C}_0) = \frac{\exp(\beta_0 \cdot \mathbb{E}_{\xi \sim \psi_0(c_i)}[r(\xi)])}{\sum_{j \in [n]} \exp(\beta_0 \cdot \mathbb{E}_{\xi \sim \psi_0(c_j)}[r(\xi)])}$
Credit assignment	$r(\xi^*) \geq r(\xi) \quad \forall \xi \in \mathcal{C}$	$\mathbb{P}(\xi^* r, \mathcal{C}) = \frac{\exp(\beta \cdot r(\xi^*))}{\sum_{\xi \in \mathcal{C}} \exp(\beta \cdot r(\xi))}$

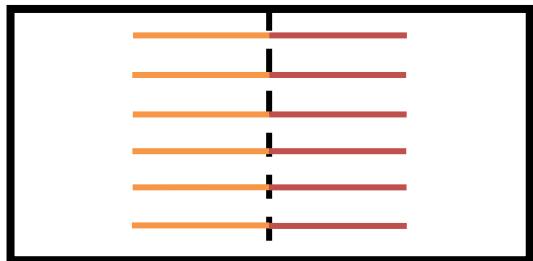
Today...

- Learning from human feedback
 - Pairwise comparisons
 - Reinforcement learning from human feedback (RLHF)
 - Comparative language feedback

Mutual information maximization

$$\max_{\xi_A, \xi_B} I(q ; w \mid \mathcal{C}, \xi_A, \xi_B)$$

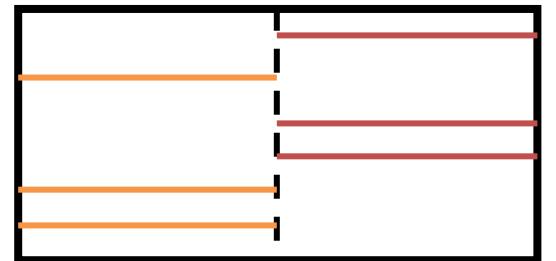
$$\max_{\xi_A, \xi_B} H(q \mid \mathcal{C}, \xi_A, \xi_B) - H(q \mid \mathcal{C}, \xi_A, \xi_B, w)$$



User Choice

Model
Uncertainty

User
Uncertainty



User Choice

Where do these trajectories
come from in the first place?

Incorporating comparisons

$$\underset{w}{\operatorname{argmax}} P(w \mid \mathcal{D}, \mathcal{C})$$

$$P(w \mid \mathcal{D}, \mathcal{C}) \propto P(w) P(\mathcal{D} \mid w) P(\mathcal{C} \mid w)$$

How do we solve this optimization problem?

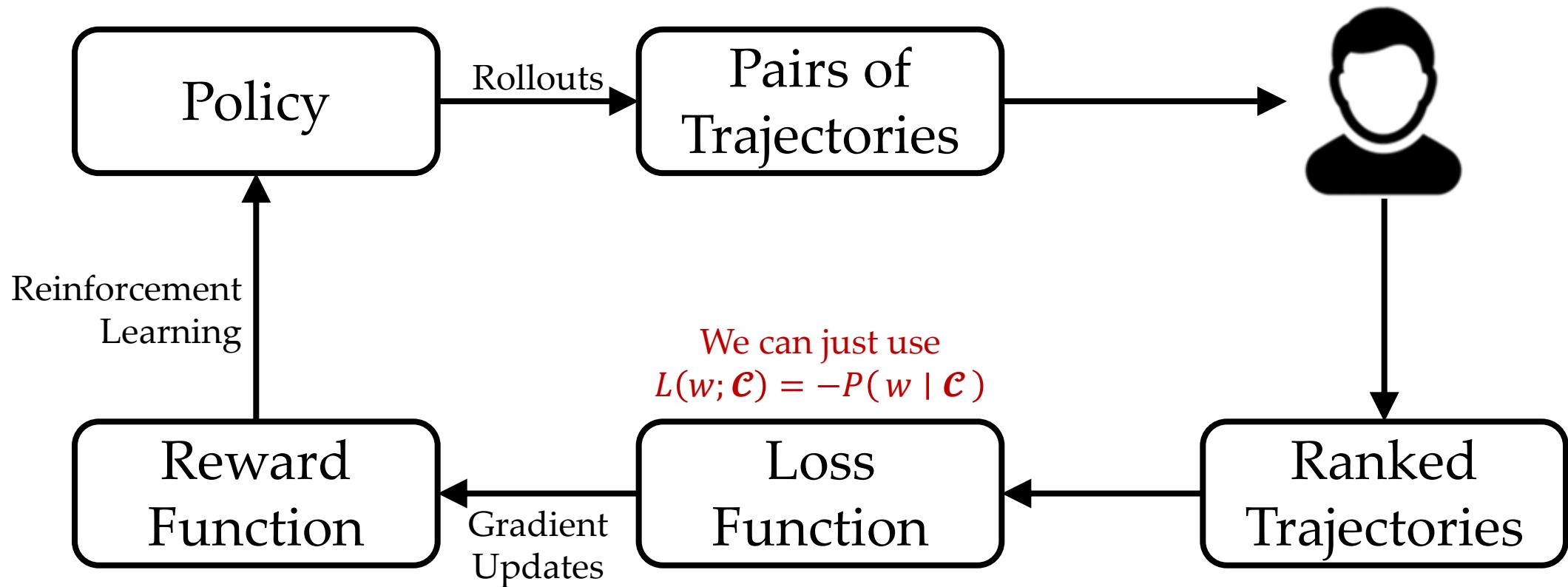
$$\begin{aligned} &= P(w) \prod_{i=1}^L P(\xi_i \mid w) \prod_{i=1}^N P(q^{(i)} \mid w, \xi_A^{(i)}, \xi_B^{(i)}) \\ &\propto P(w) \prod_{i=1}^L \exp f_w(\xi_i) \prod_{i=1}^N \frac{\exp f_w(\xi_{q^{(i)}}^{(i)})}{\exp f_w(\xi_{q^{(i)}}^{(i)}) + \exp f_w(\xi_{\neg q^{(i)}}^{(i)})} \end{aligned}$$

