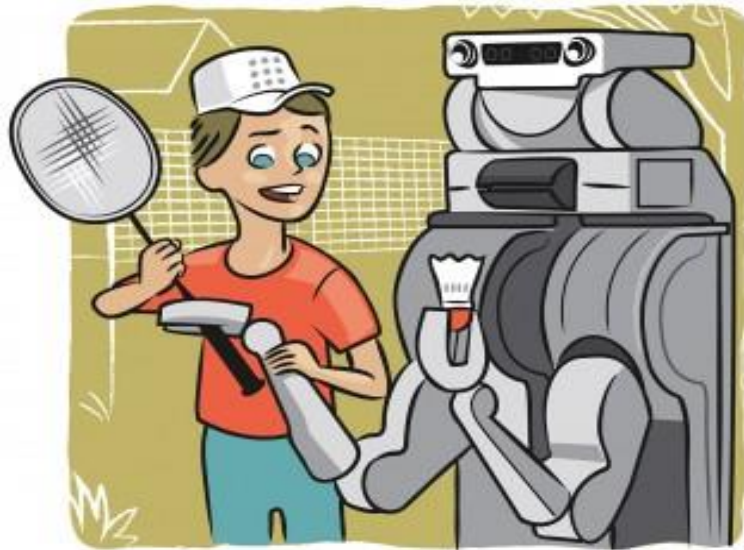


Behavioral Cloning and Interactive Imitation Learning

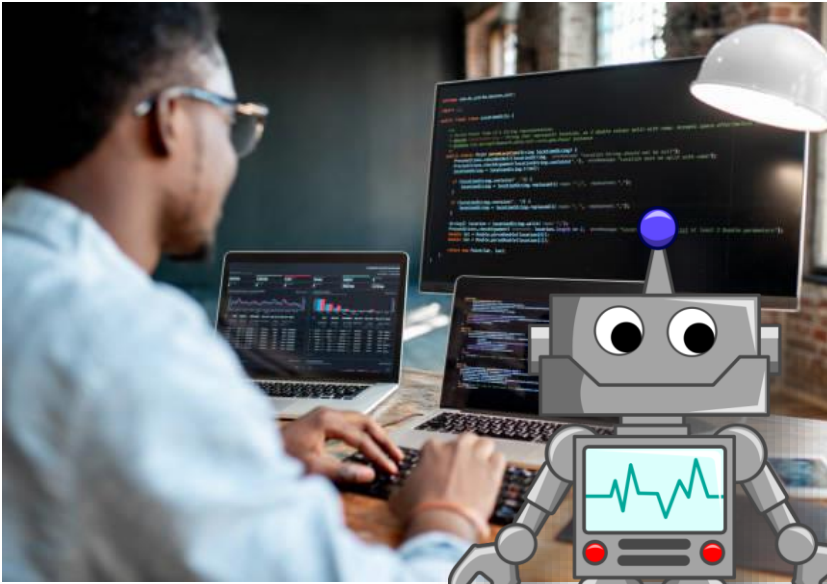


Instructor: Daniel Brown

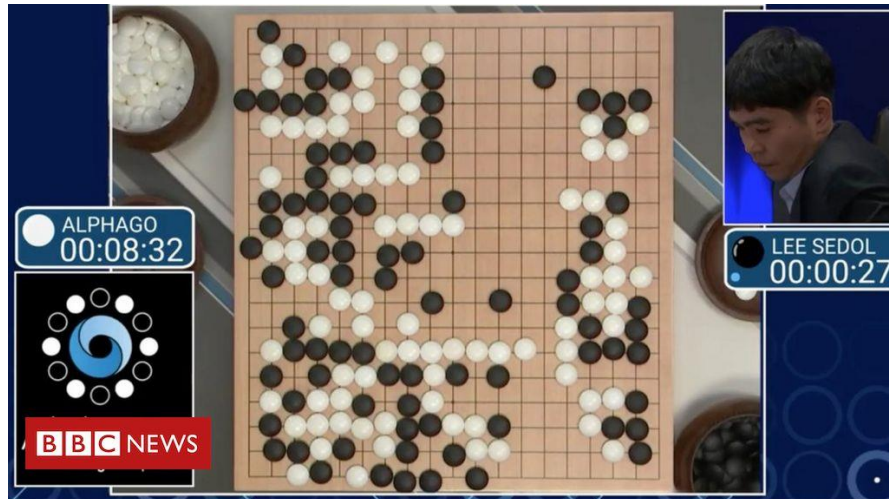
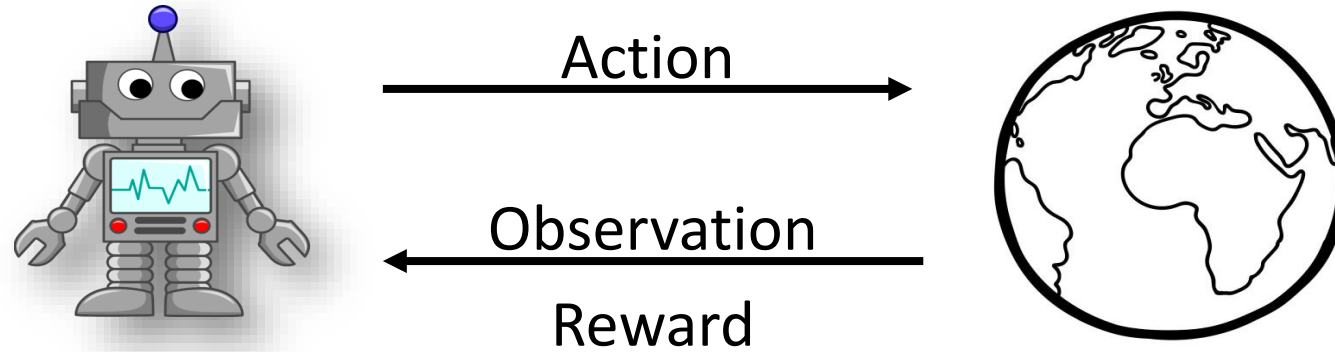
[Some slides adapted from Sergey Levine (CS 285) and Alina Vereshchaka (CSE4/510)]

Course feedback is open

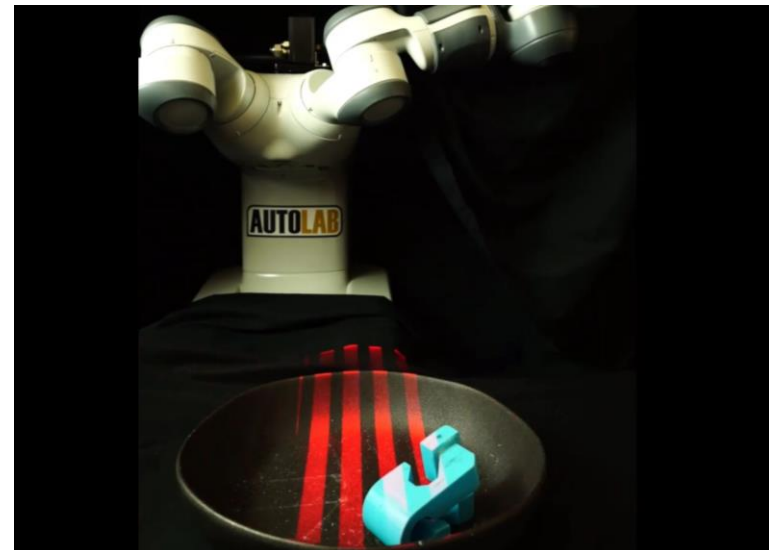
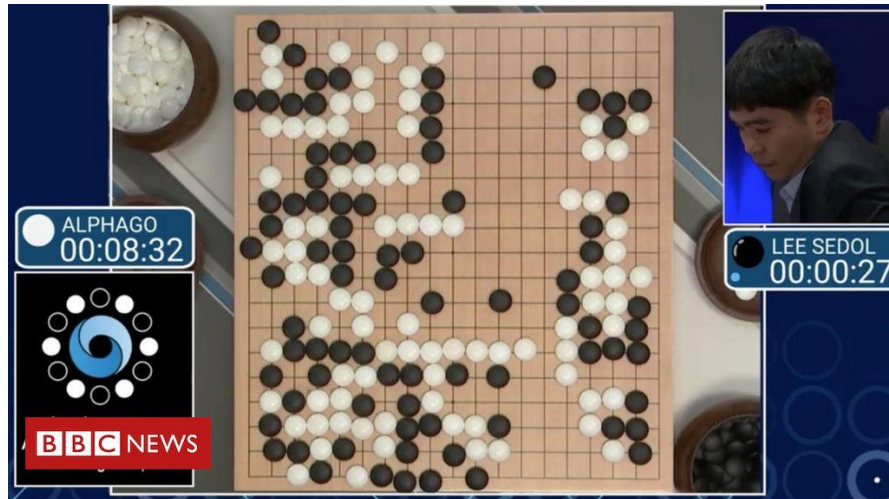
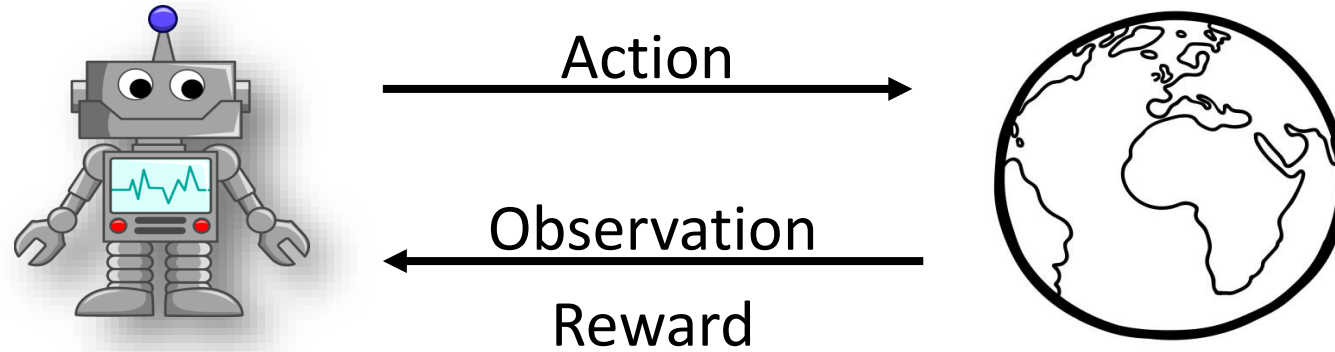
- Extra credit if class response rate is 70% or higher
 - Sliding scale if we reach 70%:
 - Extra credit points = response rate percentage / 10



