

## Why Environment Design (ED)?

- Human-robot collaboration requires the robot to **predict human intent** and **adapt to their preferences**.
- Current approaches focus on better inference algorithms but overlook an important lever for improving collaboration: ED.



## Problem Formulation

Markov Decision Process (MDP) with  $\mathbf{S}$  (state space),  $\mathbf{A}$  (action space),  $\mathbf{T}$  (transition function),  $\mathbf{R}$  (reward function), and  $\gamma$  (discount factor)

We define ED as modifying the transition function  $\mathbf{T}$  in a MDP.

$$\max_{\mathbf{T}} \max_{\pi} E\left[\sum_t \gamma^t R(s_t, a_t) \mid \pi\right]$$

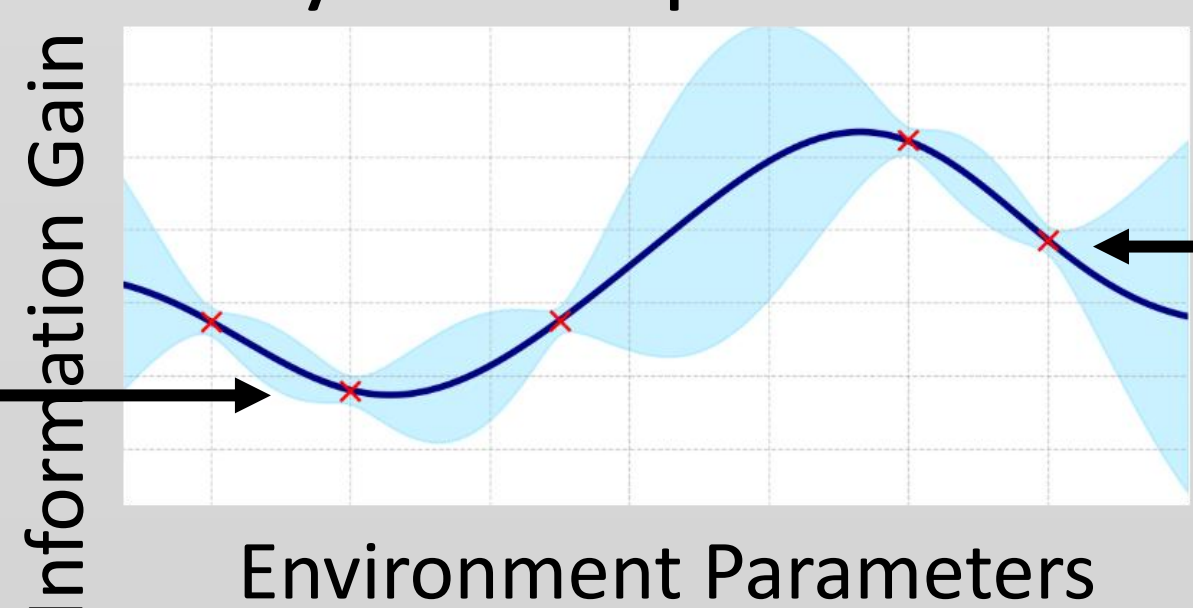
## Reward Alignment

### Key Insight

Environment design can generate informative scenarios to query the human, improving the efficiency of reward function learning.