RLHF

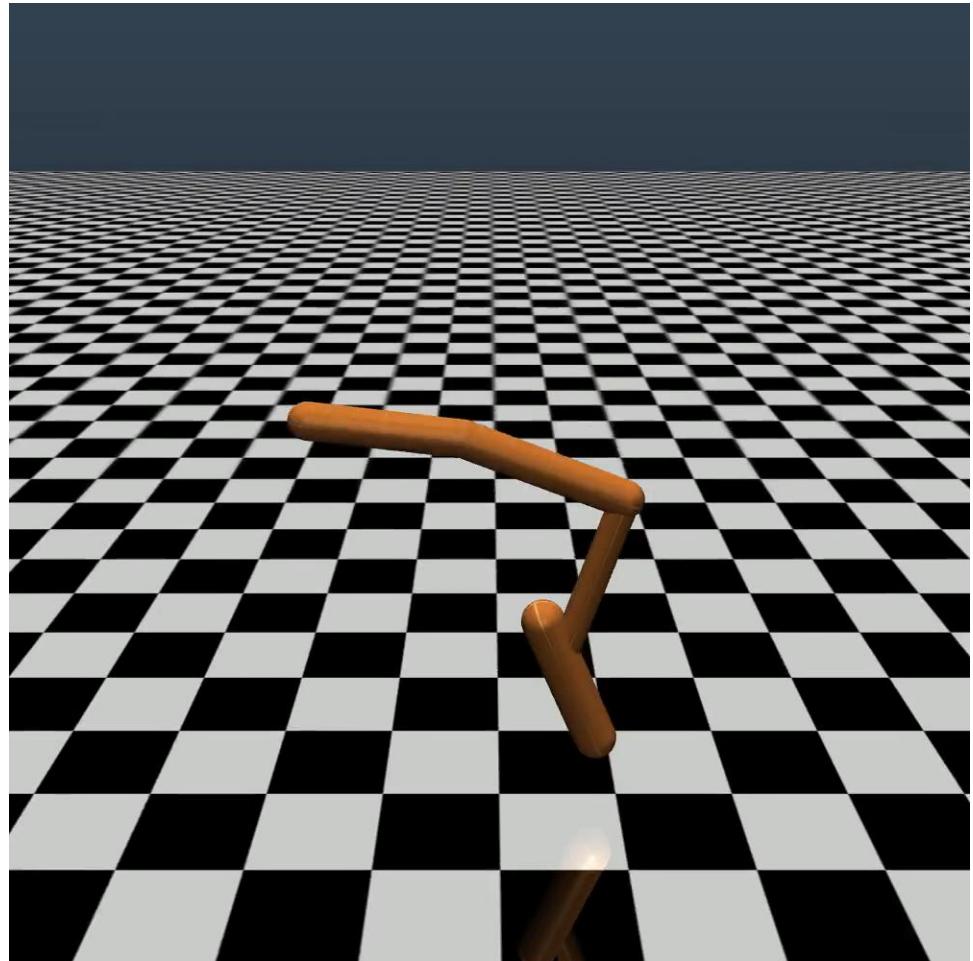
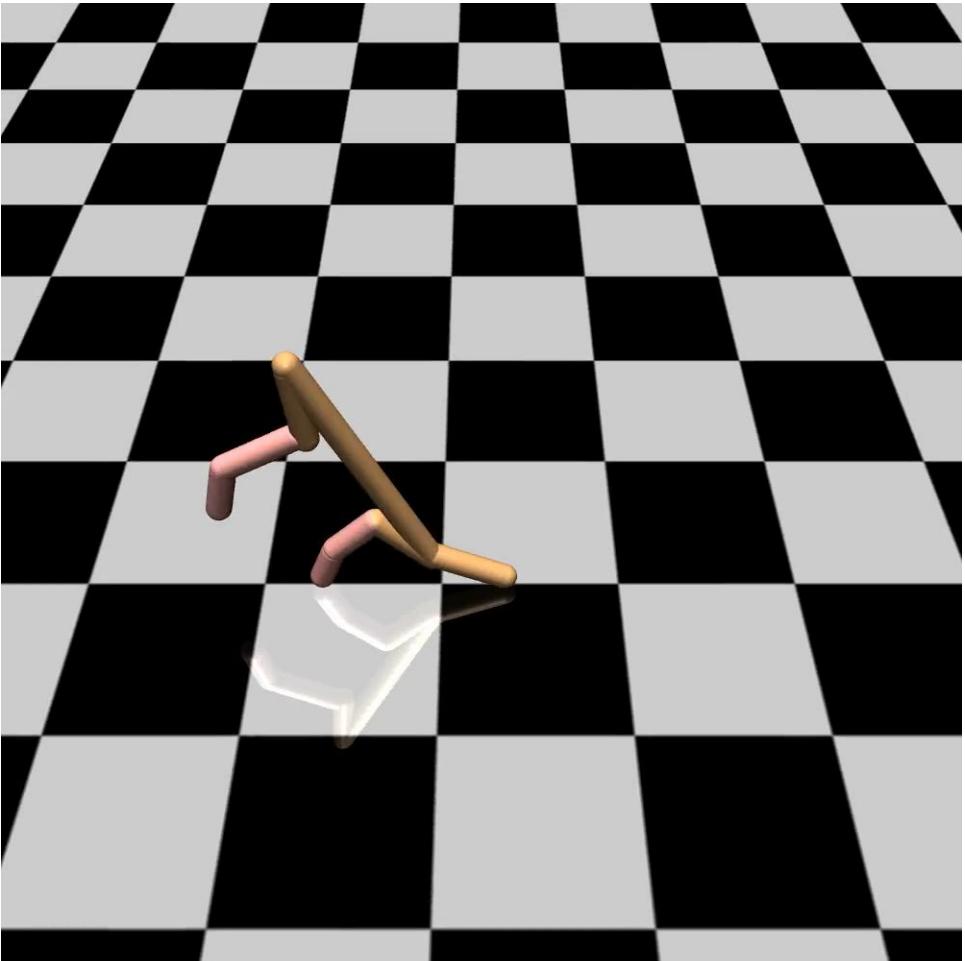
Two major changes to preference-based reward learning:

1. Instead of Bayesian learning, write a loss function and learn with gradient updates
2. After learning a reward, train a policy to generate new trajectories for the next iteration of reward learning

RLHF



RLHF



InstructGPT

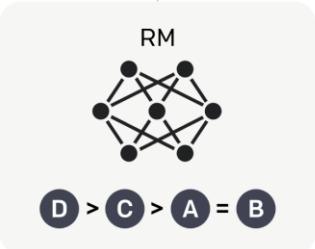
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



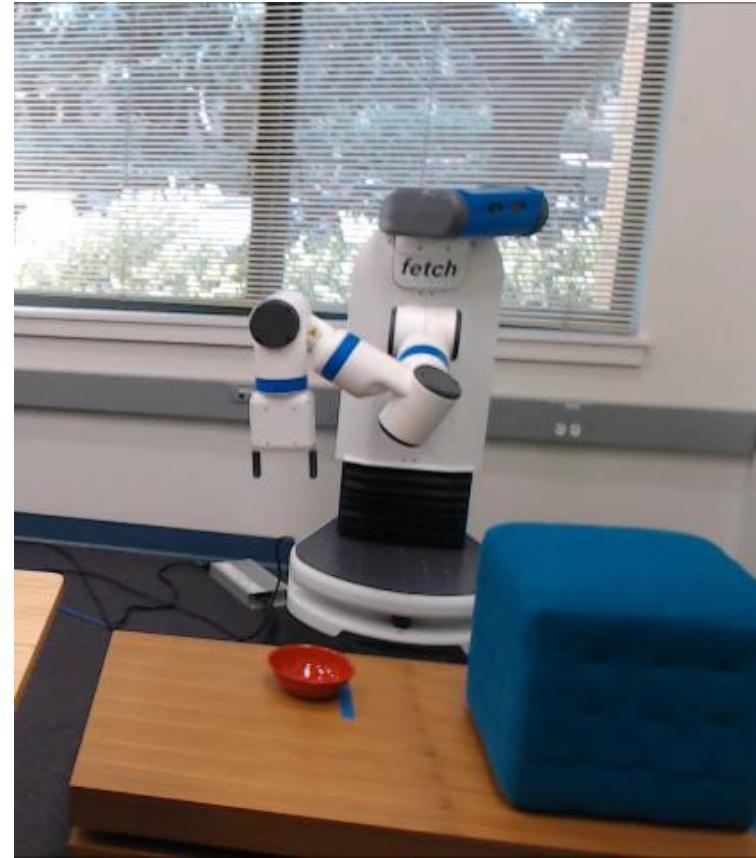
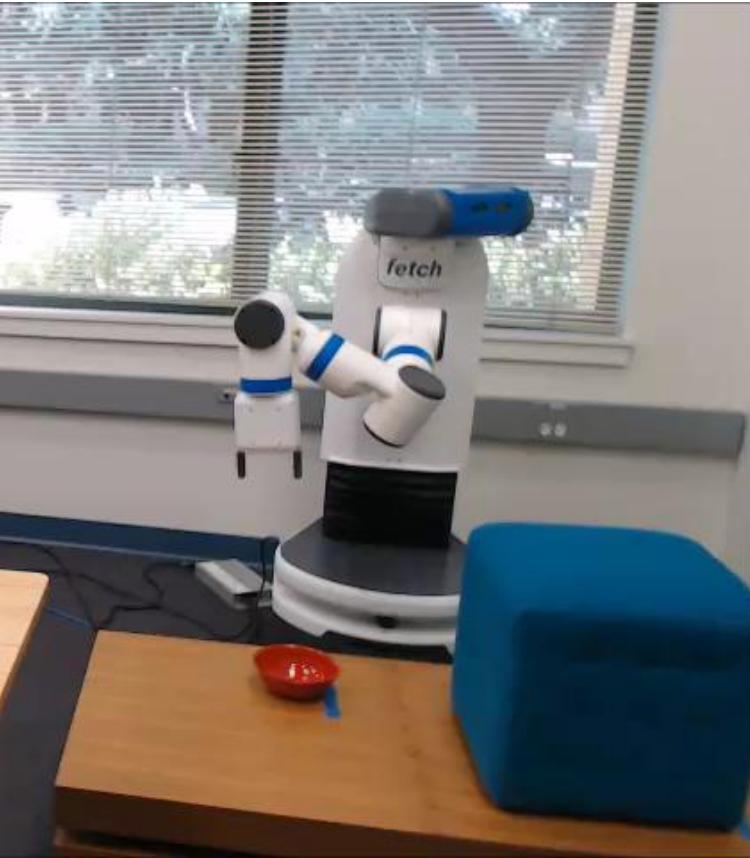
This data is used to train our reward model.