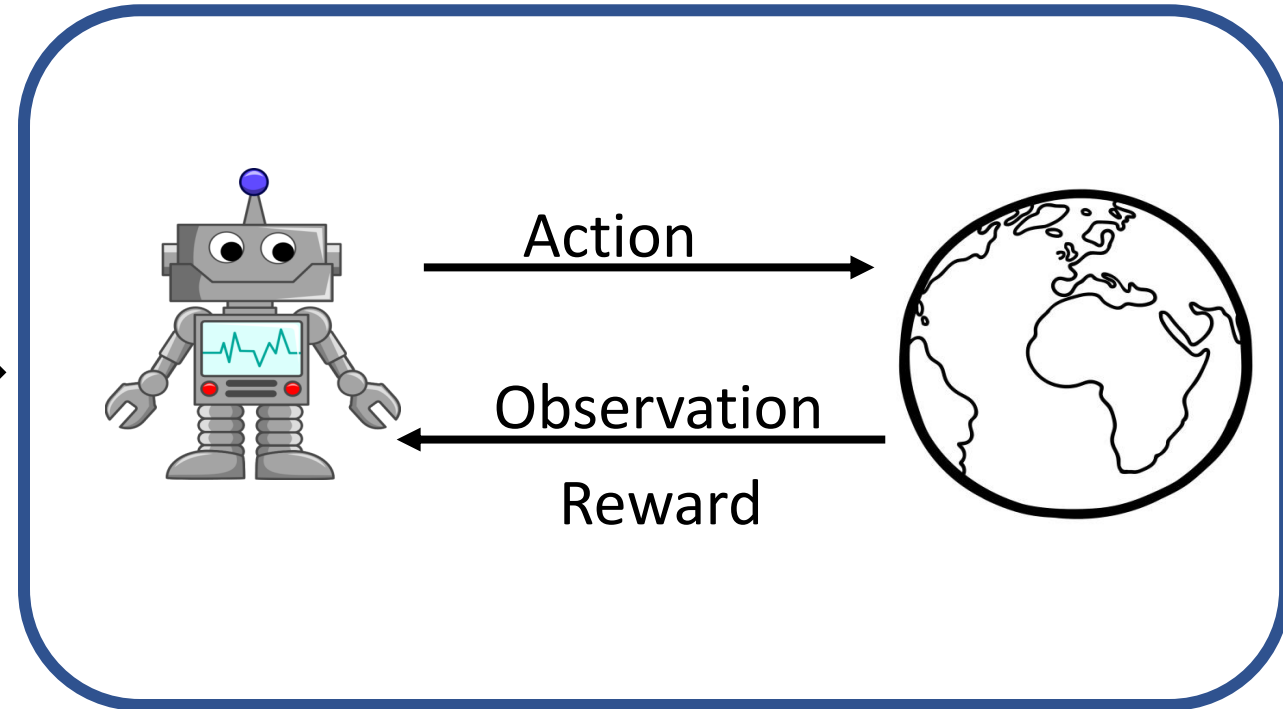
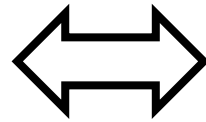
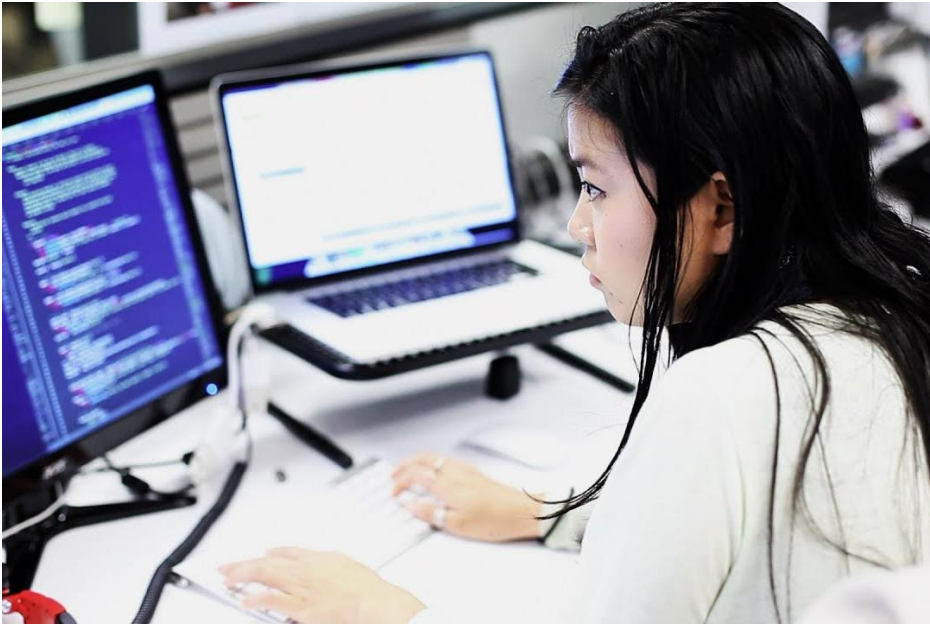
Reinforcement Learning



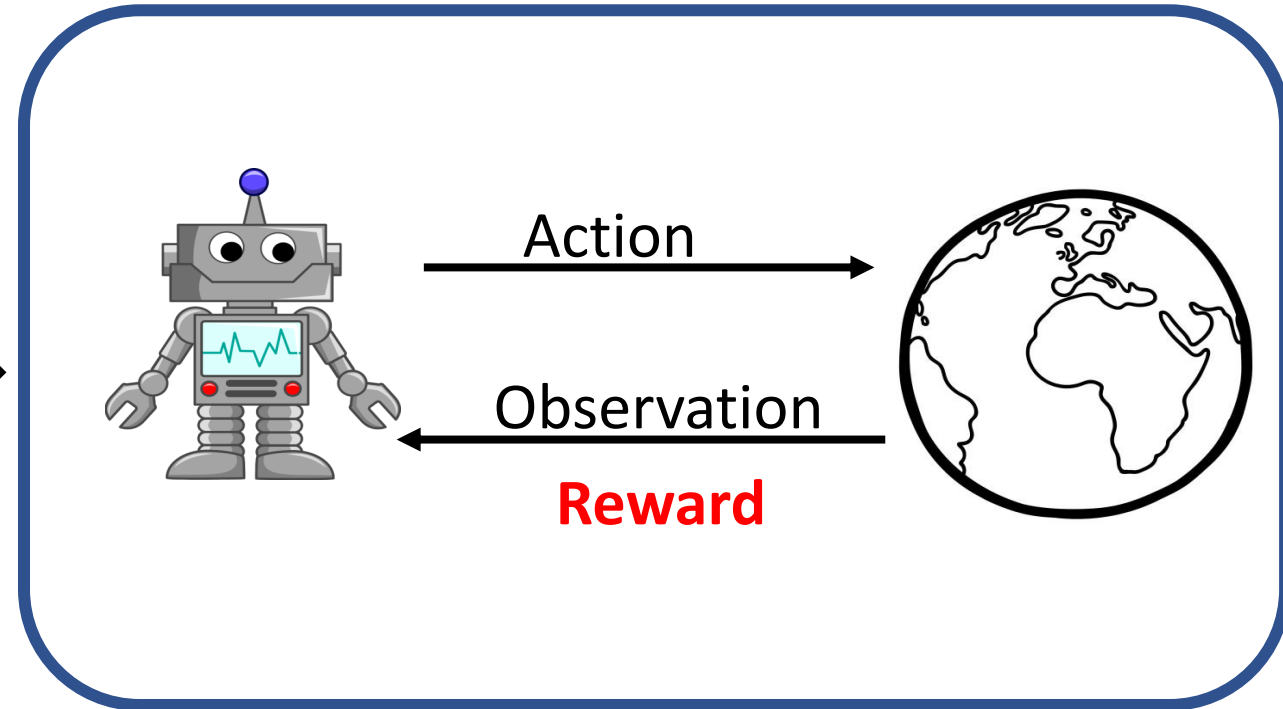
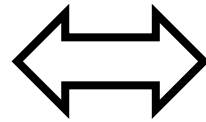
Reinforcement Learning



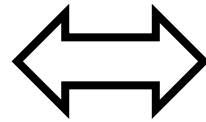
Reward engineering is hard!



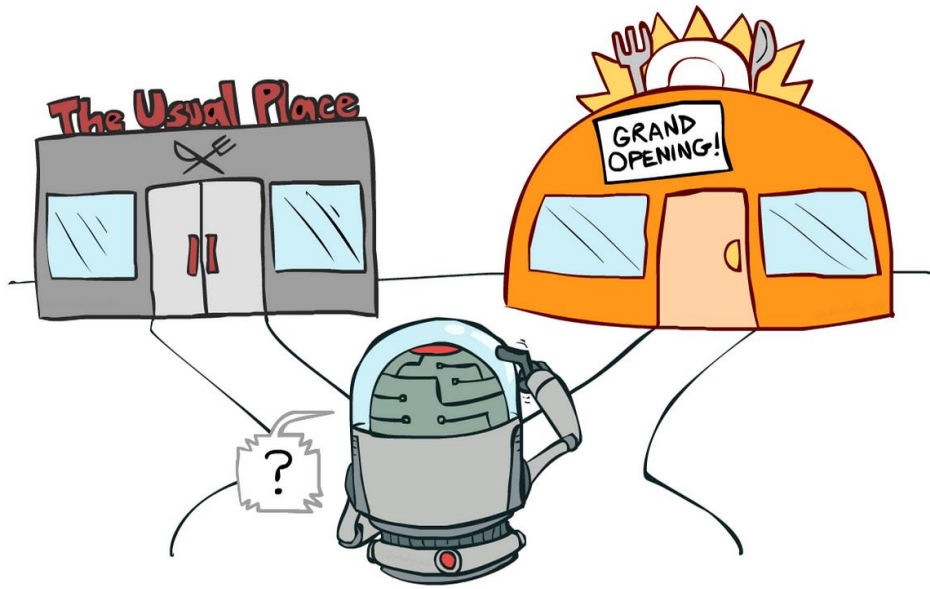
Reward engineering is hard!



Reward engineering is hard!



