# Principles of Robot Autonomy II

Human-Robot Interaction

Stanford University

iliad

intelligent and interactive autonomous systems

# Announcement

Paper presentation for 4-credit students due next Wednesday.

Turn them in on Gradescope (slides for live presenters, video recordings for non-live presenters)
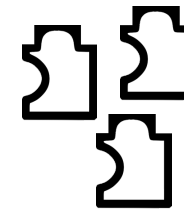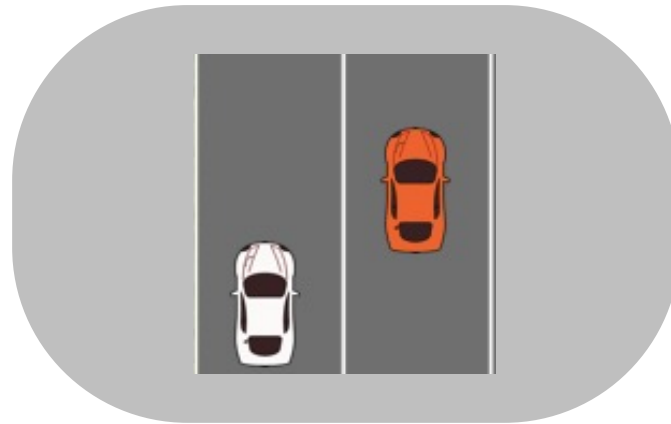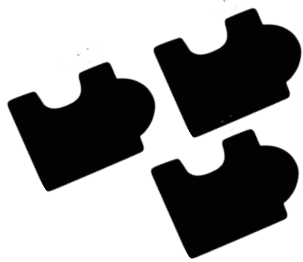
# Today's itinerary

- Game-Theoretic Views on Multi-Agent Interactions

- Partner Modeling: Active Info Gathering over Human's Intent

- Partner Modeling: Learning and Influencing Latent Intent

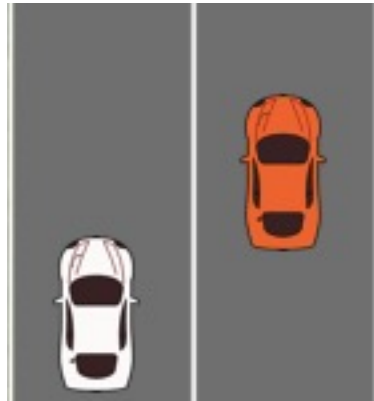- Partner Modeling: Role Assignment

# Today's itinerary

- Game-Theoretic Views on Multi-Agent Interactions

- Partner Modeling: Active Info Gathering over Human's Intent

- Partner Modeling: Learning and Influencing Latent Intent

- Partner Modeling: Role Assignment

# Nth order Theory of Mind

[*Sadigh, Sastry, Seshia, Dragan,* RSS 2016, IROS 2016, AURO 2018]

# Interaction as a Dynamical System

$$u_{\mathcal{R}}^* = \operatorname*{argmax}_{u_{\mathcal{R}}} R_{\mathcal{R}}(x, u_{\mathcal{R}}, u_{\mathcal{H}}^*(x, u_{\mathcal{R}}))$$



Find optimal actions for the robot while accounting for the human response $u_{\mathcal{H}}^*$.

Model $u_{\mathcal{H}}^*$ as optimizing the human reward function $R_{\mathcal{H}}$.

$$u_{\mathcal{H}}^*(x, u_{\mathcal{R}}) \approx \operatorname*{argmax}_{u_{\mathcal{H}}} R_{\mathcal{H}}(x, u_{\mathcal{R}}, u_{\mathcal{H}})$$

# Today's itinerary

- Game-Theoretic Views on Multi-Agent Interactions

- Partner Modeling: Active Info Gathering over Human's Intent

- Partner Modeling: Learning and Influencing Latent Intent

- Partner Modeling: Role Assignment

# Modeling Intent Inference using POMDPs



[Javdani et al.]
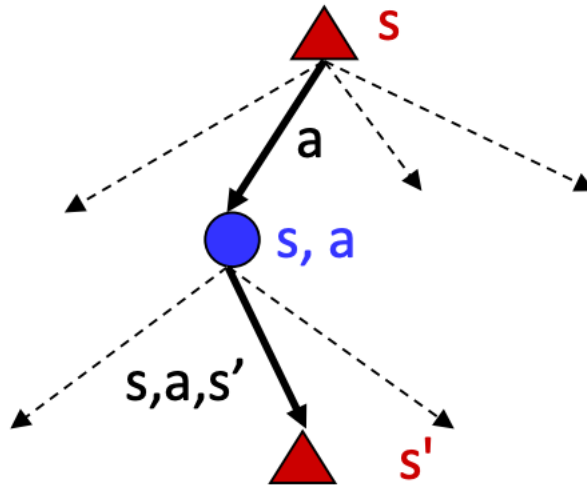
# POMDP Formulation

**MDPs have:**

States $S$

Actions $A$

Transition Function $P(s'|s, a)$

Reward $R(s, a, s')$

**POMDPs add:**

Observations $O$

Observation Function $P(o|s)$

# Tiger Example

+10    -100

?

-1

**Actions** $a = \{0, : \text{listen } 1: \text{open left}, 2: \text{open right}\}$
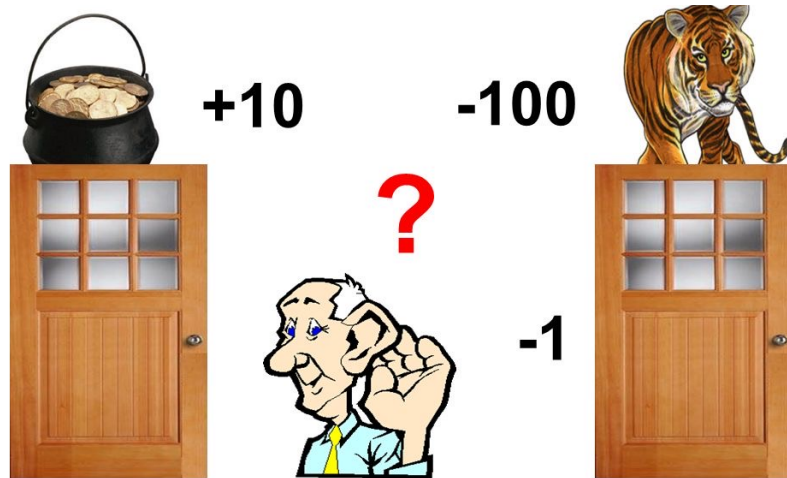
**Reward Function:**
- Penalty for wrong opening: -100
- Reward for correct opening: +10
- Cost of listening: -1

**Observations:**
- To hear the tiger on the left
- To hear the tiger on the right

# Tiger Example

**+10**  **-100**

**-1**

Belief update based on observations:

$$b_1(s_i) \propto p(o|s_i, a) \sum_{s_j \in S} p(s_i|s_j, a) \cdot b_0(s_j)$$

*Immediate return*   *Discounted future return*

Value Iteration over Beliefs

$$V^*(b) = \max_{a \in A} [\sum_{s \in S} b(s) \cdot R(s, a) + \gamma \sum_{o \in O} P(o|b, a) \cdot V^*(b_o^a)]$$

Hard to compute continuous space MDPs -> Approximation

# Tiger Example

Value Iteration over Beliefs

*Immediate return*　　*Discounted future return*

$$V^*(b) = \max_{a \in A}[\sum_{s \in S} b(s) \cdot R(s,a) + \gamma \sum_{o \in O} P(o|b,a) \cdot V^*(b_o^a)]$$

Hard to compute continuous space MDPs -> Approximation

Q-MDP Approximation

$$V^*(b) = \mathbb{E}_s[V^*(s)] = \sum_s b(s) \cdot V^*(s)$$
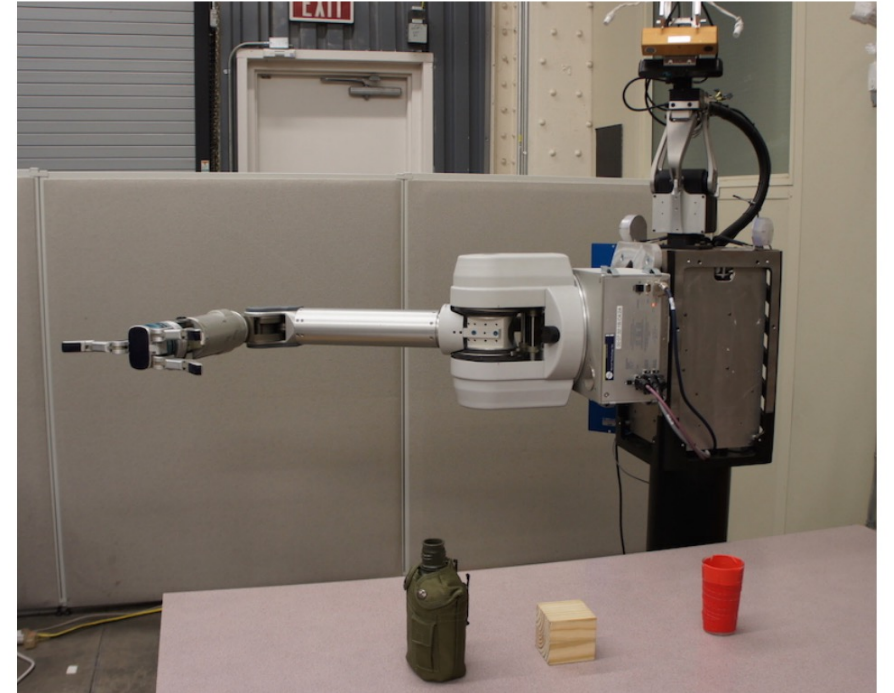
# Intent Inference