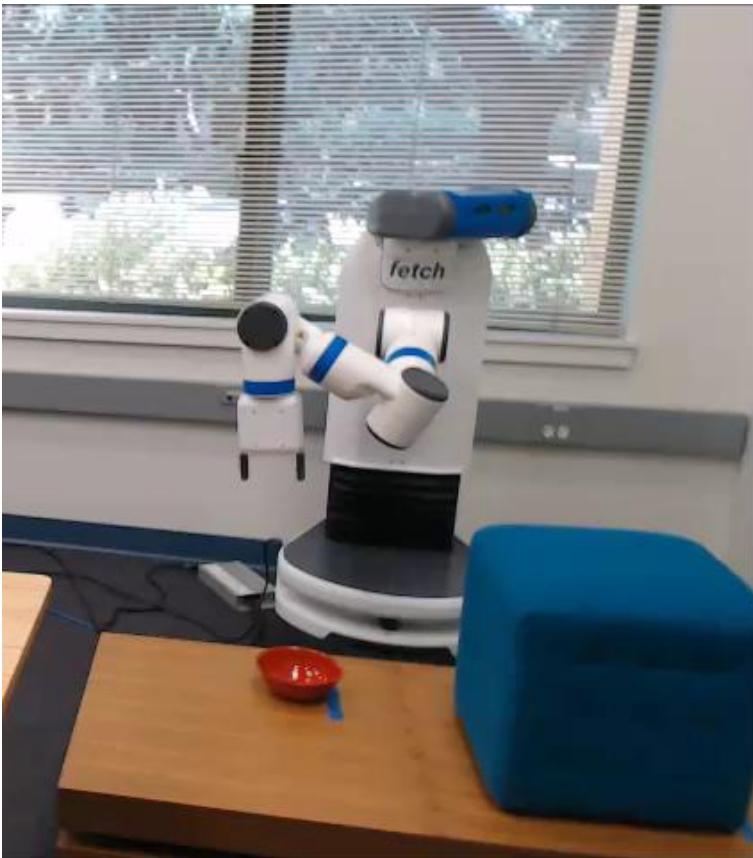
Today...

- Learning from human feedback
 - Pairwise comparisons
 - Reinforcement learning from human feedback (RLHF)
 - Comparative language feedback

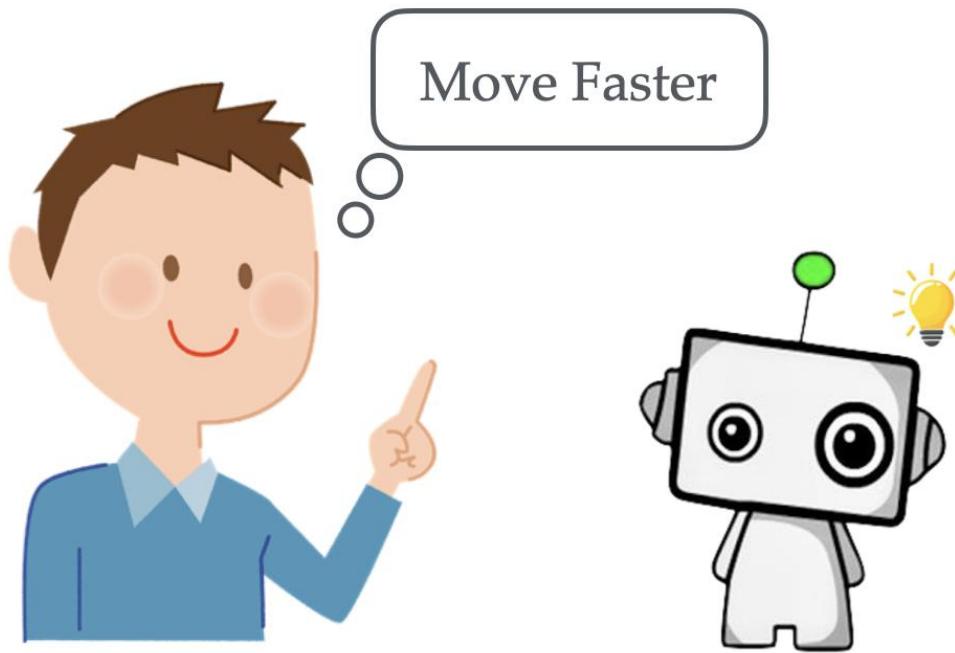
Comparisons take too long



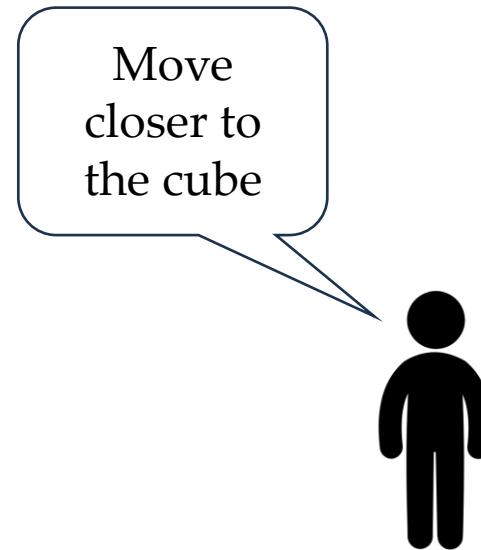
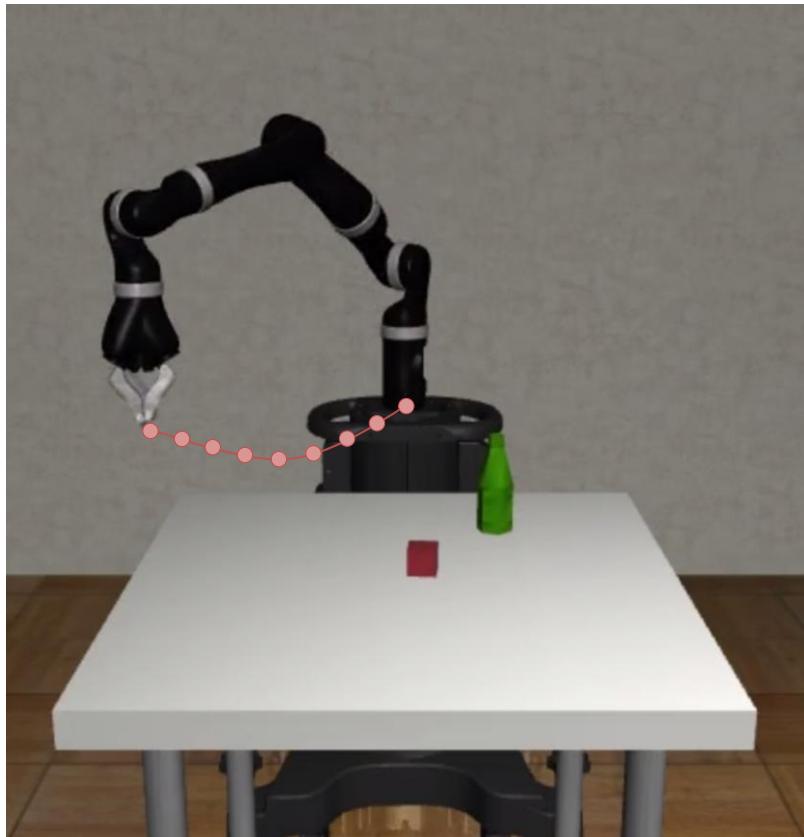
Why do you prefer that?