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



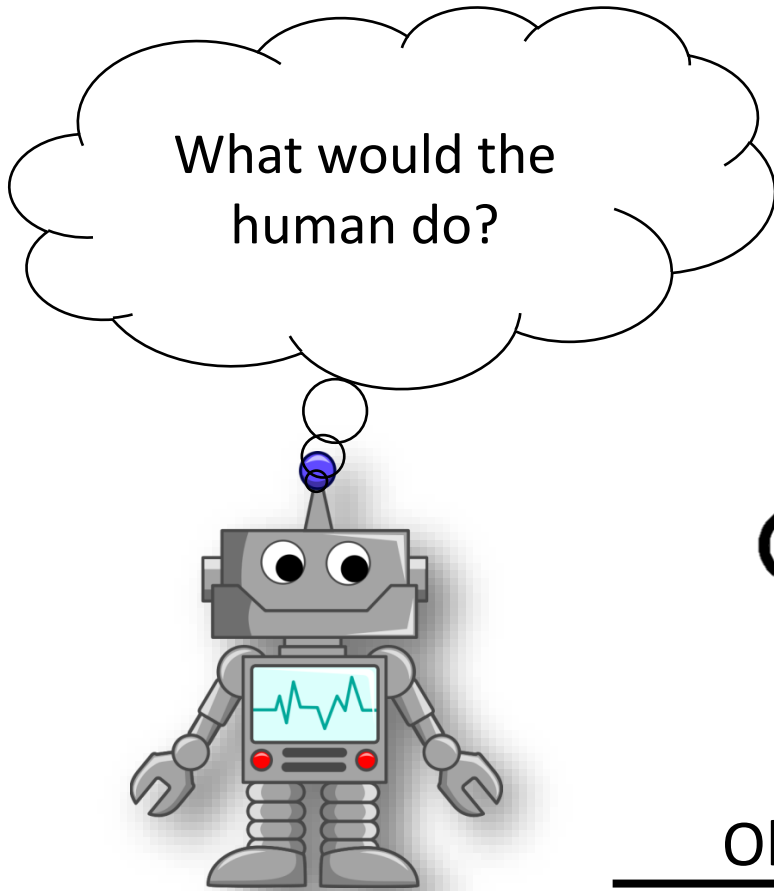
Imitation Learning:

Learn a policy from examples of good behavior.



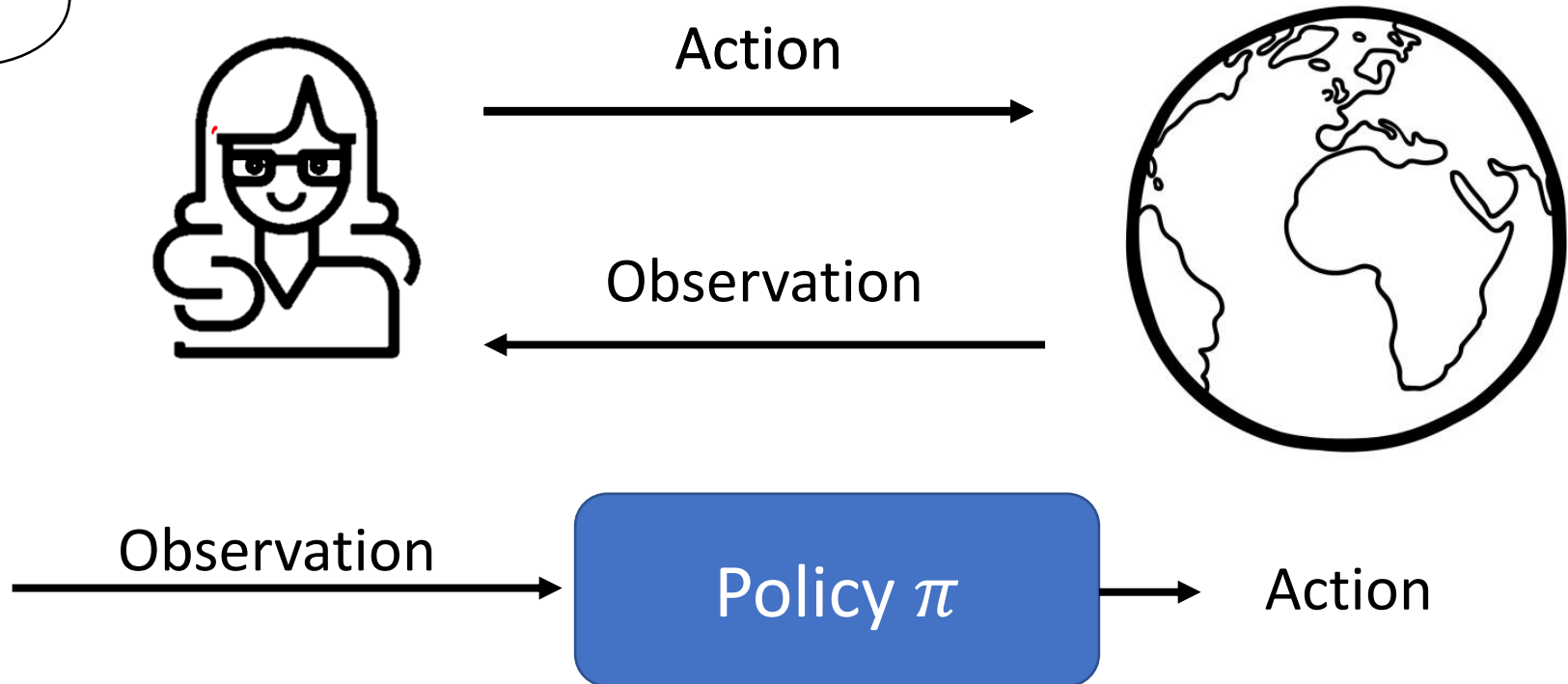
- Often showing is easier than telling.
- Alleviates problem of exploration.

Behavioral Cloning



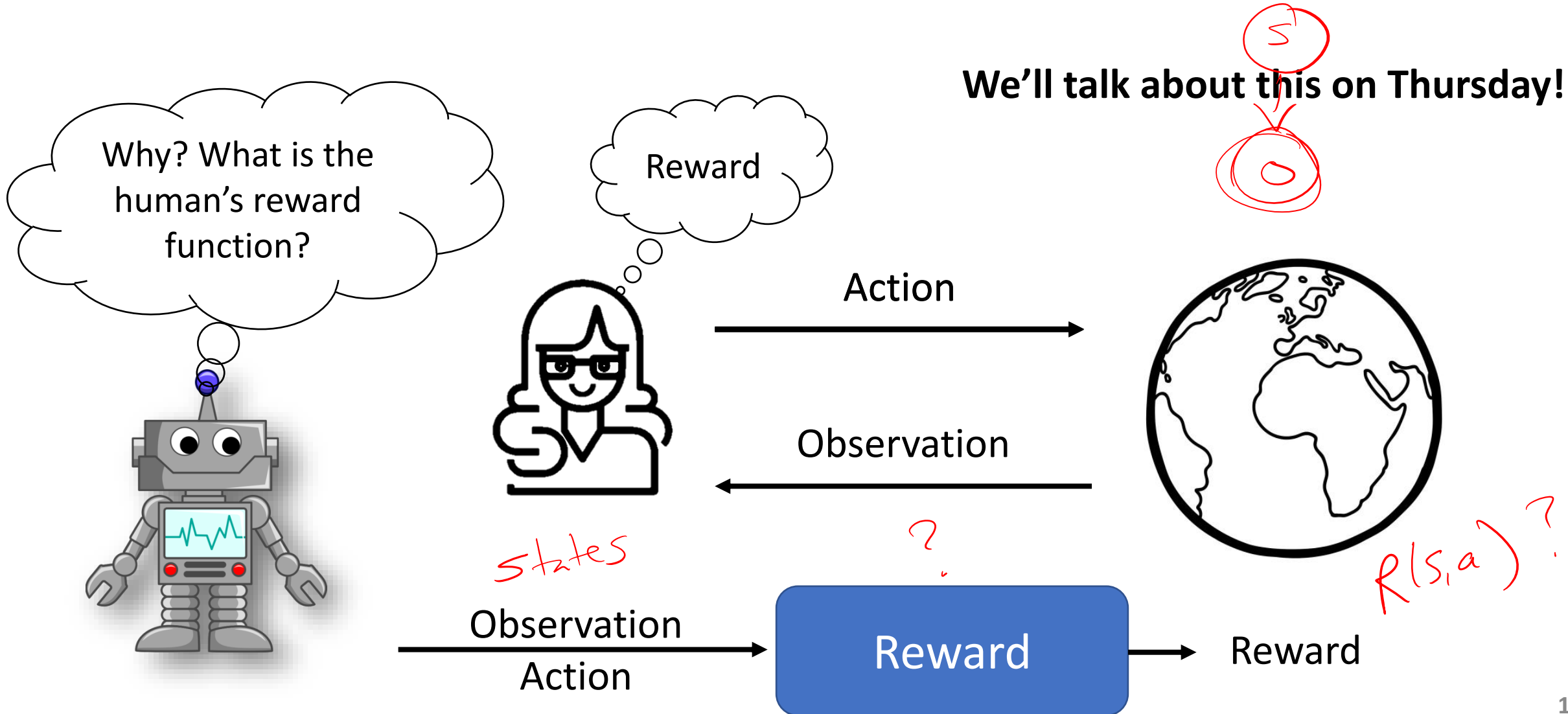
$$D = \{ (s_0, a_0) (s_1, a_1) \dots \}$$

Handwritten notes:
- $s_0 \rightarrow \boxed{} \rightarrow a_0$ with "Classifier, regression" written below the box.
- "disc actions" with an arrow pointing to the first pair (s_0, a_0) .
- "cont. actions" with an arrow pointing to the second pair (s_1, a_1) .