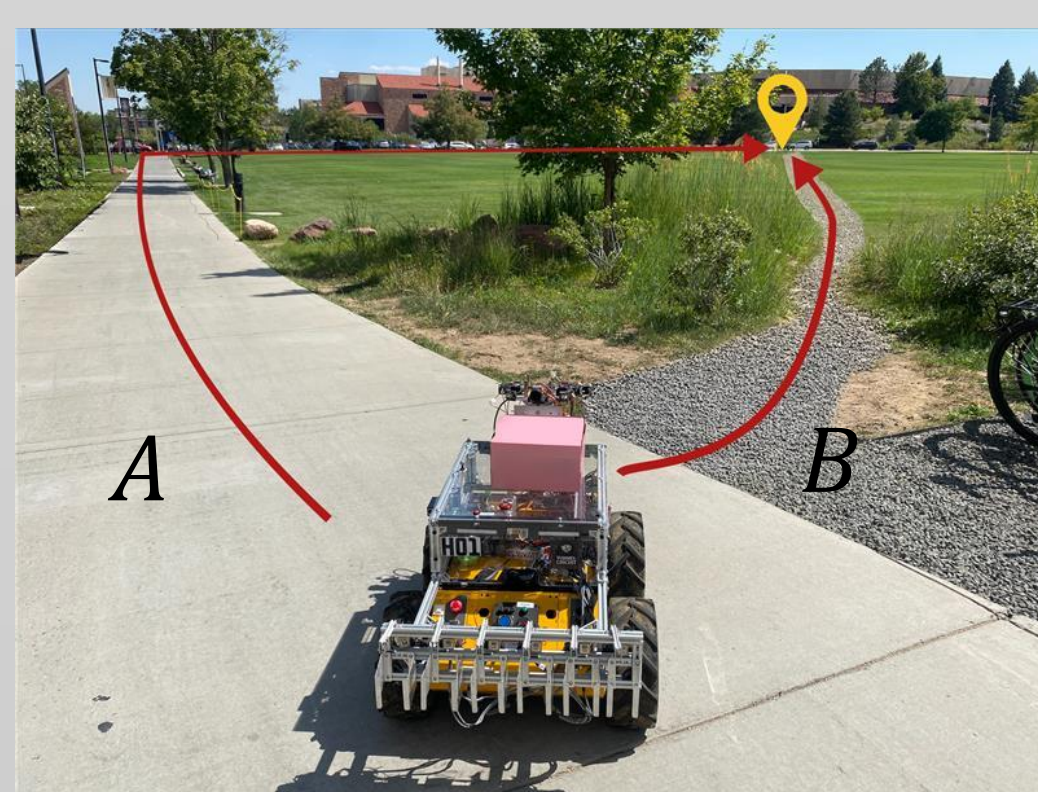
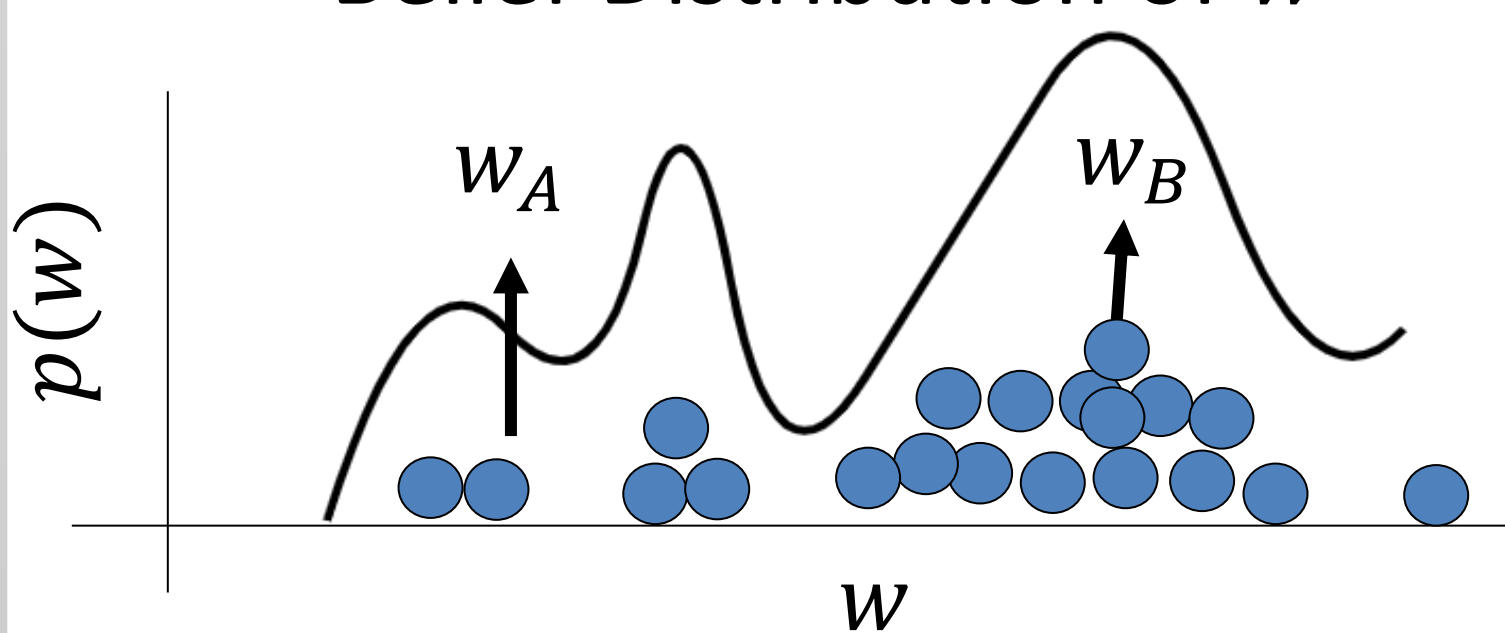
### 1. Generate Environments