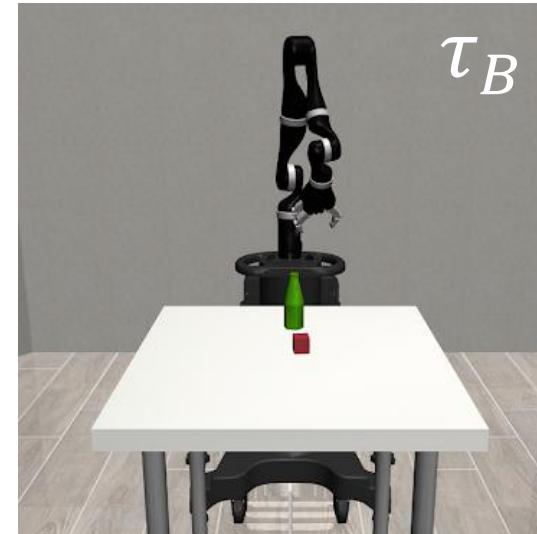
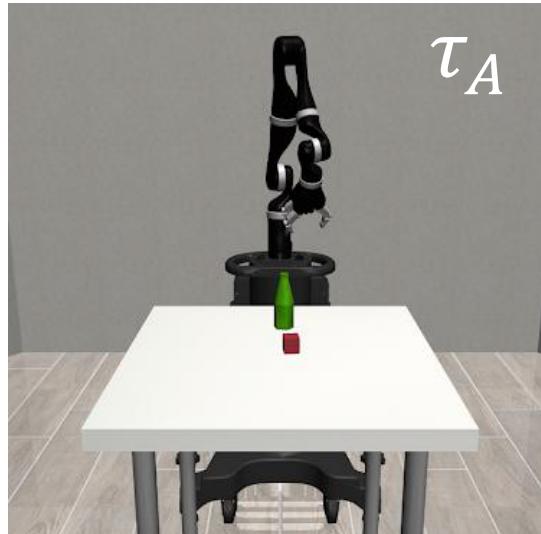
Comparative Language Feedback



Comparative Language Feedback



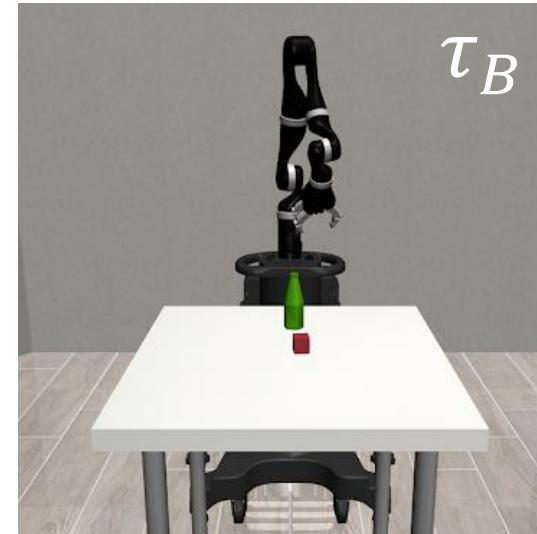
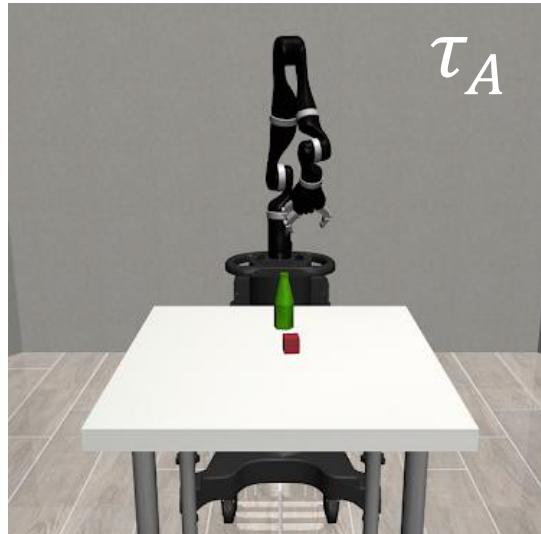
Shared Latent Space



τ^B moves faster than τ^A

Just describes a difference
between two trajectories!