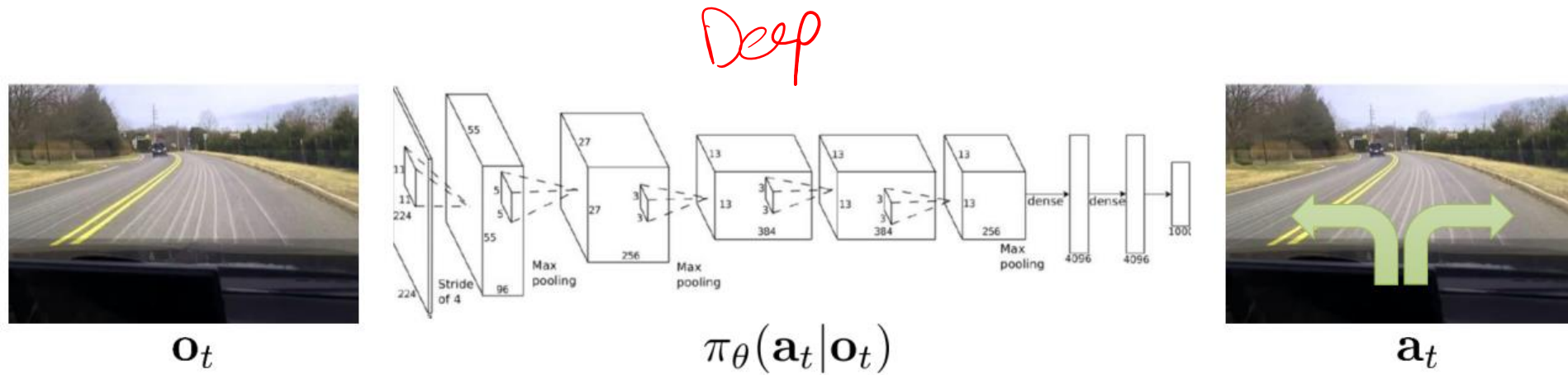


Inverse Reinforcement Learning

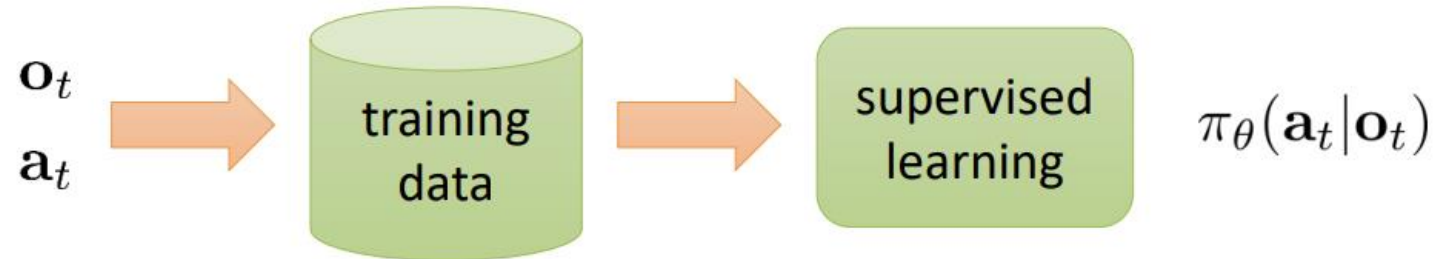
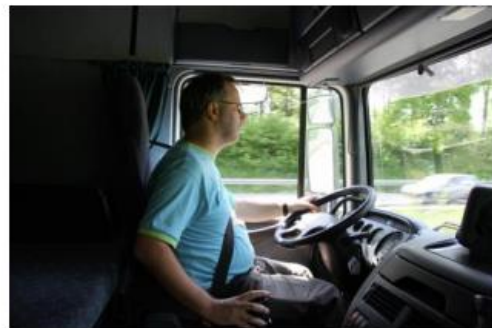
We'll talk about this on Thursday!



Imitation Learning via Behavioral Cloning



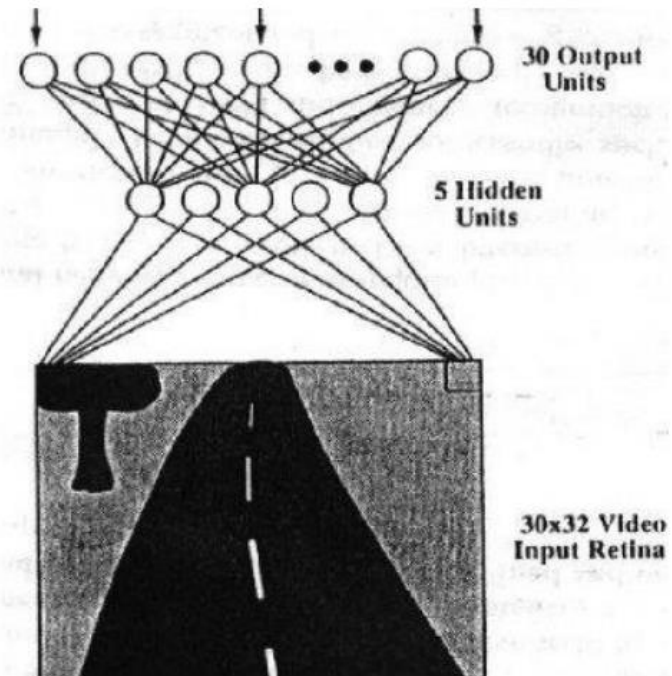
Supervised machine learning
~~input output~~
 x, y $f(x) = y$



ALVINN: One of the first imitation learning systems

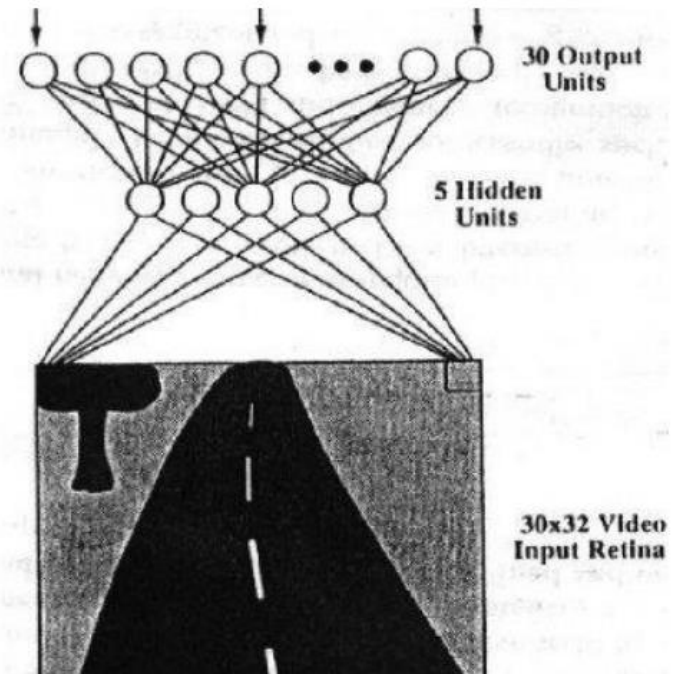
ALVINN: Autonomous Land Vehicle In a Neural Network

1989

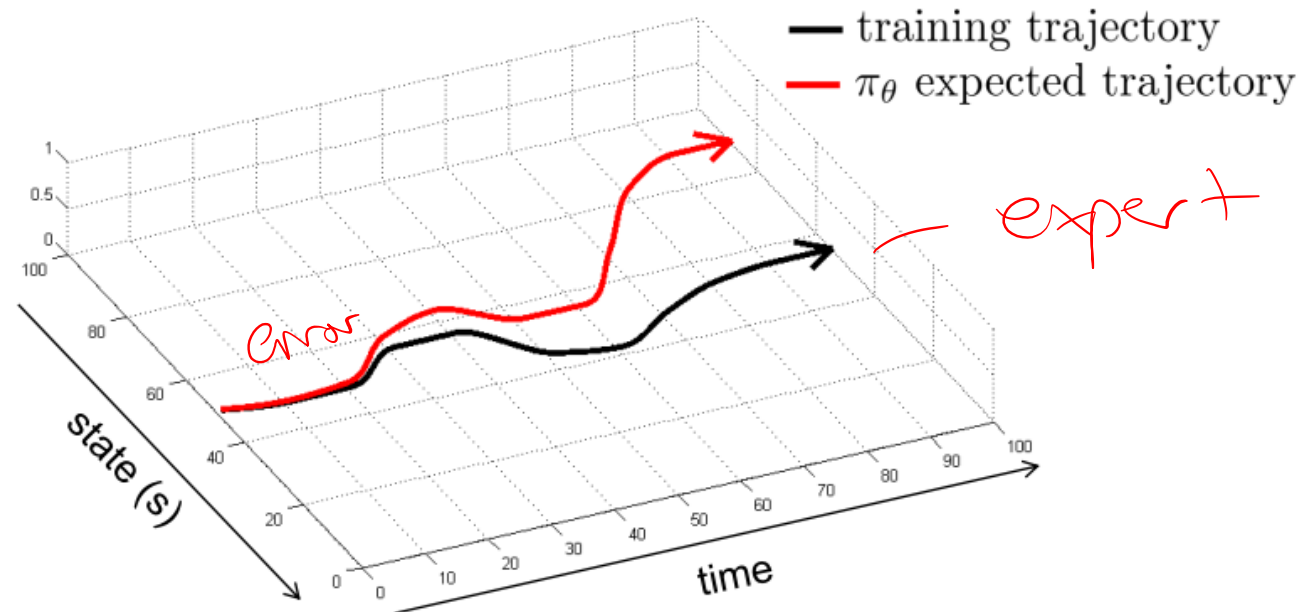
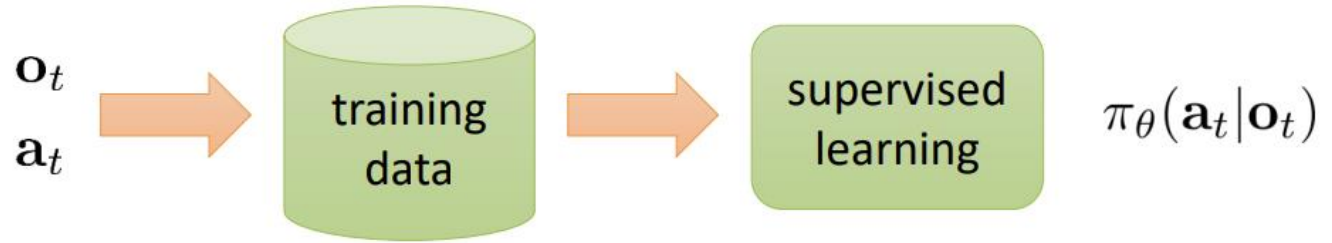
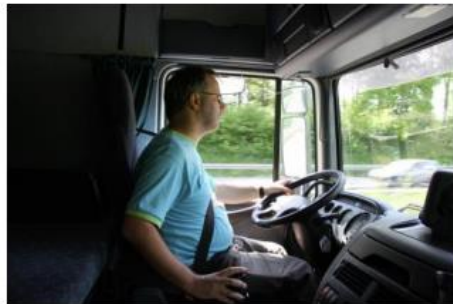


ALVINN: One of the first imitation learning systems

ALVINN: Autonomous Land Vehicle In a Neural Network
1989

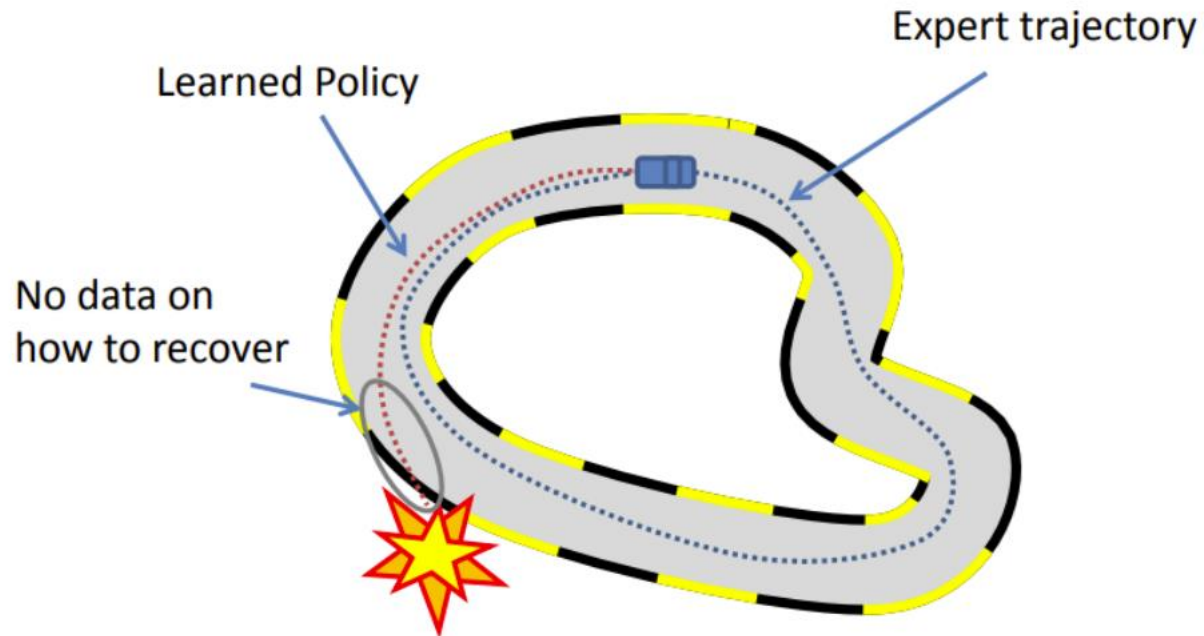


What could go wrong?



