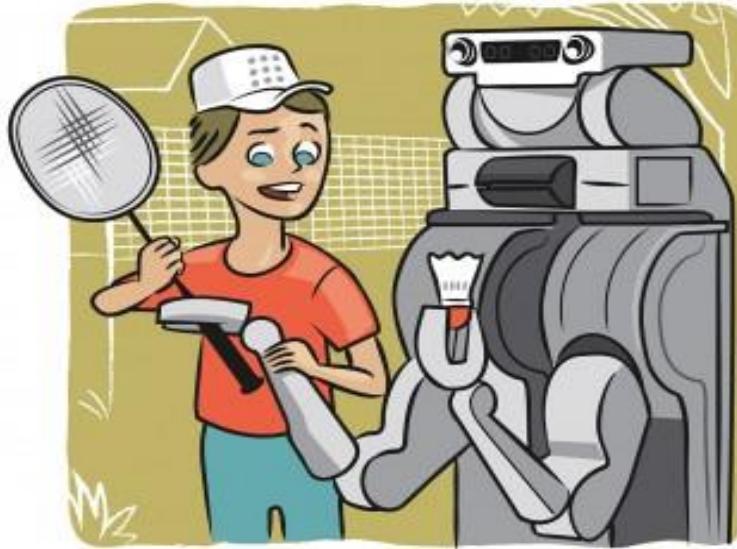


Behavioral Cloning and Interactive Imitation Learning



Instructor: Daniel Brown

[Some slides adapted from Sergey Levine (CS 285) and Alina Vereshchaka (CSE4/510)]



Brief Machine Learning Refresher

There are roughly 3 main branches of machine learning

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

Supervised Learning

- **Setting/Assumptions:** In supervised learning, the model is trained on labeled data, where the input data is paired with the correct output (i.e., the "ground truth").
- **Goal:** To learn a mapping from inputs to outputs so that the model can predict the output for new, unseen inputs.
- **Common Use Cases:**
 - Classification (e.g., spam email detection, image recognition).
 - Regression (e.g., predicting house prices, stock market trends).
- **Example models:**
 - Linear regression, decision trees, support vector machines, and neural networks.

Classification

$$\text{Cross-Entropy Loss} = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

Regression

$$\text{MSE Loss} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

Unsupervised Learning

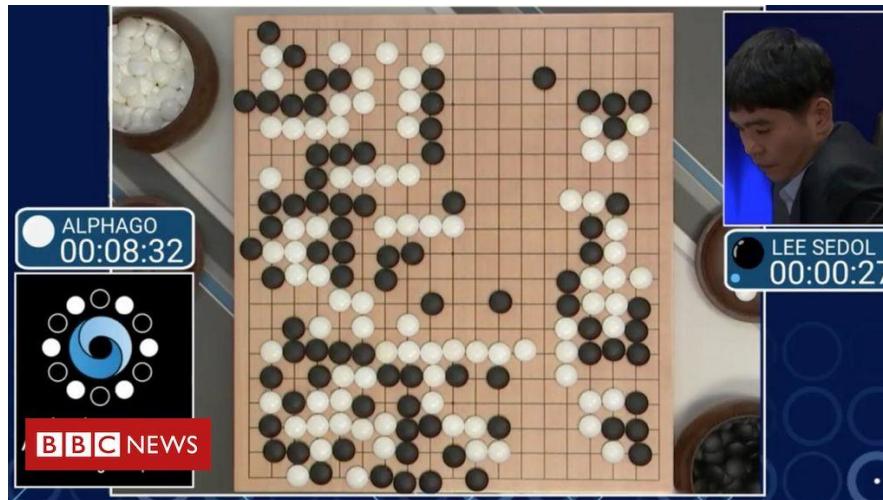
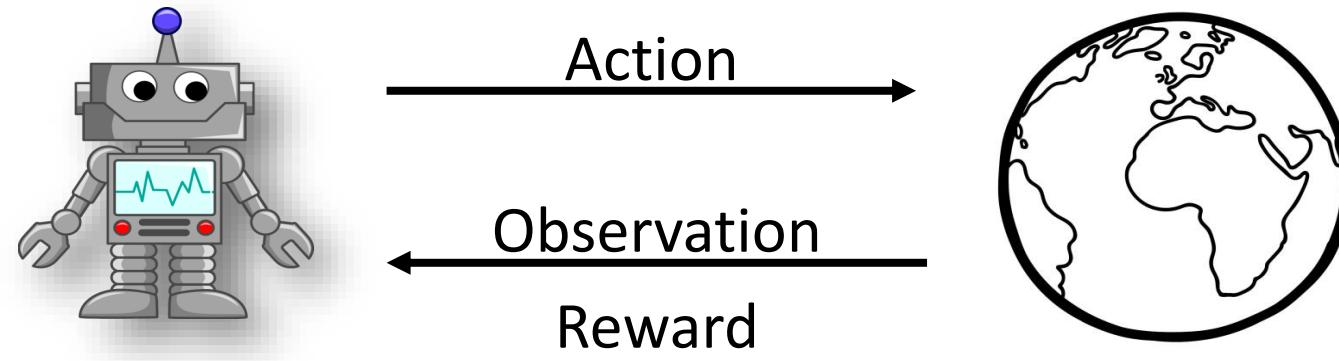
- **Setting/Assumptions:** In unsupervised learning, the model is trained on data without labeled outputs. It seeks to find patterns, structures, or relationships in the data. No “ground truth” labels.
- **Goal:** To explore the data and identify meaningful clusters, associations, or representations.
- **Common Use Cases:**
 - Clustering (e.g., customer segmentation).
 - Dimensionality reduction (e.g., PCA for visualization).
 - Anomaly detection (e.g., fraud detection).
- **Example models:**
 - K-means clustering, hierarchical clustering, and autoencoders.

Auto-Encoders

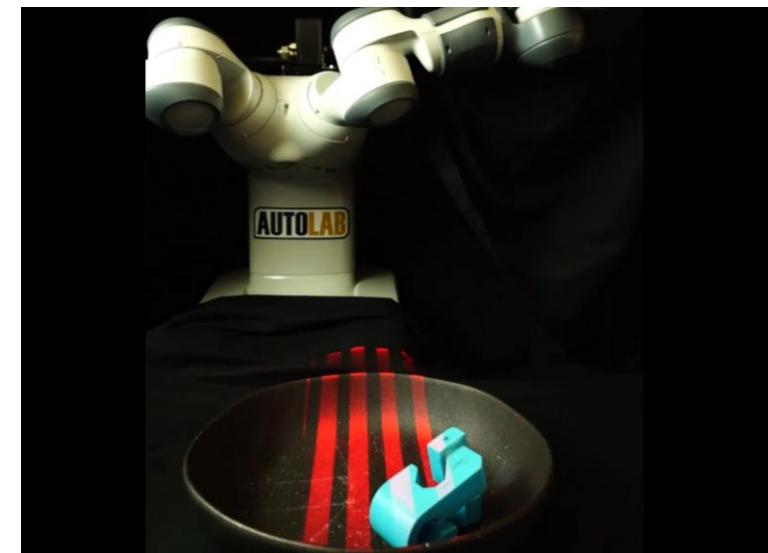
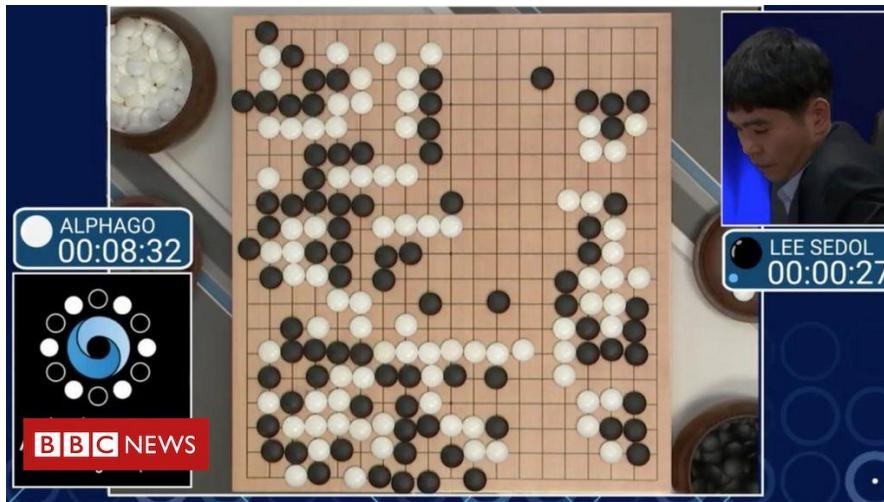
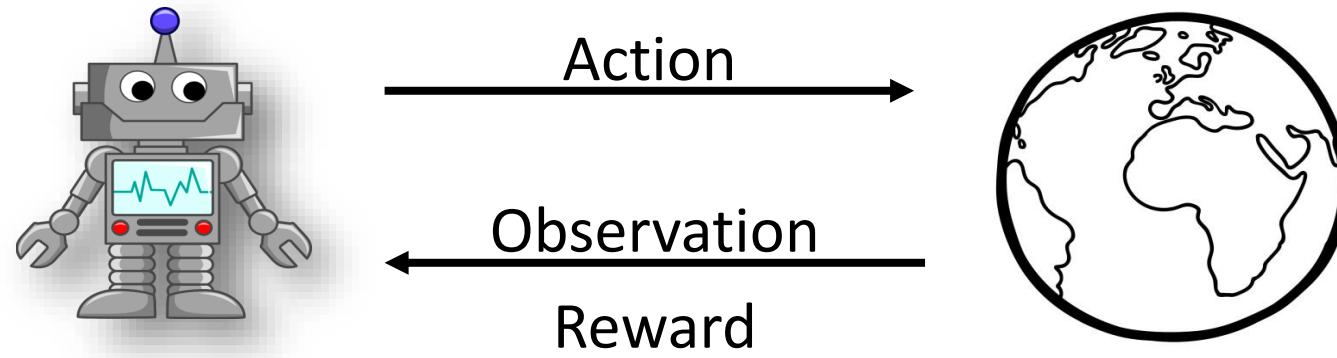
Reinforcement Learning

- **Setting/Assumptions:** Reinforcement learning (RL) involves training an agent to make decisions by interacting with an environment. The agent learns through trial and error (receiving rewards and penalties), optimizing its behavior to maximize cumulative rewards.
- **Goal:** To learn a policy that maps states of the environment to actions that achieve the highest reward.
- **Common Use Cases:**
 - Game-playing AI (e.g., AlphaGo, chess-playing bots).
 - Robotics (e.g., autonomous navigation).
 - Dynamic resource allocation (e.g., in networking or traffic management).
- **Examples:**
 - Q-learning, Deep Q-Networks (DQN), and Proximal Policy Optimization (PPO).

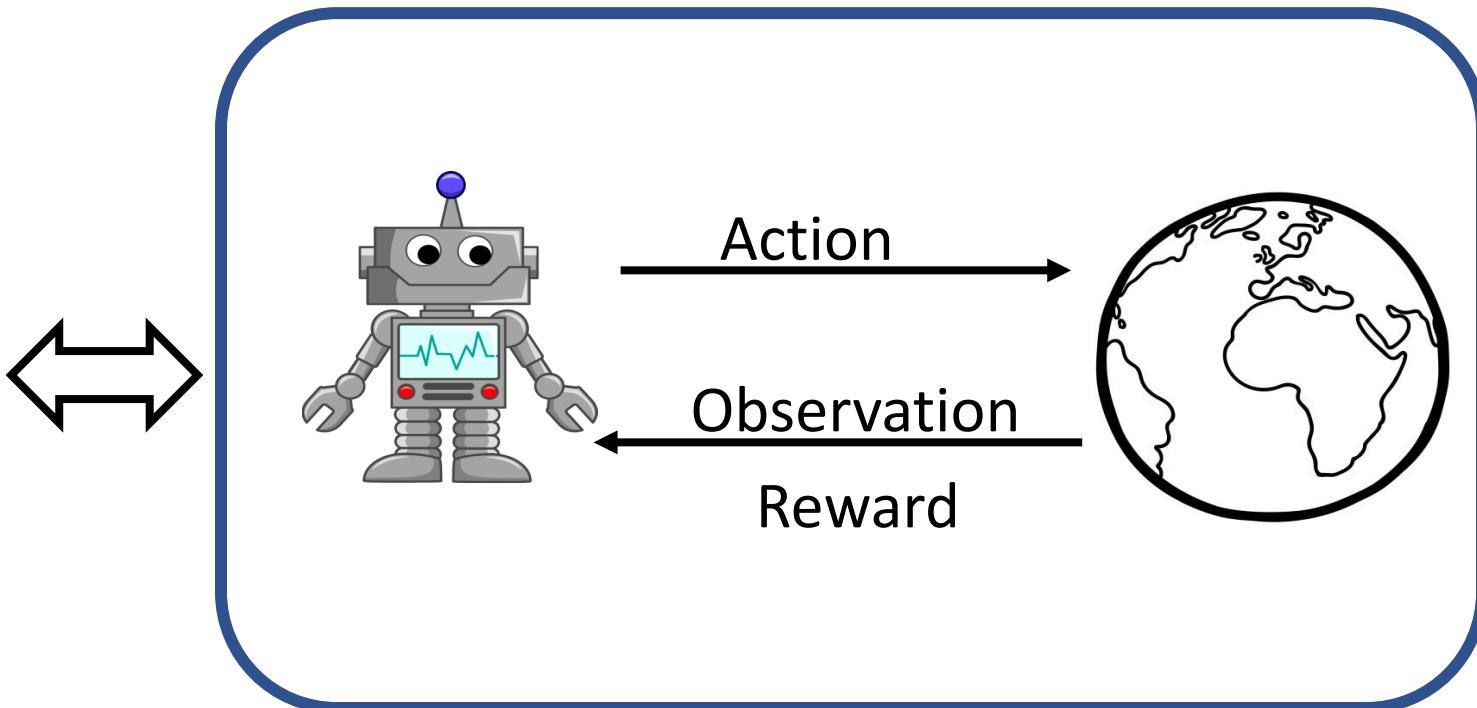
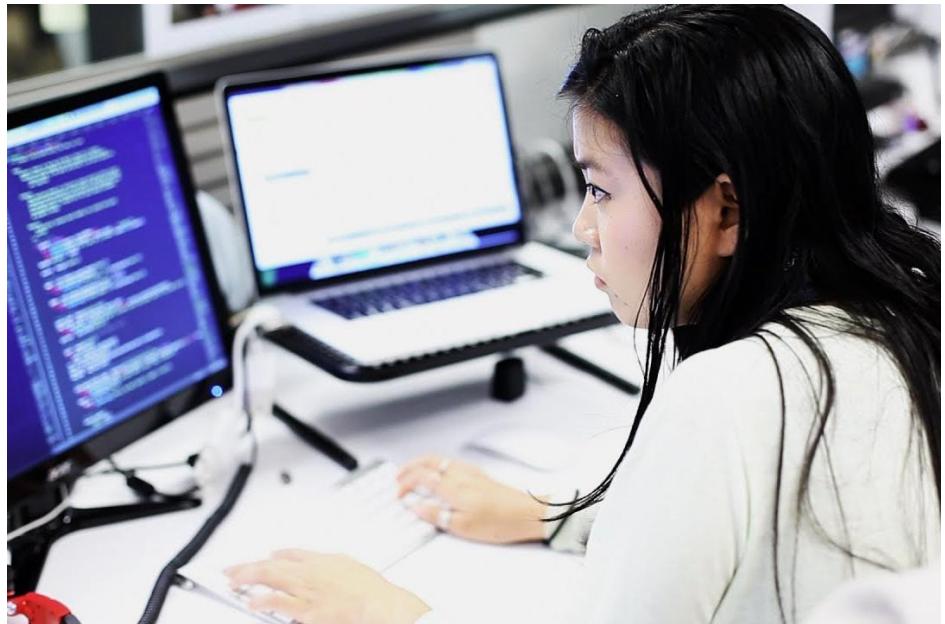
Reinforcement Learning



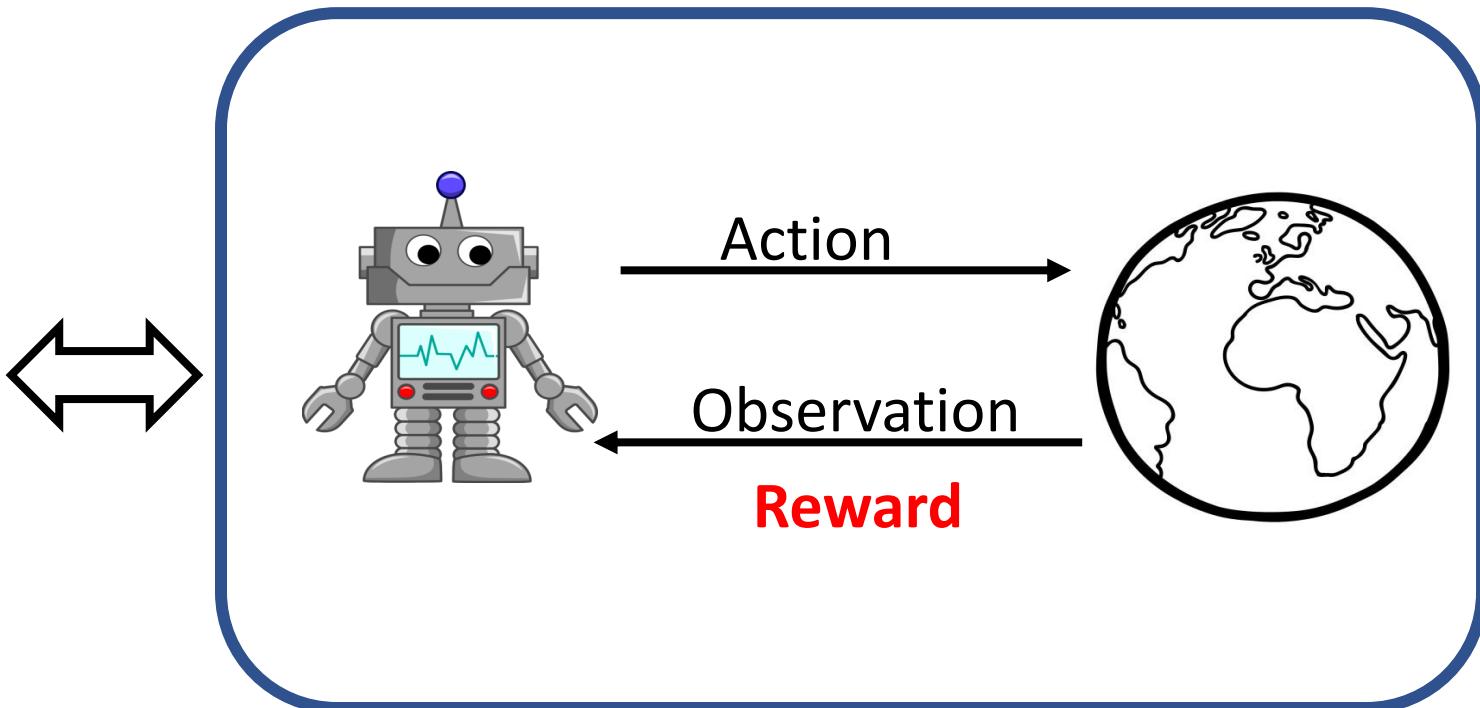
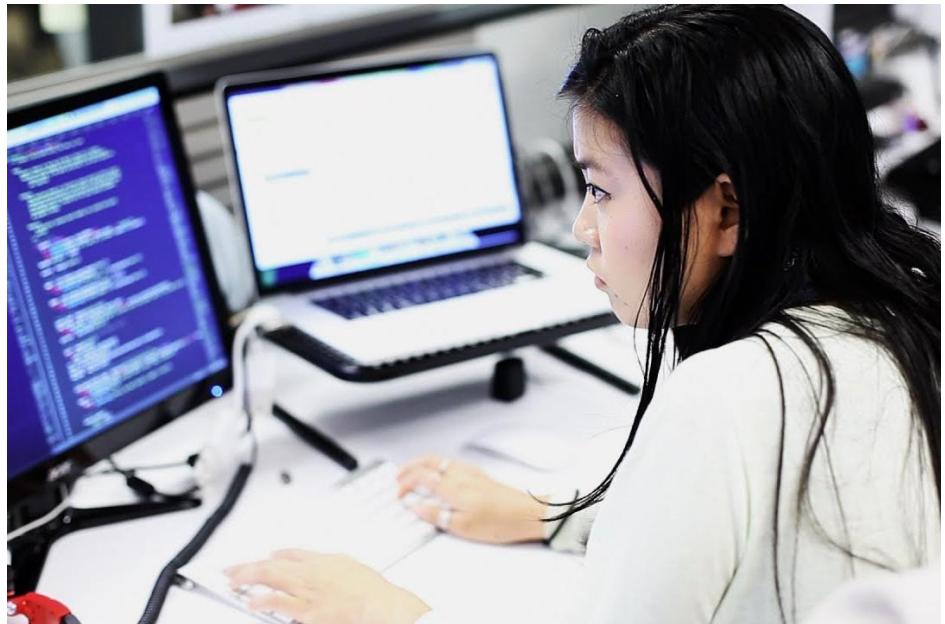
Reinforcement Learning



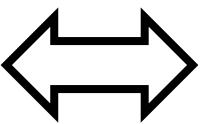
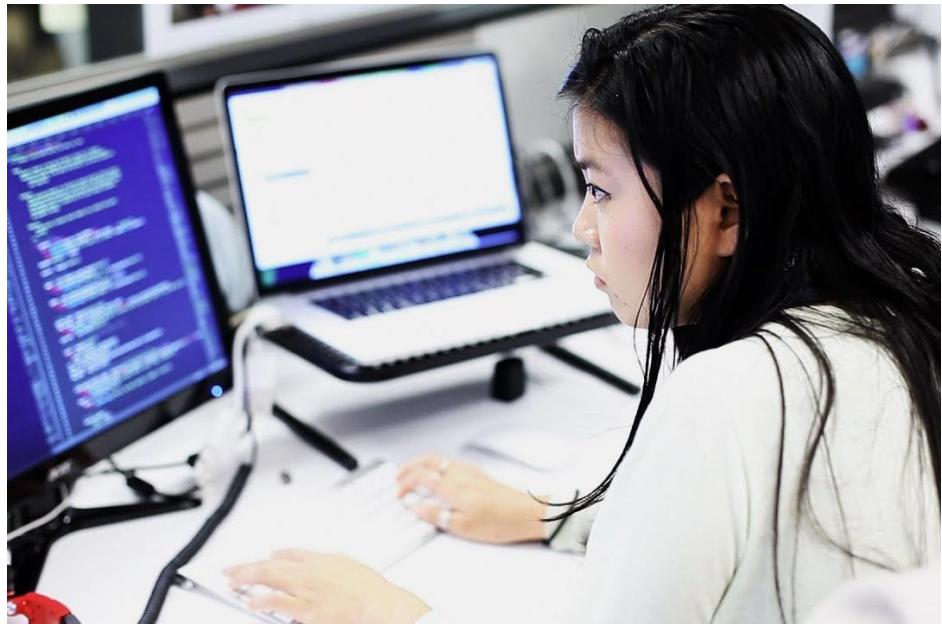
Reward engineering is hard!



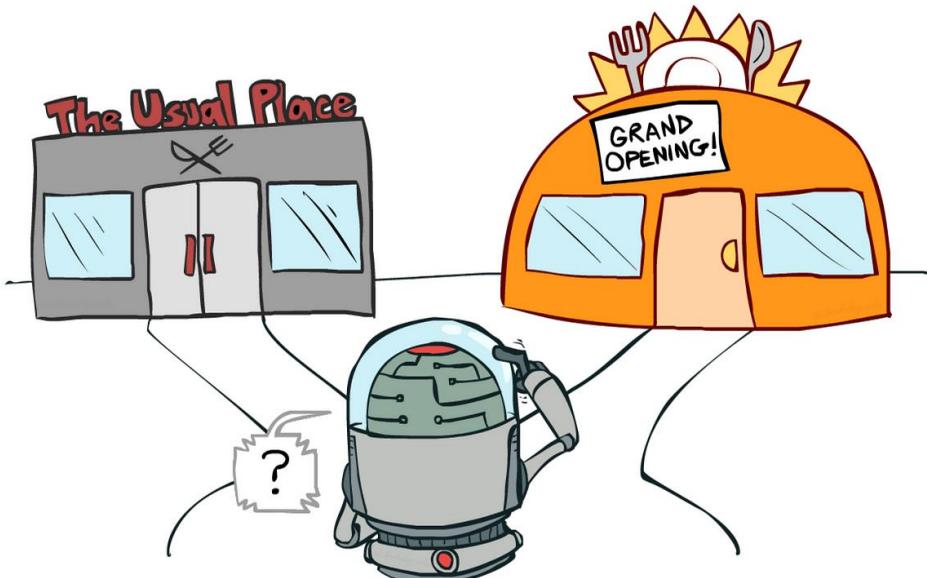
Reward engineering is hard!