#### Bayesian Optimization



### 2. Create Counterfactual Trajectories

#### Belief Distribution of $w$



We use the Bradley-Terry Model of preferences where trajectories with higher rewards are exponentially more likely.

$$p(A > B) = \frac{\exp(R(A))}{\exp(R(A)) + \exp(R(B))}$$

$$R(A) = \phi(A)^T w$$

$R(A)$  – reward of trajectory  $A$

$\phi(A)$  – features of trajectory  $A$

We maintain a Bayesian belief of the human's reward functions by using Monte Carlo Markov Chain (MCMC).

**Goal:** Find trajectories to ask the human for preferences such that the uncertainty of our belief over the human's reward function is minimized.

$$\max_{A, B} H(w) - E[H(w|I)]$$

$H(w)$  – entropy of our belief  $w$

$I$  – human's input

## Research Questions

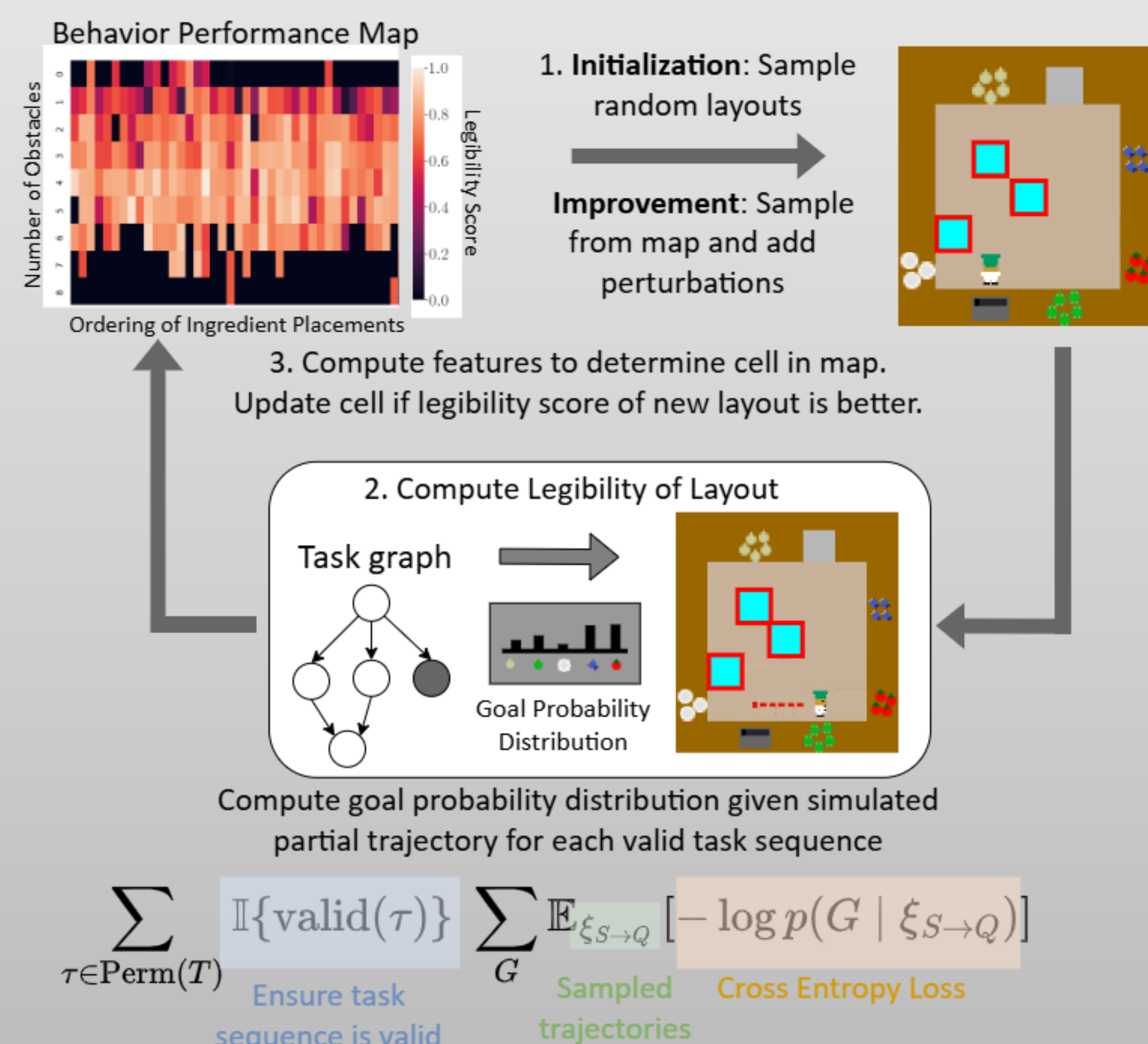
**RQ 1:** What aspects of the environment can we modify to improve the overall task performance and fluency of human-robot interaction?