Distribution Shift

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



Demonstrator / Expert policy π^*

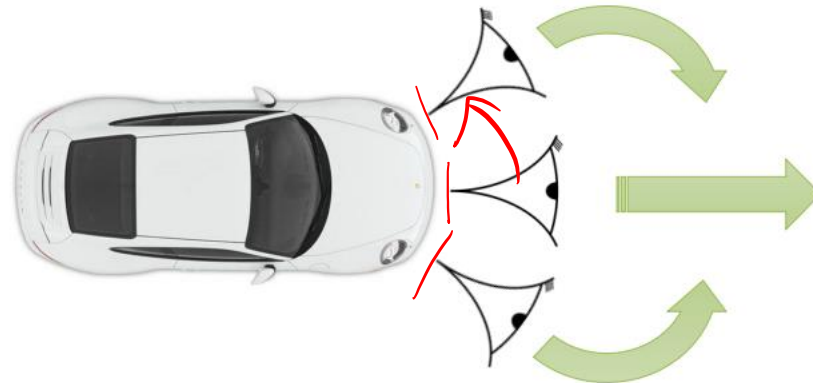
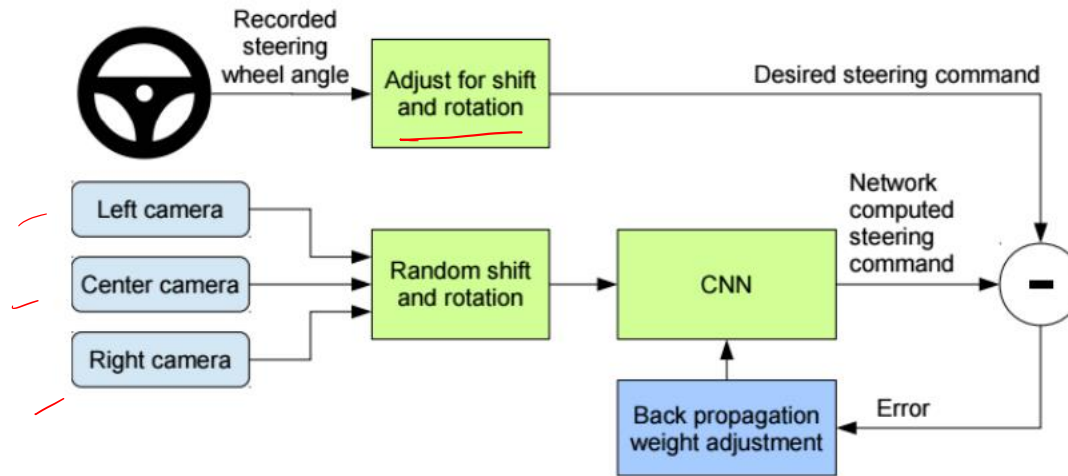
	Supervised Learning	Supervised Learning + Control
Train	$(x, y) \sim D$	$s \sim T(s, a, \pi^*(s))$
Test	$(x, y) \sim D$	$s \sim T(s, a, \pi(s))$

learned policy

But it still can work in practice...

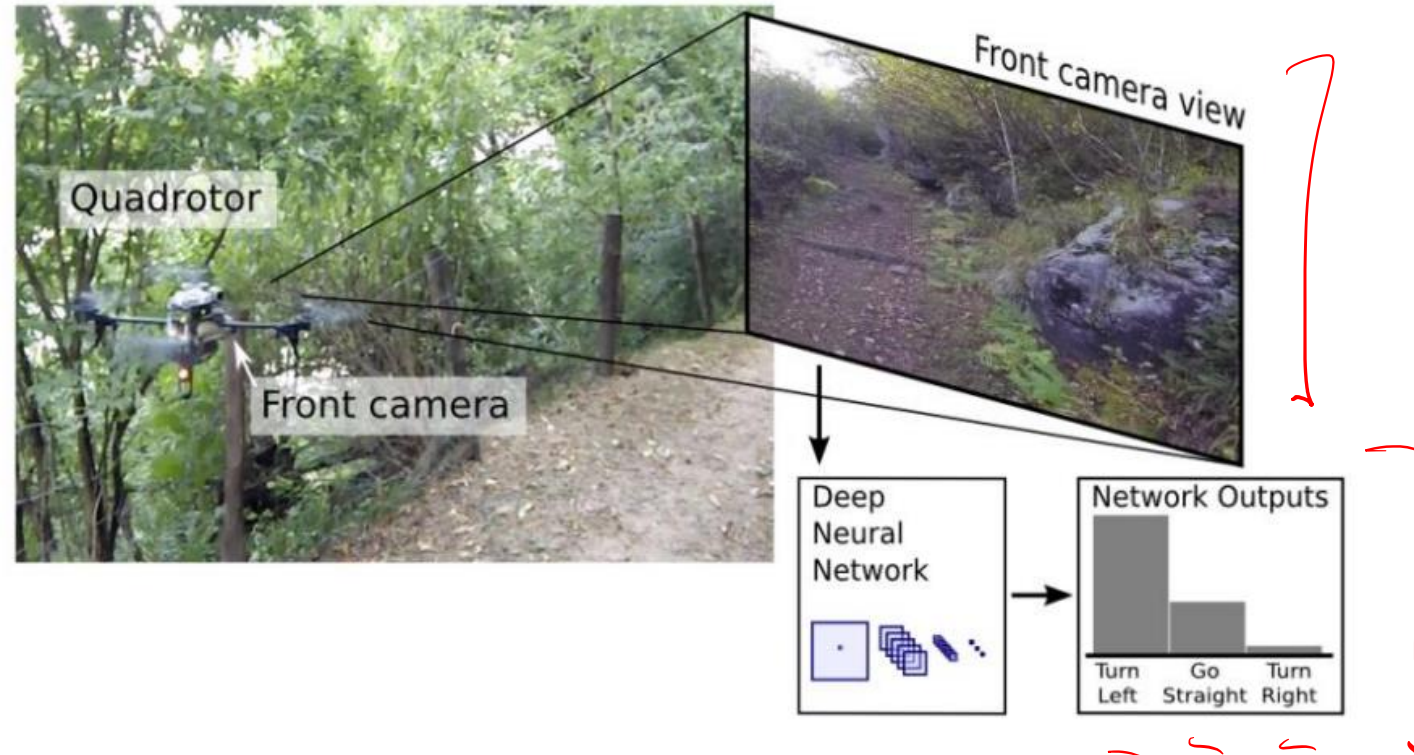


How?