Reward engineering is hard!

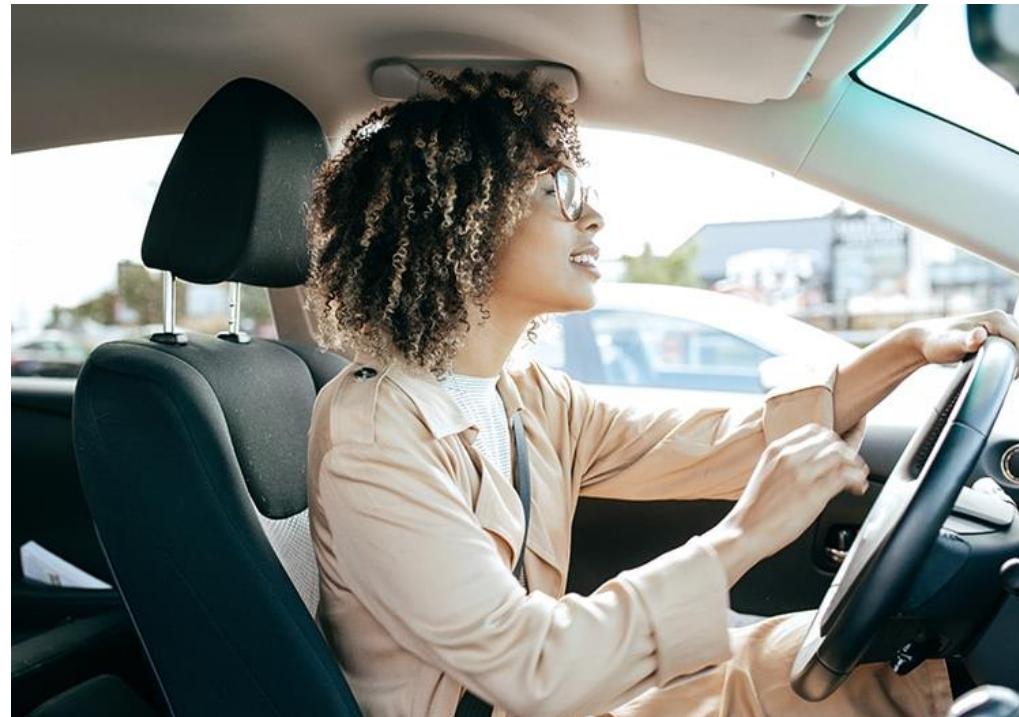


Reinforcement learning is hard...even with a reward function!



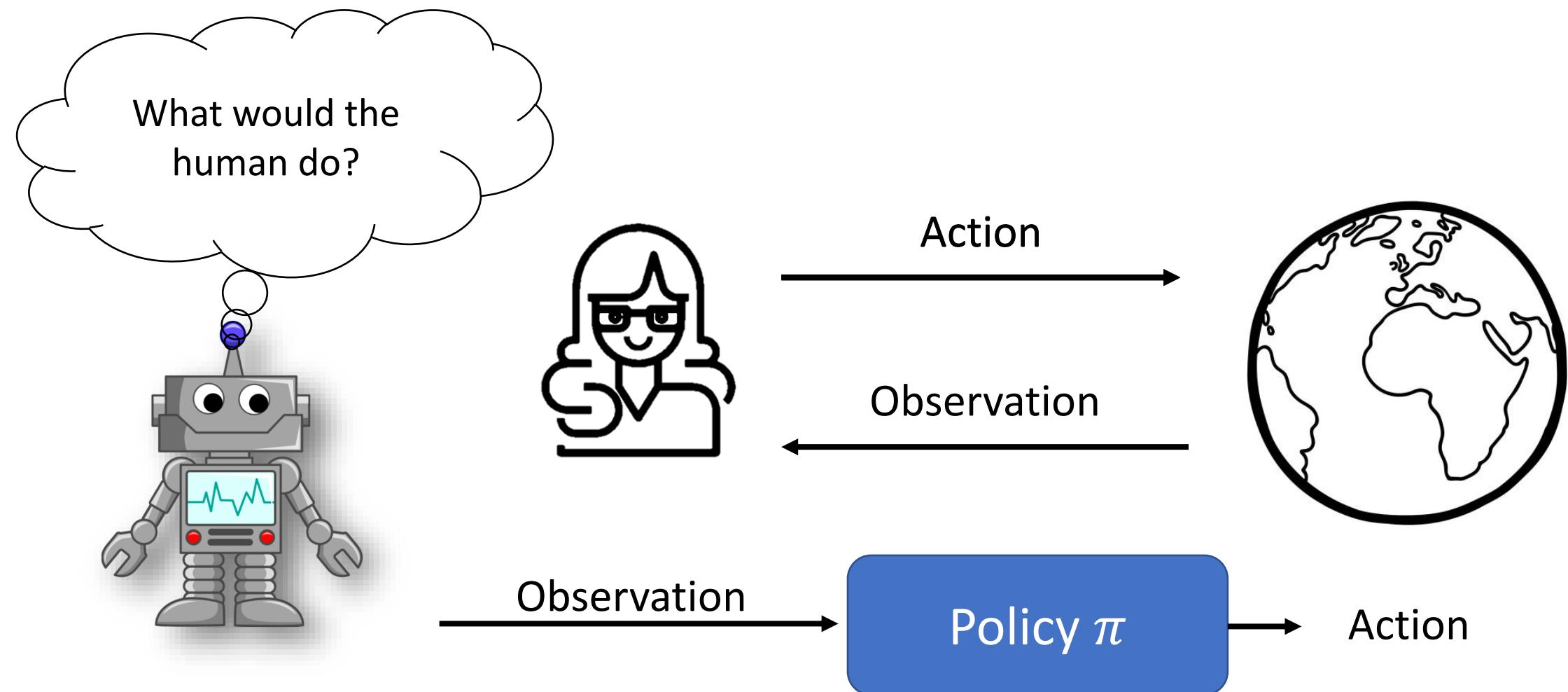
Imitation Learning:

Learn a policy from examples of good behavior.



- Often showing is easier than telling.
- Alleviates problem of exploration.

Behavioral Cloning



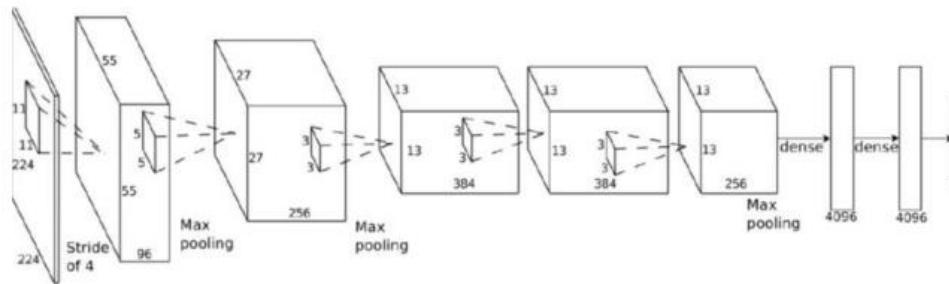
Inverse Reinforcement Learning



Imitation Learning via Behavioral Cloning



\mathbf{o}_t



\mathbf{a}_t



\mathbf{o}_t
 \mathbf{a}_t

training
data

supervised
learning

$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$

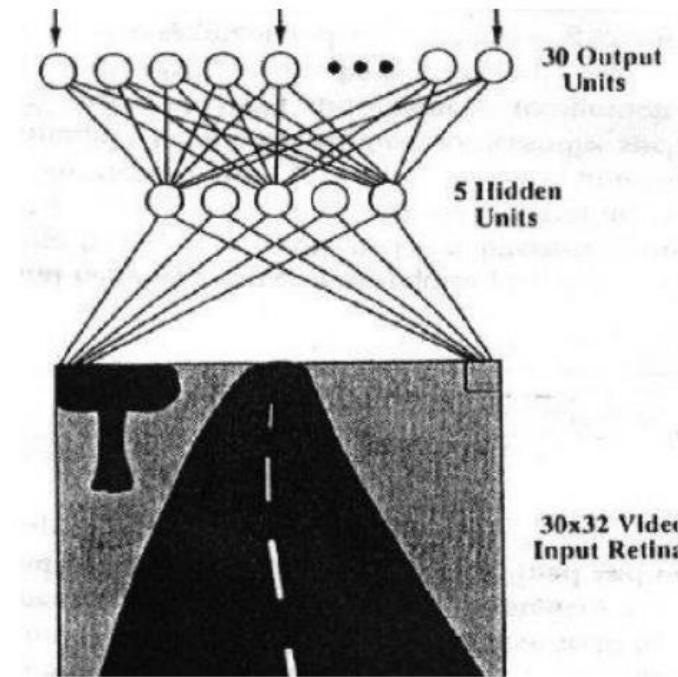
Live demo

`python test_gym.py`

`python mountain_car_bc.py --num_demos 1`

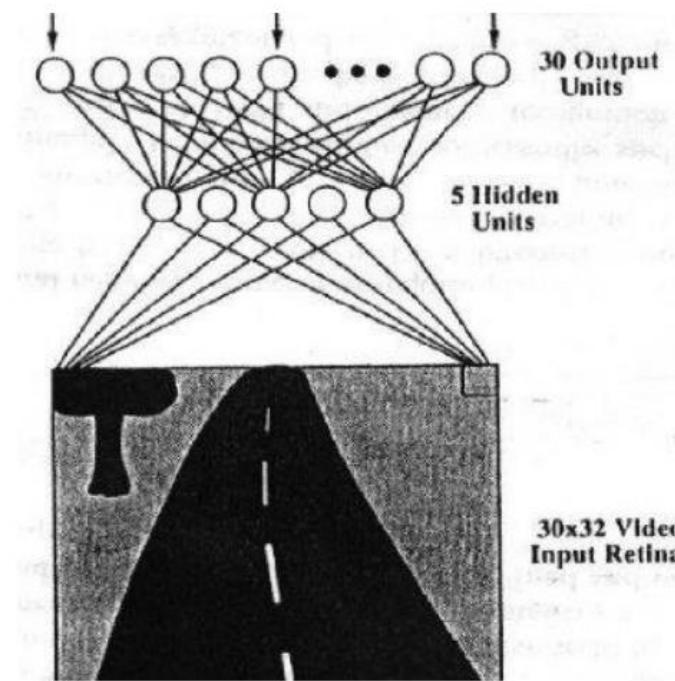
ALVINN: One of the first imitation learning systems

ALVINN: Autonomous Land Vehicle In a Neural Network
1989



ALVINN: One of the first imitation learning systems

ALVINN: Autonomous Land Vehicle In a Neural Network
1989

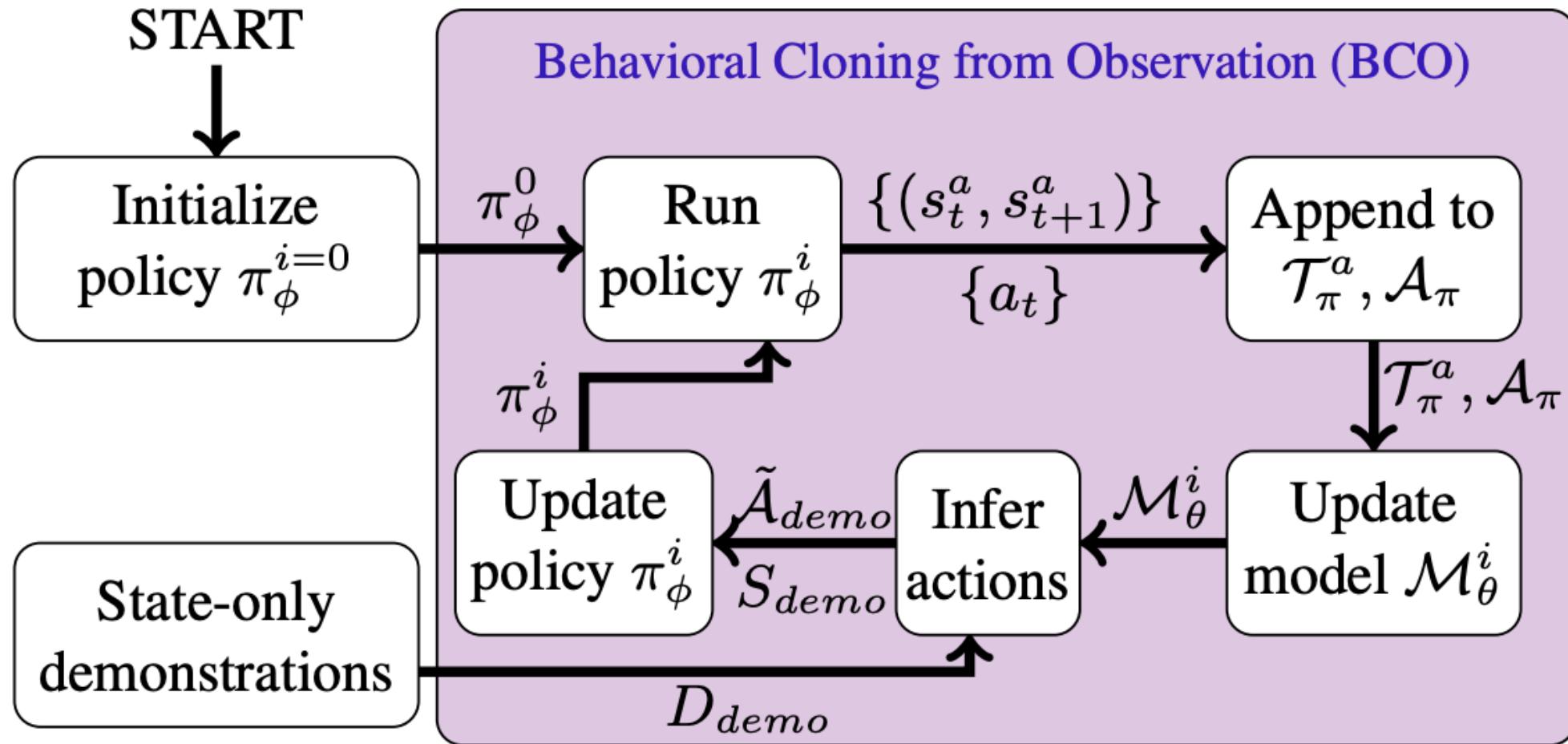


What if you don't have actions?

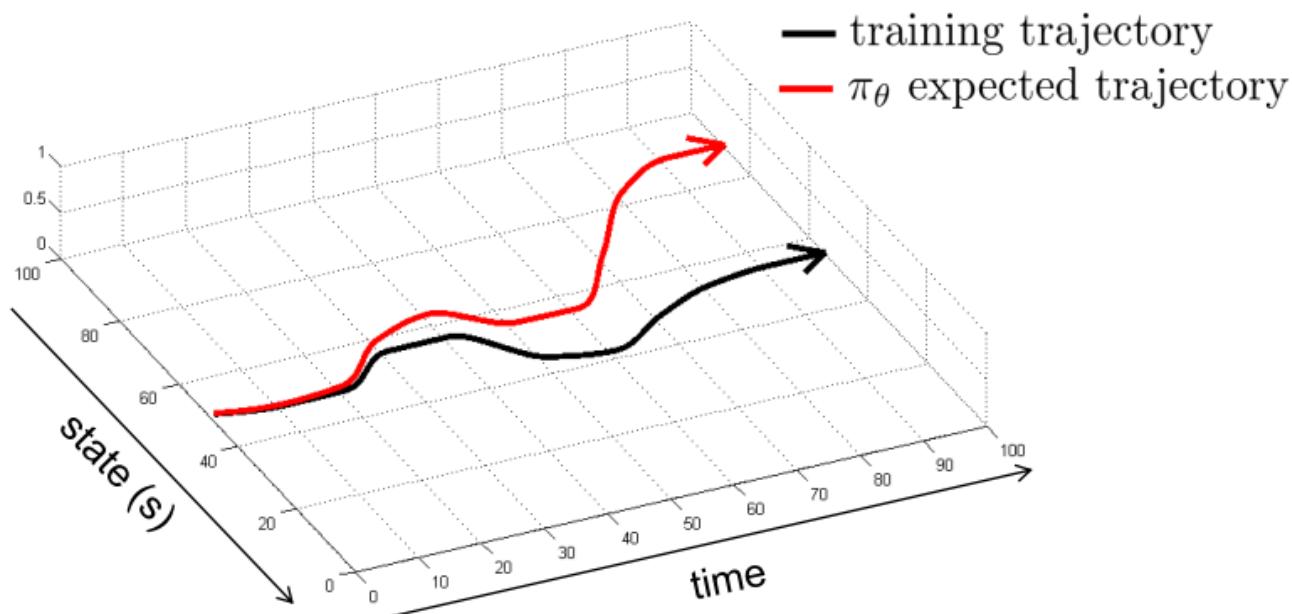
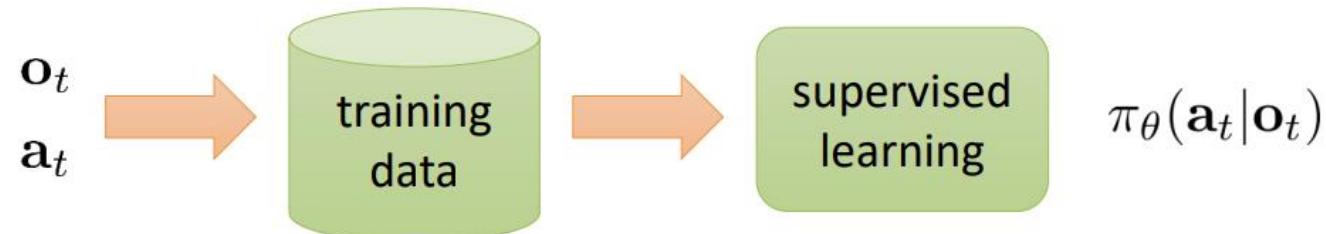
The screenshot shows a YouTube search results page for the query "tire change". The main video thumbnail features a man in a white shirt changing a tire on a silver car. Below the video, the title is "How to Change a Tire | Change a flat car tire step by step" and it has 3.5M views and is 15 years old. The video is from the channel "Howdini" with 714K subscribers. The video player shows the progress bar at 3:17 / 5:34. To the right of the main video, there is an advertisement for the TV show "Ahsoka" on Disney+. Below the ad, several related video thumbnails are displayed:

- "How to change a tire | Dad, how do I?" by Pushing Pistons (13:24, 886K views, 3 years ago)
- "How to Replace your Flat Tire" by Pushing Pistons (0:57, 117K views, 2 years ago)
- "Steer Tire Change" by The Tire Doctor (1:01, 571K views, 1 year ago)
- "Winter Tire Swap" by Your Home Garage (12:47, 48K views, 2 years ago)
- "How to Plug a Flat Tire (easily)" by ChrisFix (1:00, 922K views, 9 months ago)
- "POV Tire Change with your Dad #shorts" by Charlie Berens (7.4M views, < 6 months ago)

Behavioral Cloning from Observation (Torabi et al. 2018)

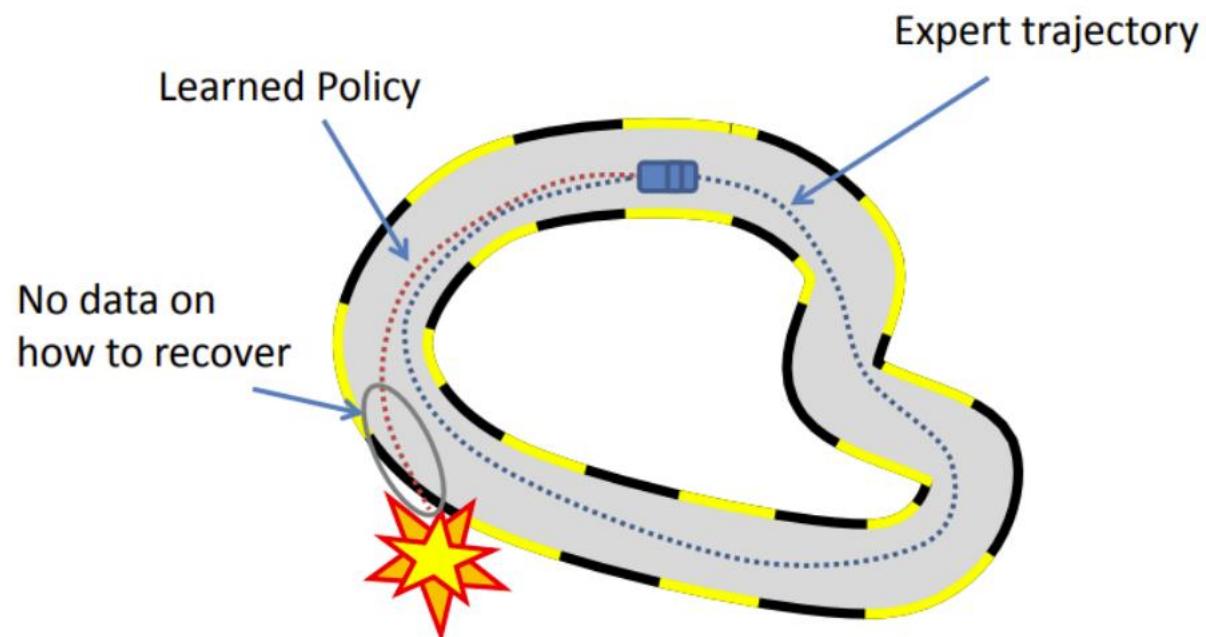


What could go wrong?



Distribution Shift

$$p_{\pi^*}(o_t) \neq p_{\pi_\theta}(o_t)$$

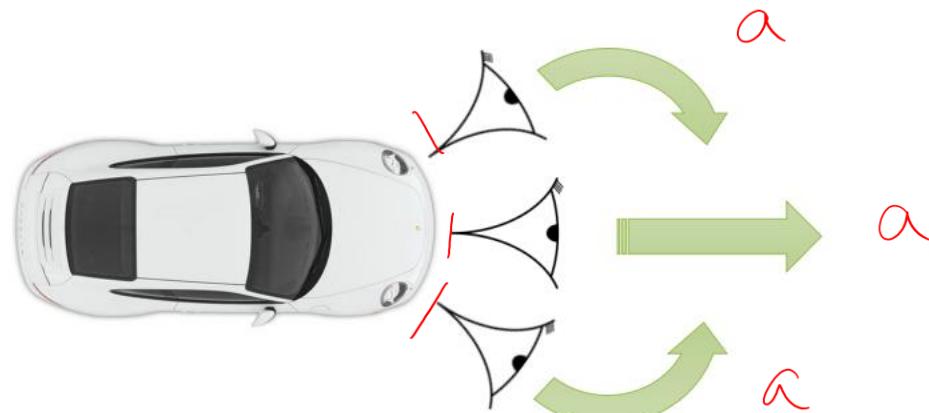
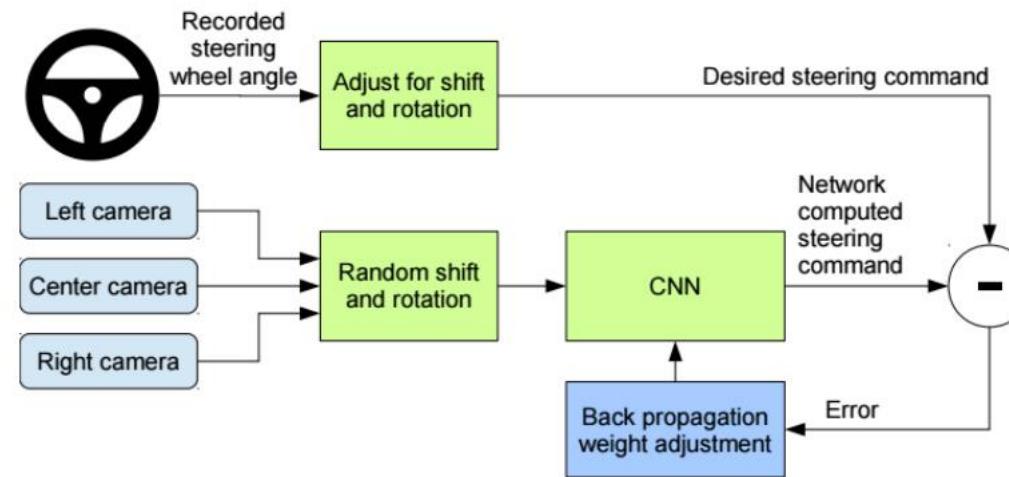


	Supervised Learning	Supervised Learning + Control
Train	$(x, y) \sim D$	$s \sim P(\cdot s, \pi^*(s))$
Test	$(x, y) \sim D$	$s \sim P(\cdot s, \pi(s))$

But it still can work in practice...