**RQ 2:** How can we efficiently optimize over the high-dimensional space of possible environment layouts?

## Human Goal Prediction

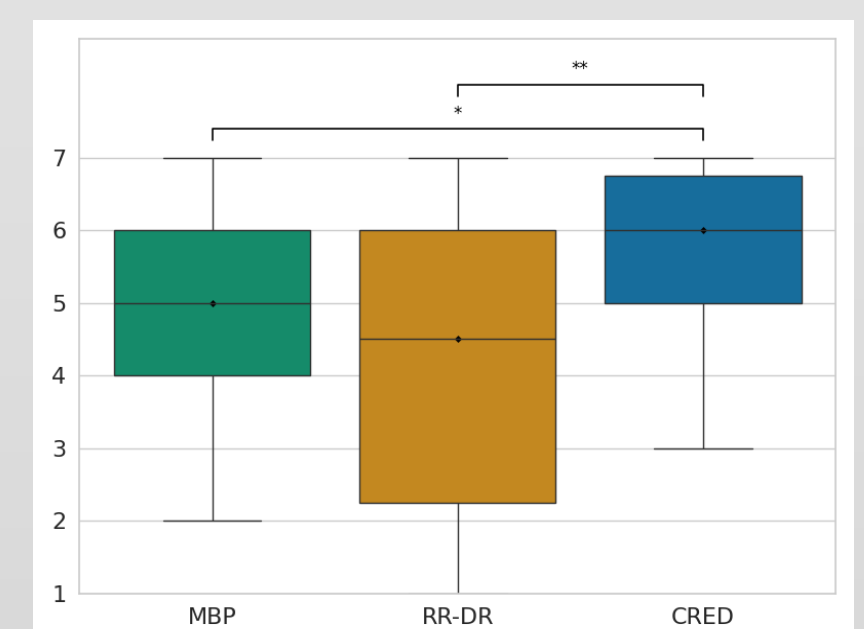