Shared Latent Space



τ^B moves faster than τ^A



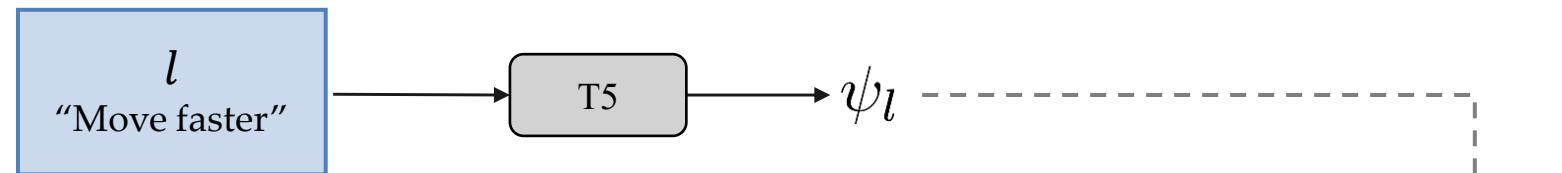
τ^B moves more quickly than τ^A



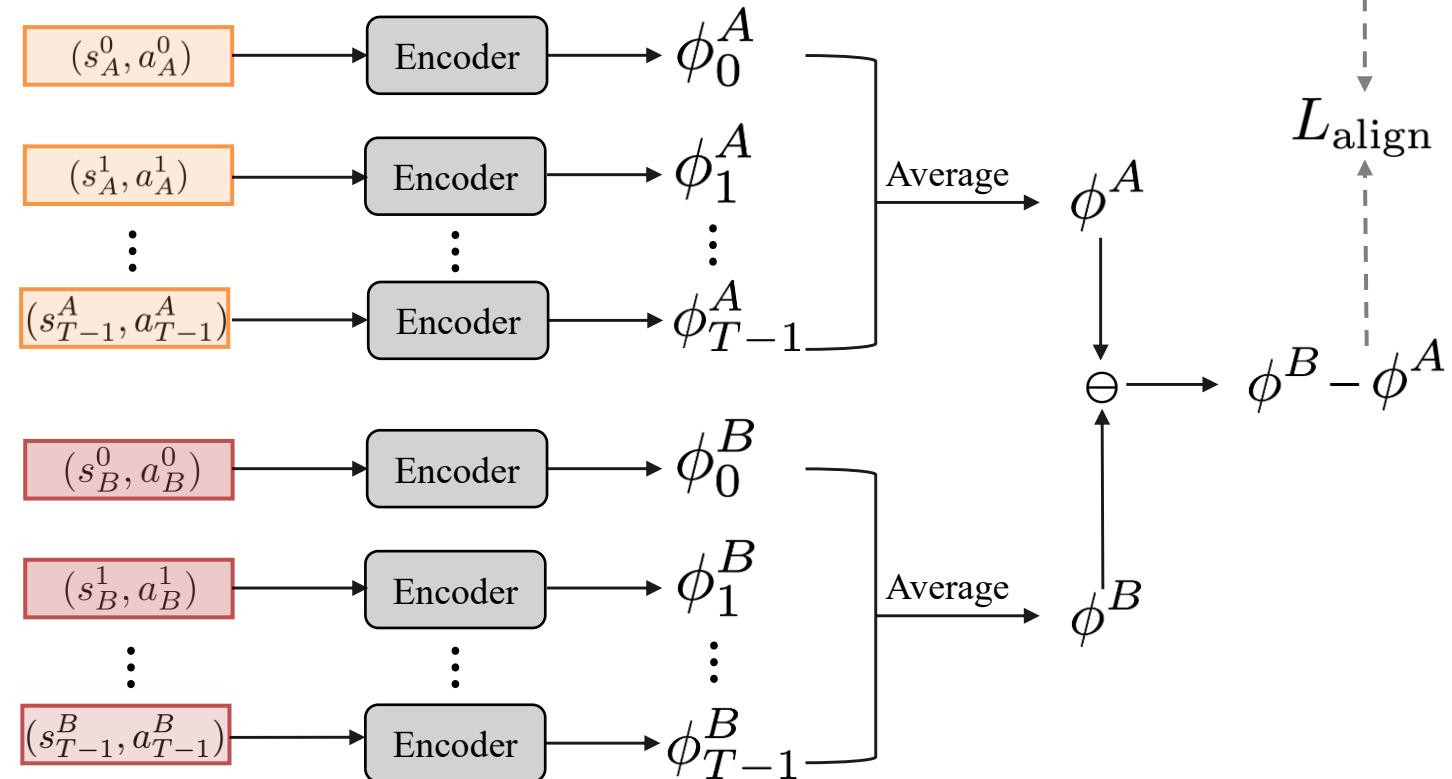
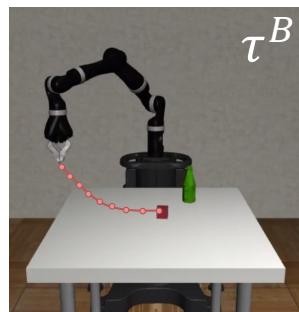
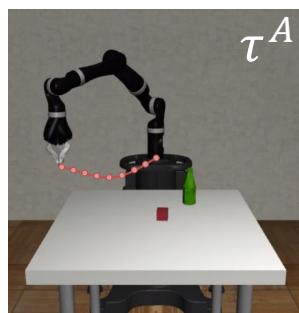
τ^A moves slower than τ^B

Just describes a difference
between two trajectories!