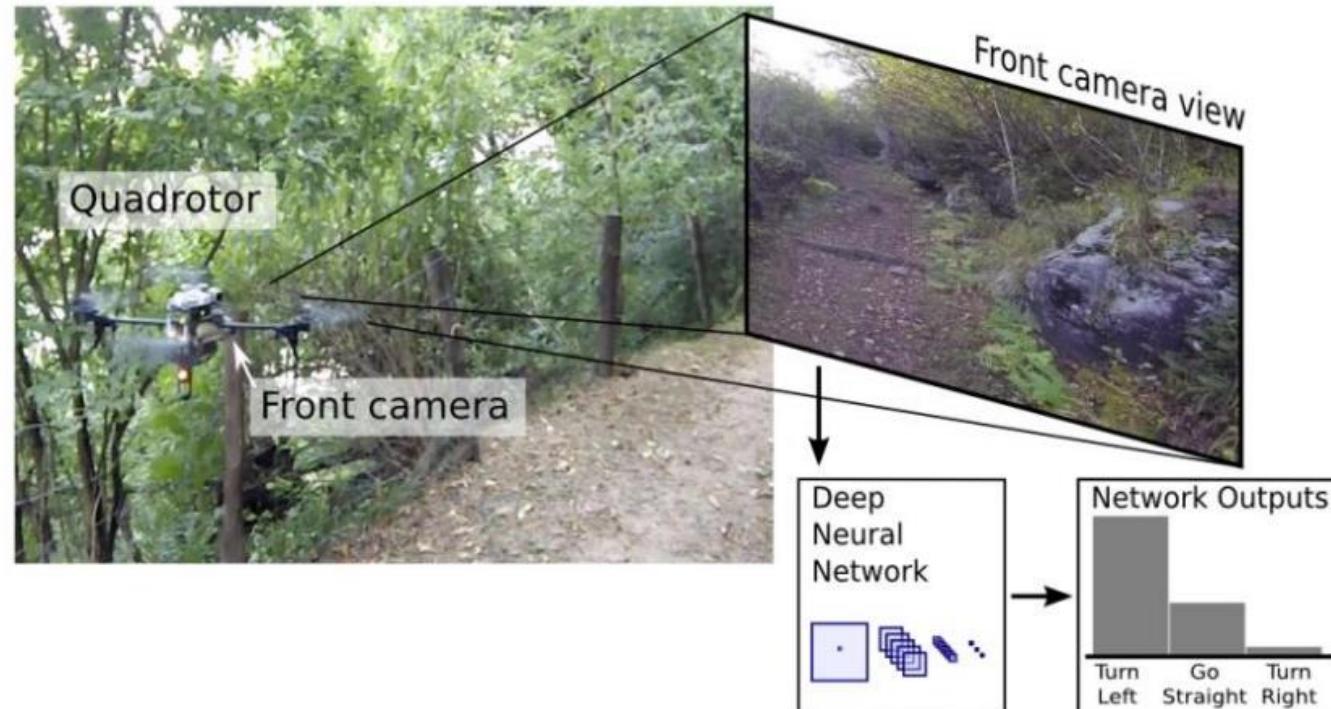


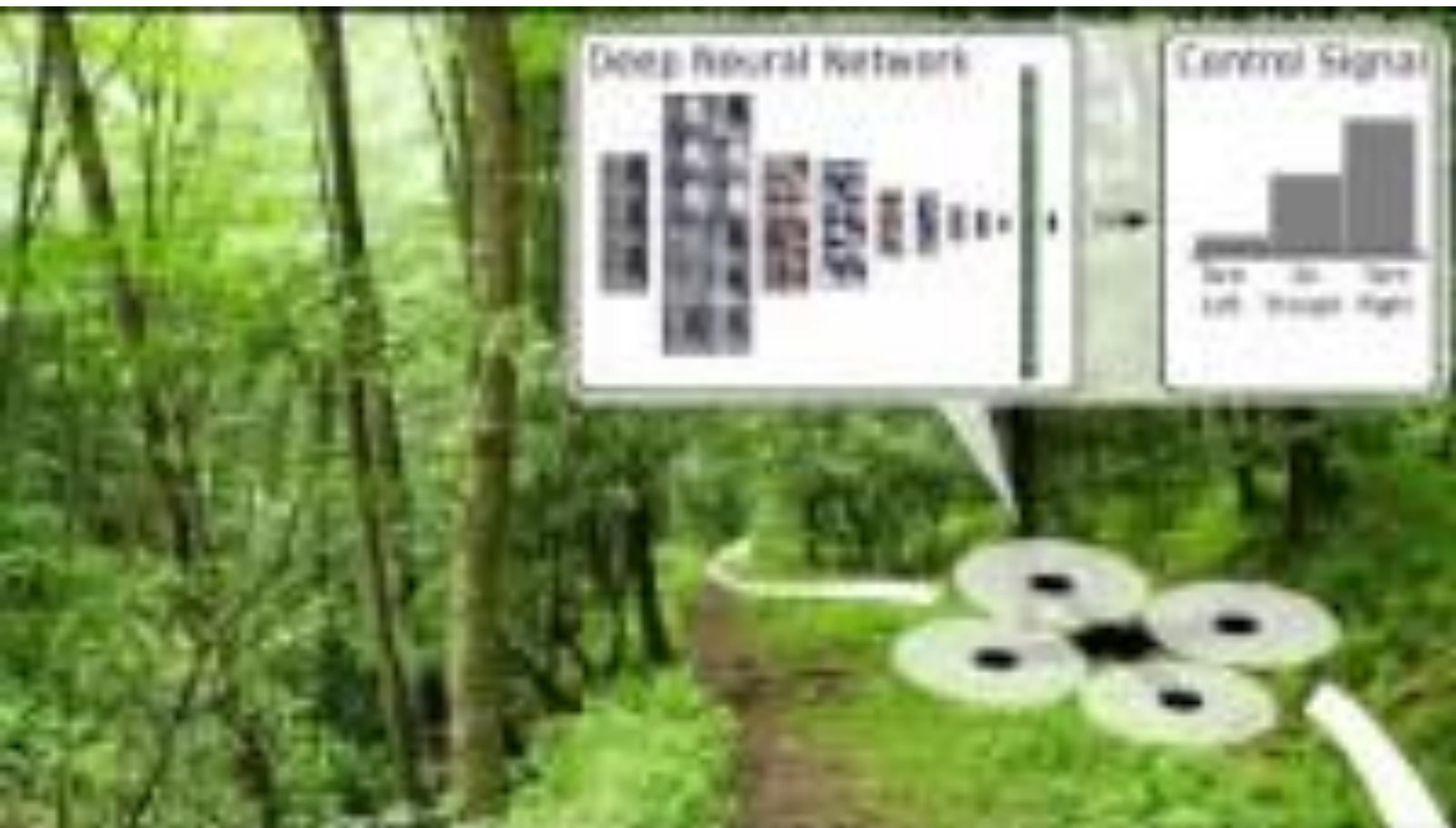
How?



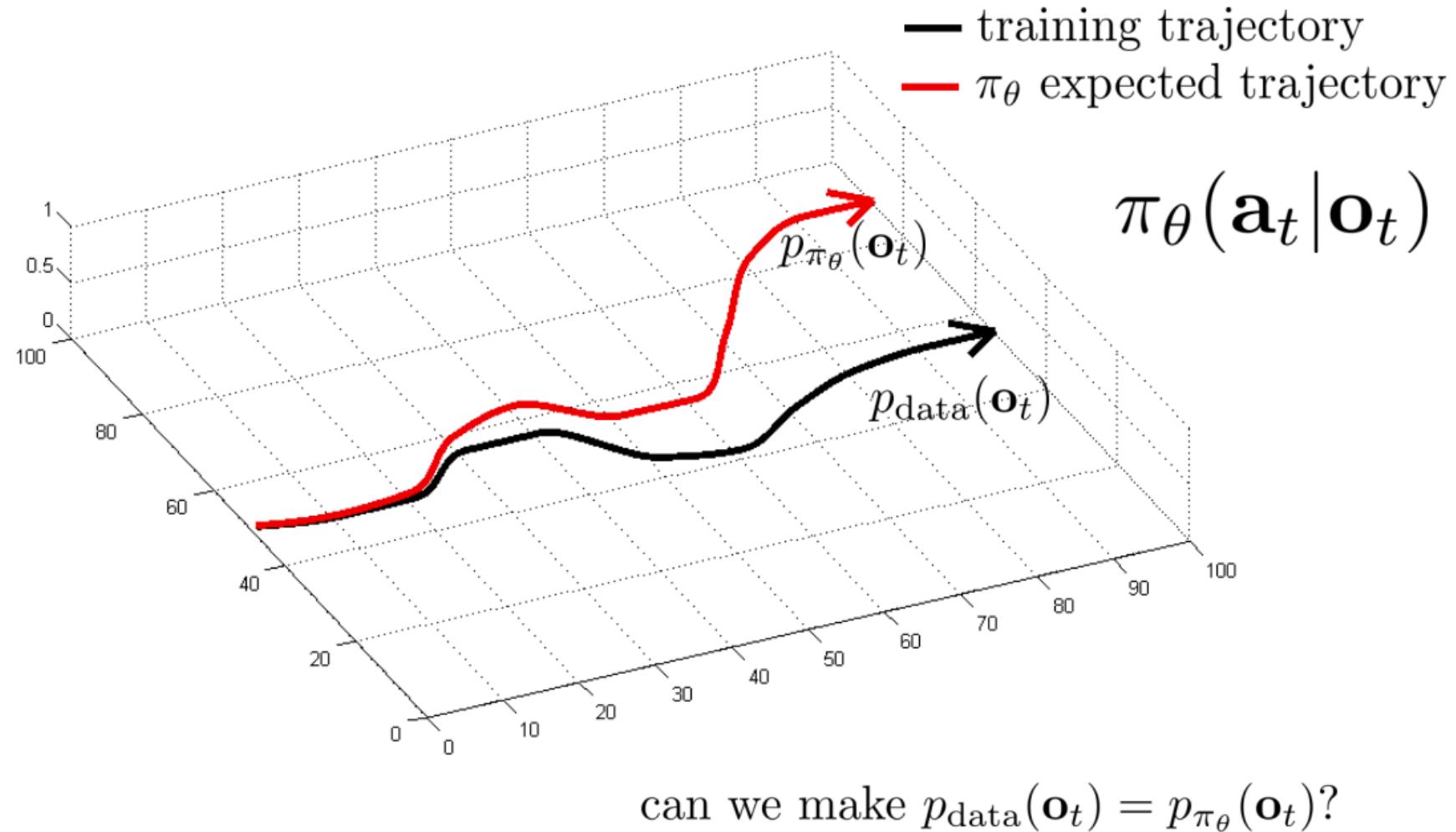
A Machine Learning Approach to Visual Perception of Forest Trails for Mobile Robots

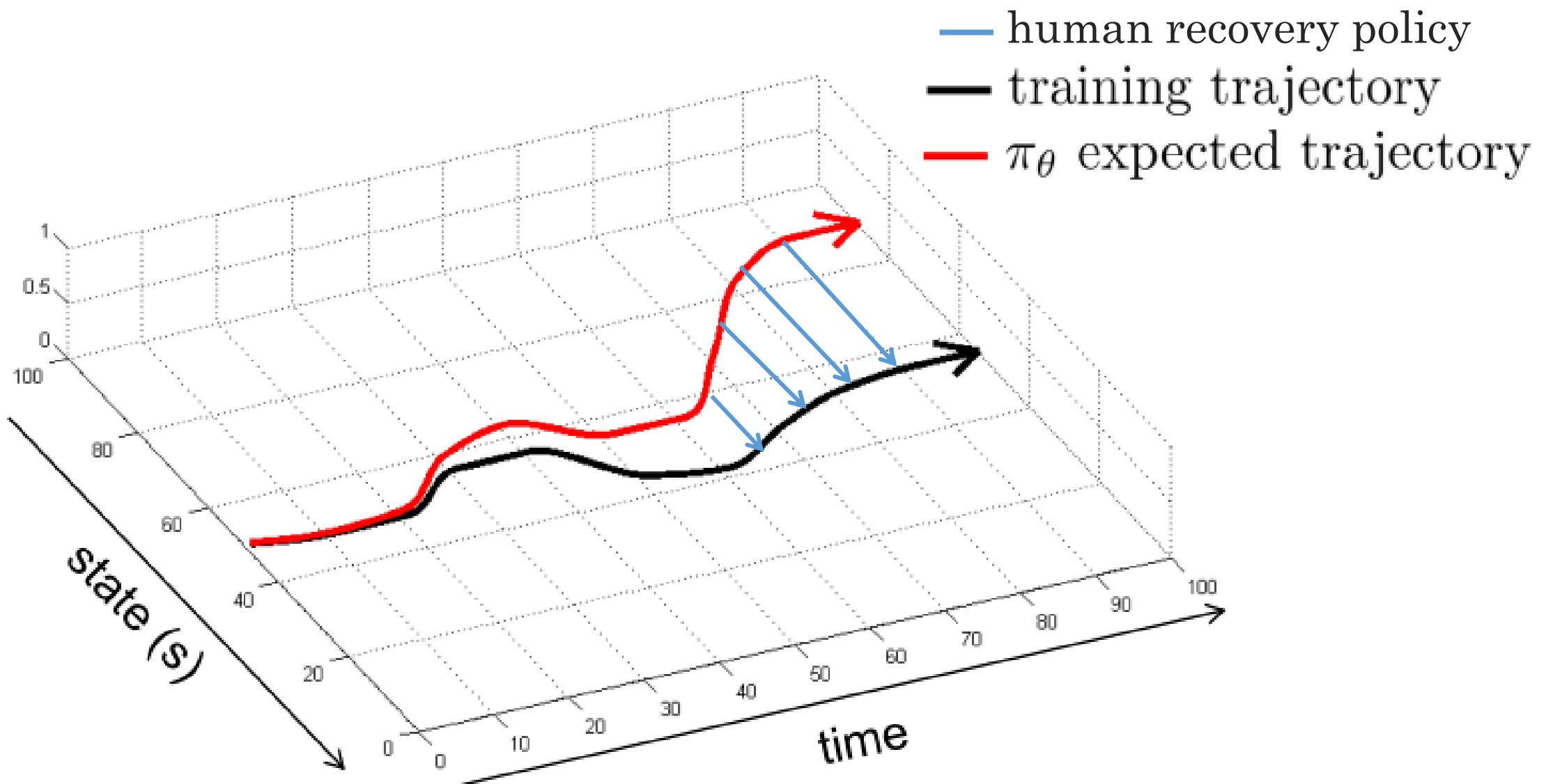
Alessandro Giusti¹, Jérôme Guzzi¹, Dan C. Cireşan¹, Fang-Lin He¹, Juan P. Rodríguez¹
Flavio Fontana², Matthias Faessler², Christian Forster²
Jürgen Schmidhuber¹, Gianni Di Caro¹, Davide Scaramuzza², Luca M. Gambardella¹





Can we make it work more often?





DAgger

can we make $p_{\text{data}}(\mathbf{o}_t) = p_{\pi_\theta}(\mathbf{o}_t)$?

idea: instead of being clever about $p_{\pi_\theta}(\mathbf{o}_t)$, be clever about $p_{\text{data}}(\mathbf{o}_t)$!

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

how? just run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$

but need labels \mathbf{a}_t !

- 
1. train $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

DAgger has very nice theoretical guarantees.

Why might it be **hard** to implement in practice?

DAgger: Dataset Aggregation

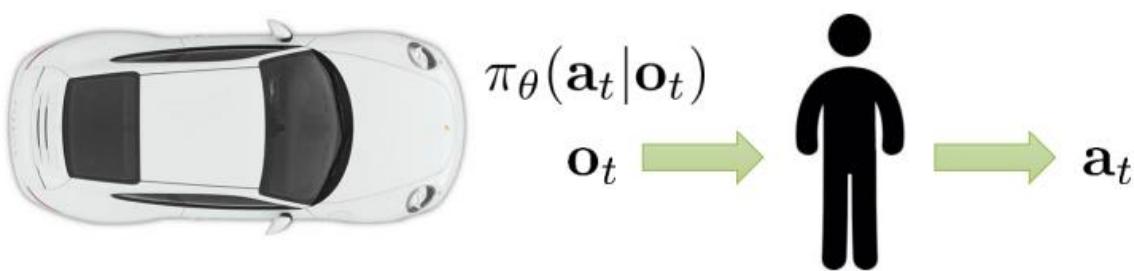
goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

how? just run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$

but need labels \mathbf{a}_t !

- 
1. train $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t | \mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$



DAgger has very nice theoretical guarantees.

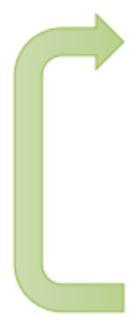
Why might it be **easy** to implement in practice?

DAgger: Dataset Aggregation

goal: collect training data from $p_{\pi_\theta}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

how? just run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$

but need labels \mathbf{a}_t !

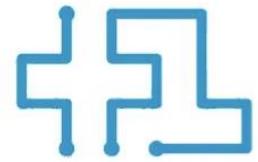
- 
1. train $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_\theta(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_\pi = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 3. Ask human to label \mathcal{D}_π with actions \mathbf{a}_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

Learn from an Algorithmic Supervisor!



But we don't always have access to an algorithmic supervisor...

Can we make DAgger more practical when dealing with real human labeling?