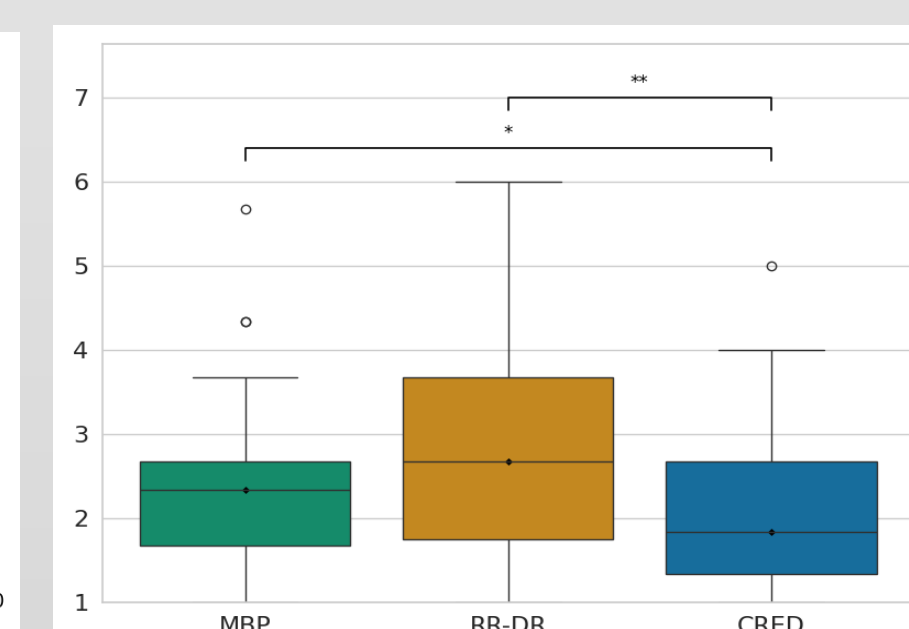
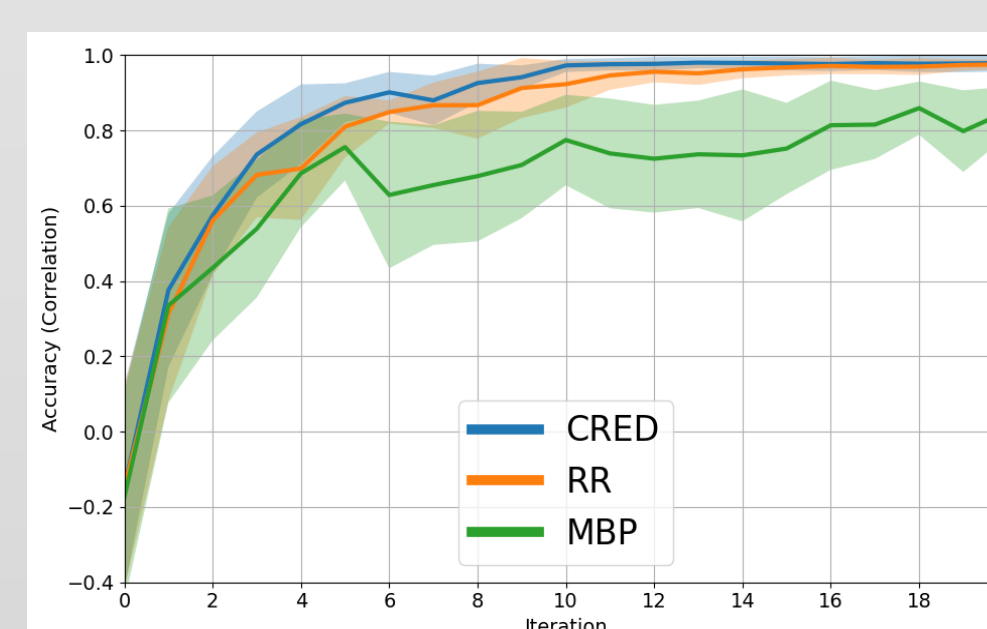
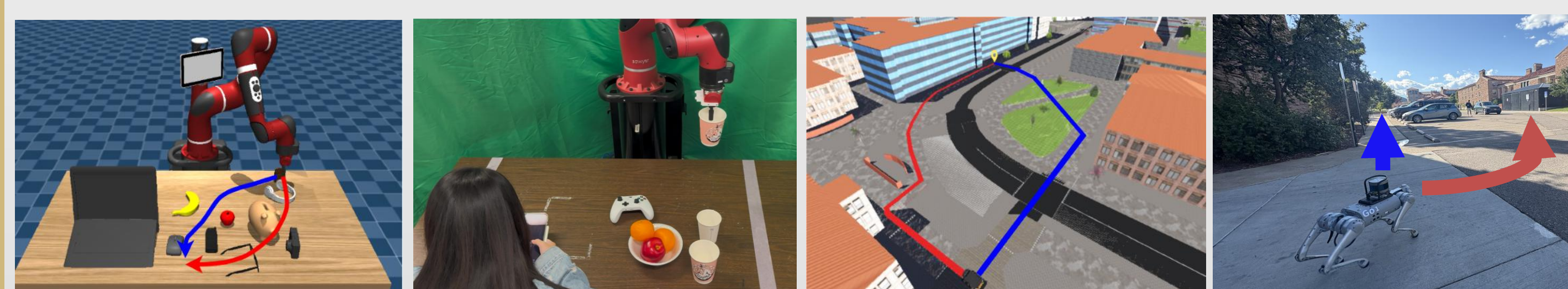
### Key Insight