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



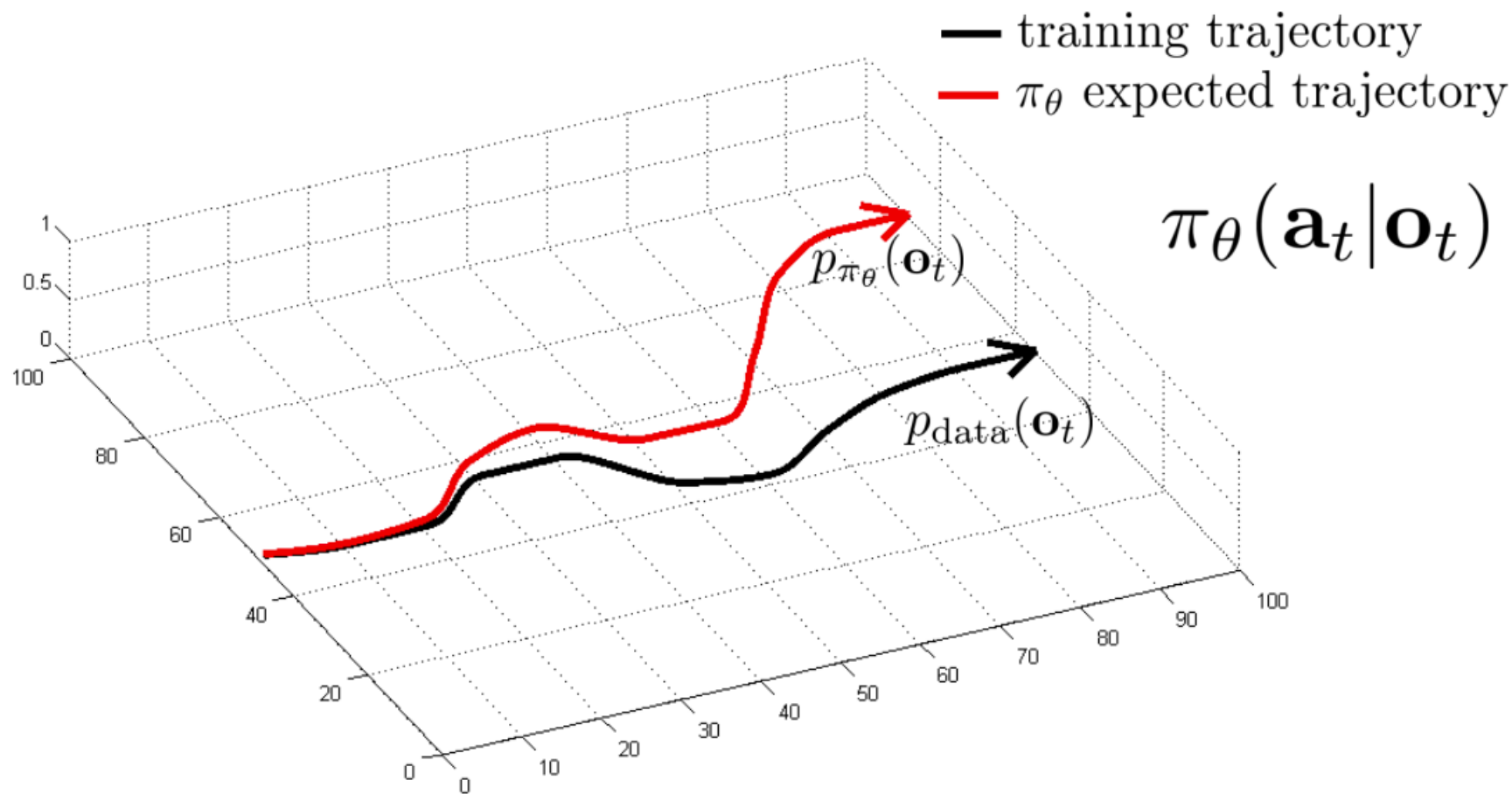
Deep Neural Network



Control Signal

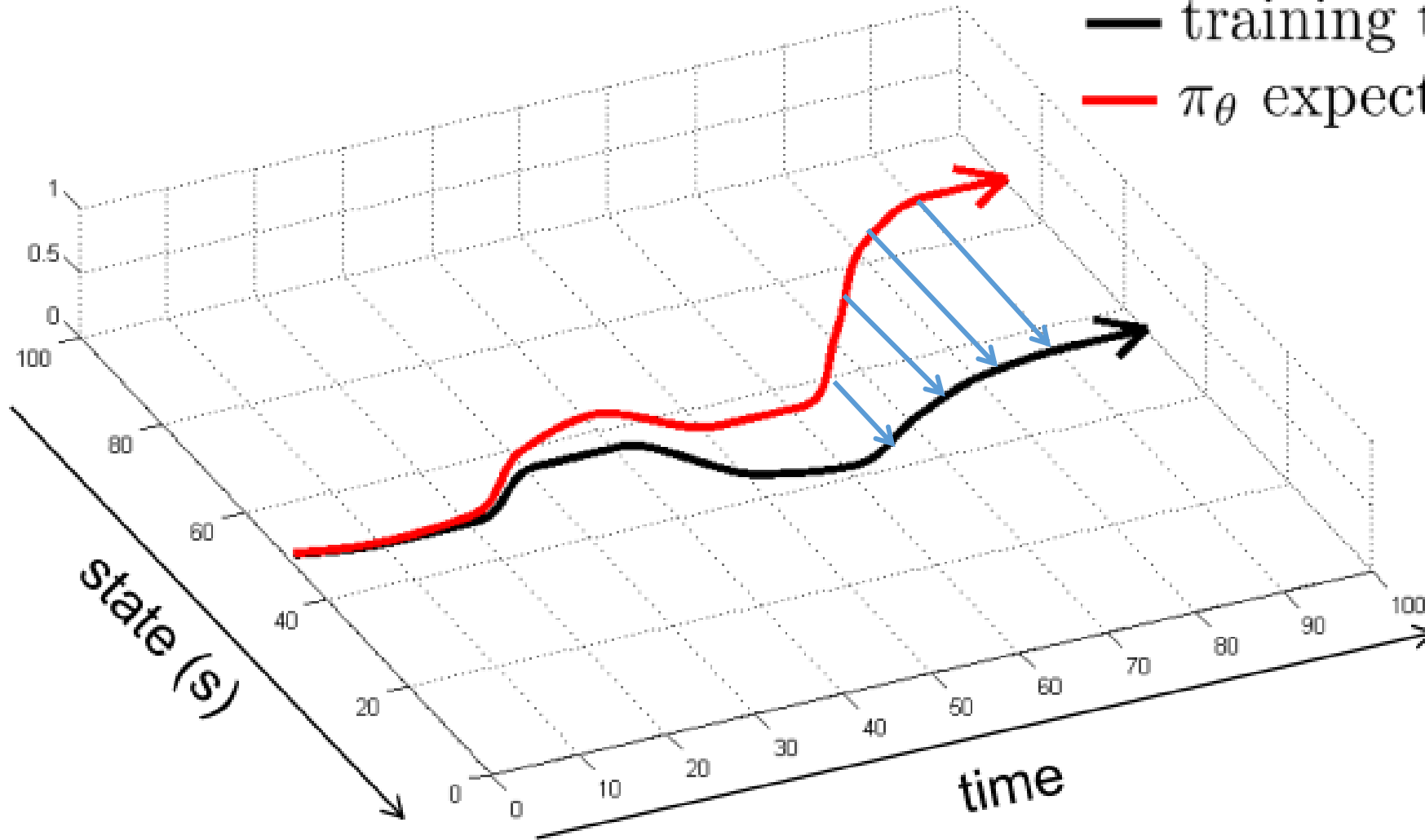


Can we make it work more often?



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

- human recovery policy
- training trajectory
- π_θ expected trajectory



Dagger

π_θ learned policy

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

$\mathcal{O}_{t=0} \dots \mathcal{O}_{t=T}$

but need labels \mathbf{a}_t !

$\mathbf{a}_0 \dots \mathbf{a}_T$

- 
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: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

only new

Dagger has very nice theoretical guarantees.


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

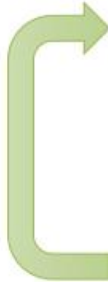
Dagger: Dataset Aggregation

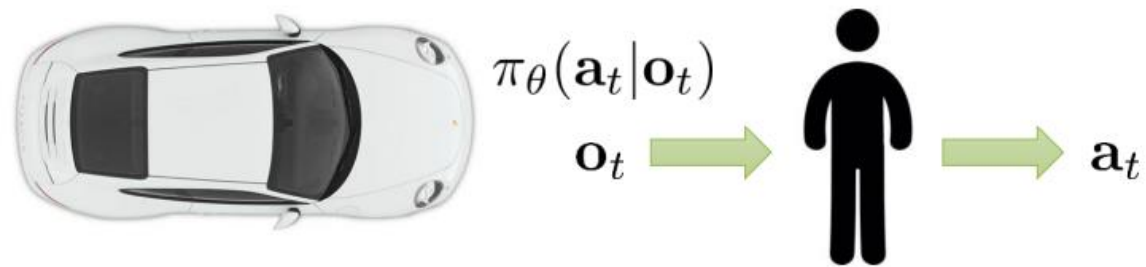
goal: collect training data from $p_{\pi_{\theta}}(\mathbf{o}_t)$ instead of $p_{\text{data}}(\mathbf{o}_t)$

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 !

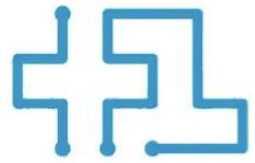
1. train $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ from human data $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
 2. run $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$ to get dataset $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
 - ~~3. Ask human to label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