PLUS ONE
ROBOTICS



ZOOX

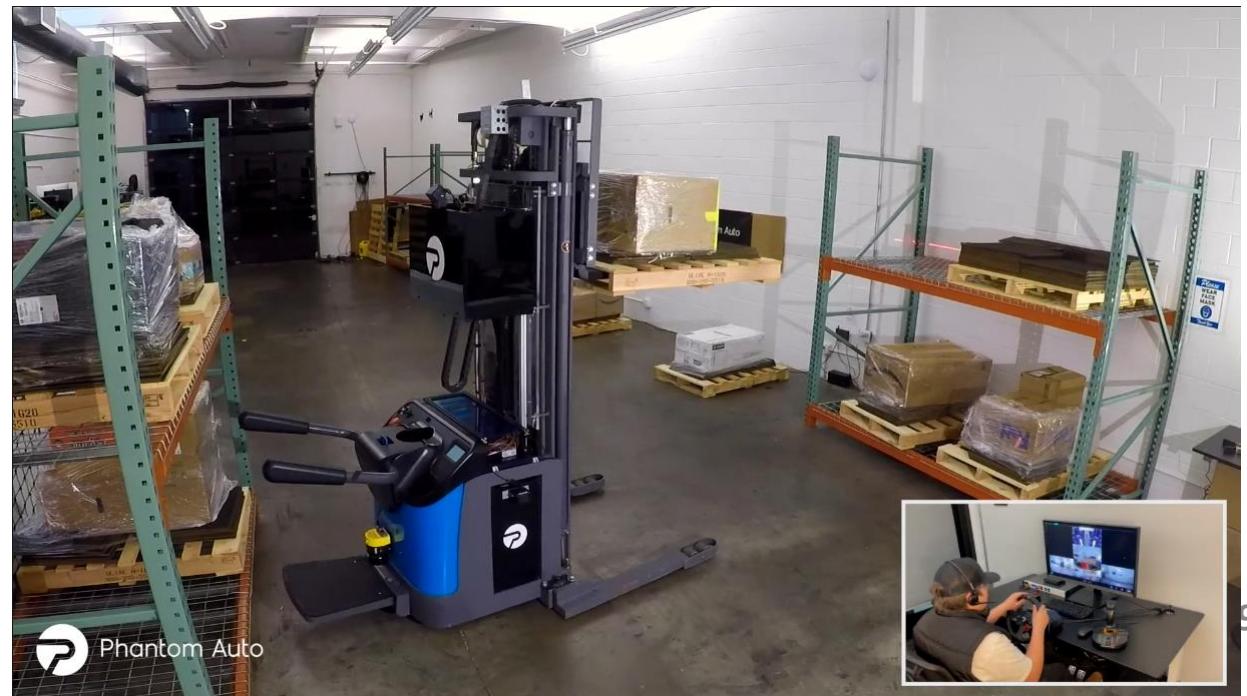


nimble

WAYMO

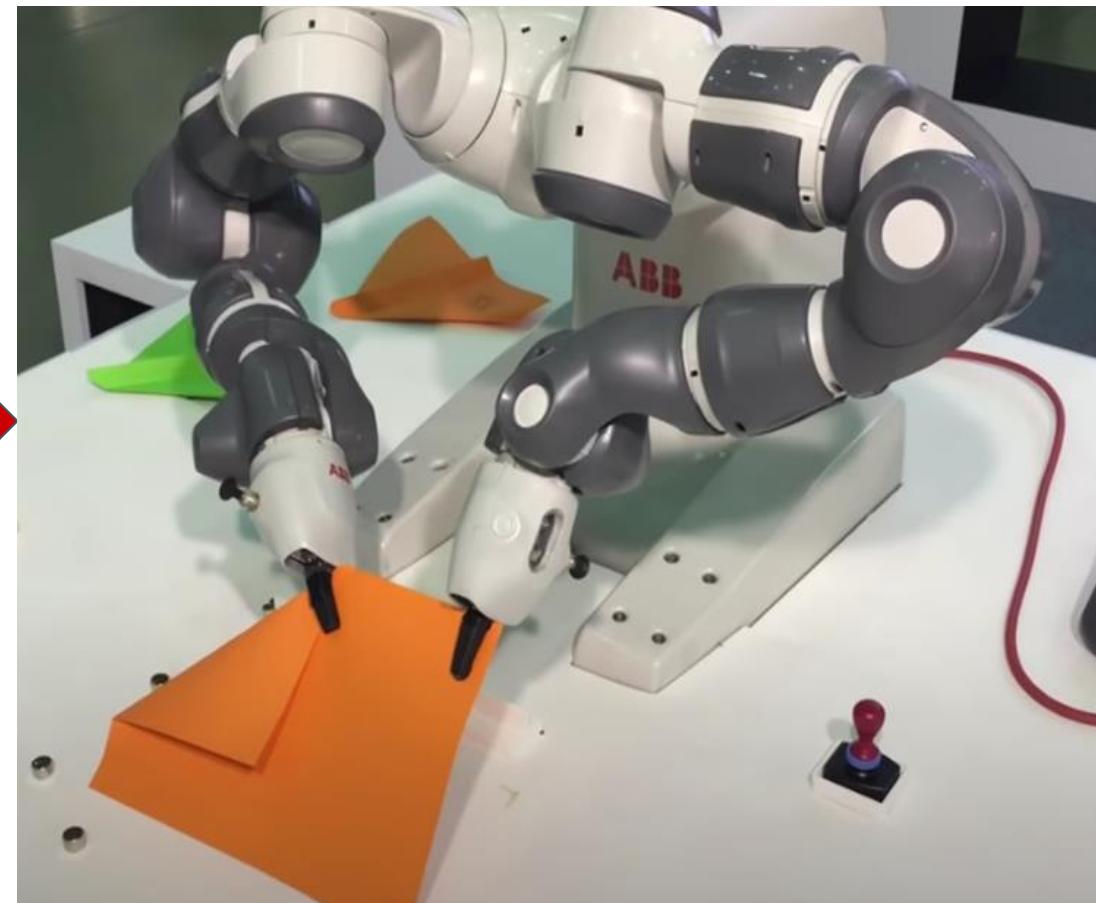


Phantom Auto



Phantom Auto

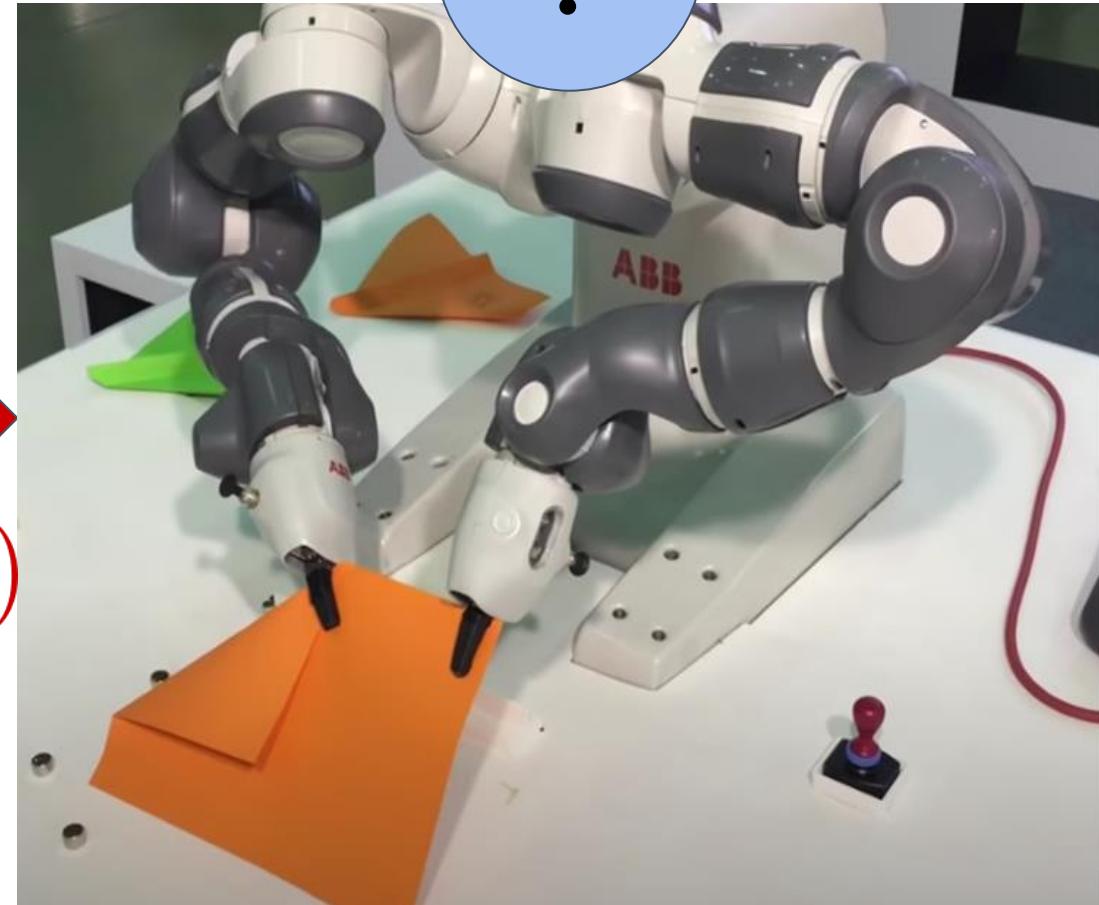
Interactive IL



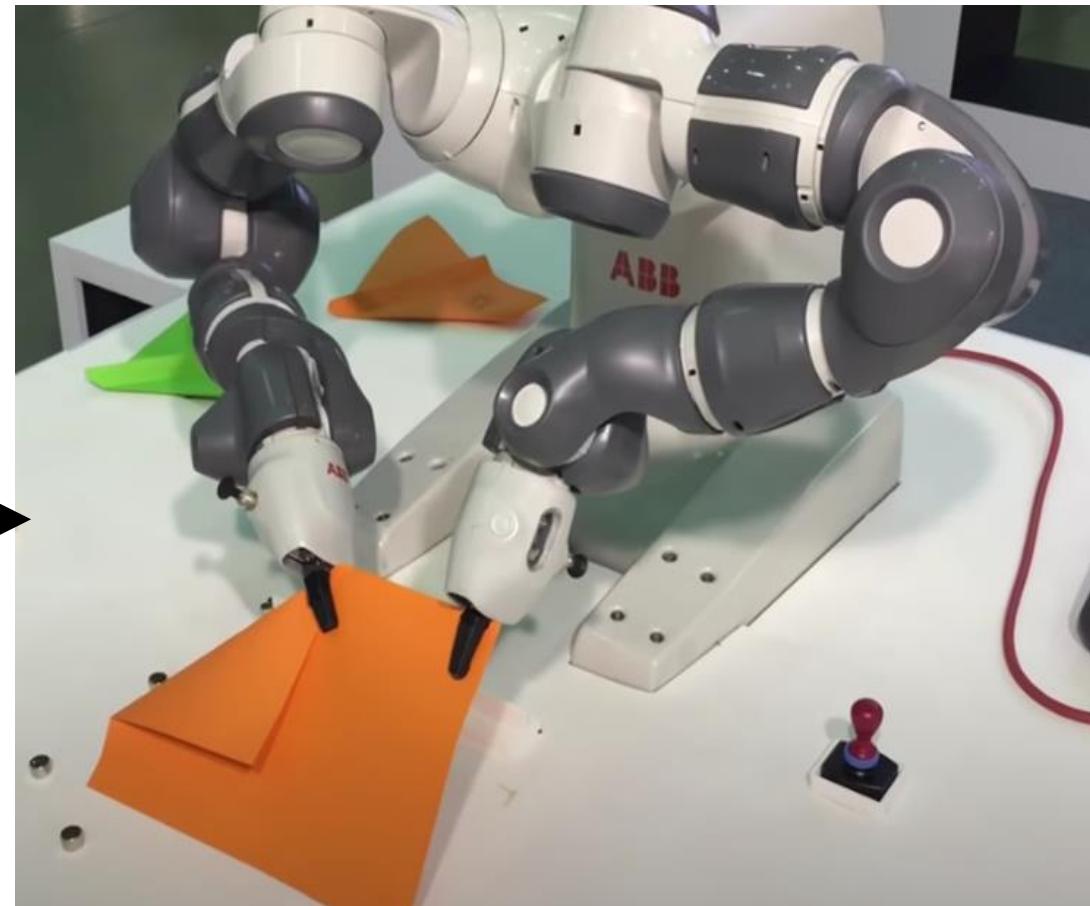
Interactive IL

 $\pi_H(s)$

↔
 $\pi_{\text{meta}}(s)$
???

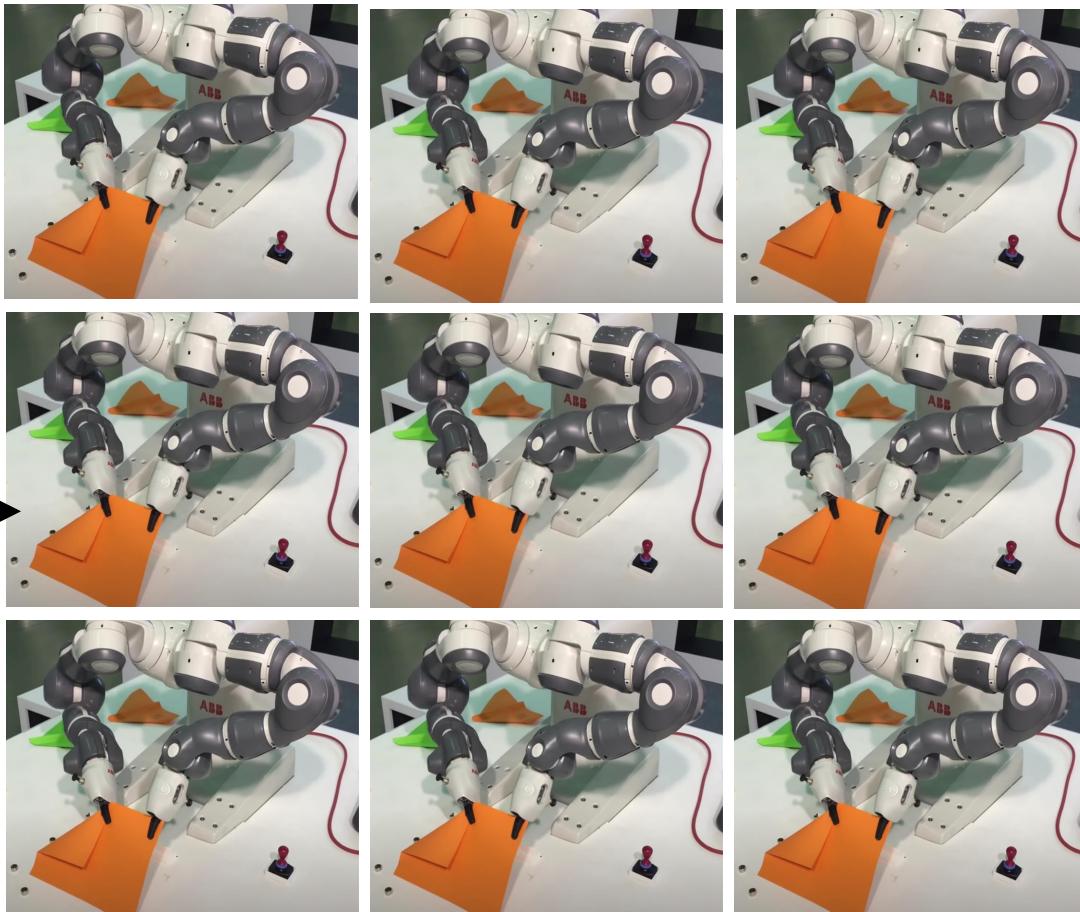
 $\pi_R(s)$

Human-Gated Interactive IL



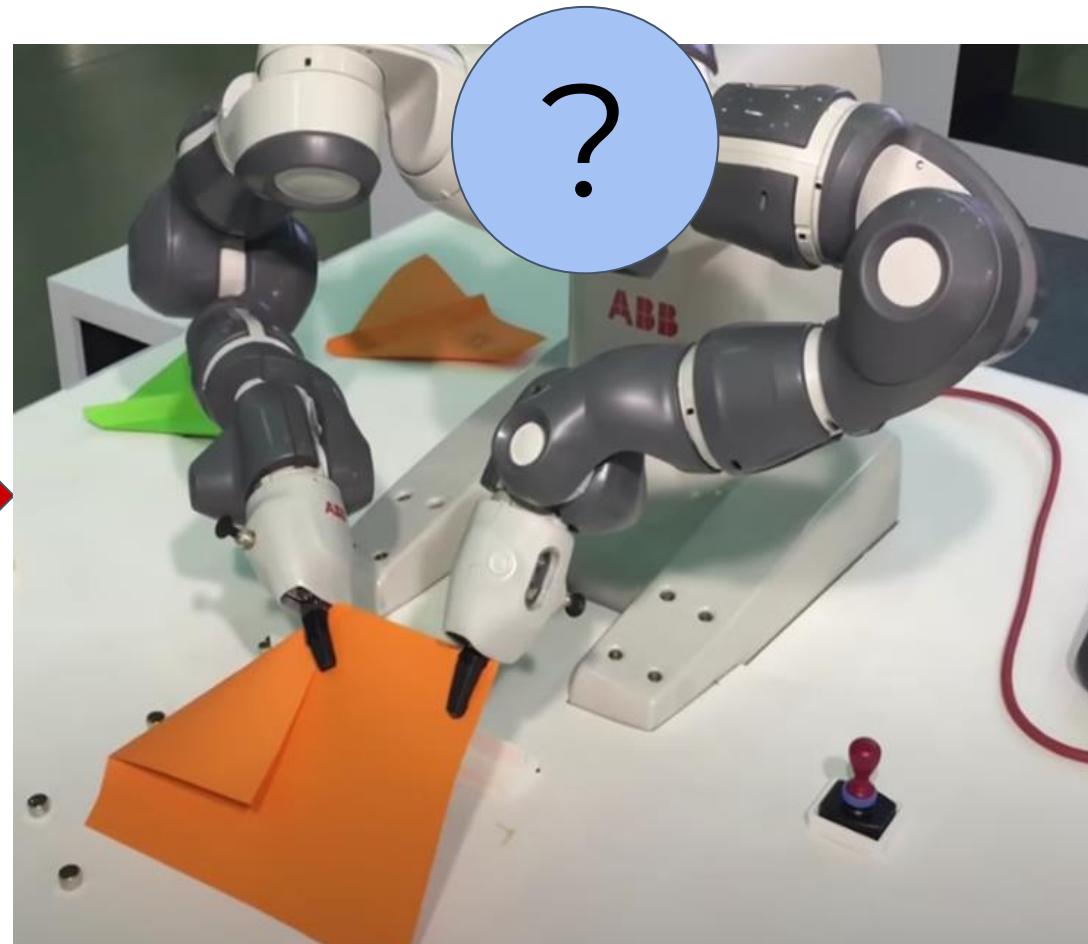
[3] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. HG-Dagger: Interactive Imitation Learning with Human Experts. ICRA 2019.

Human-Gated Interactive IL



[3] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. HG-Dagger: Interactive Imitation Learning with Human Experts. ICRA 2019.

Robot-Gated Interactive IL



- [4] J. Zhang, K. Cho. Query-Efficient Imitation Learning for End-to-End Autonomous Driving. AAAI 2017.
- [5] K. Menda, K. Driggs-Campbell, M. Kochenderfer. EnsembleDagger: A Bayesian Approach to Safe Imitation Learning. IROS 2019.

Minimizing Supervisor Burden

- C = Number of context switches
- L = Latency of context switching
- I = Expected number of supervisor actions per intervention

$$B(\pi) \triangleq C(\pi) \cdot (L + I(\pi))$$

Ideally, we want

$$\begin{aligned} \pi &= \arg \min_{\pi' \in \Pi} L(\pi'_r) \\ \text{s.t. } B(\pi') &\leq \Gamma_b \end{aligned}$$

Minimizing Supervisor Burden

- C = number of context switches
- L = Latency of context switching
- I = expected number of supervisor actions per intervention

$$B(\pi) \triangleq \textcolor{orange}{C}(\pi) \cdot (L + I(\pi))$$

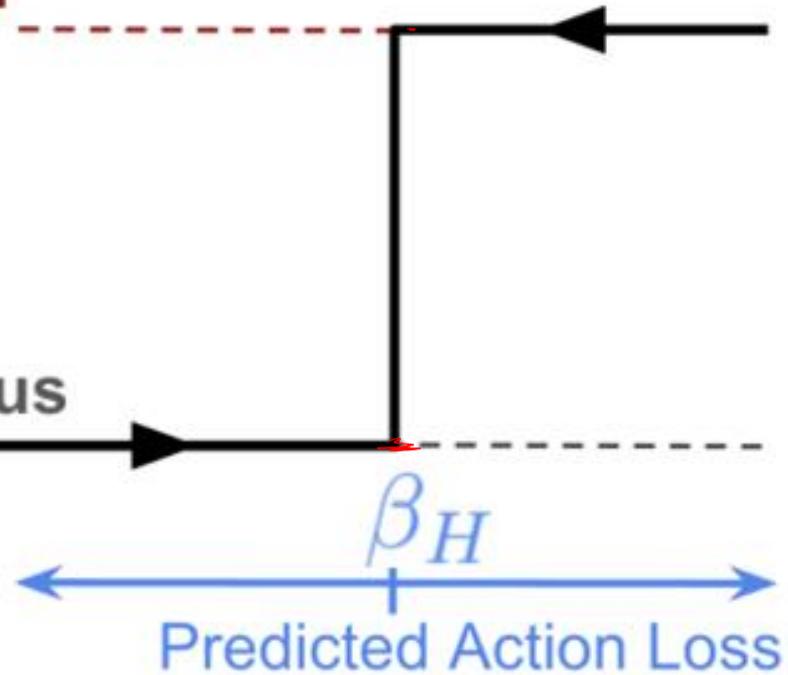
In practice, we approximate this by focusing on limiting the number of interventions (number of context switches)

$$\begin{aligned}\pi &= \arg \min_{\pi' \in \Pi} L(\pi'_r) \\ s.t. \quad B(\pi') &\leq \Gamma_b\end{aligned}$$

SafeDAgger

**Supervisor
Mode**

**Autonomous
Mode**



Predicted action loss = predicted difference between human and robot action.

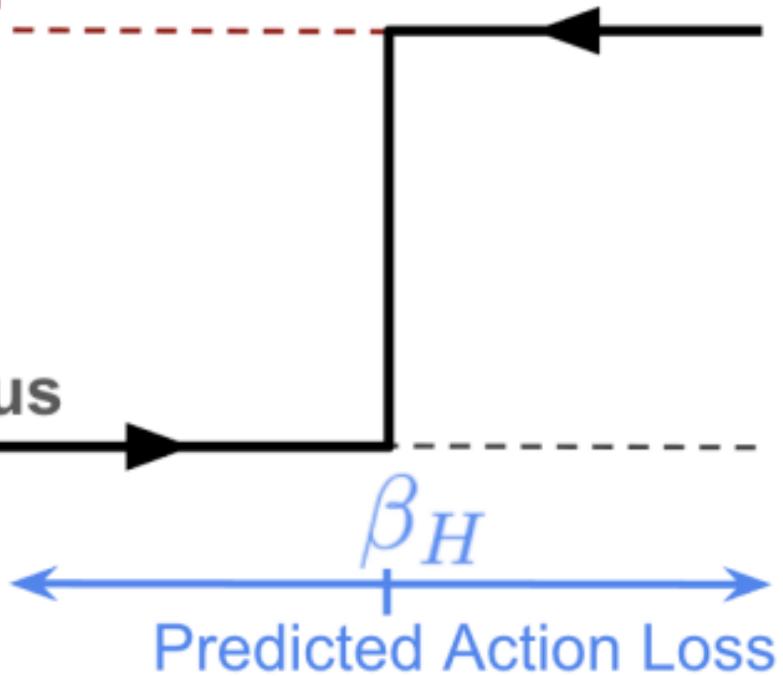
Trained using held-out set of data from human.