Language feedback



Trajectory pair



Accuracy

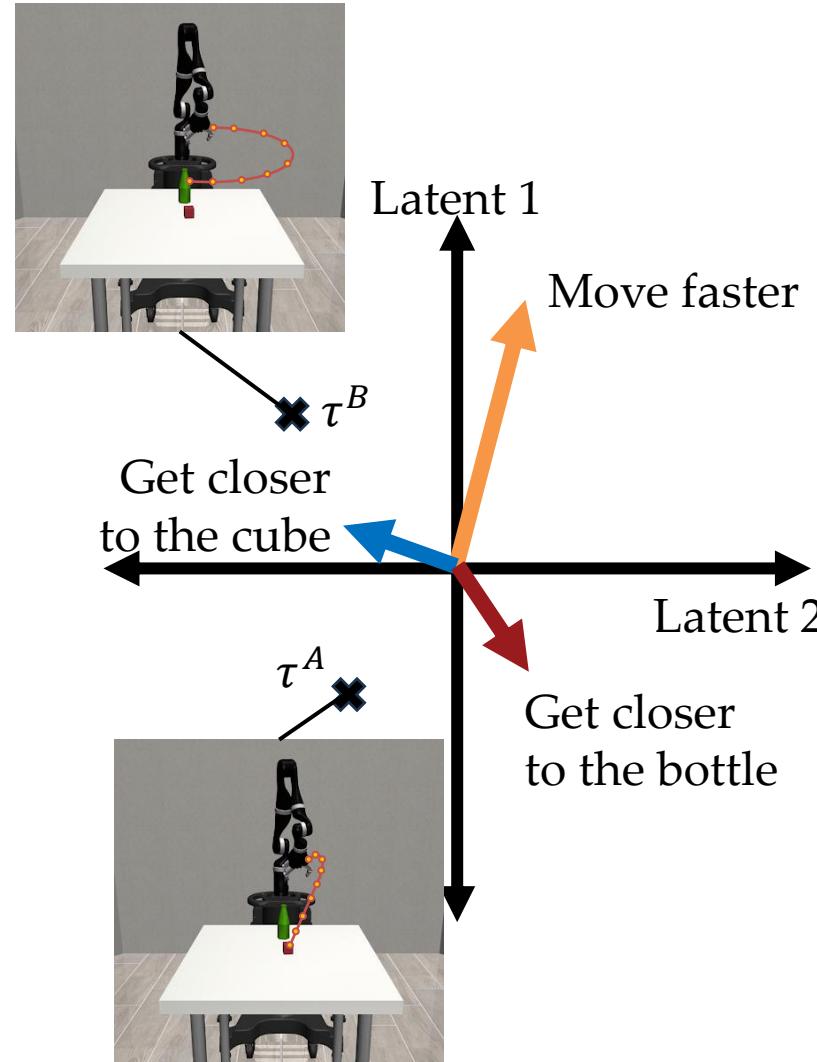
Robosuite

84.9%

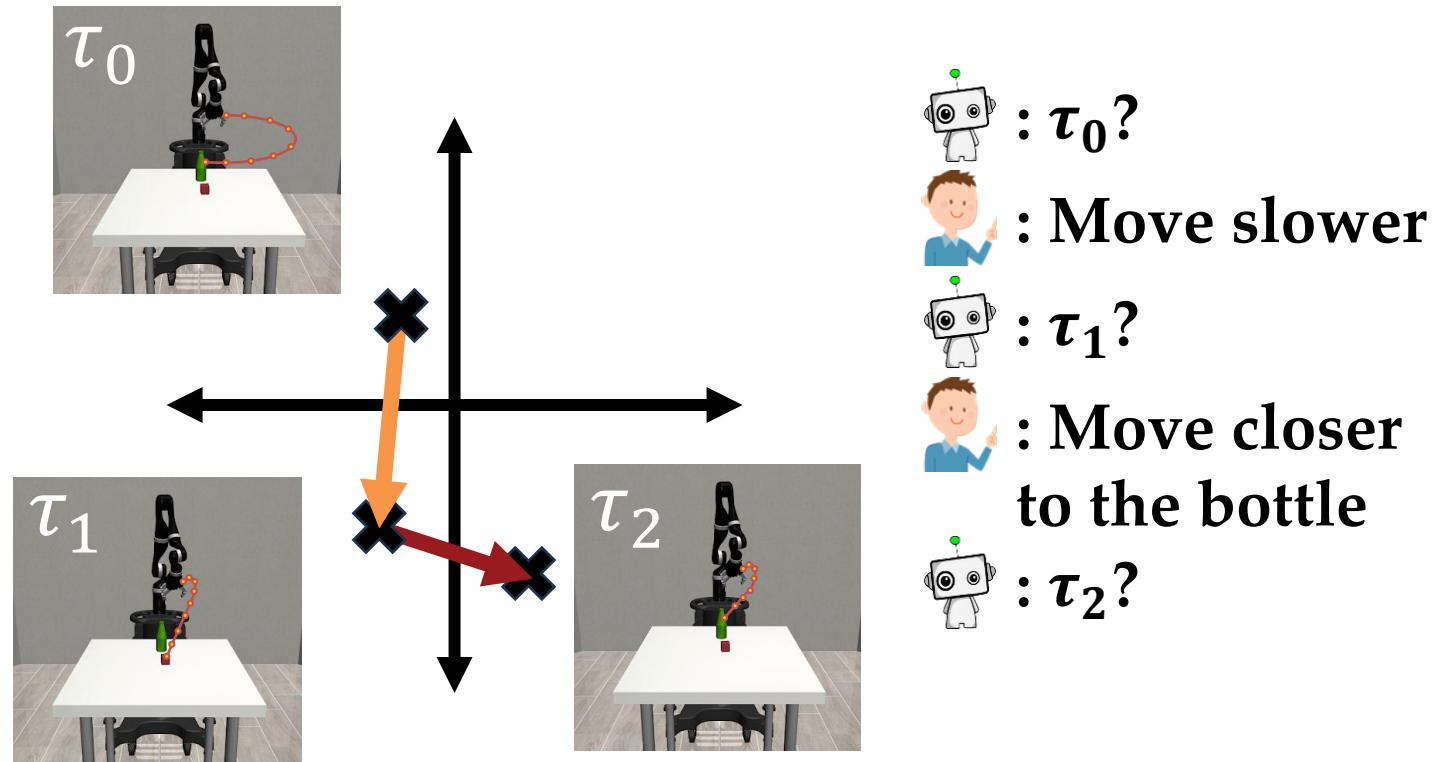
Meta-World

82.9%

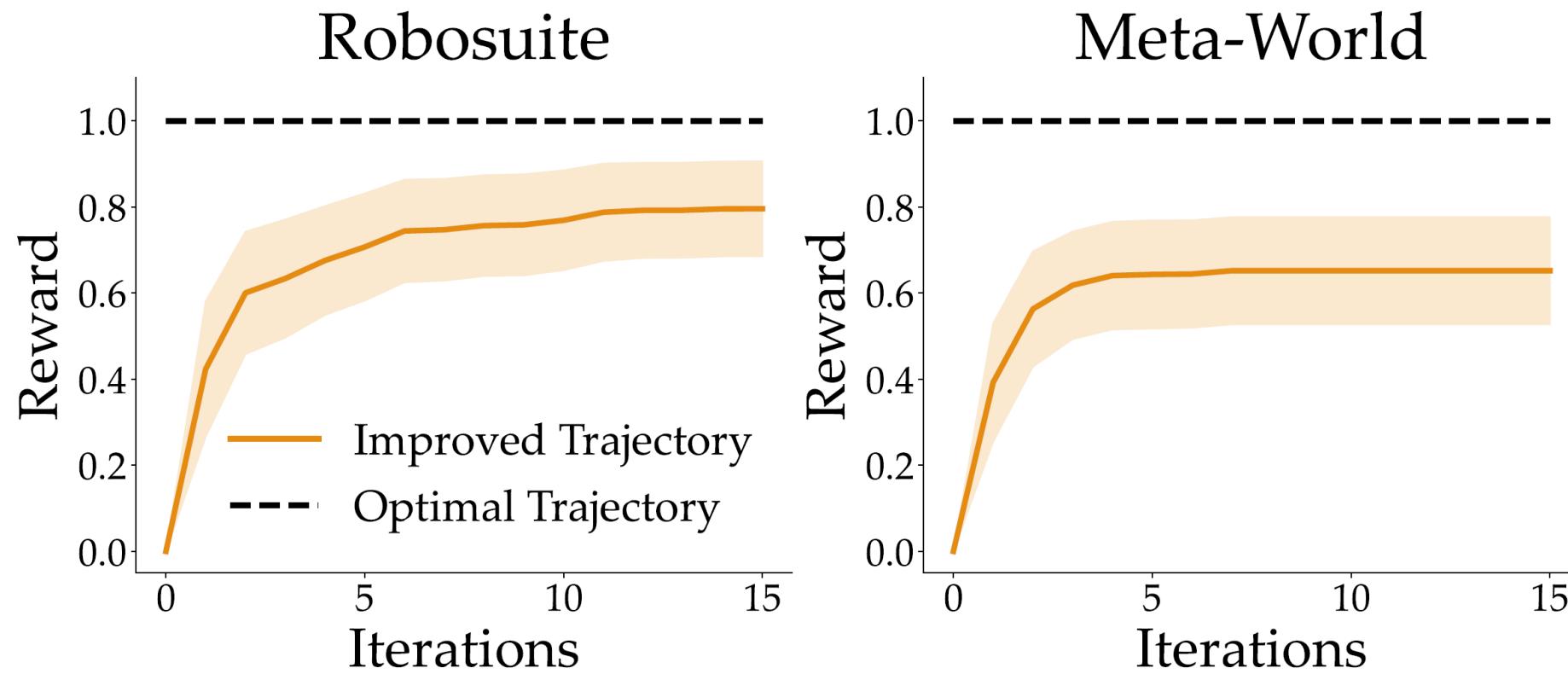
Shared Latent Space



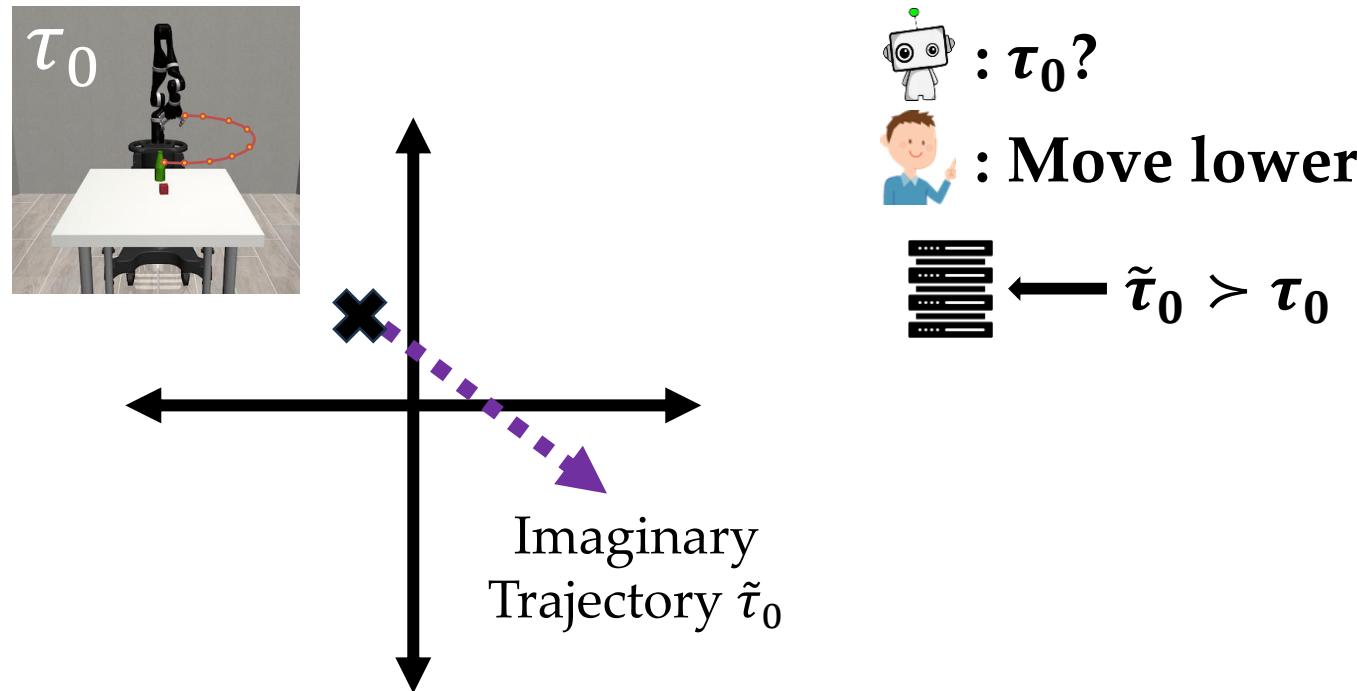
Using the Latent Space