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



WAYMO

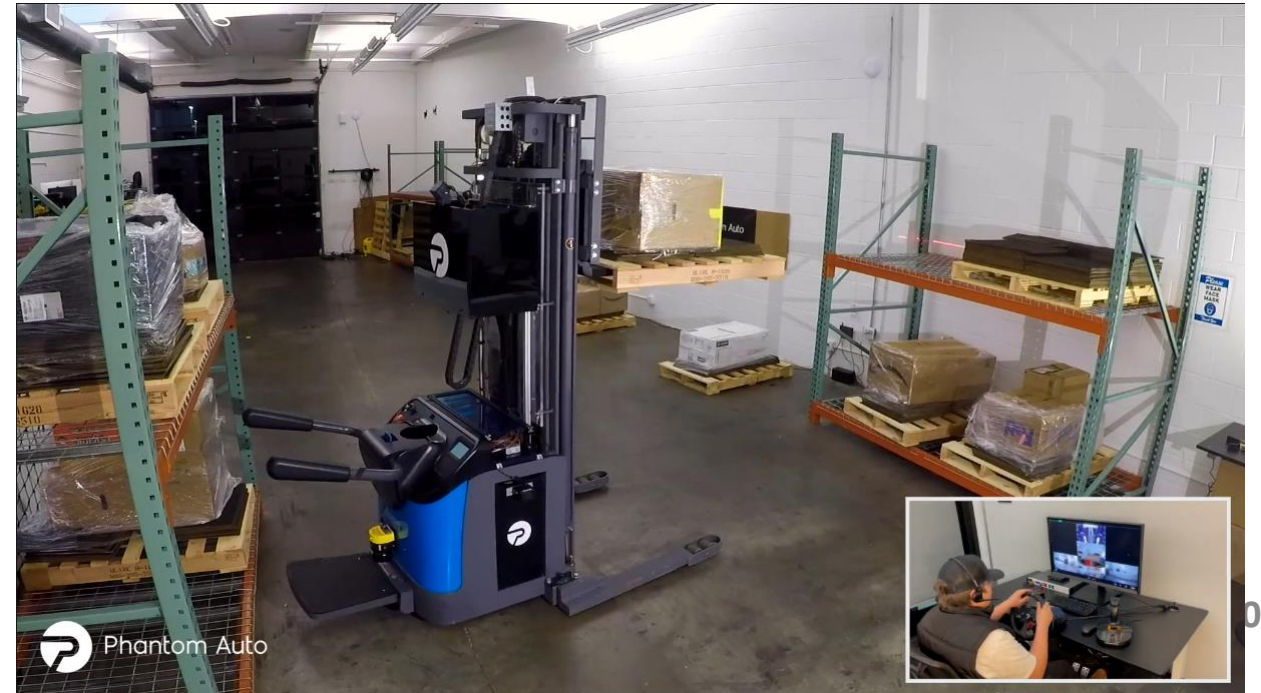
ZOOX



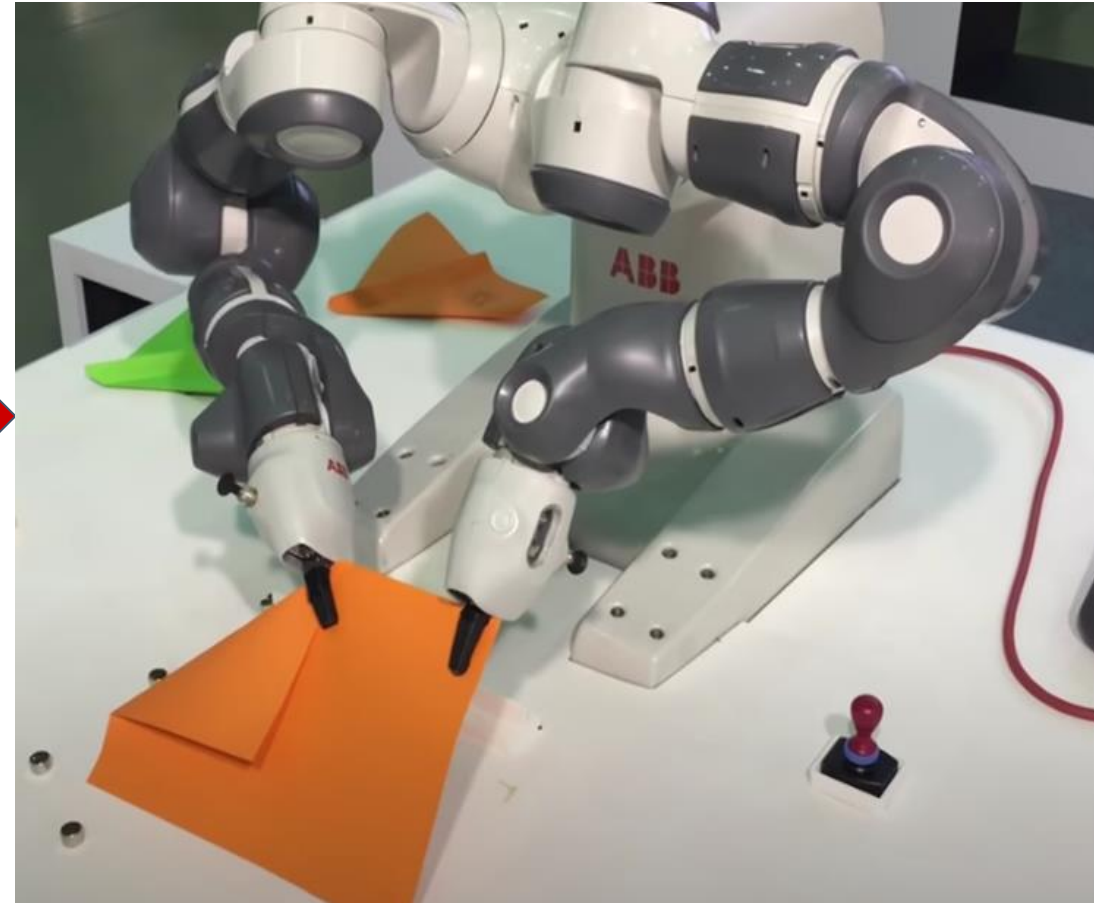
nimble



Phantom Auto



Interactive IL



Interactive IL

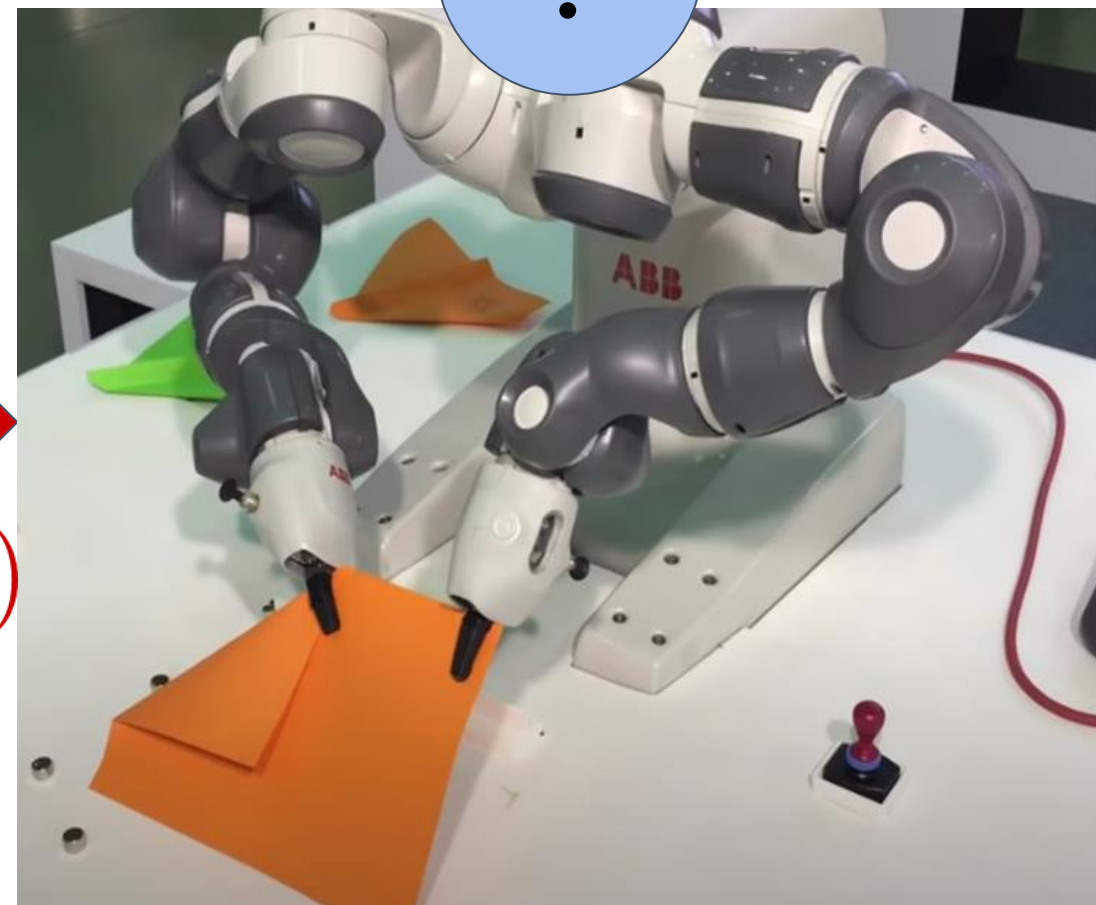


$\pi_H(s)$



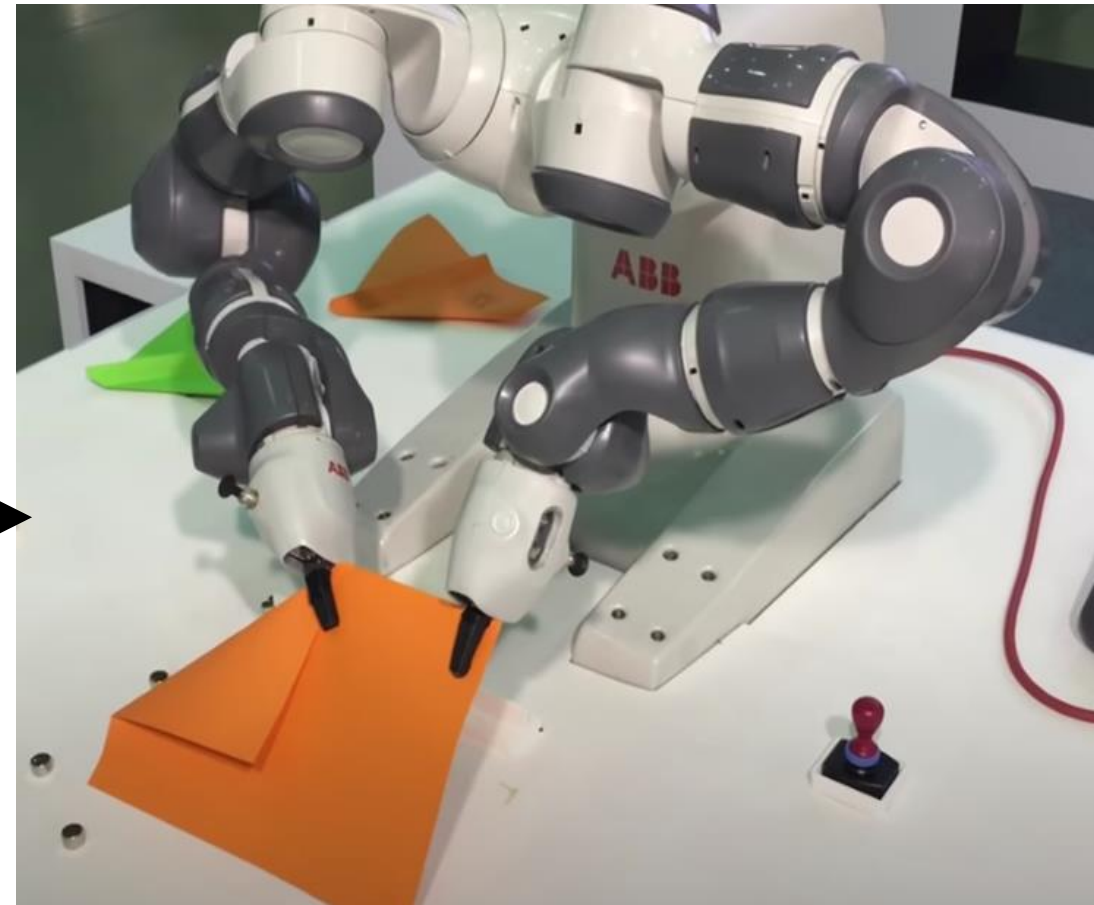
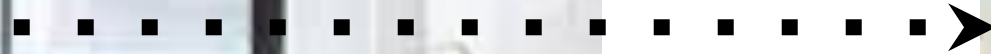
$\pi_{\text{meta}}(s)$

???