SafeDAgger

Supervisor
Mode

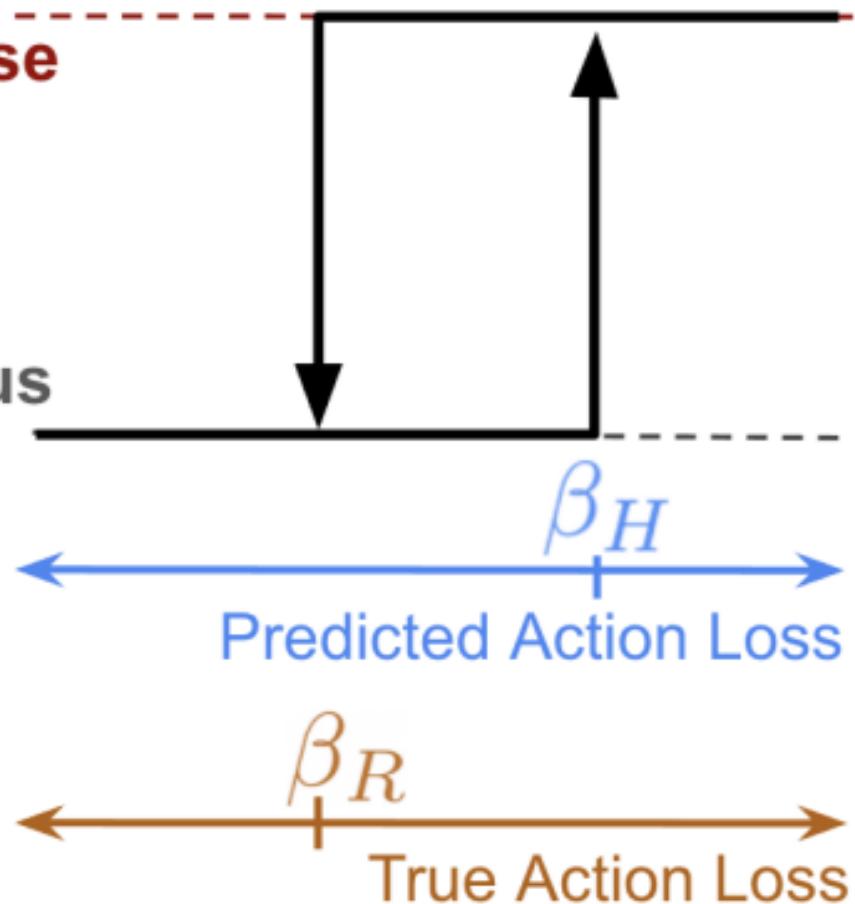
Autonomous
Mode



LazyDAgger

Supervisor
Mode + Noise

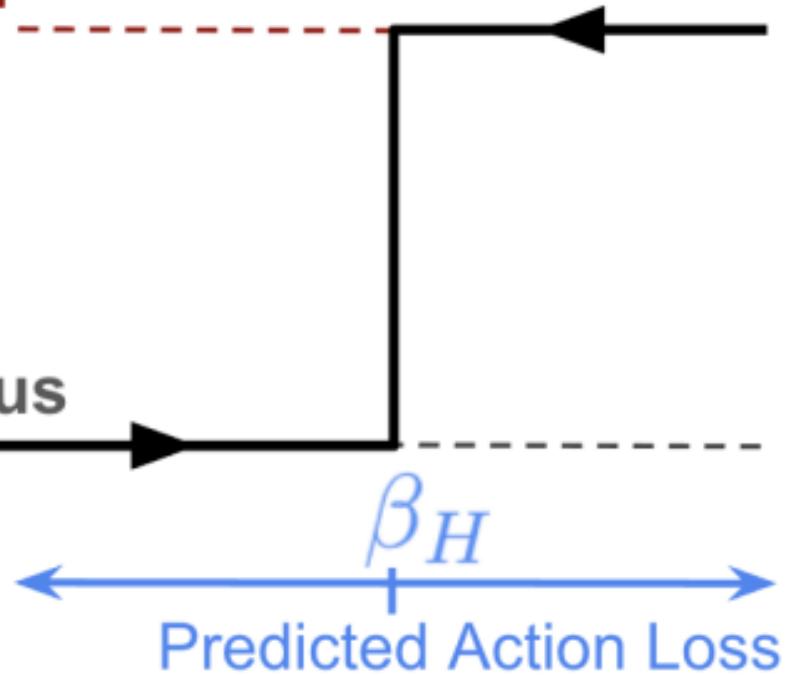
Autonomous
Mode



SafeDAgger

Supervisor
Mode

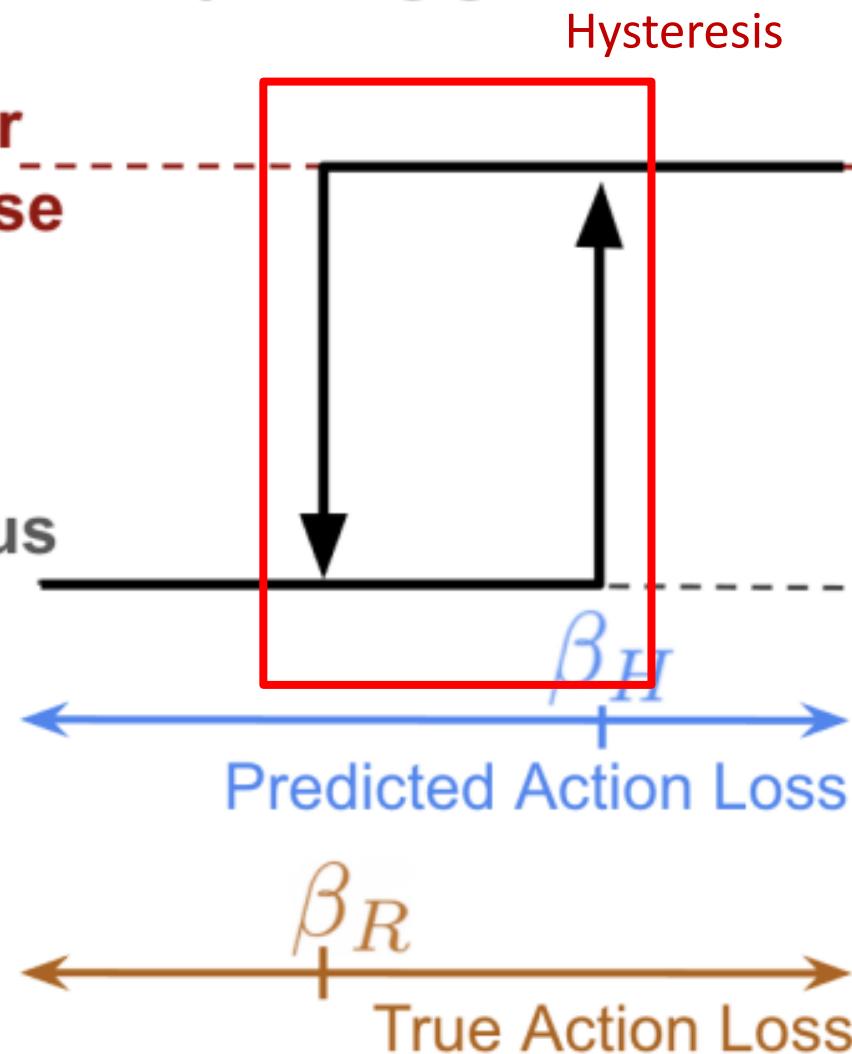
Autonomous
Mode



LazyDAgger

Supervisor
Mode + Noise

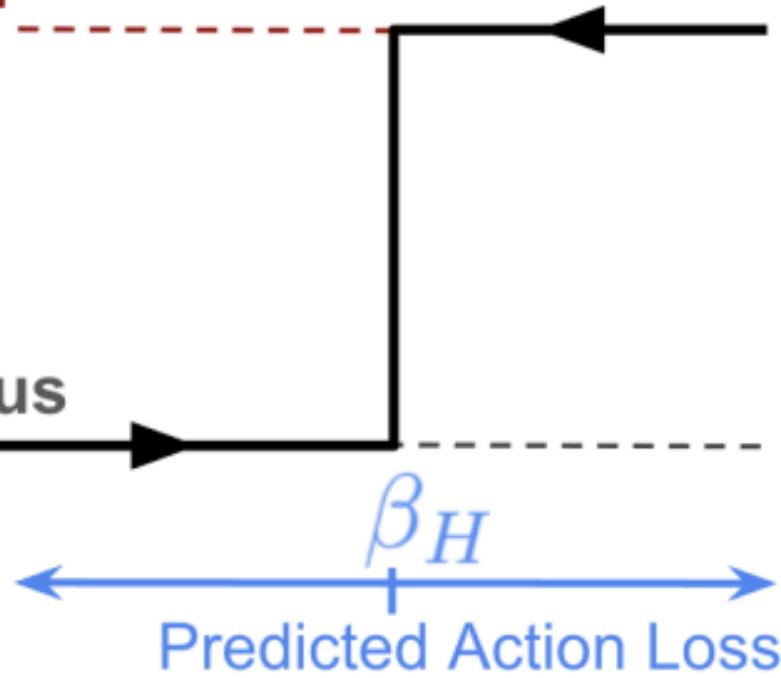
Autonomous
Mode



SafeDAgger

Supervisor
Mode

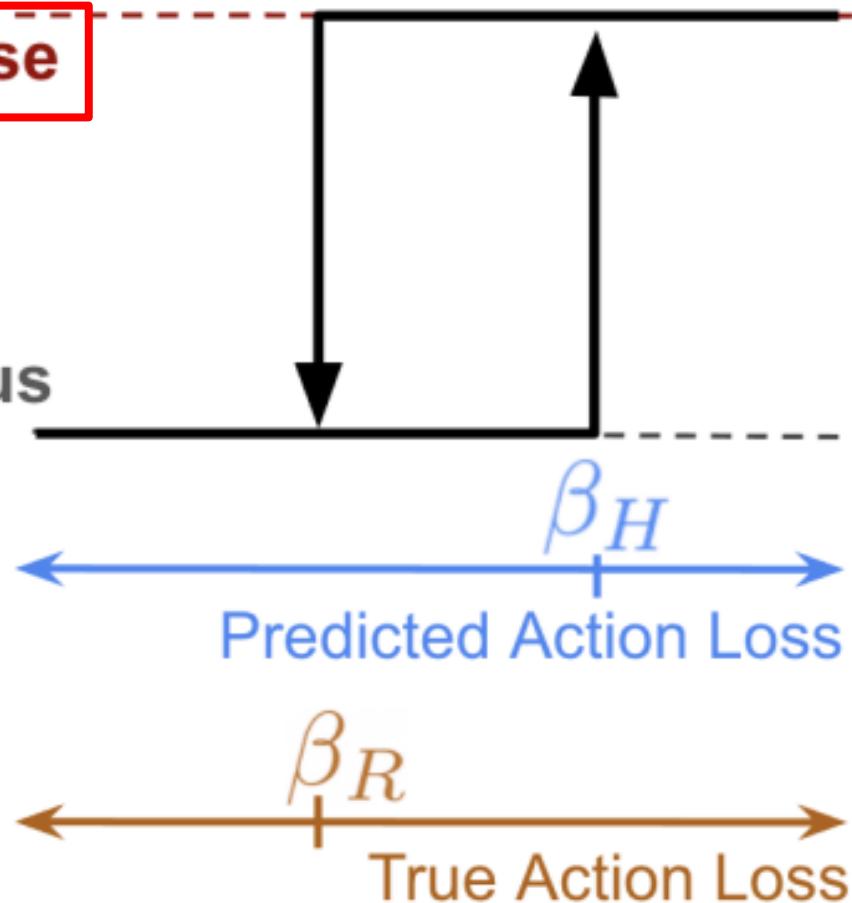
Autonomous
Mode



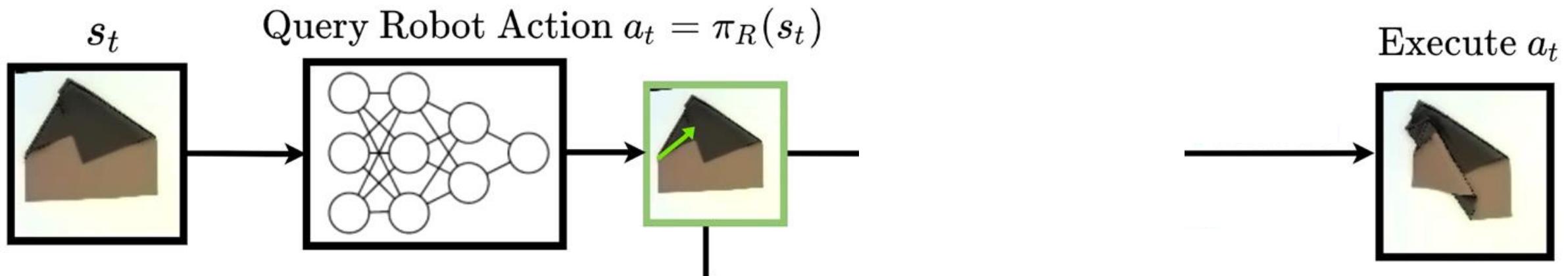
LazyDAgger

Supervisor
Mode + Noise

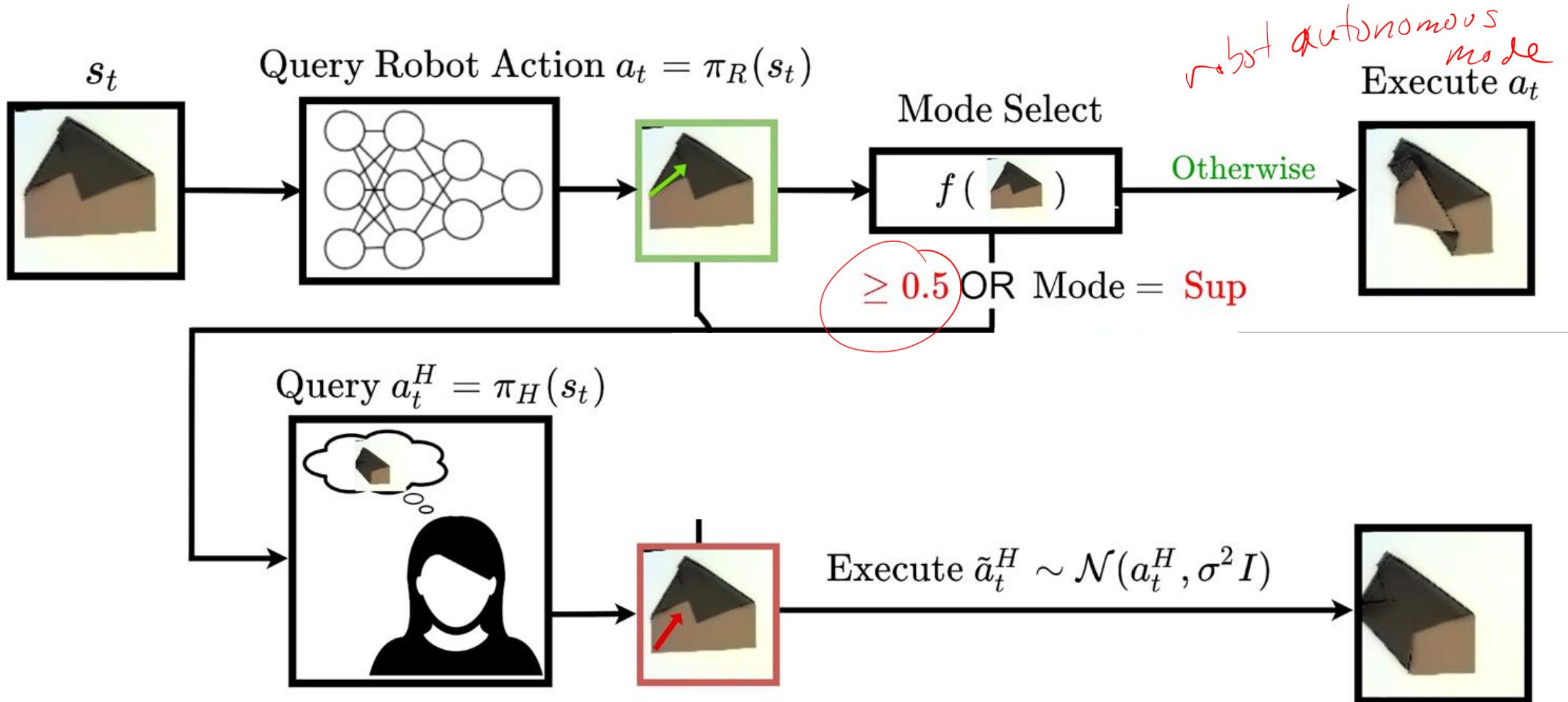
Autonomous
Mode



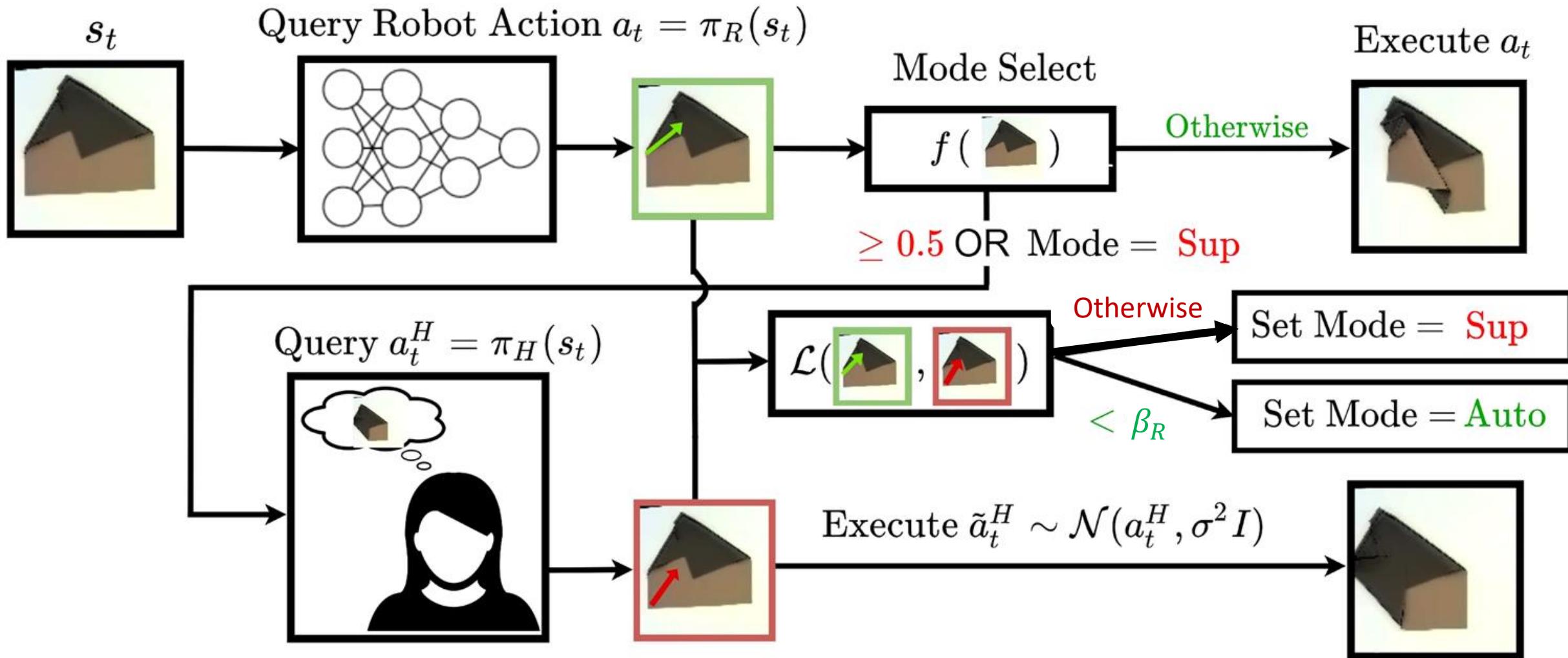
LazyDAgger



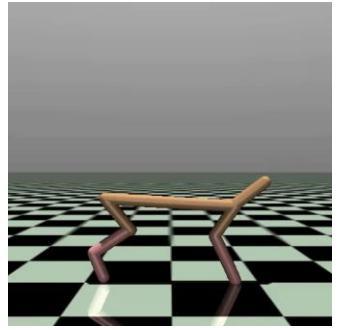
LazyDAgger



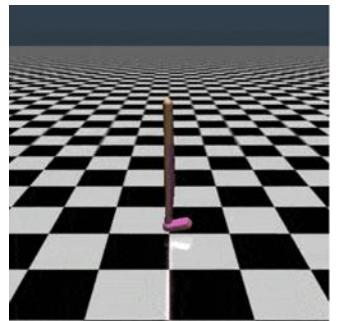
LazyDAgger



Simulation Experiments

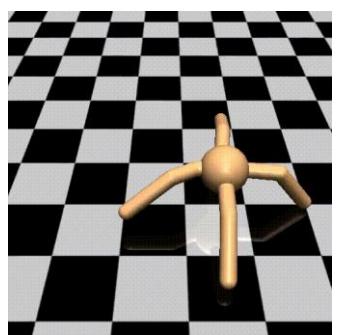


Number of Context Switches

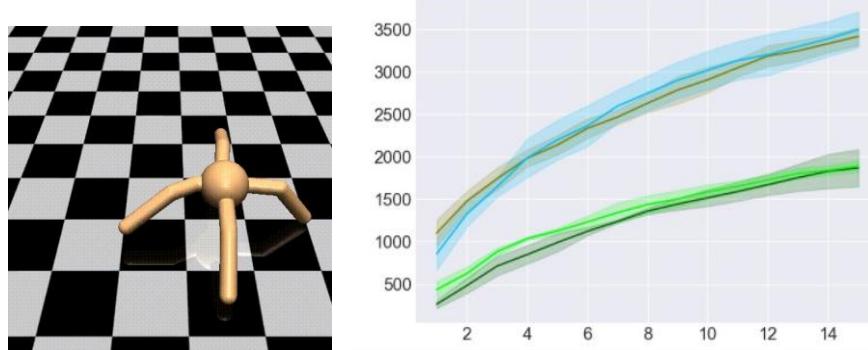
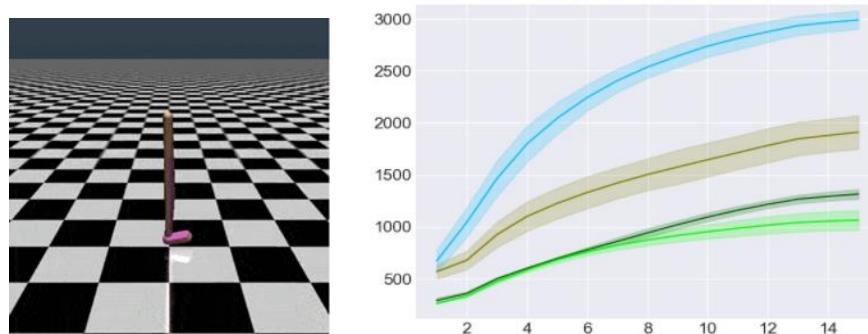
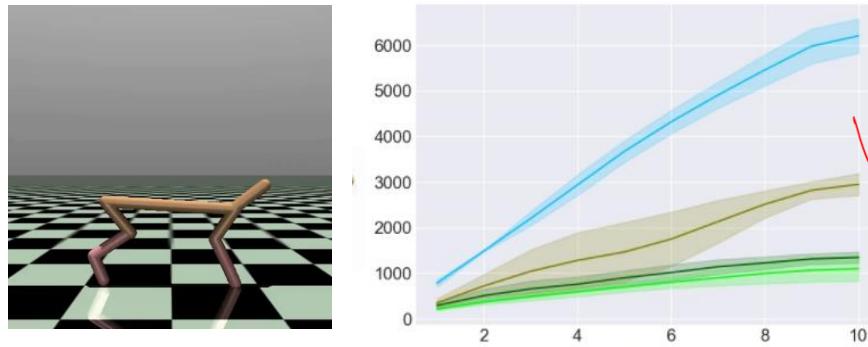


Context Switching Reduction

79%

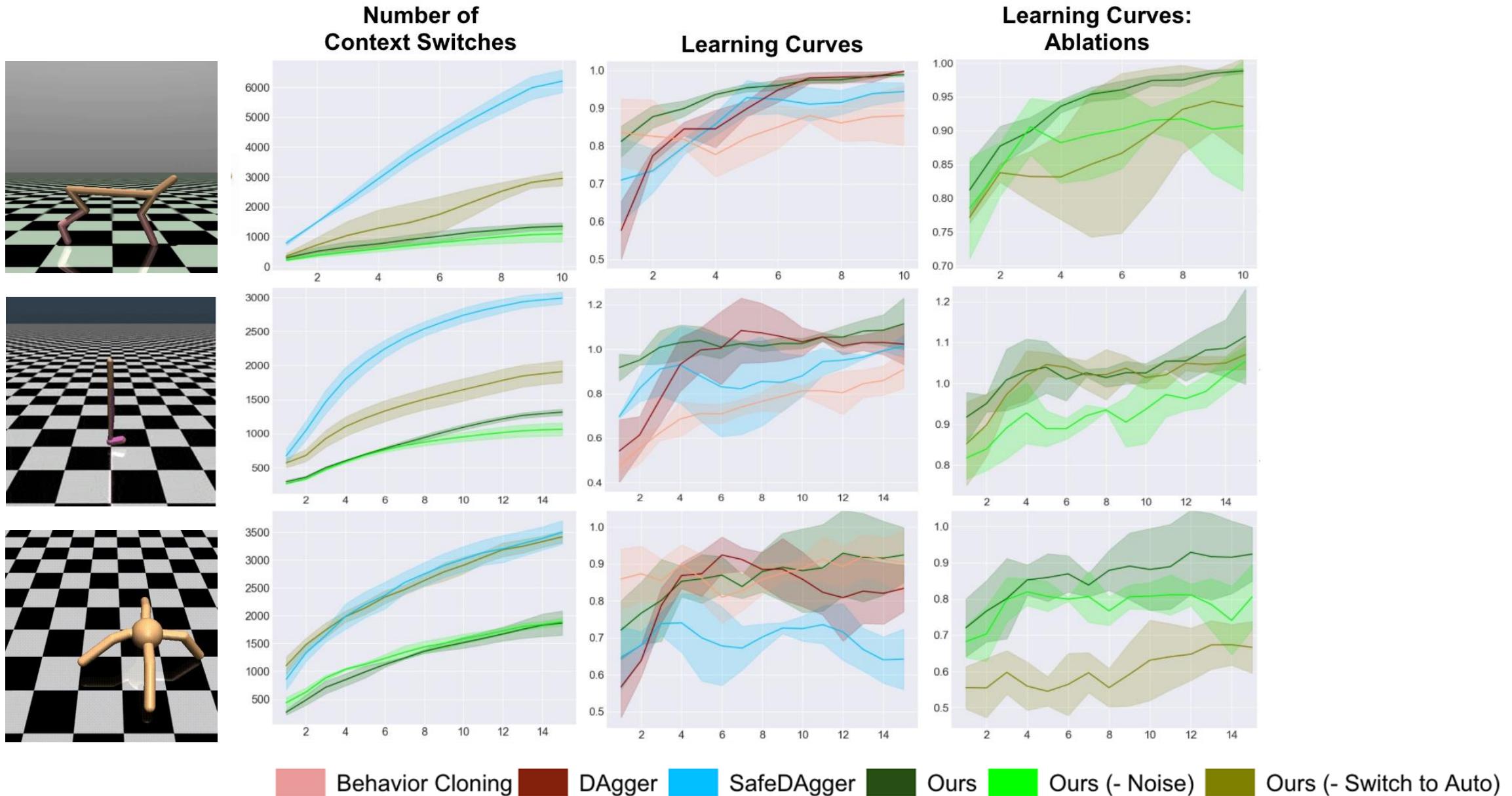


56%

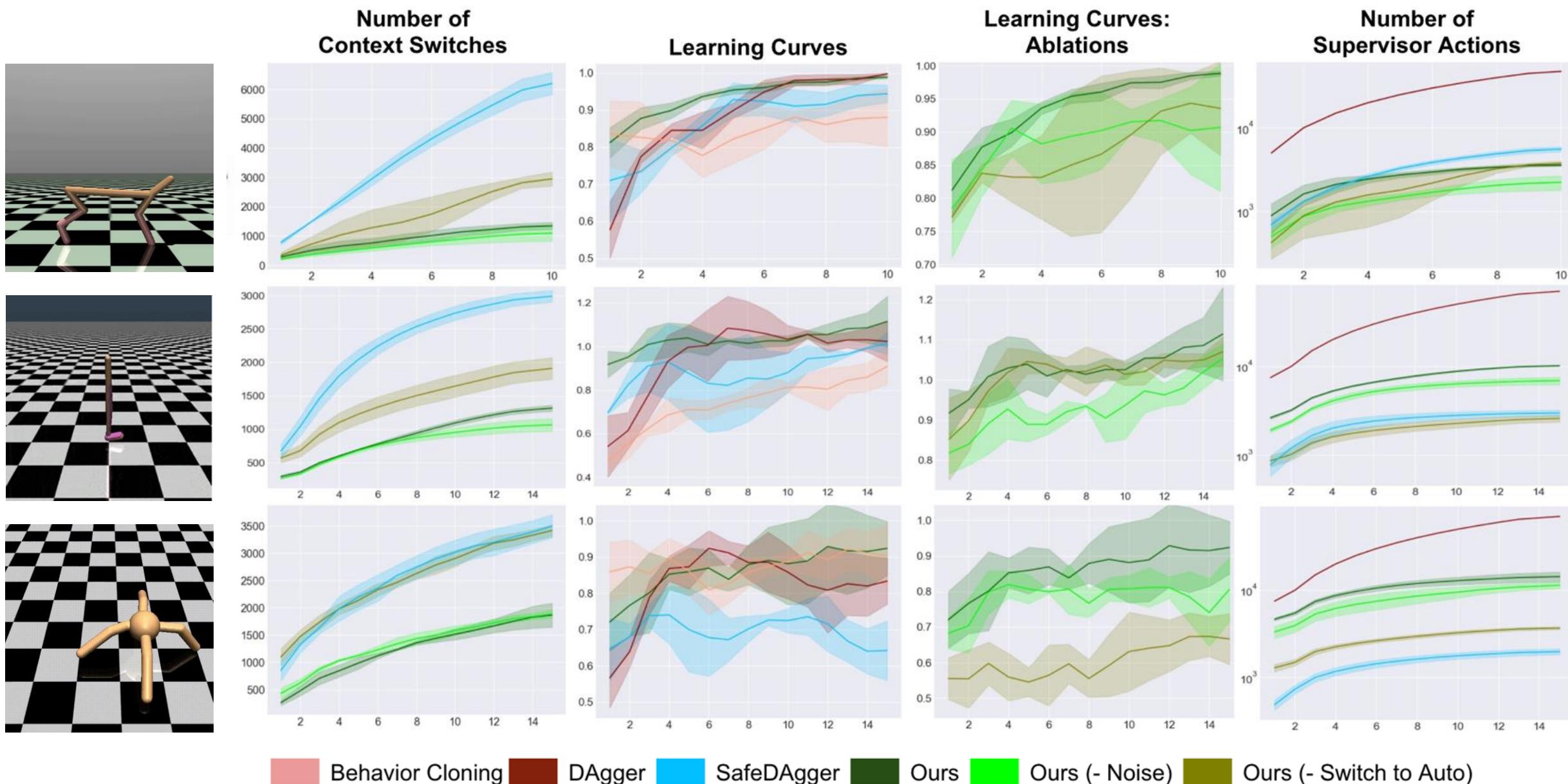


Behavior Cloning DAgger SafeDAgger Ours Ours (- Noise) Ours (- Switch to Auto)