$X$      Robot States

$A$      Robot Actions

$T: X \times A \to X$      Transition function

$u \in U$      Human continuous input

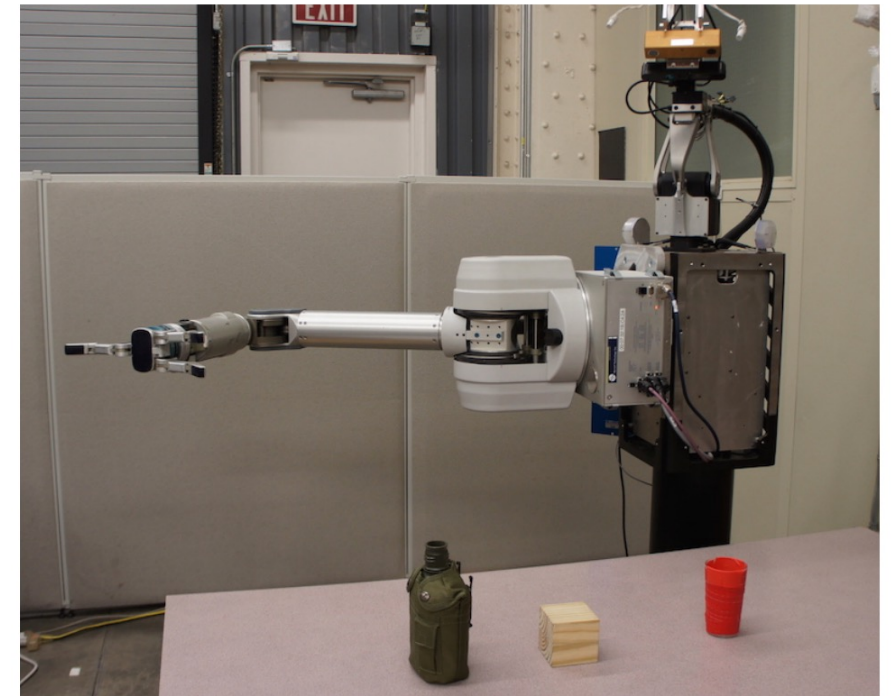$D: U \to A$      Mapping between human input and robot actions

# User's Policy is Learned from IRL

$$\pi_g^{usr}(x) = p(u|x, g)$$  We learn a policy for each goal

$$p(\xi|g) \propto \exp(-C_g^{usr}(\xi))$$

$$p(g|\xi) \propto p(\xi|g) \cdot p(g)$$  Bayes Rule

POMDP Observation Model

# Hindsight Optimization (Q-MDP)

Estimate cost-to-go of the belief by assuming full observability will be obtained at the next time step.

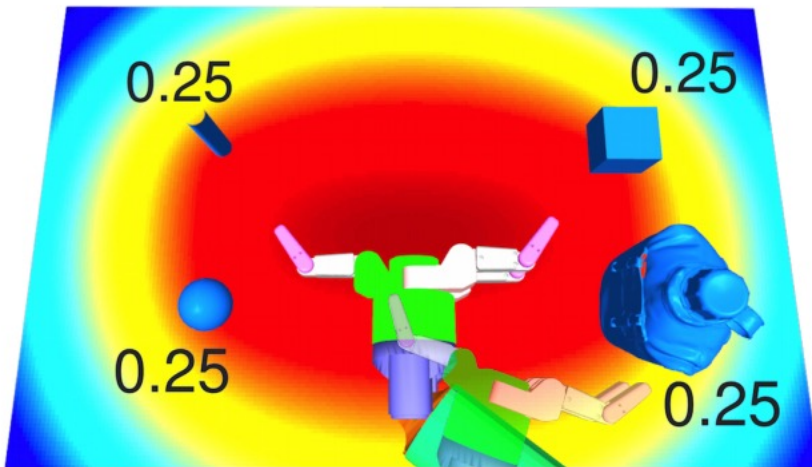You never gather information, but can plan efficiently in deterministic subproblems.

$$b(s) = b(g) = p(g|\xi) \qquad \text{Uncertainty is only over goals}$$

$$Q(b, a, u) = \sum_g b(g) \cdot Q_g(x, a, u)$$

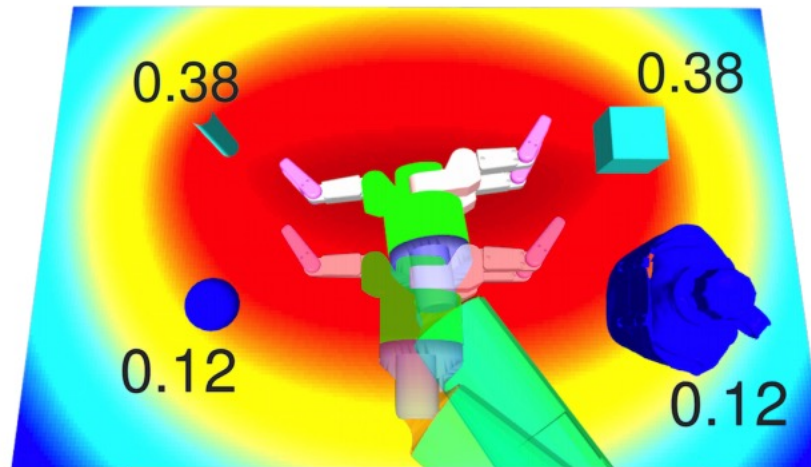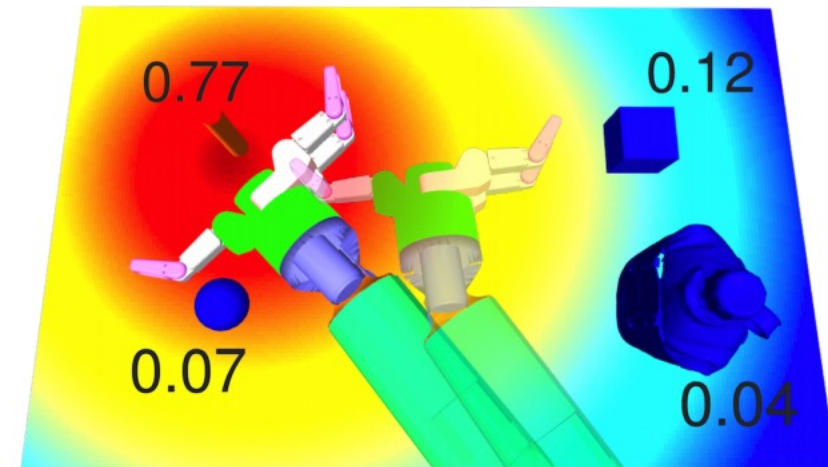Action-Value function of the POMDP          Cost-to-Go of Acting optimally and going towards goal $g$
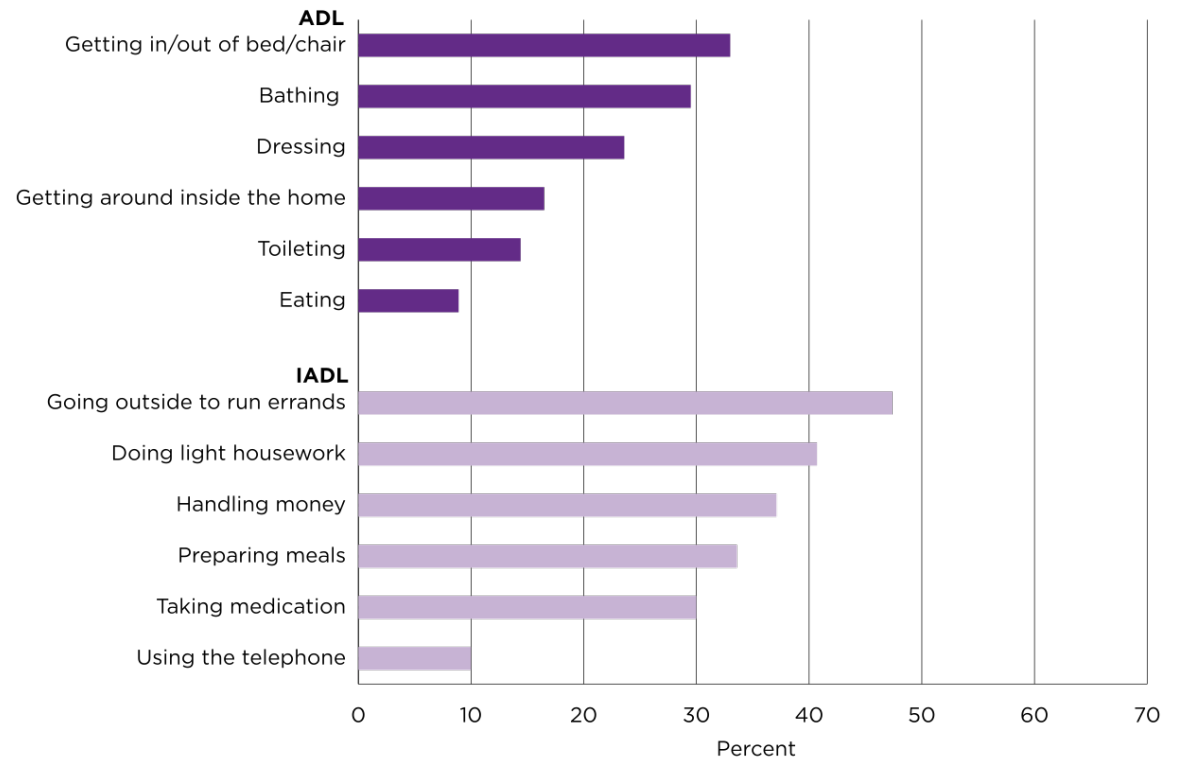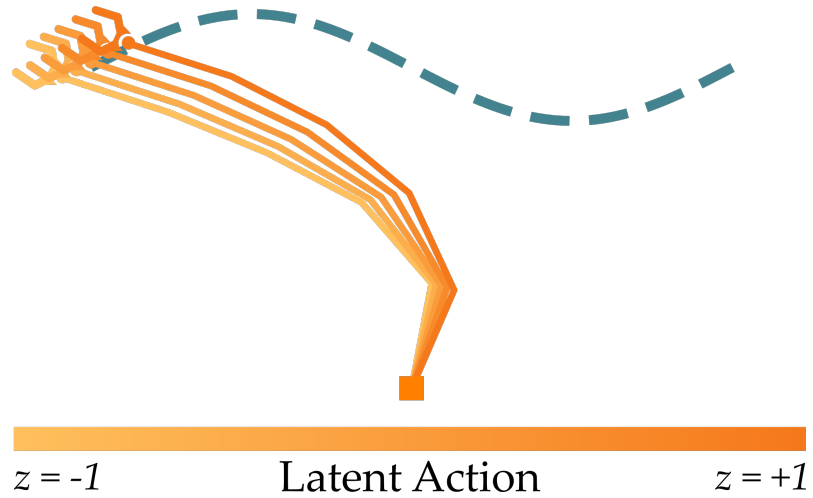
# Shared Autonomy with Hindsight Optimization