By shaping the workspace through object placement and AR obstacles, the robot can make human motion more legible and improve goal prediction.



## Experiments / Results

### 1. Reward Alignment

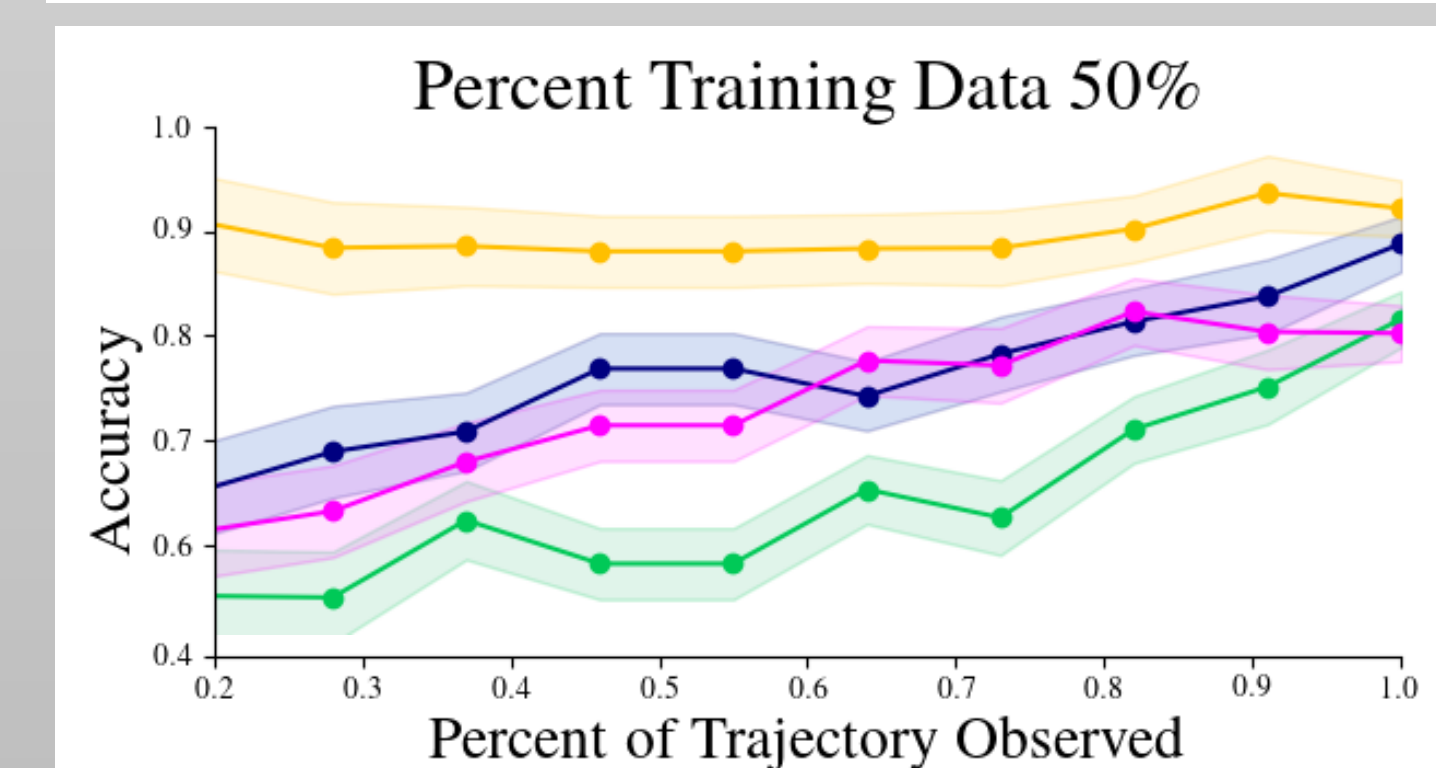
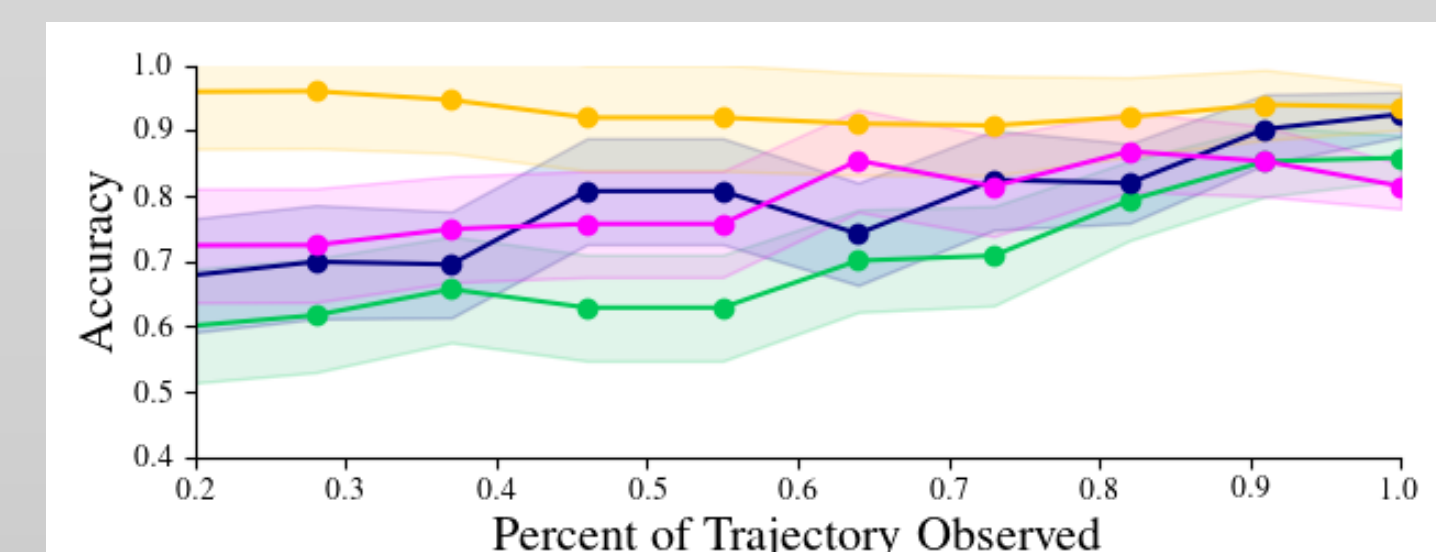
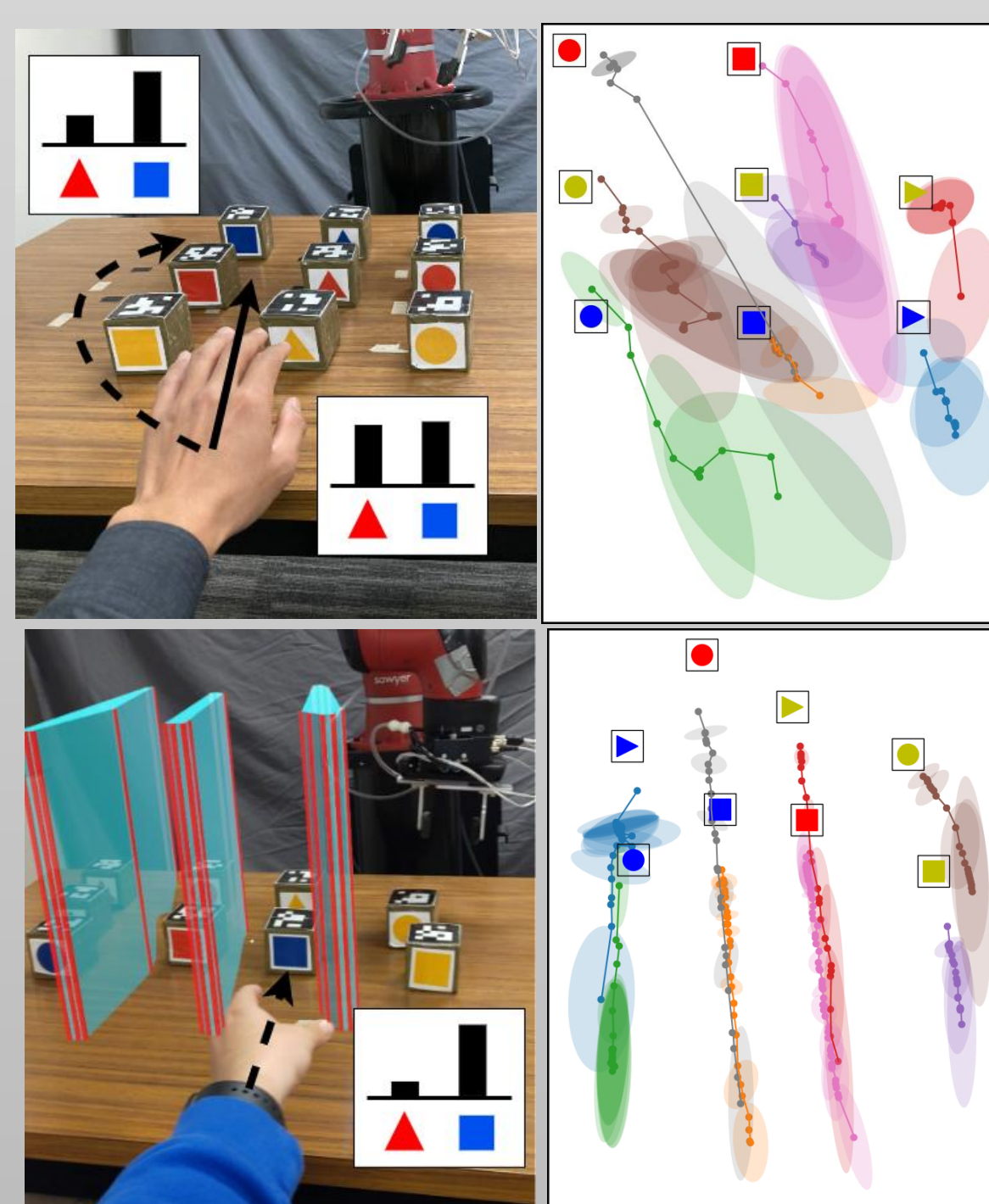


Reward Accuracy

NASA-TLX

Ease of Choice

### 2. Human Goal Prediction



Legend: Baseline (green), Placement Optimized (blue), Virtual Obstacle Optimized (magenta), Both Optimized (yellow)