Simulation Experiments



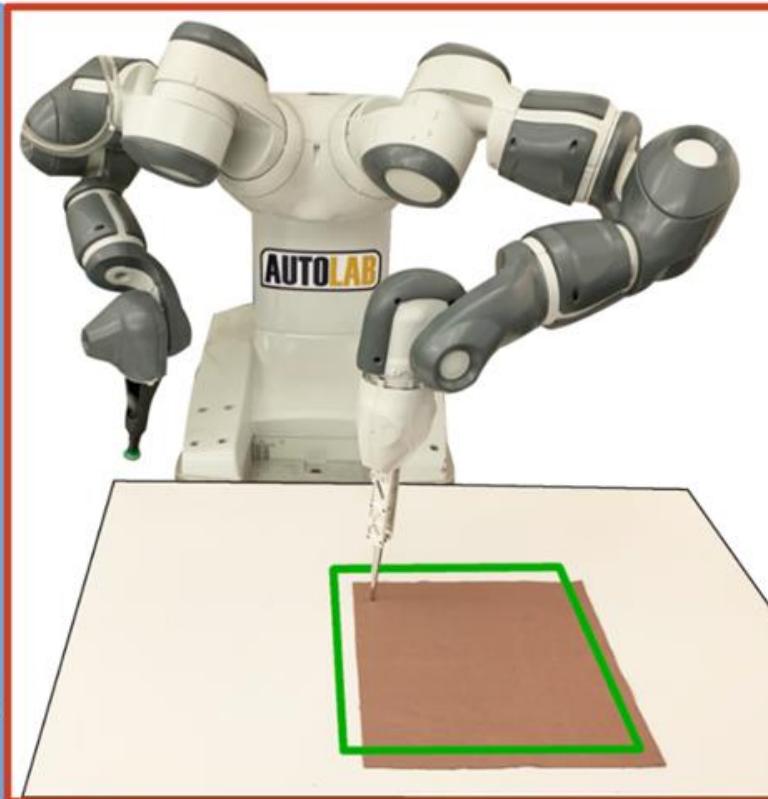
Simulation Experiments



(1) SMOOTH



(2) ALIGN

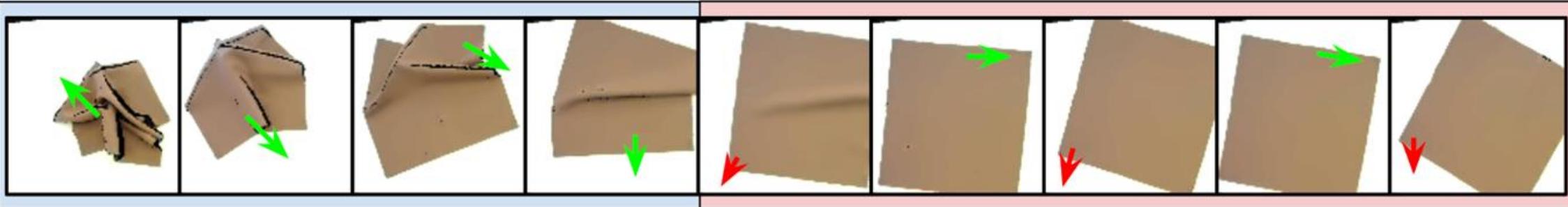


(3) FOLD

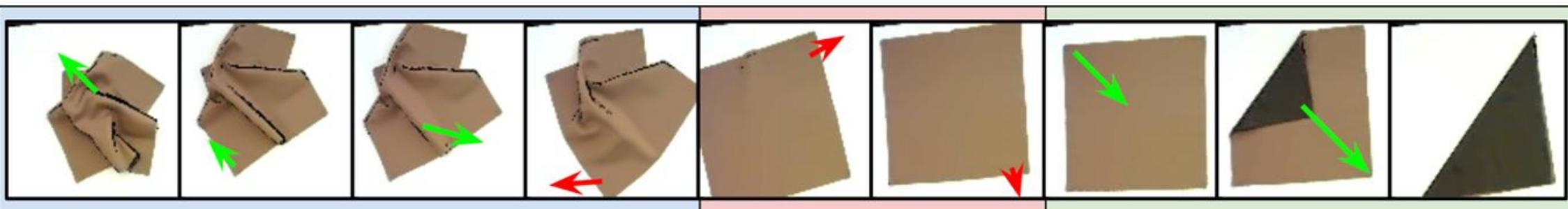




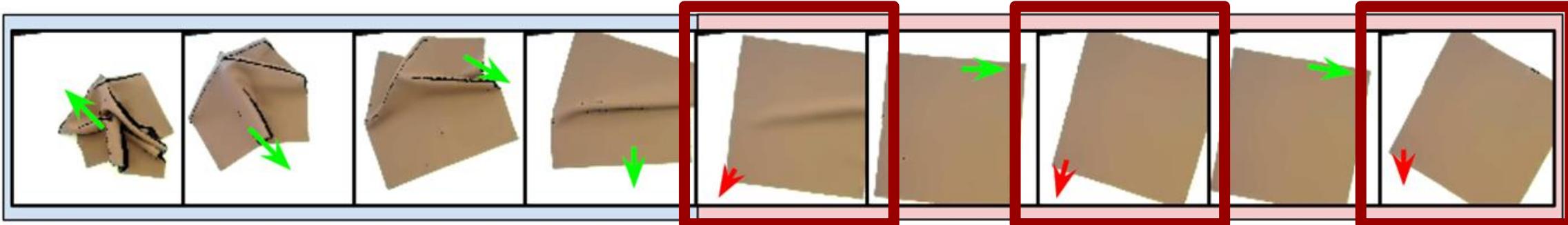
**Safe
Dagger**



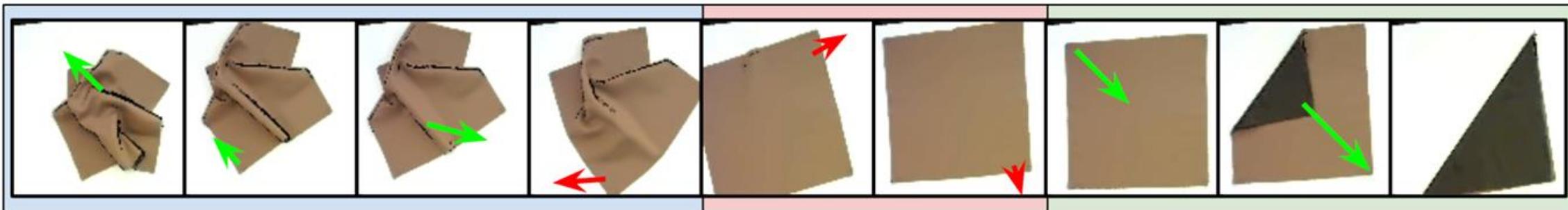
**Lazy
Dagger**



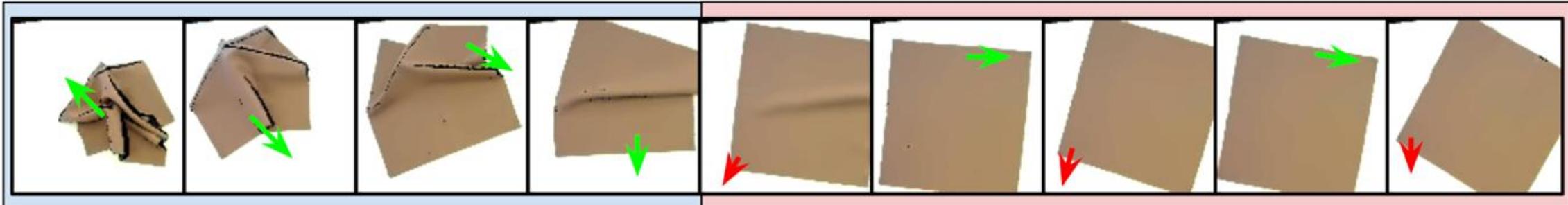
**Safe
Dagger**



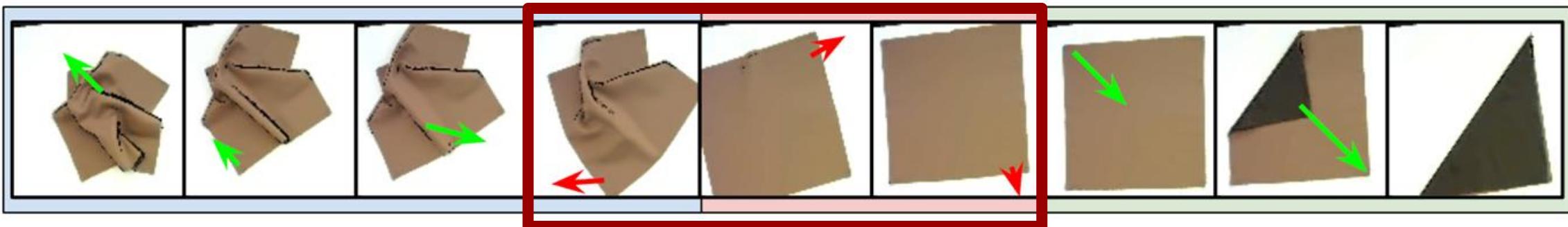
**Lazy
Dagger**



**Safe
DAgger**

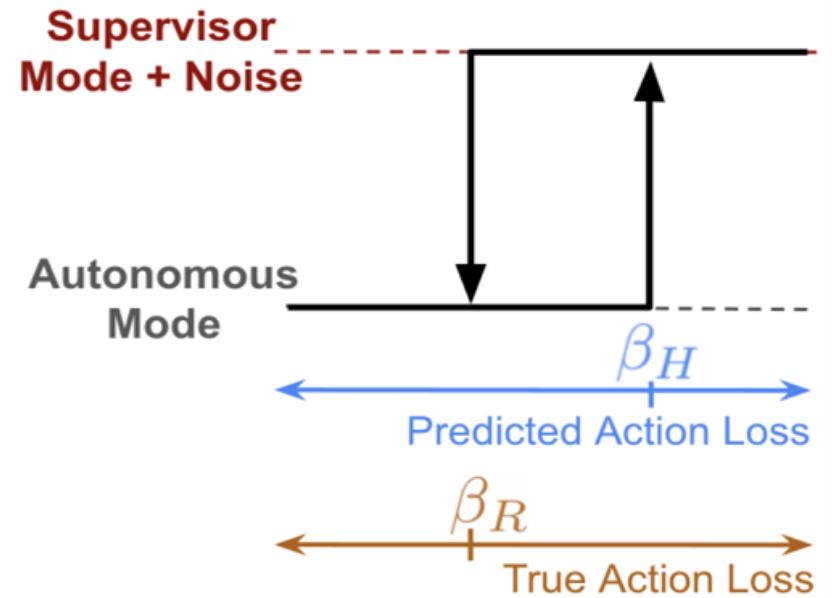


**Lazy
DAgger**



Limitations

- Parameter tuning
- Hard to know how many interventions will be requested.
- One human managing one robot.



When should a robot ask for help?



Novel (and risky)

When should a robot ask for help?

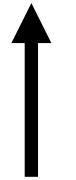
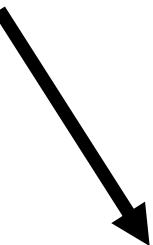
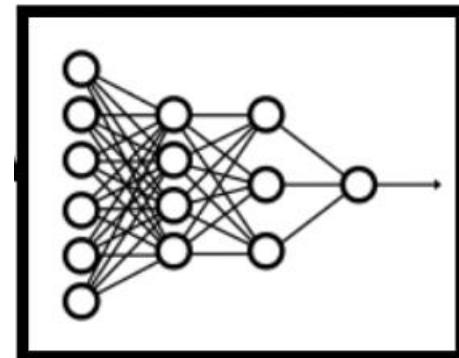
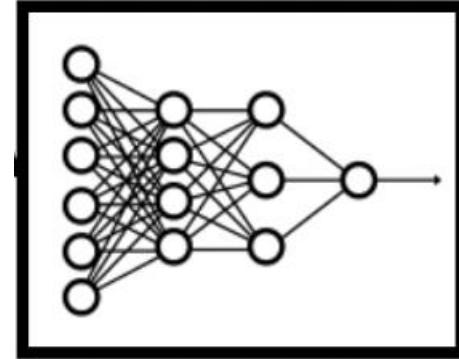
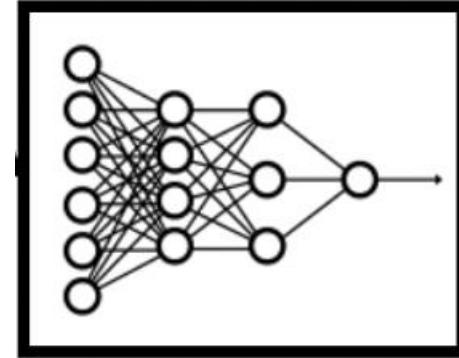
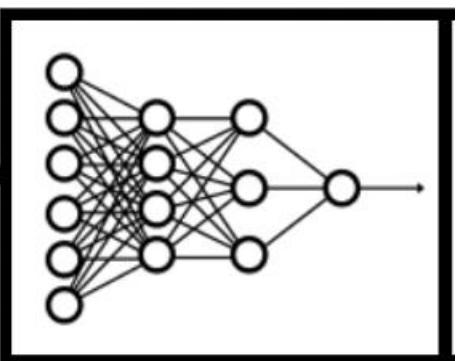
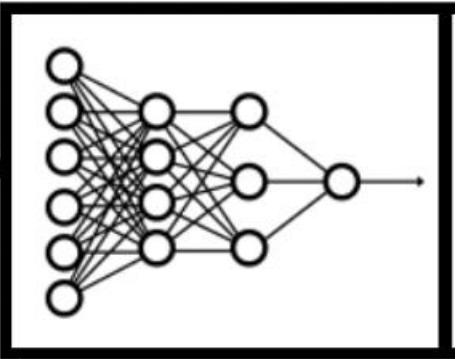
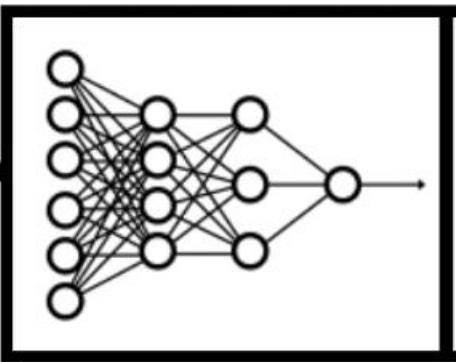


Novel (and risky)

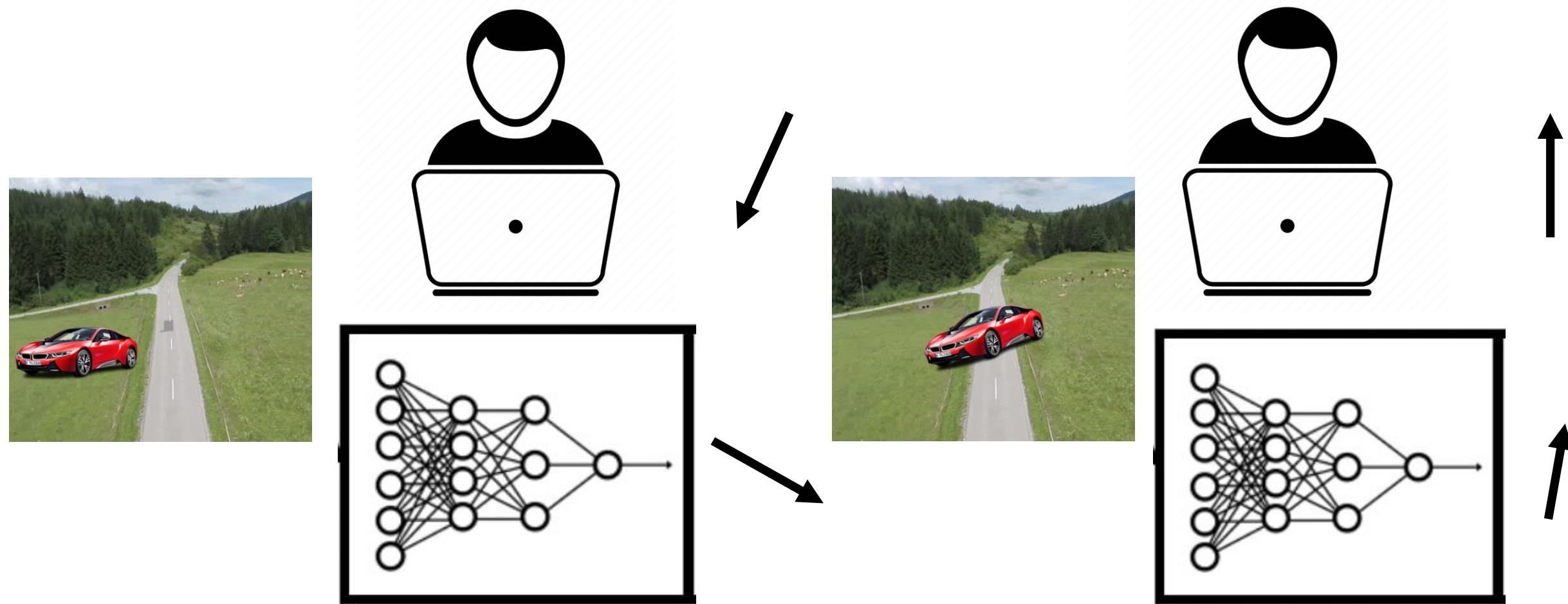


Risky (but not novel)

Novelty Estimation



Novelty Estimation: Supervisor Mode



Risk Estimation

$$Q_{\mathcal{G}}^{\pi_r}(s_t, a_t) = \mathbb{E}_{\pi_r} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} \mathbf{1}_{\mathcal{G}}(s'_t) | s_t, a_t \right]$$

Risk Estimation

$$Q_{\mathcal{G}}^{\pi_r}(s_t, a_t) = \mathbb{E}_{\pi_r} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} \mathbf{1}_{\mathcal{G}}(s'_t) | s_t, a_t \right]$$

$$\text{Risk}^{\pi_r}(s, a) = 1 - \hat{Q}_{\phi, \mathcal{G}}^{\pi_r}(s, a)$$

Risk Estimation

$$Q_{\mathcal{G}}^{\pi_r}(s_t, a_t) = \mathbb{E}_{\pi_r} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} \mathbb{1}_{\mathcal{G}}(s'_t) | s_t, a_t \right]$$

$$\text{Risk}^{\pi_r}(s, a) = 1 - \hat{Q}_{\phi, \mathcal{G}}^{\pi_r}(s, a)$$

$$\begin{aligned} J_{\mathcal{G}}^Q(s_t, a_t, s_{t+1}; \phi) = \\ \frac{1}{2} \left(\hat{Q}_{\phi, \mathcal{G}}^{\pi_r}(s_t, a_t) - (\mathbb{1}_{\mathcal{G}}(s_t) + (1 - \mathbb{1}_{\mathcal{G}}(s_t))\gamma \hat{Q}_{\phi, \mathcal{G}}^{\pi_r}(s_{t+1}, \pi_r(s_{t+1}))) \right)^2 \end{aligned}$$

Putting it all together...

**AUTONOMOUS
MODE**

$$\begin{array}{c} \text{Novelty}(s_t) > \delta_h \\ \text{OR} \\ \text{Risk}^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h \end{array} \longrightarrow$$

Switch to
**SUPERVIS
OR MODE**

Putting it all together...

**AUTONOMOUS
MODE**

Novelty(s_t) > δ_h
OR
Risk $^{\pi_r}(s_t, \pi_r(s_t))$ > β_h

Switch to
**SUPERVISOR
MODE**

**SUPERVISOR
MODE**

$||\pi_r(s_t) - \pi_h(s_t)||_2^2 < \delta_r$
AND
Risk $^{\pi_r}(s_t, \pi_r(s_t)) < \beta_r$

Switch to
**AUTONOMOUS
MODE**

Putting it all together...

**AUTONOMOUS
MODE**

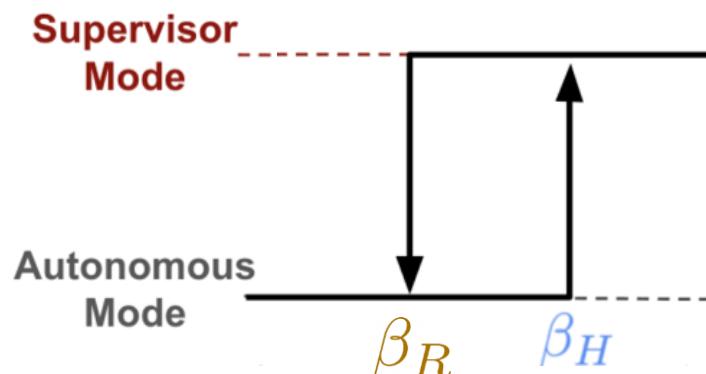
Novelty(s_t) > δ_h
OR
Risk $^{\pi_r}(s_t, \pi_r(s_t))$ > β_h

Switch to
**SUPERVISOR
MODE**

**SUPERVISOR
MODE**

$||\pi_r(s_t) - \pi_h(s_t)||_2^2 < \delta_r$
AND
Risk $^{\pi_r}(s_t, \pi_r(s_t)) < \beta_r$

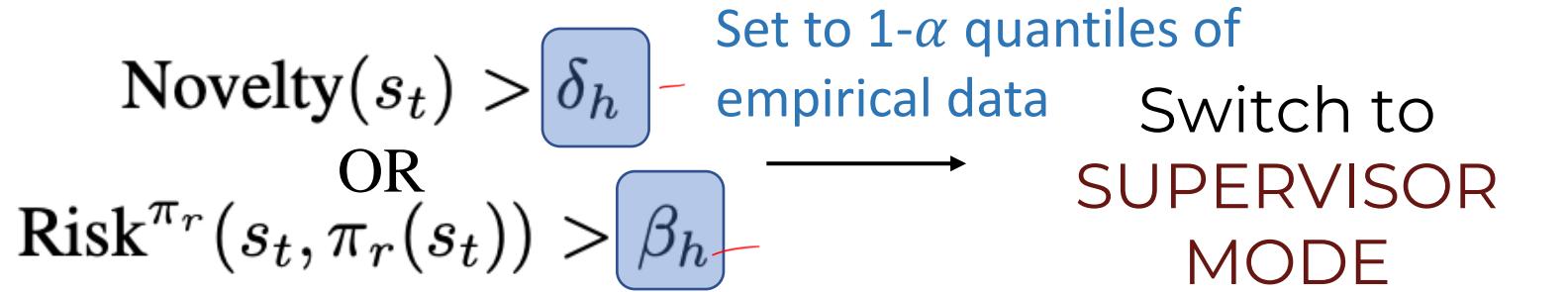
Switch to
**AUTONOMOUS
MODE**



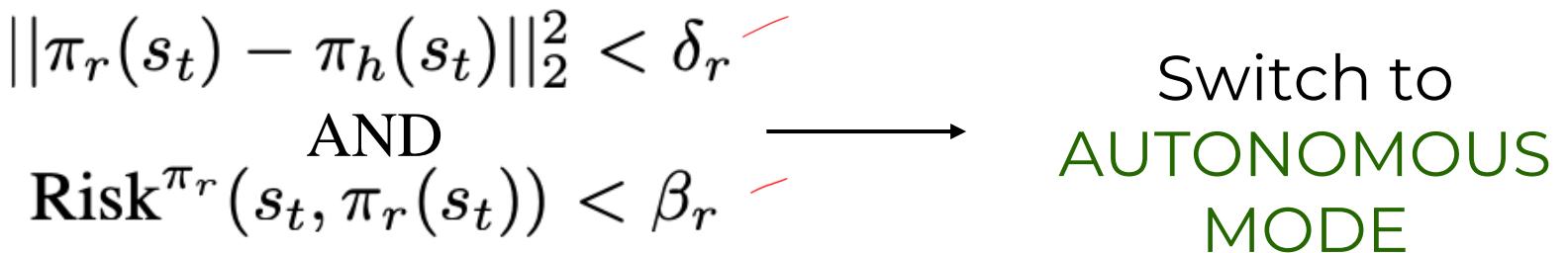
Wait, didn't we just double
the number of
hyperparameters?