**Prevalence of Difficulty Performing ADLs and IADLs in Adults 18 Years and Older With One or More Selected Symptoms That Interfere With Everyday Activities: 2014**

**ADL**
- Getting in/out of bed/chair
- Bathing
- Dressing
- Getting around inside the home
- Toileting
- Eating

**IADL**
- Going outside to run errands
- Doing light housework
- Handling money
- Preparing meals
- Taking medication
- Using the telephone
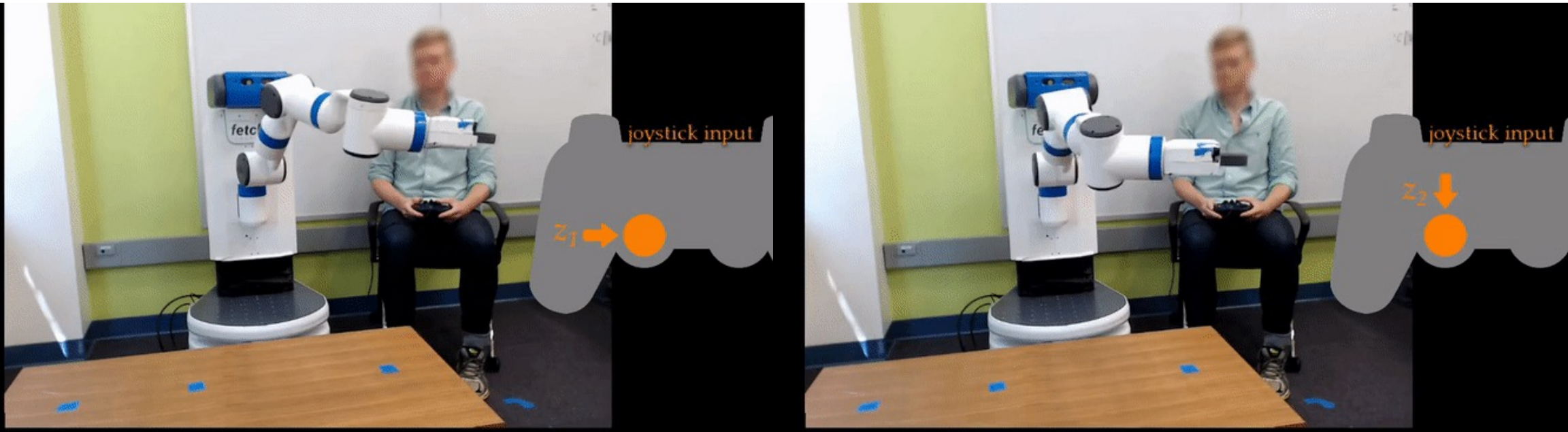
Percent: 0 10 20 30 40 50 60 70

- Assistive robotic arms are *dexterous*
- This dexterity makes it hard for users to *control* the robot

- How can robots *learn* low-dimensional representations that make controlling the robot intuitive?
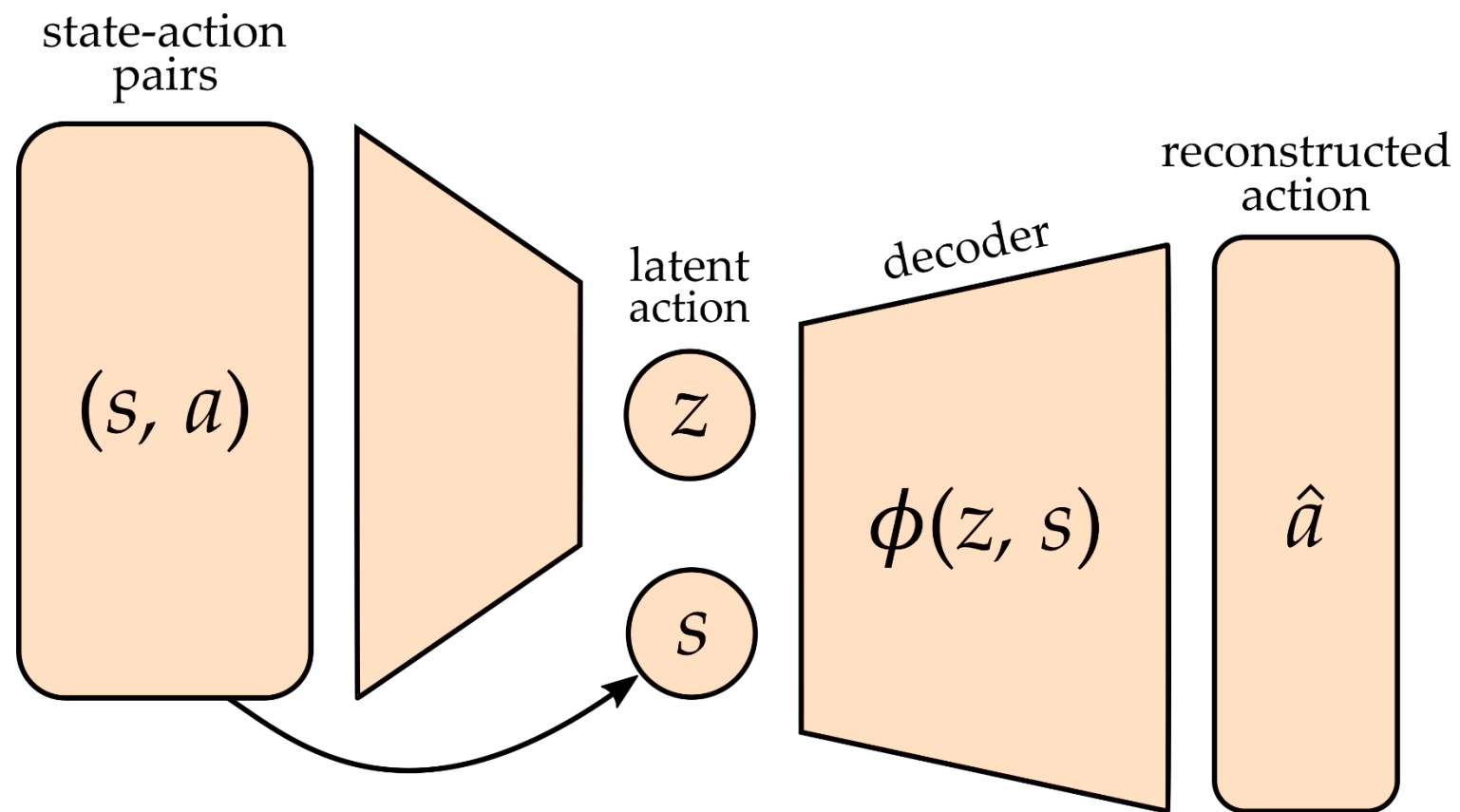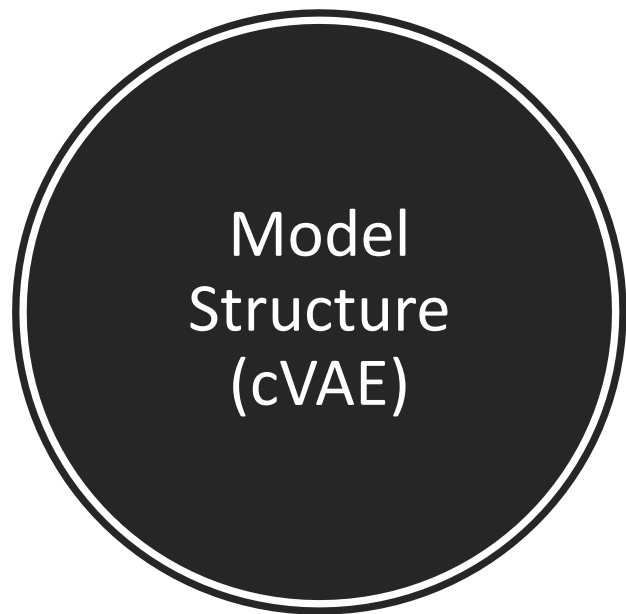
# Our Vision