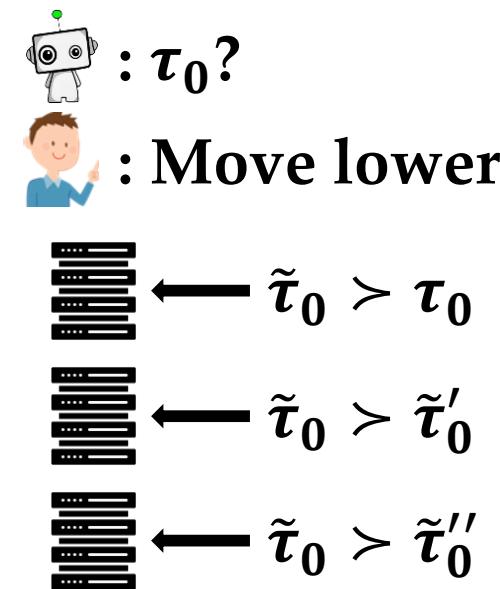
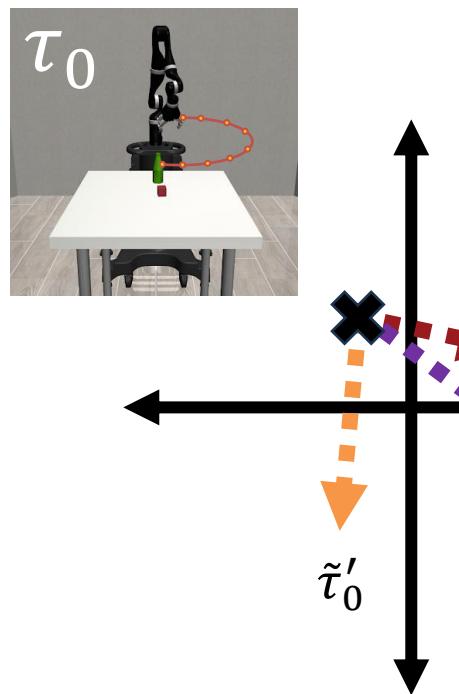
Using the Latent Space



Using the Latent Space

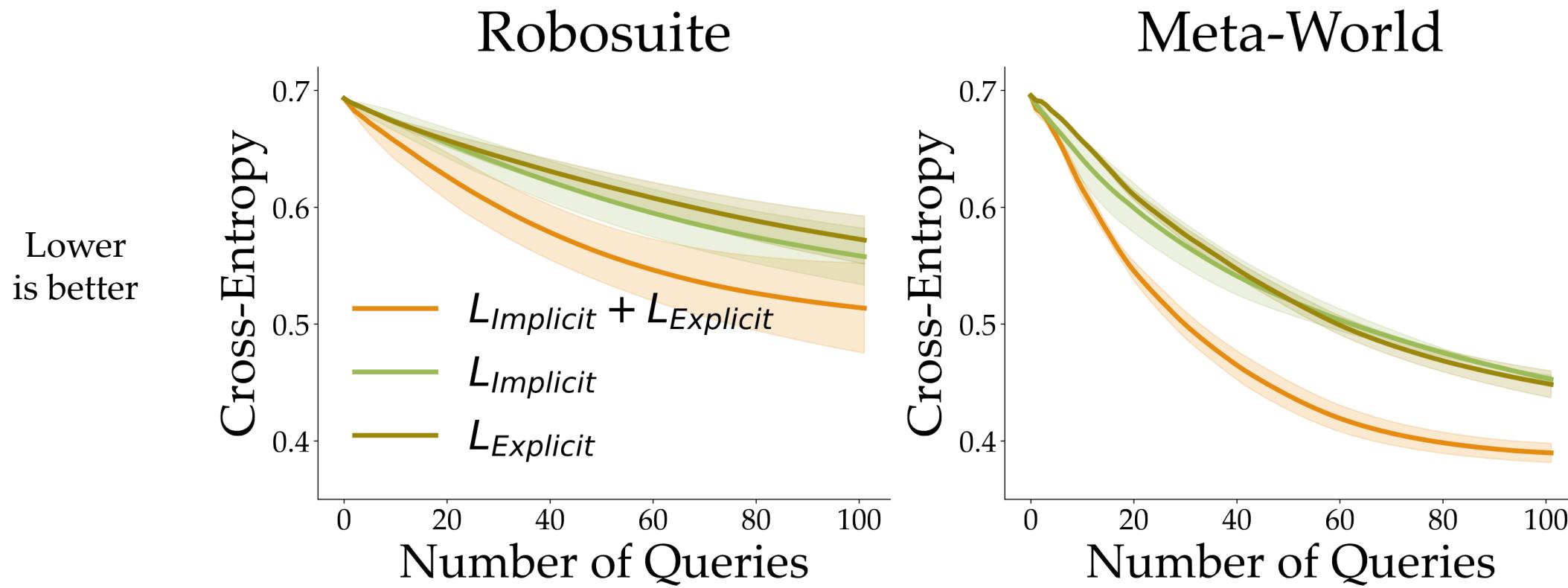


Using the Latent Space



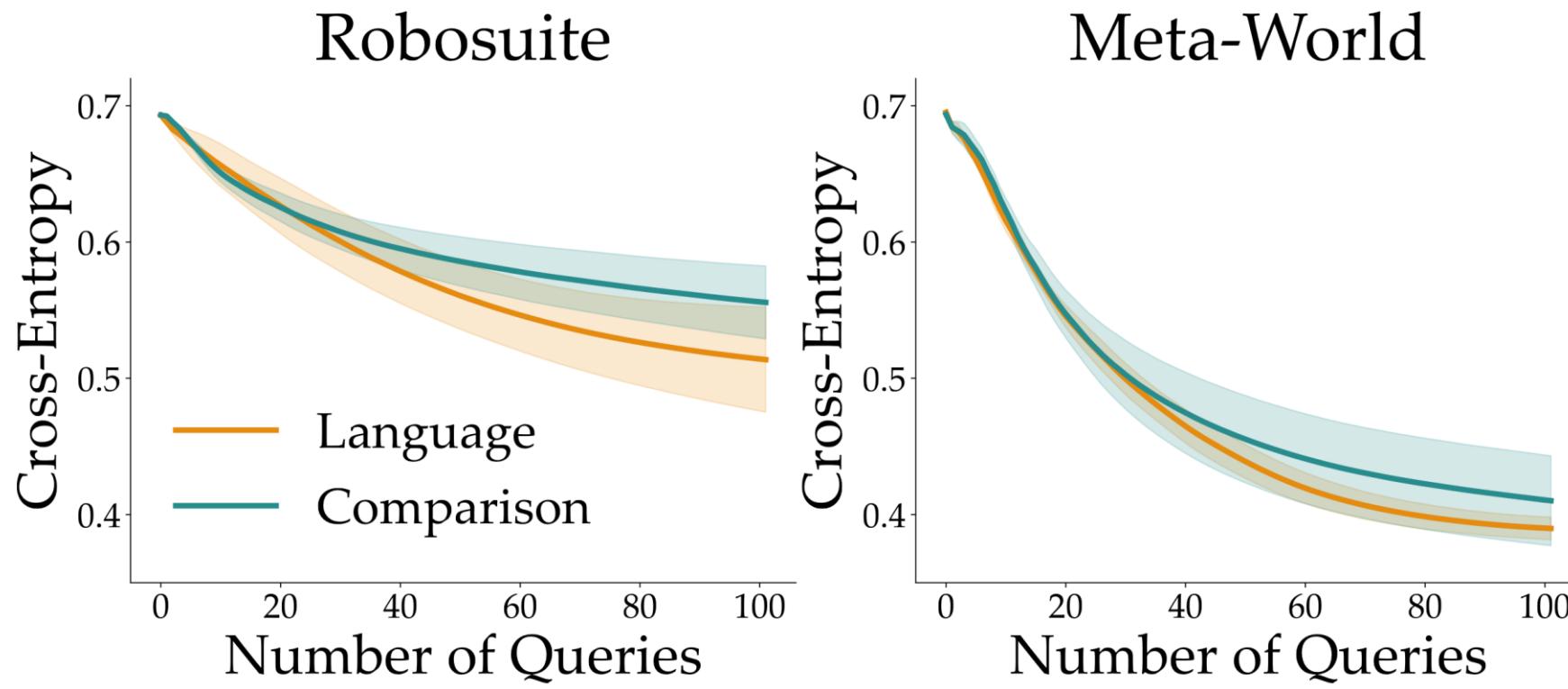
The user intentionally chose to say “move lower”