Putting it all together...

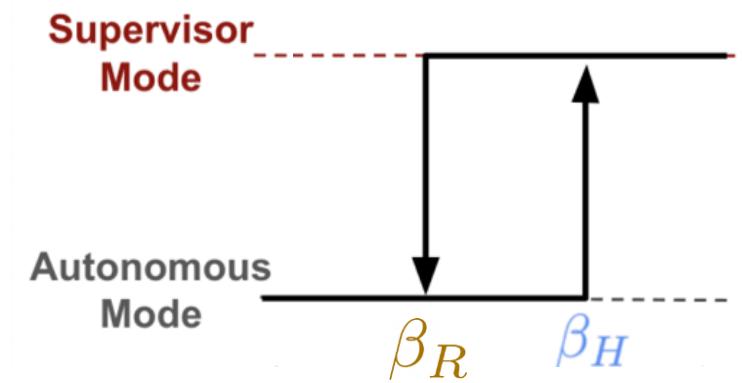
**AUTONOMOUS
MODE**



**SUPERVISOR
MODE**

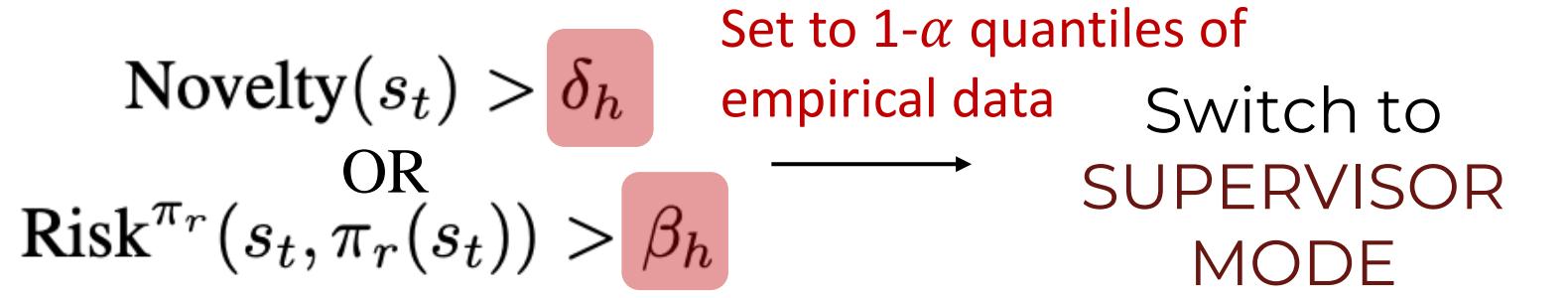


$$\alpha = \frac{\text{desired } \# \text{ interventions}}{\# \text{ robot actions}}$$

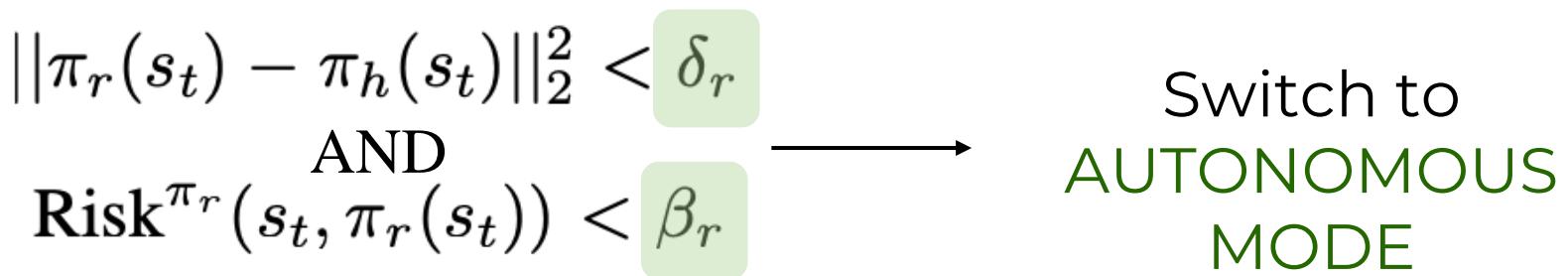


Putting it all together...

**AUTONOMOUS
MODE**

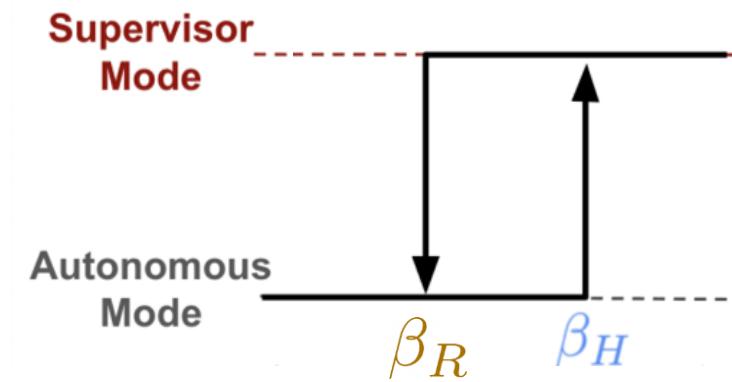


**SUPERVISOR
MODE**



$$\alpha = \frac{\text{desired } \# \text{ interventions}}{\# \text{ robot actions}}$$

Set to medians of empirical data



Putting it all together...

**AUTONOMOUS
MODE**

Novelty(s_t) > δ_h
OR
Risk $^{\pi_r}(s_t, \pi_r(s_t))$ > β_h

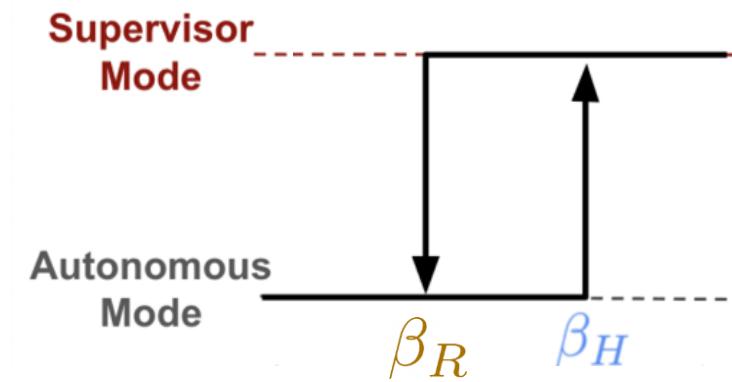
Switch to
**SUPERVISOR
MODE**

**SUPERVISOR
MODE**

$||\pi_r(s_t) - \pi_h(s_t)||_2^2 < \delta_r$
AND
Risk $^{\pi_r}(s_t, \pi_r(s_t)) < \beta_r$

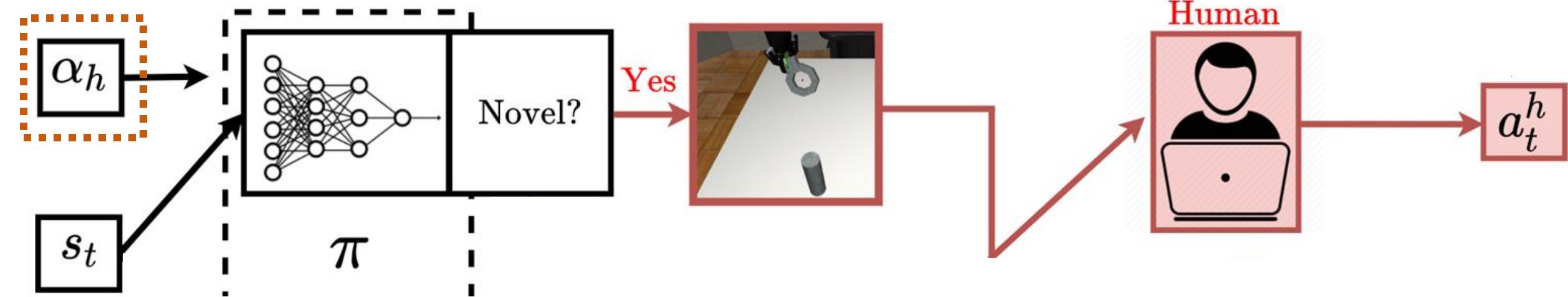
Switch to
**AUTONOMOUS
MODE**

$$\alpha = \frac{\text{desired } \# \text{ interventions}}{\# \text{ robot actions}}$$

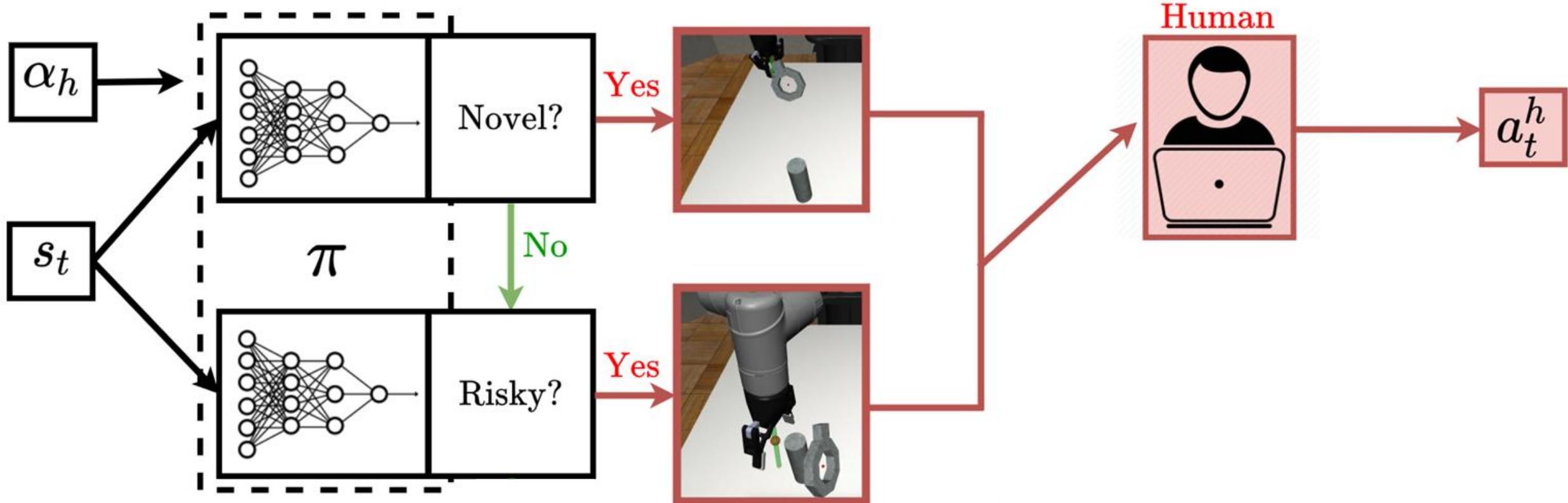


Target percent of time human wants to give interventions.

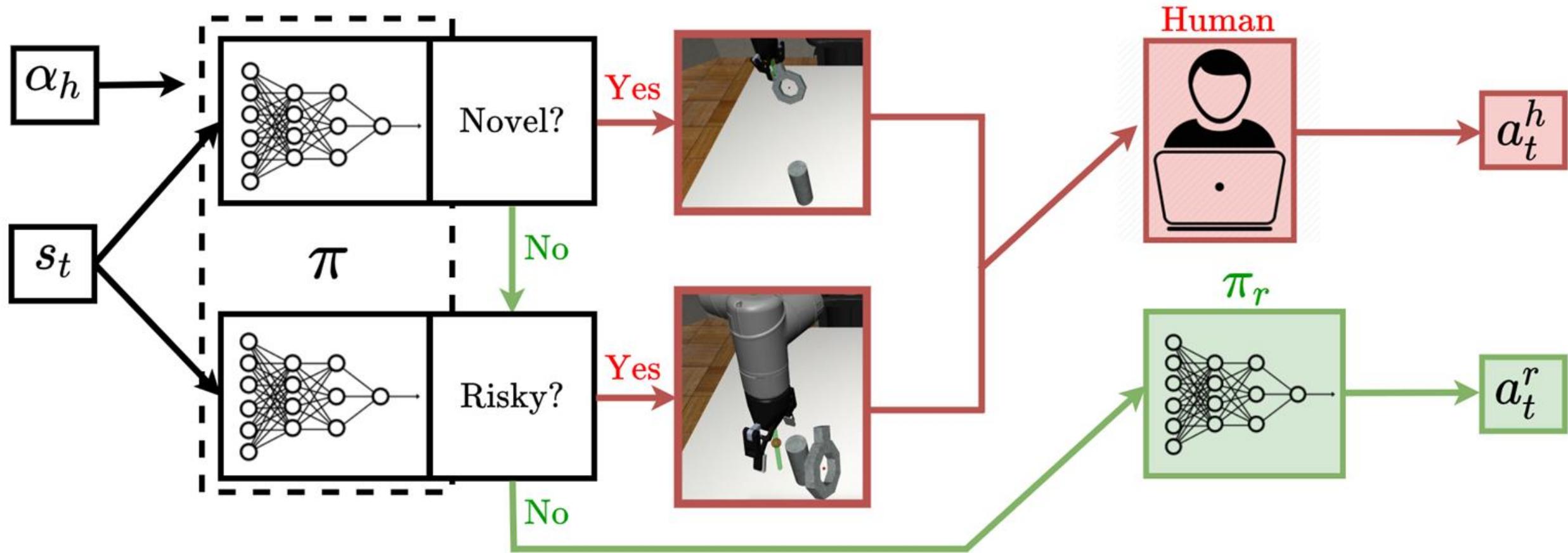
ThriftyDAgger

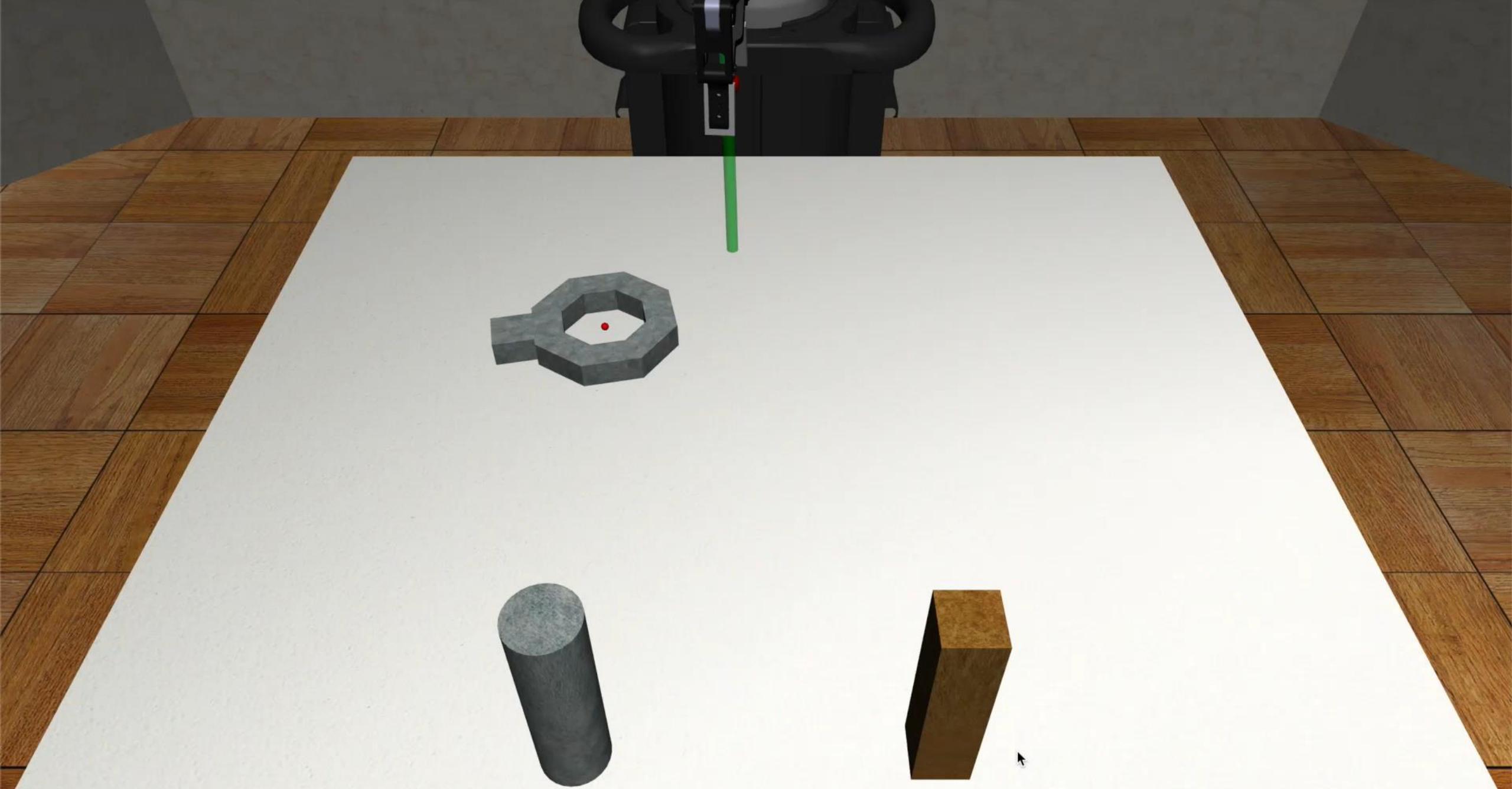


ThriftyDAgger

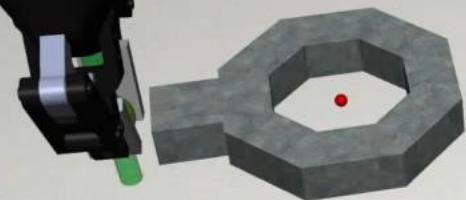


ThriftyDAgger

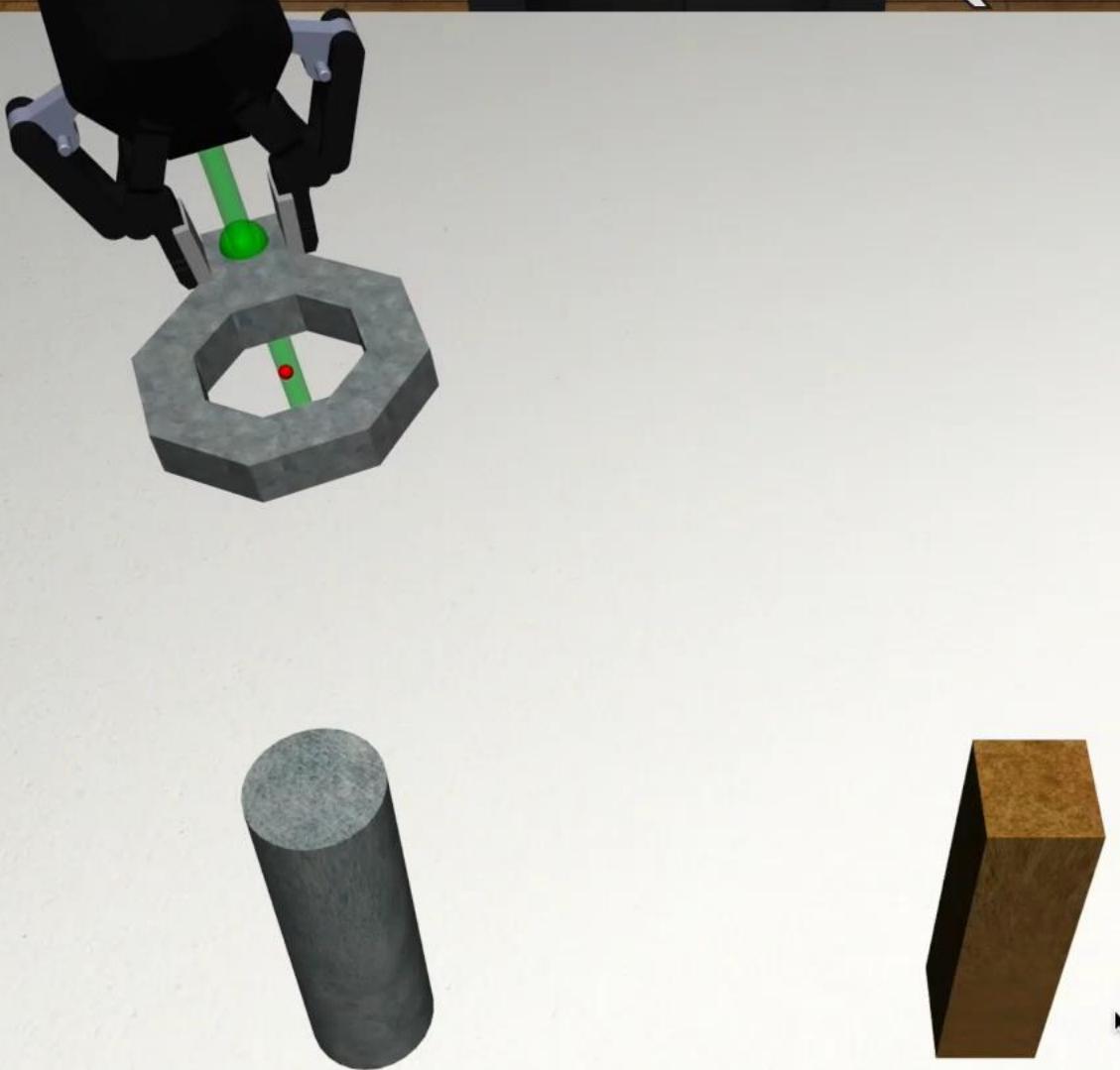


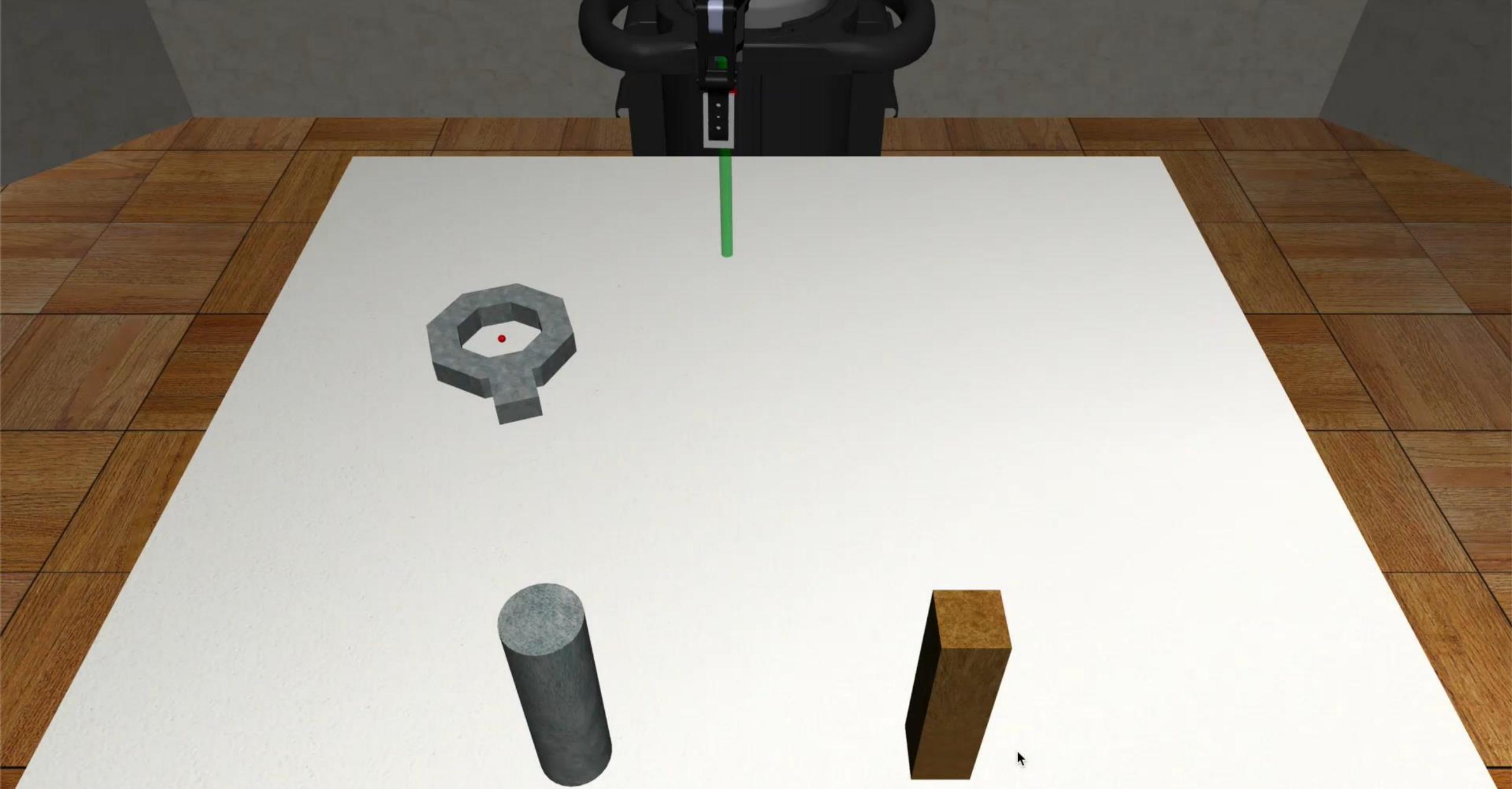


Autonomous Mode

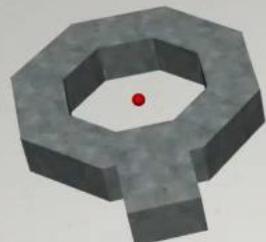


Supervisor Mode (Novel)