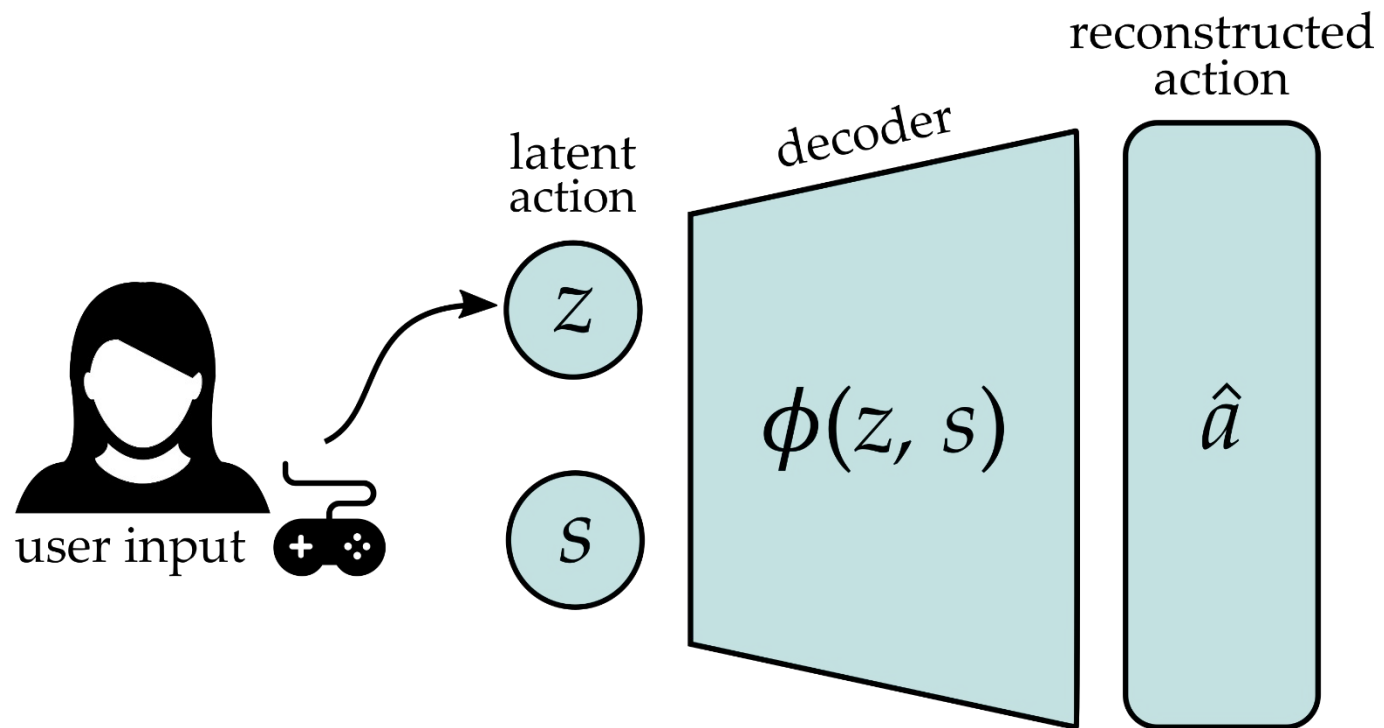Offline, expert demonstrations of *high-dimensional* motions
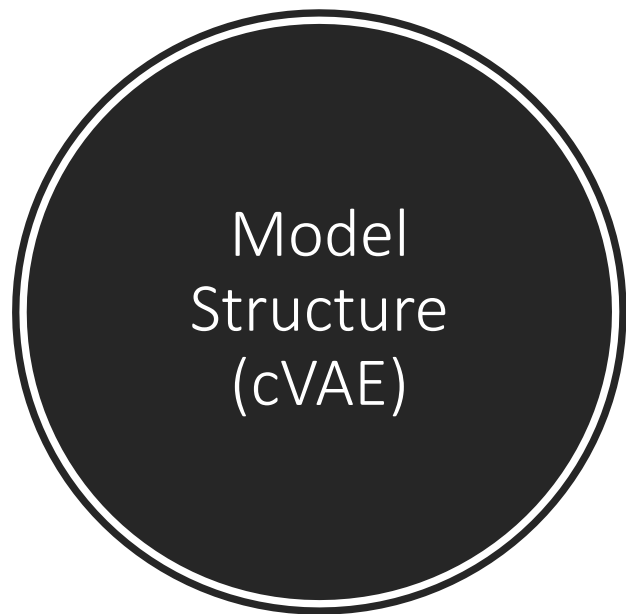
# Our Vision



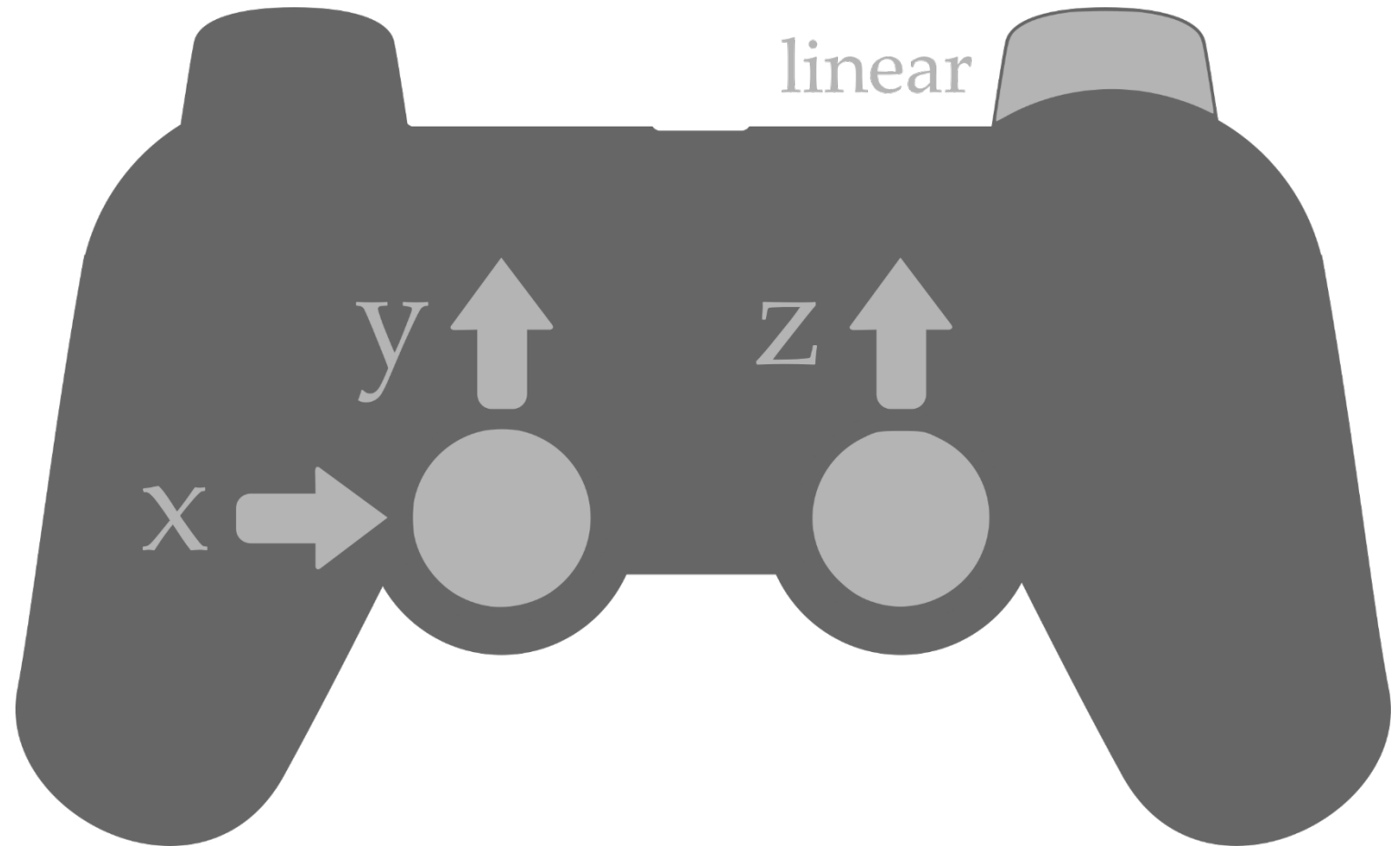Learn *low-dimensional* latent representations for online control