Using the Latent Space



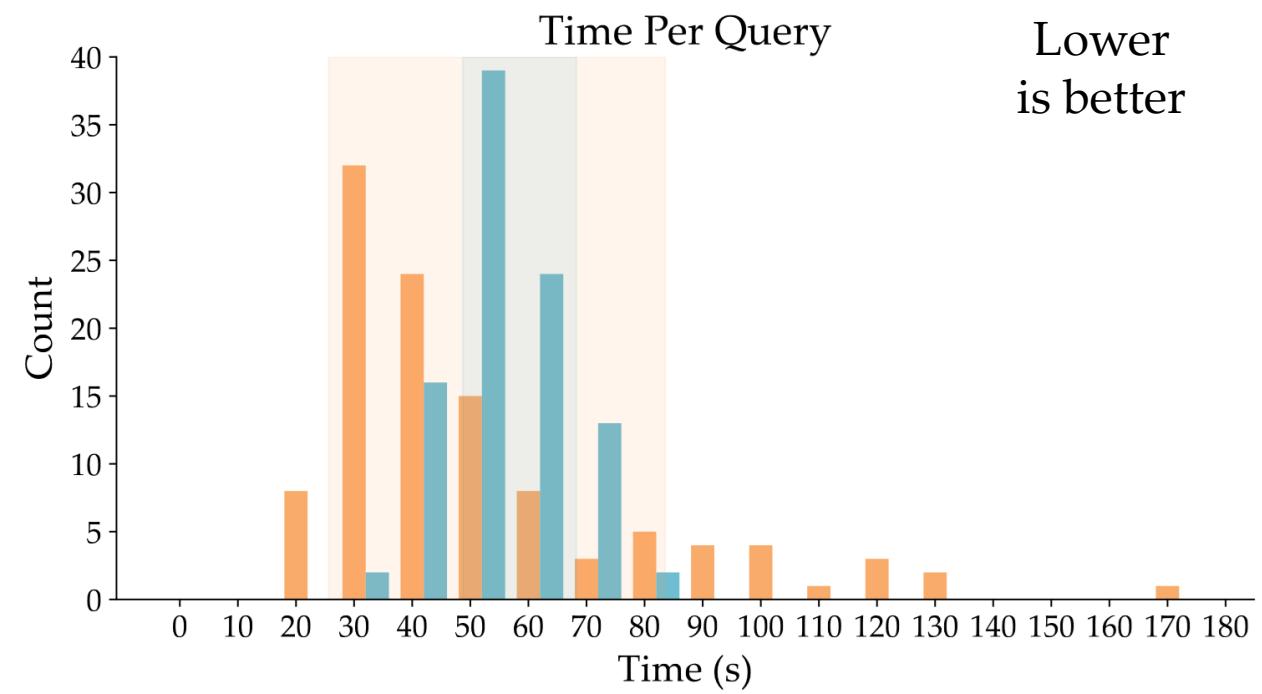
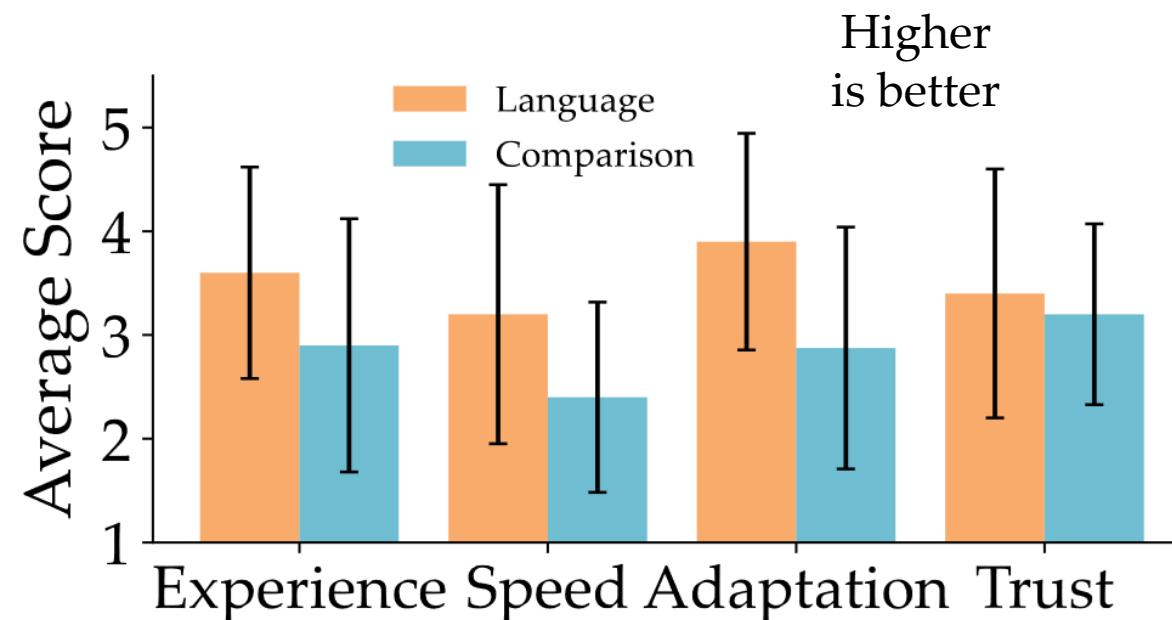
Using the Latent Space

Lower
is better





User Studies



Key Takeaway:

Comparative language feedback brings the benefits of pairwise comparisons and is preferred by humans.

Comparative language for learning



The first law of thermodynamics, also known as the law of energy conservation, states that energy cannot be created or destroyed in an isolated system. Instead, the total energy of a system and its surroundings remains constant. It can be transformed from one form to another or transferred between the system and its surroundings. Mathematically, it is expressed as $\Delta U = Q - W$, where ΔU is the change in internal energy of the system, Q is the heat added to the system, and W is the work done by the system. This principle underscores that all energy changes in a system must account for energy conservation.

Make it less technical



The first law of thermodynamics says that energy in a system can't be created or destroyed, only changed from one form to another or moved around. For example, when you heat water on a stove, the heat energy from the stove is transferred to the water, raising its temperature. The total amount of energy stays the same; it just changes its form or location.

Today...

- Learning from human feedback
 - Pairwise comparisons
 - Reinforcement learning from human feedback (RLHF)
 - Comparative language feedback