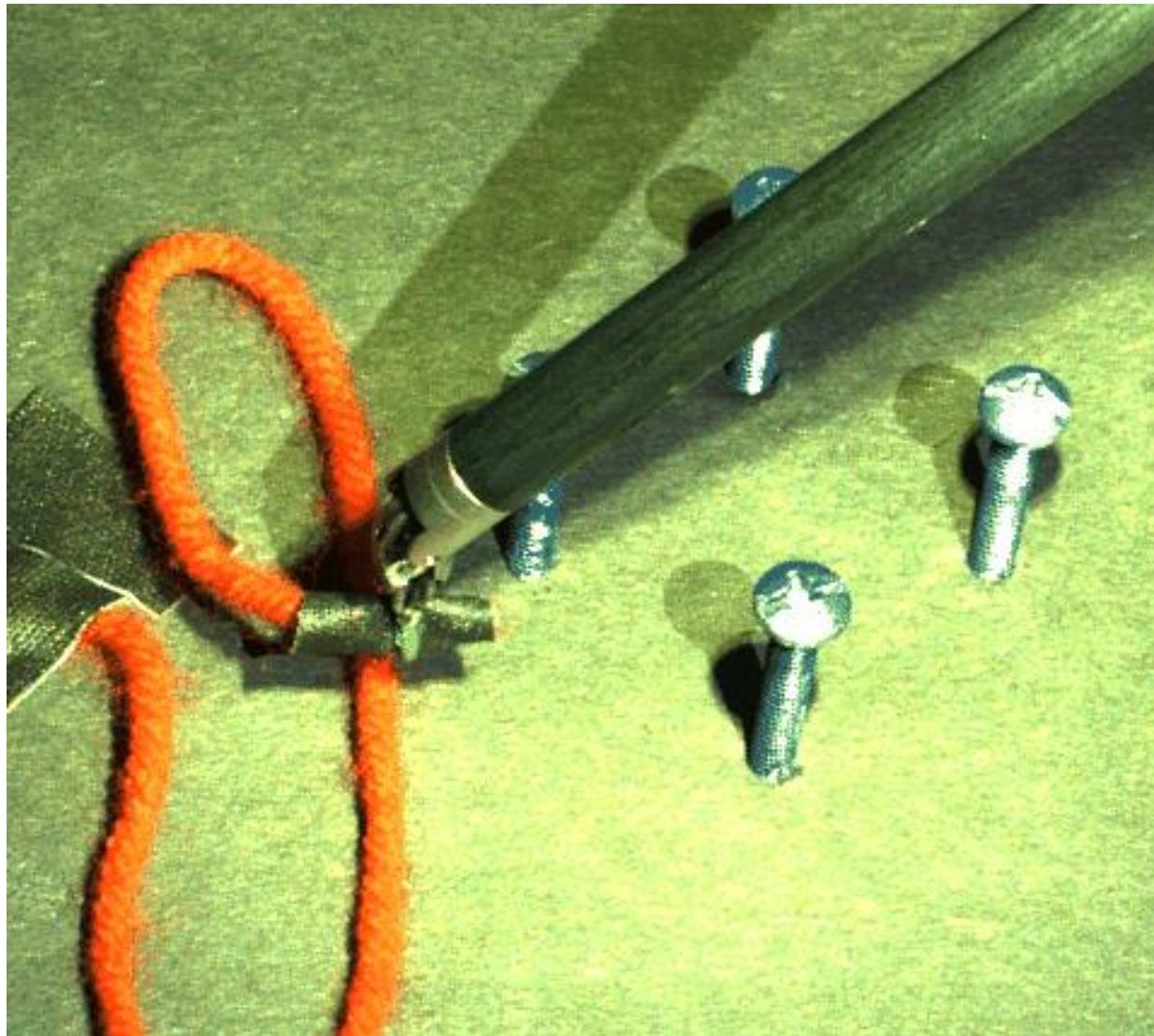
Supervisor Mode (Risk)



Supervisor Mode (Risk)



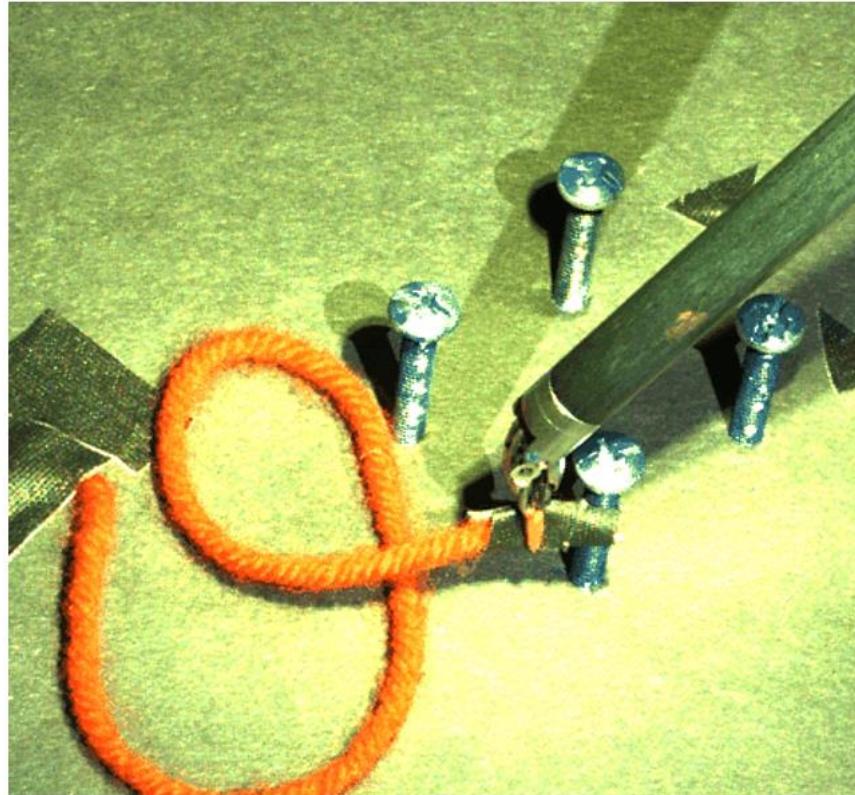
Human Demonstration



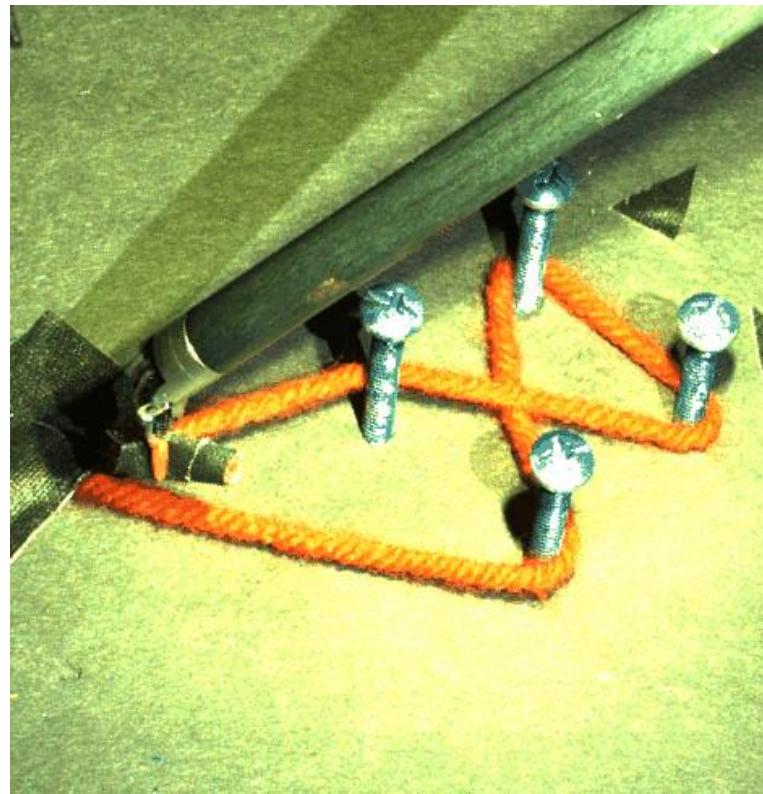
Behavior Cloning



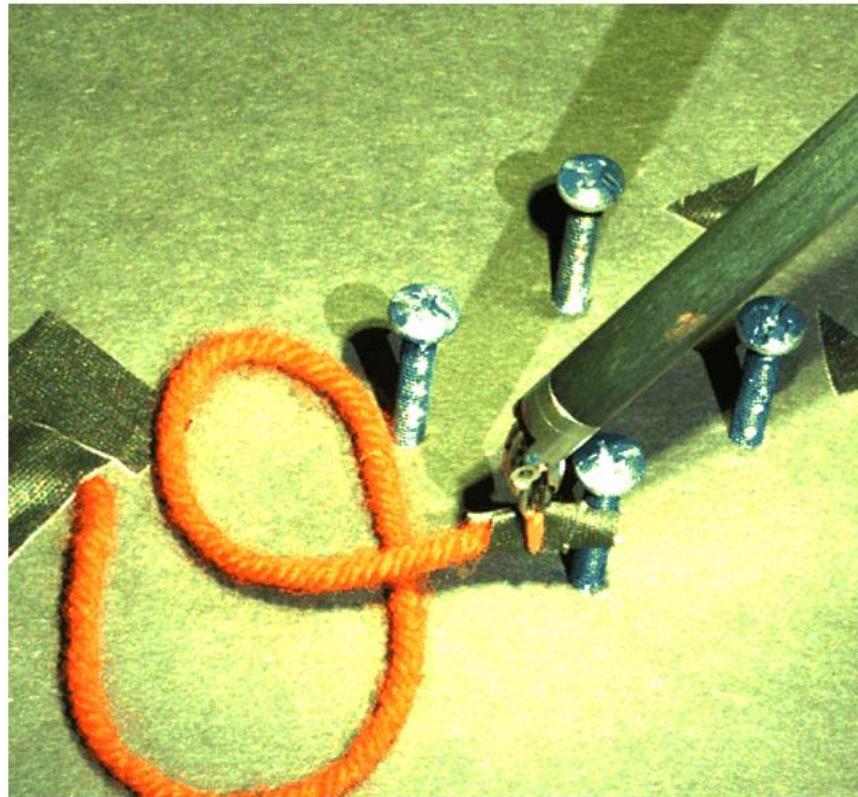
Behavior Cloning



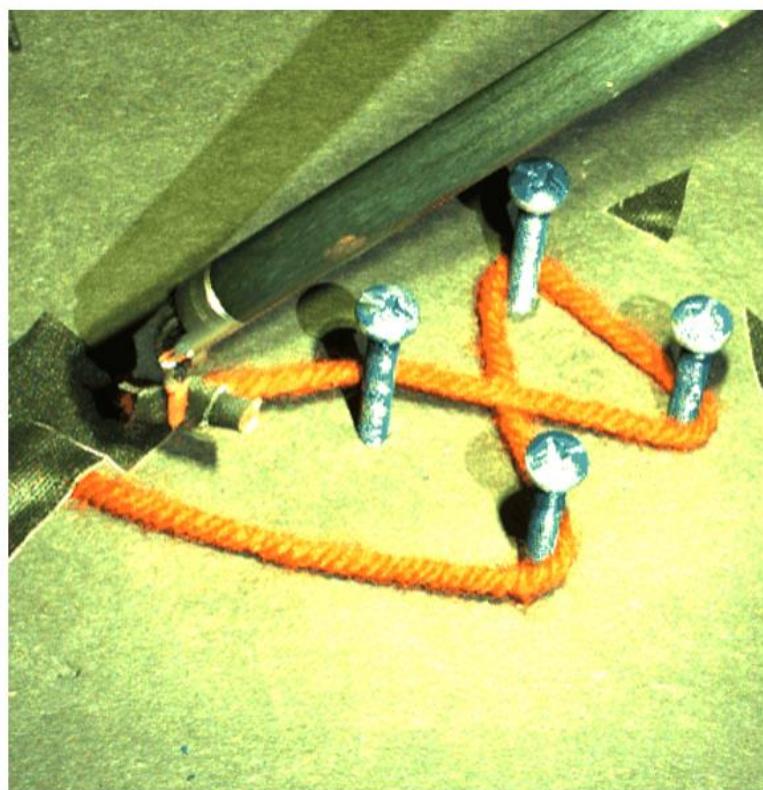
ThriftyDAgger (autonomous)



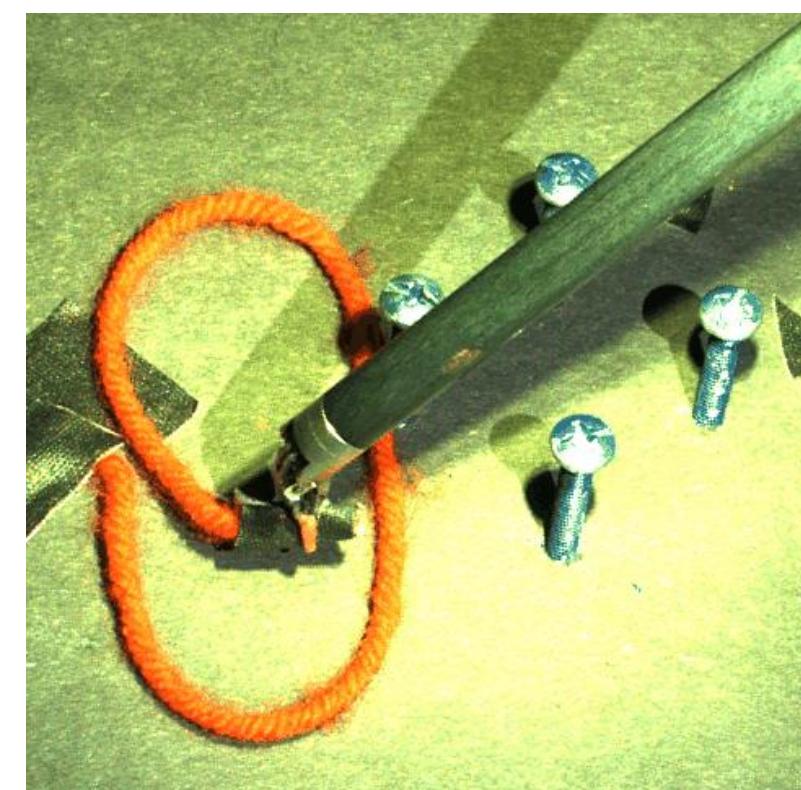
Behavior Cloning



ThriftyDAgger (autonomous)



ThriftyDAgger (+human)



User Study

N=10 subjects each control 3 robots in simulation.

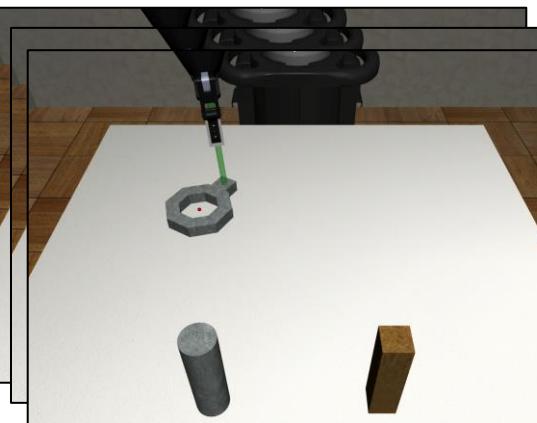
Robot-Gated

Memory: Non-Match

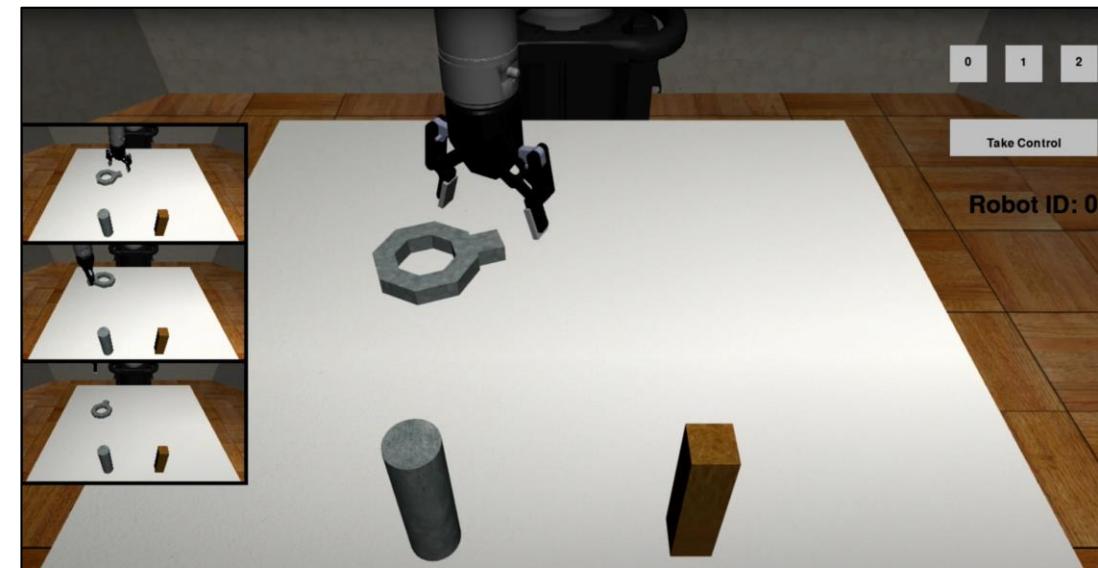
H	H	H	H	H
H				H
H	H	H	H	H
H	H	H	H	H

Memory: Match

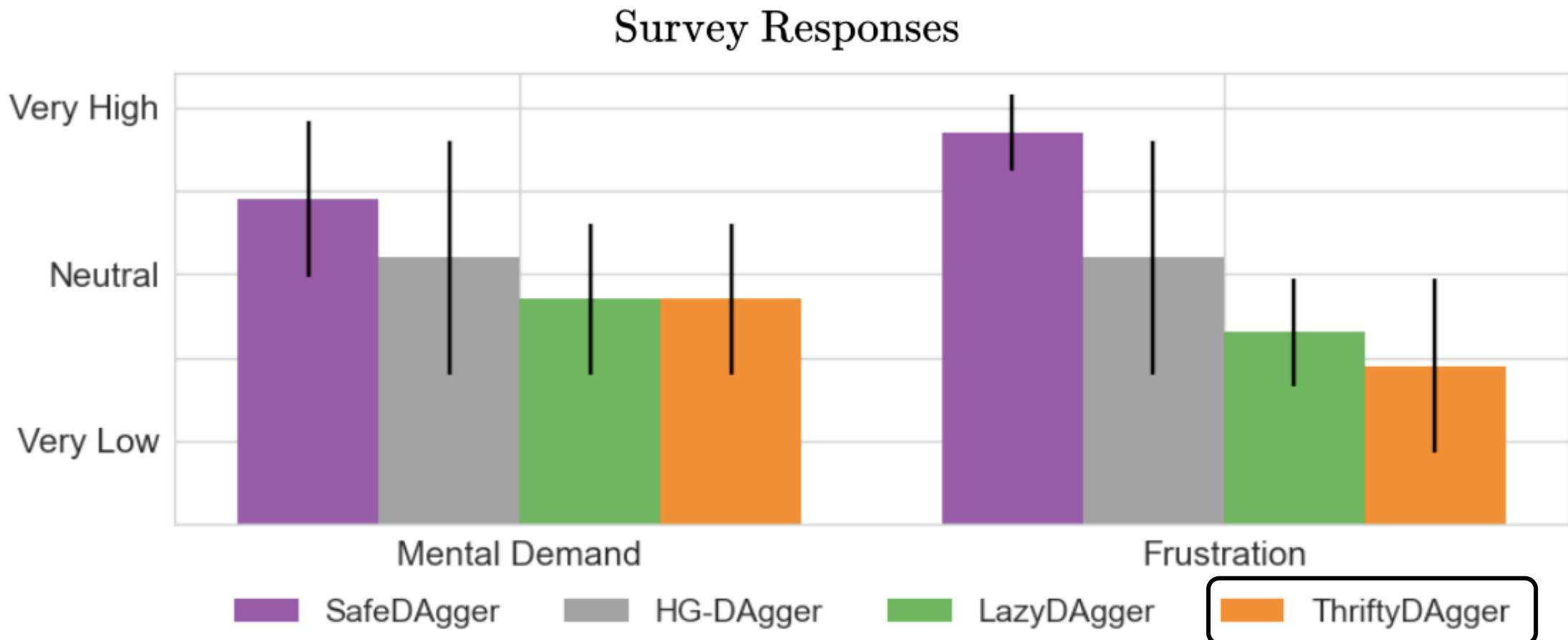
H	H	H	H	H
H	H	H		H
	H	H	H	H
H	H	H	H	H



Human-Gated



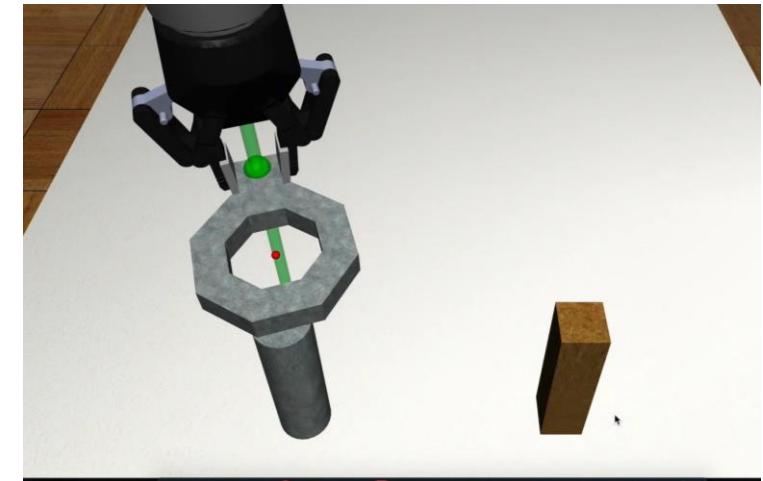
ThriftyDAgger Qualitative Results



User Study Quantitative Results

ThriftyDAgger had

- 21% fewer human interventions
- 57% more concentration pairs found
- 80% more throughput



Scalable and safe robot fleets are possible when robots ask for help in ways that minimize human supervisor burden.