We make it easier to control *high-dimensional* robots by *embedding* the robot's actions into a *low-dimensional* latent space.
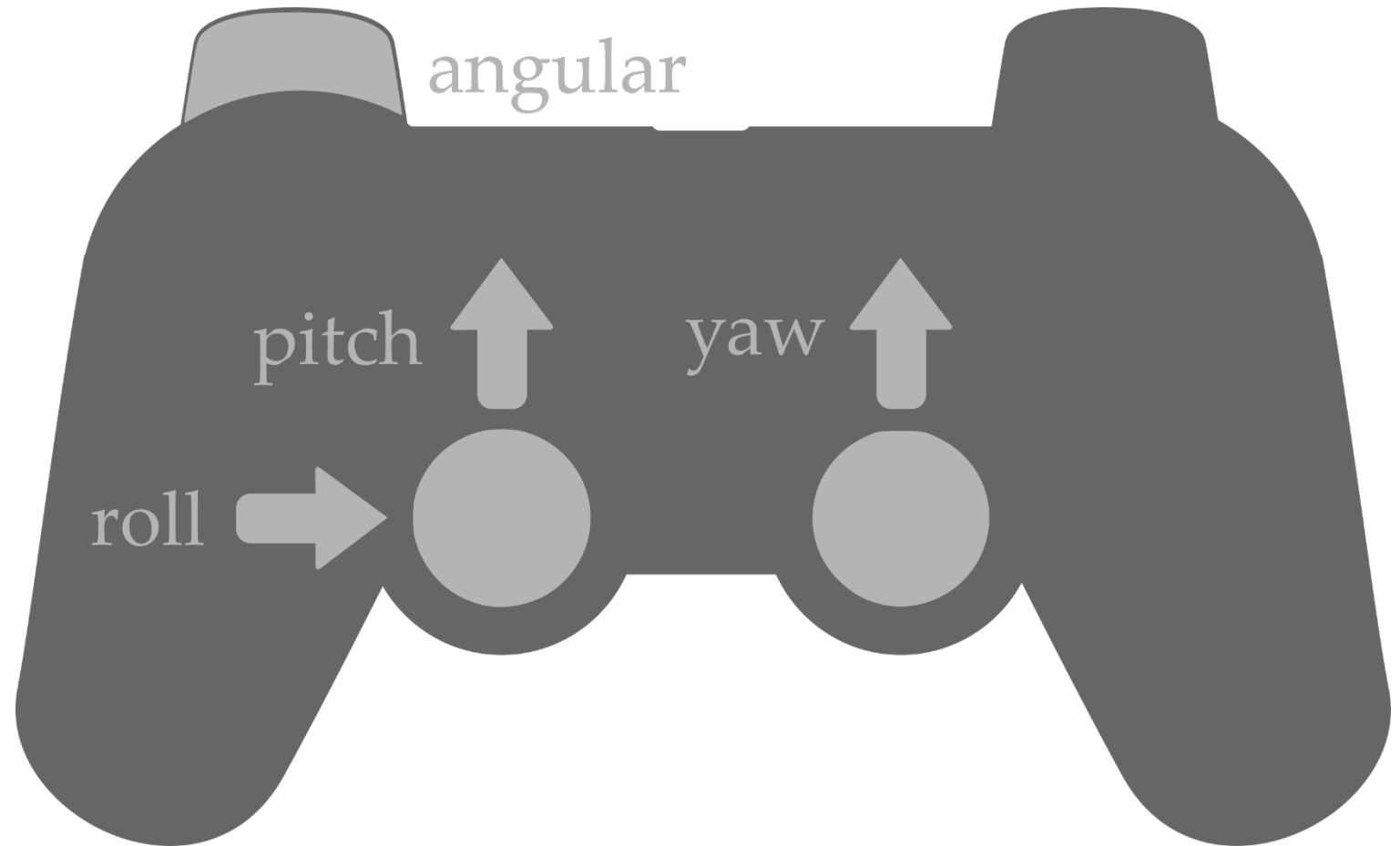
Model
Structure
(cVAE)

state-action
pairs

$(s, a)$

latent
action

$z$

$s$

decoder

$\phi(z, s)$

reconstructed
action

$\hat{a}$

Model Structure (cVAE)

latent action

user input

$z$

$s$

decoder

$\phi(z, s)$

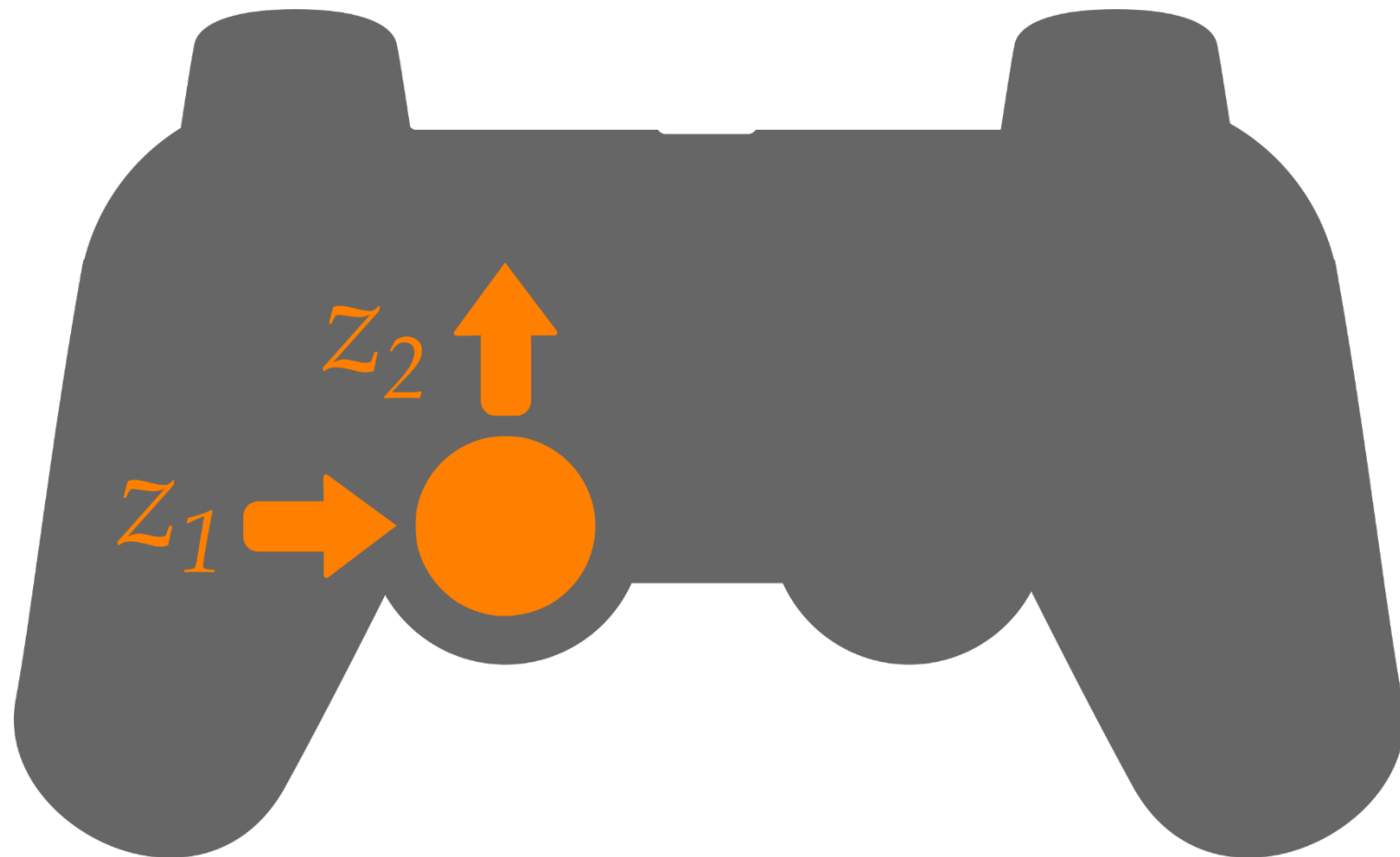reconstructed action

$\hat{a}$

# User Study

- We trained on less than **7 minutes** of kinesthetic demonstrations

- Demonstrations consisted of moving between shelves, pouring, stirring, and reaching motions

- We compared our *Latent Action* to the current method for assistive robotic arms (*End-Effector*)
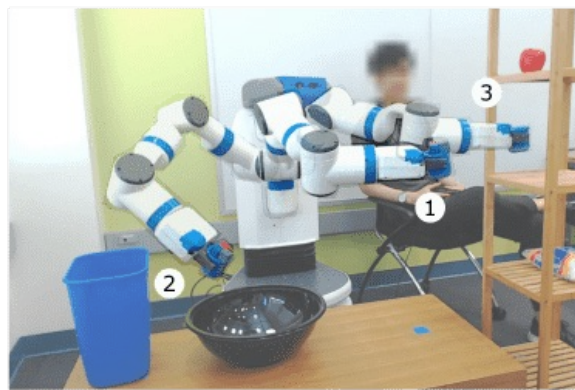
4x Speed

(1) add eggs            (1) add eggs

End-Effector            Latent Action
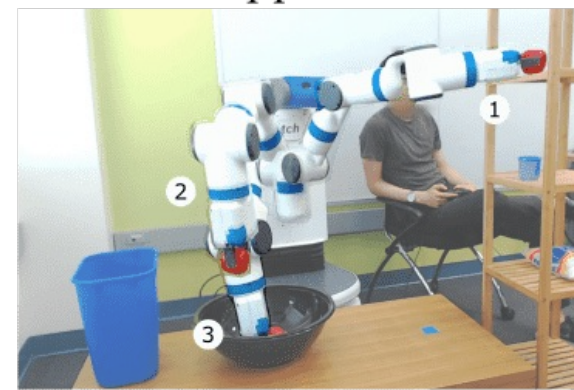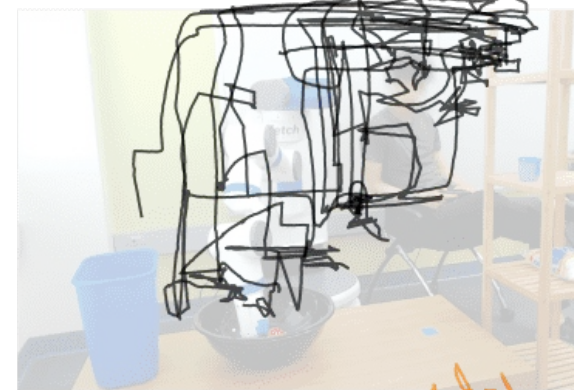
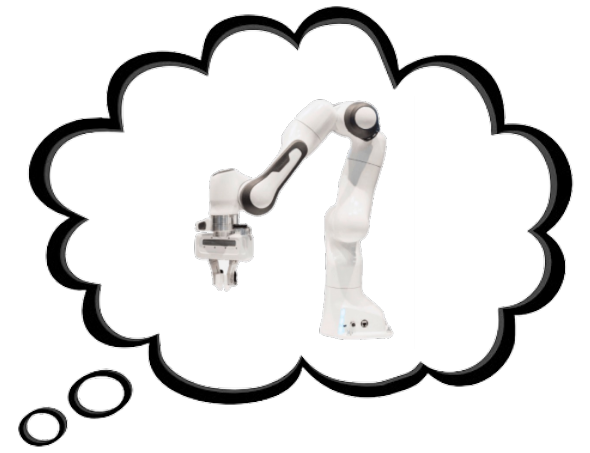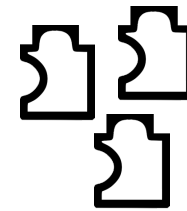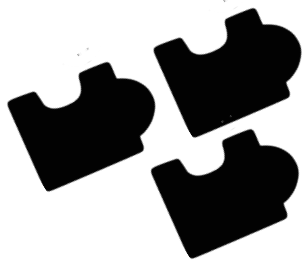| Add Eggs & Recycle | Add Flour & Return | Add Apple and Stir |

# Today's itinerary

- Game-Theoretic Views on Multi-Agent Interactions

- Partner Modeling: Active Info Gathering over Human's Intent

- Partner Modeling: Learning and Influencing Latent Intent

- Partner Modeling: Role Assignment
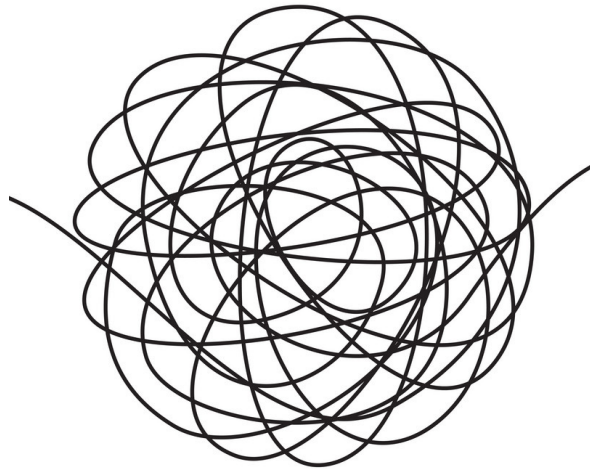
# Nth order Theory of Mind

**Most interactive tasks are not the same as playing chess!**

… low-dimensional shared representation
that captures the interaction and can change over time.

Other agents are often non-stationary:
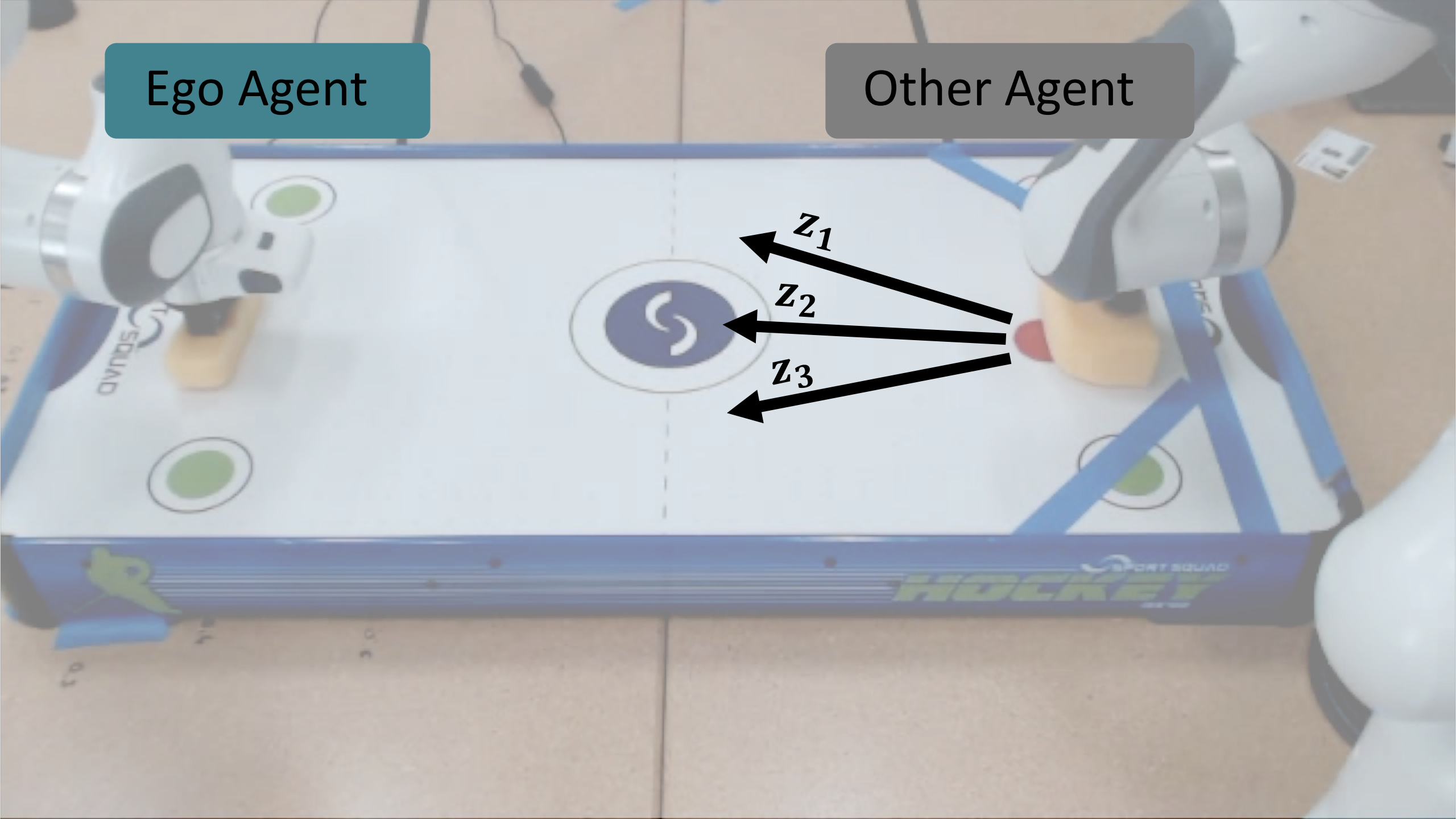They update their behavior in response to the robot.

$a \in \mathbb{R}^7$

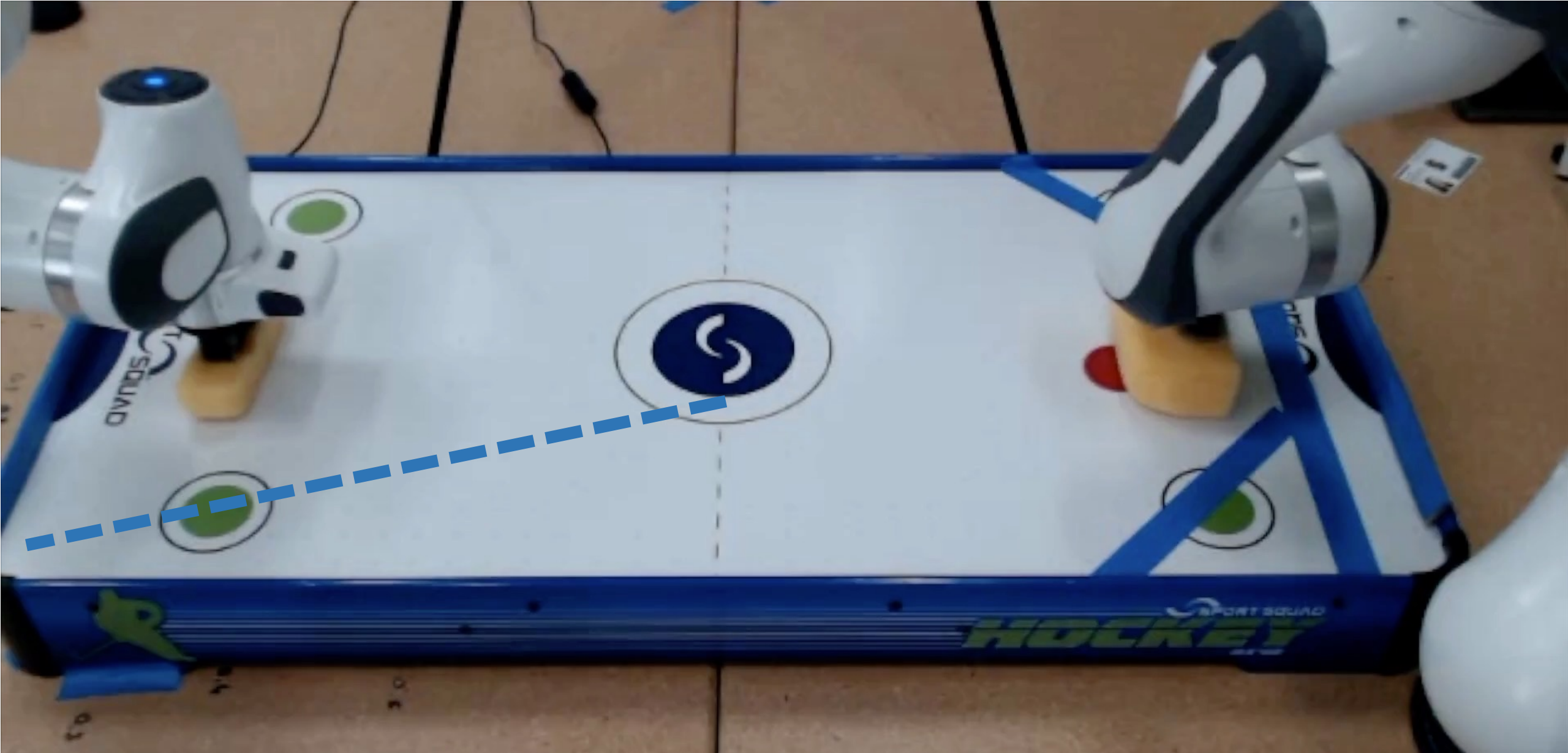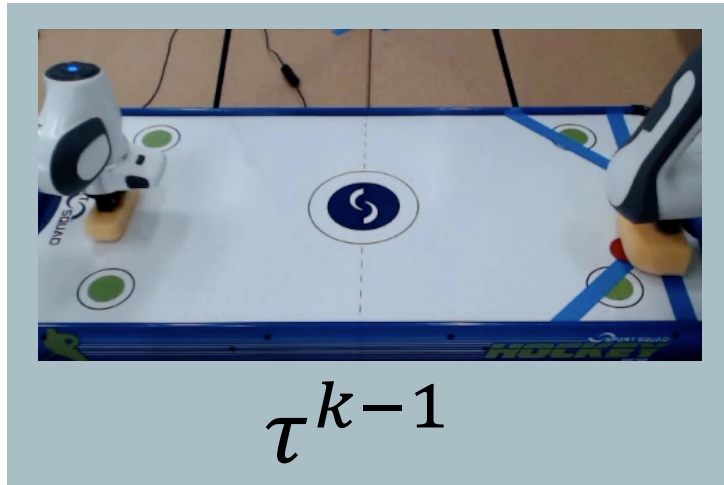$$\tau^i = \{(s_1, a_1, r_1), ..., (s_H, a_H, r_H)\}$$

$$z^{i+1} \sim f(\cdot \mid z^i, \tau^i)$$

# Modeling Other Agent's Behavior

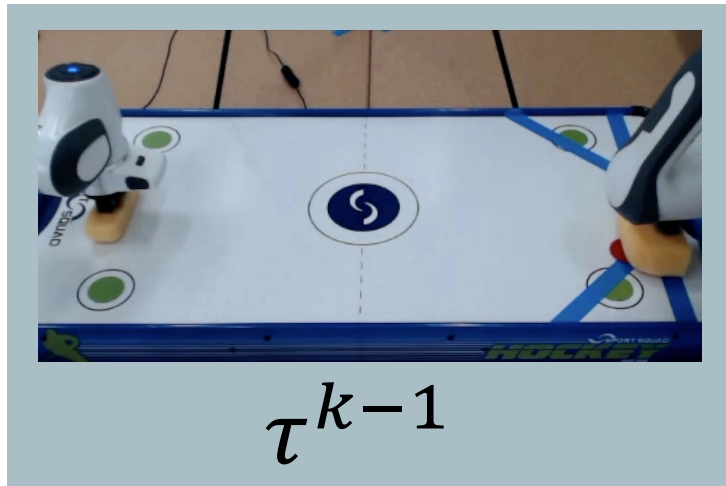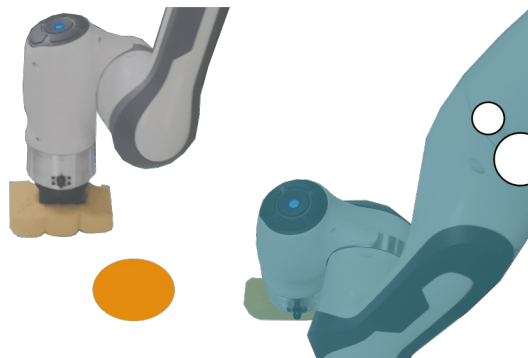# Modeling Other Agent's Behavior



$\tau^{k-1}$

$\mathcal{E}_\phi$  $z^k$

Learning objective: $\max\limits_{\phi,\psi} \sum\limits_{i=2}^{N} \sum\limits_{t=1}^{H} \log p_{\phi,\psi}(s_{t+1}^i, r_t^i \mid s_t^i, a_t^i, \tau^{i-1})$

*Representation Learning*

$\tau^{k-1}, \tau^k$

Experience
Buffer

$\tau^{k-1}$   $\mathcal{E}$   $z^k$   $\mathcal{D}$   $\hat{\tau}^k$

# Learning and Influencing Latent Intent

Maximize expected return
*within* an interaction

$$\max_{\theta} \quad \mathbb{E}_{\pi_\theta(a|s,z^i)}\left[\sum_{t=1}^{H} R(s,z^i)\right]$$

to *react to* the other agent

**Representation Learning**



$\tau^{k-1}, \tau^k$

$\tau^{k-1} \quad \mathcal{E} \quad z^k \quad \mathcal{D} \quad \hat{\tau}^k$

Experience
Buffer

[*Xie, Losey, Tolsma, Finn, Sadigh,* CoRL 2020]

Air Hockey Results

Air Hockey Results

Air Hockey Results

Air Hockey Results

Ego Agent

Other Agent

Air Hockey Results

2x speed

SAC: initial policy

2x speed

SAC: 2 hours of training

2x speed

SAC: 4 hours of training

# Air Hockey Results

2x speed

LILI: 4 hours of training

# Air Hockey Results



*always block right*

*always block left*

Random    SAC    LILI

# Reacting to Other Agents

Maximize expected return
*within* an interaction
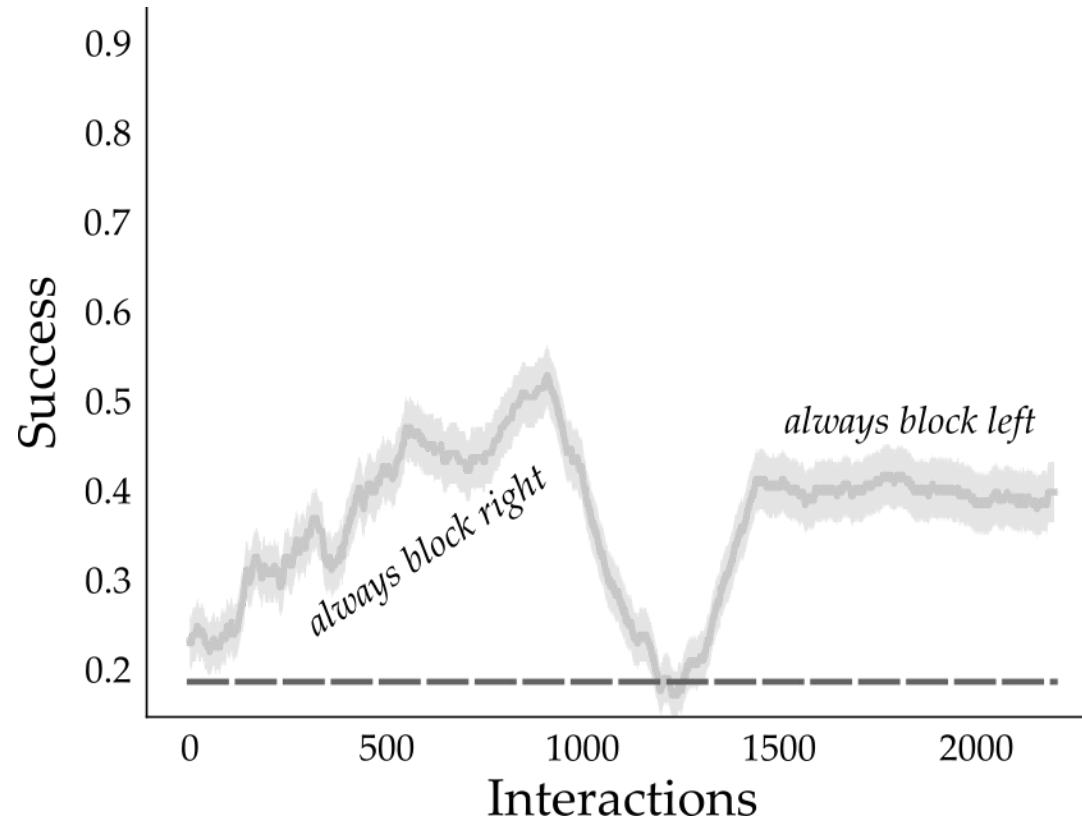
$$\max_{\theta} \quad \mathbb{E}_{\pi_\theta(a|s,z^i)}\left[\sum_{t=1}^{H} R(s,z^i)\right]$$

to *react to* the other agent

**Representation Learning**

$\tau^{k-1}, \tau^k$

Experience
Buffer

$\tau^{k-1} \quad \mathcal{E} \quad z^k \quad \mathcal{D} \quad \hat{\tau}^k$

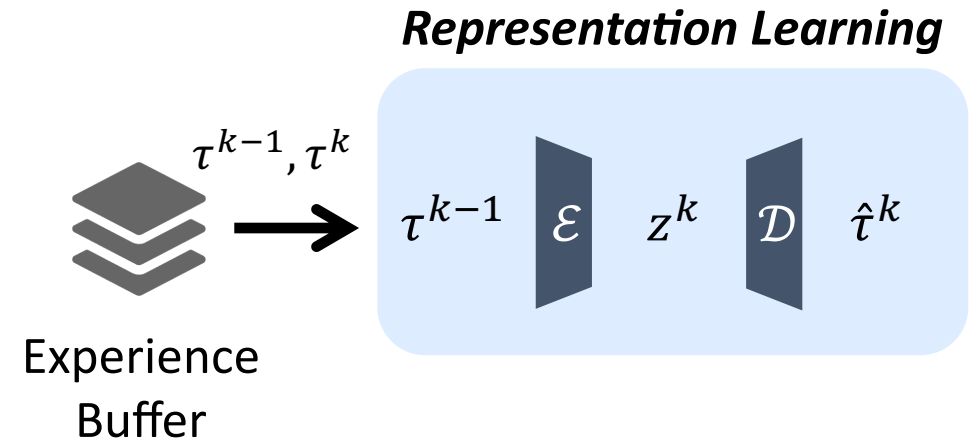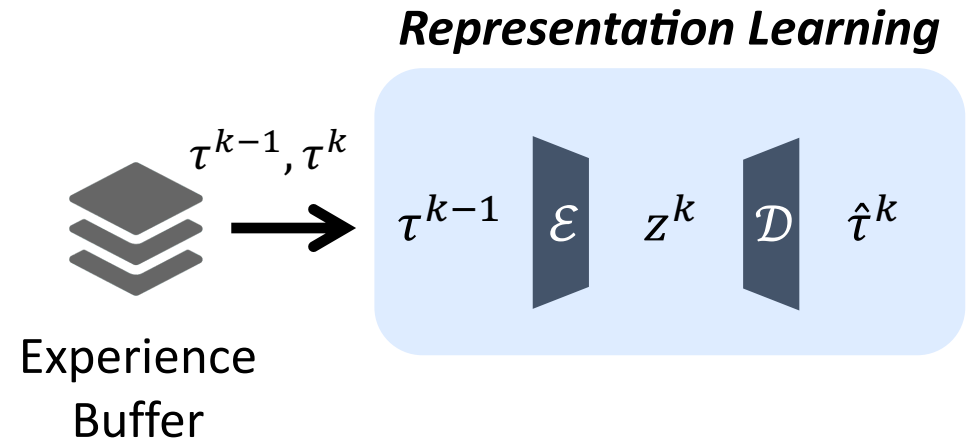# Influencing Other Agents

Maximize expected return *across* interactions

$$\max_\theta \sum_{i=1}^\infty \gamma^i \mathbb{E}_{\pi_\theta(a|s,z^i)} \left[ \sum_{t=1}^H R(s, z^i) \right]$$

to *influence* the other agent

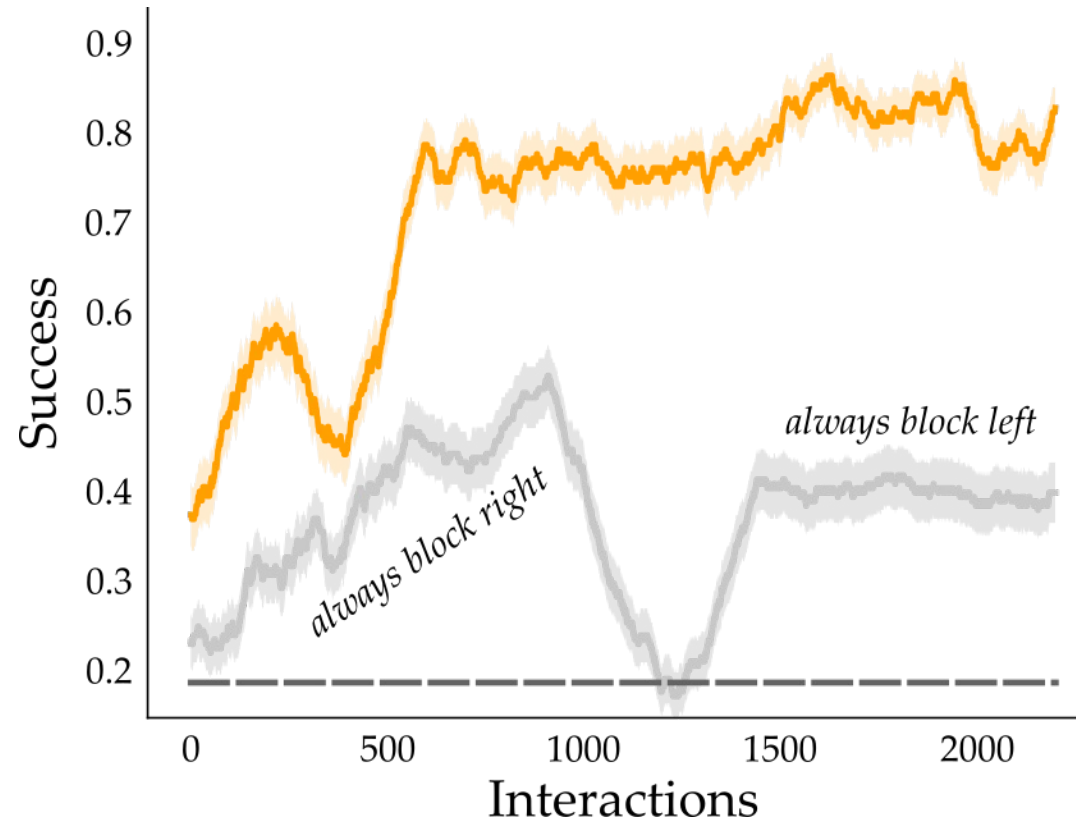$\tau^{k-1}, \tau^k$

$\tau^{k-1}$ $\mathcal{E}$ $z^k$ $\mathcal{D}$ $\hat{\tau}^k$

Experience
Buffer

2x speed

LILI (with influence): 4 hours of training

# Air Hockey Results

# Air Hockey Results



Legend: Random, SAC, LILI (no influence), LILI (ours)

# Air Hockey Results



[*Xie, Losey, Tolsma, Finn, Sadigh,* CoRL 2020]

Playing with a Human Expert

SAC: 45% success

Playing with a Human Expert

LILI : 73% success

# Key Takeaways

Human partners are often non-stationary –
which can be represented by low-dimensional latent strategies.

LILI *anticipates* the partner's policies using latent strategies to *react* and *influence* the other agent.

# Today's itinerary

- Game-Theoretic Views on Multi-Agent Interactions

- Partner Modeling: Active Info Gathering over Human's Intent

- Partner Modeling: Learning and Influencing Latent Intent

- Partner Modeling: Role Assignment