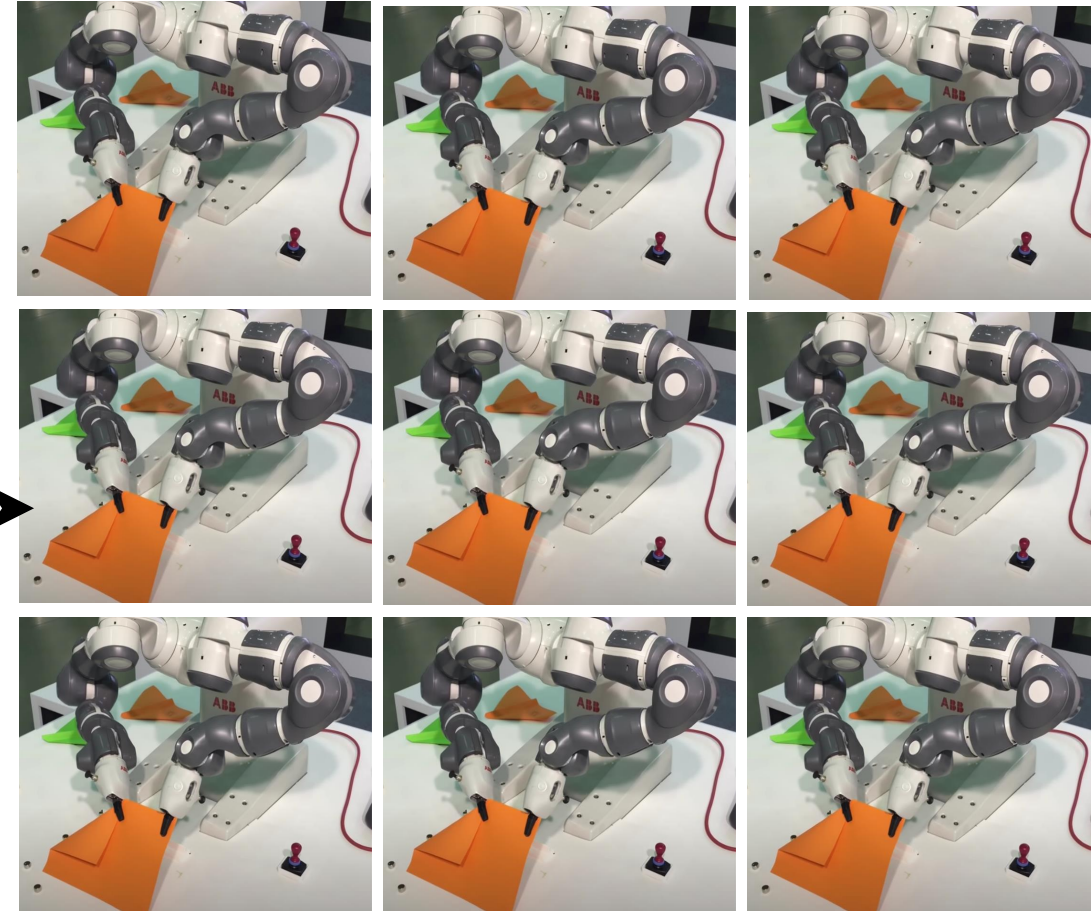
$\pi_R(s)$

Human-Gated Interactive IL



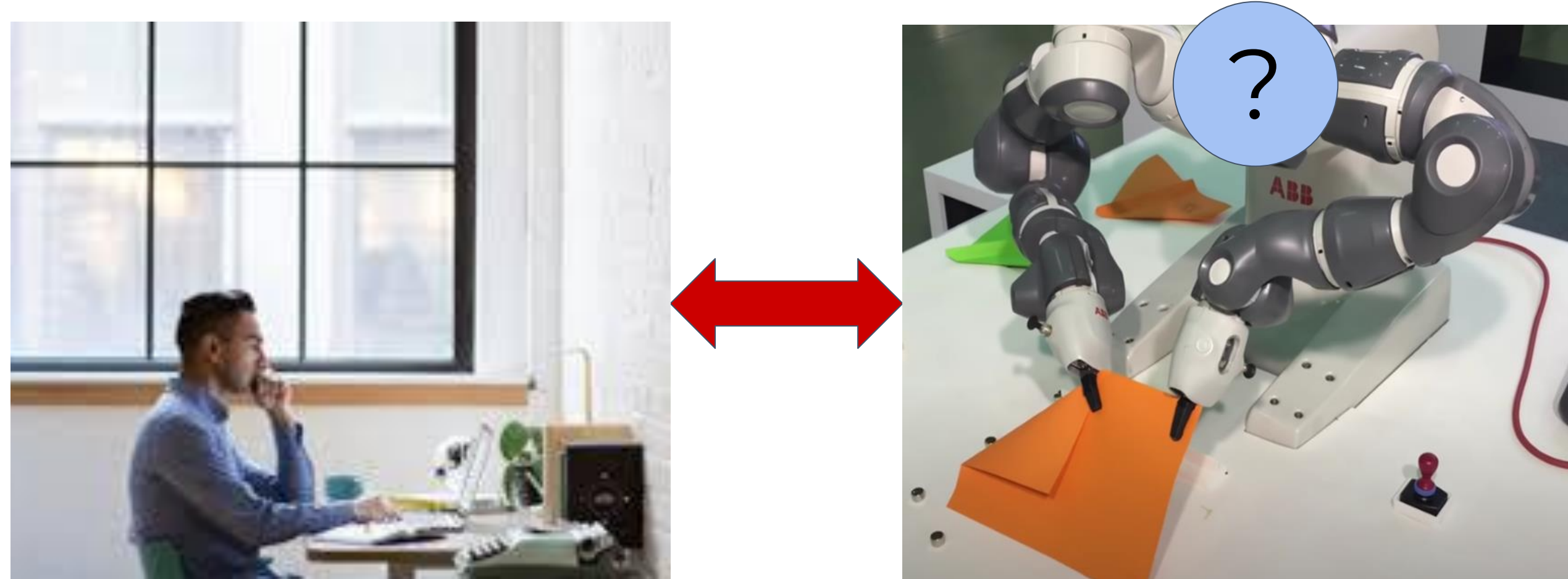
[3] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. HG-Dagger: Interactive Imitation Learning with Human Experts. ICRA 2019.

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. AAI 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

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

Loss of policy, policy error rate

Minimizing Supervisor Burden

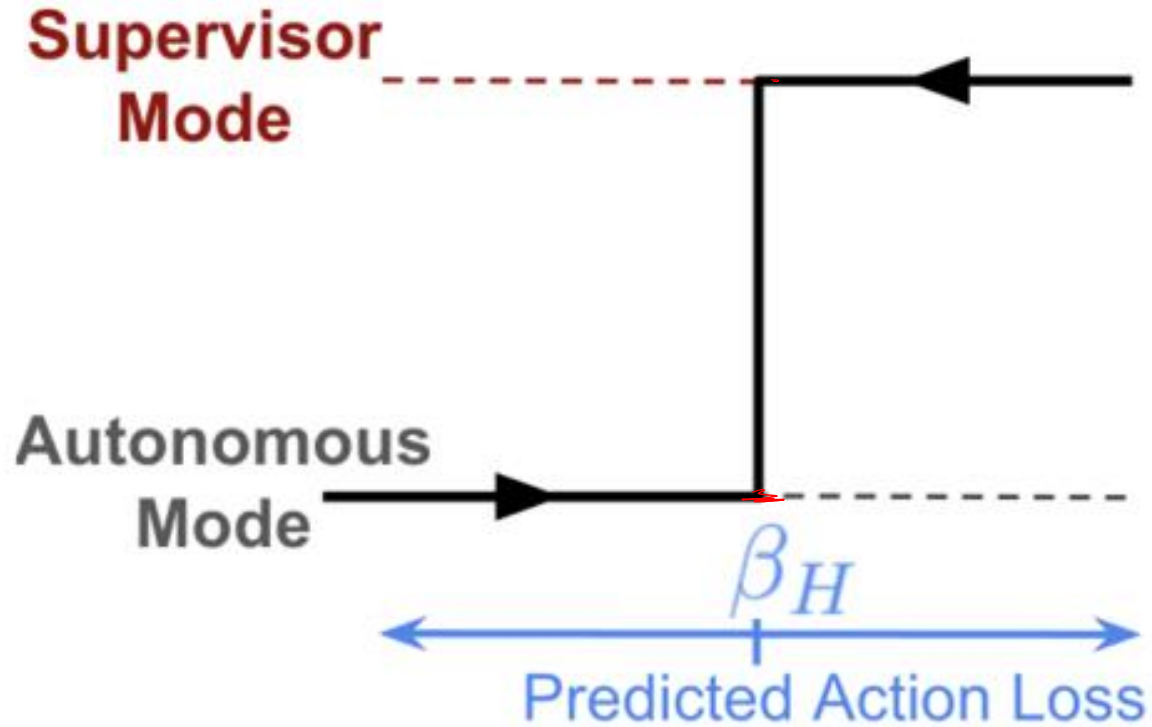
- C = number of context switches
- L = Latency of context switching
- I = expected number of supervisor actions per intervention

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

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

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

SafeDAgger

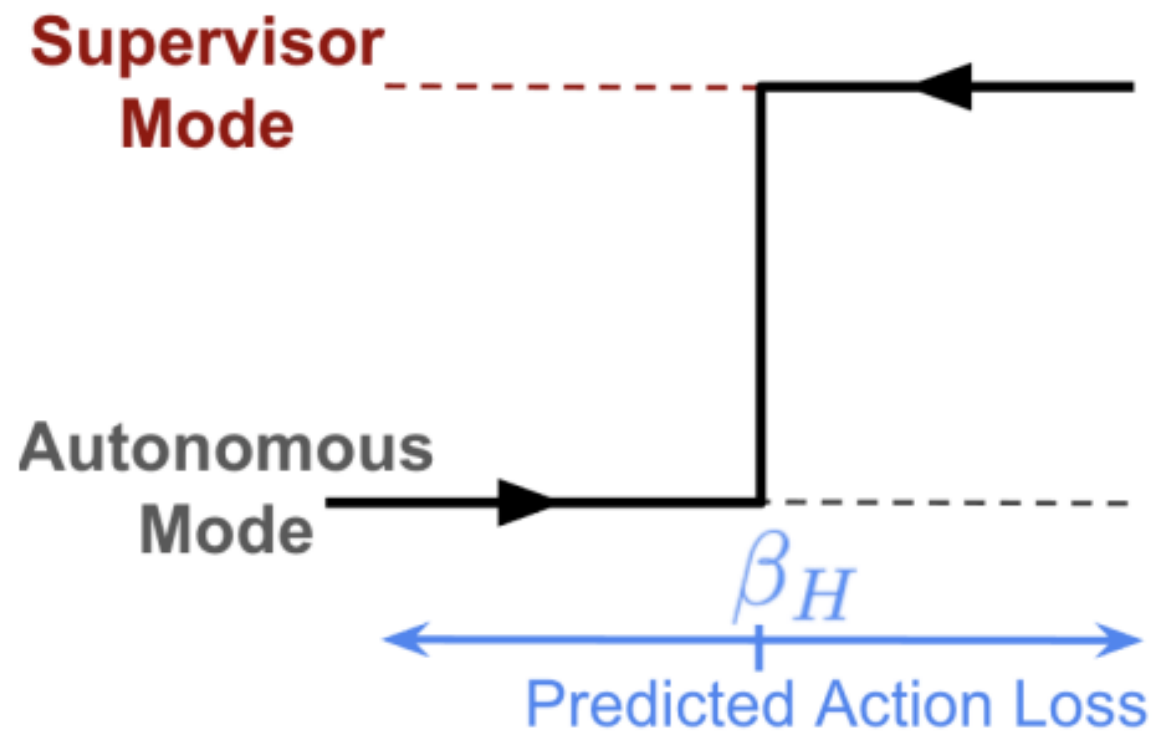


Predicted action loss = predicted difference between human and robot action.

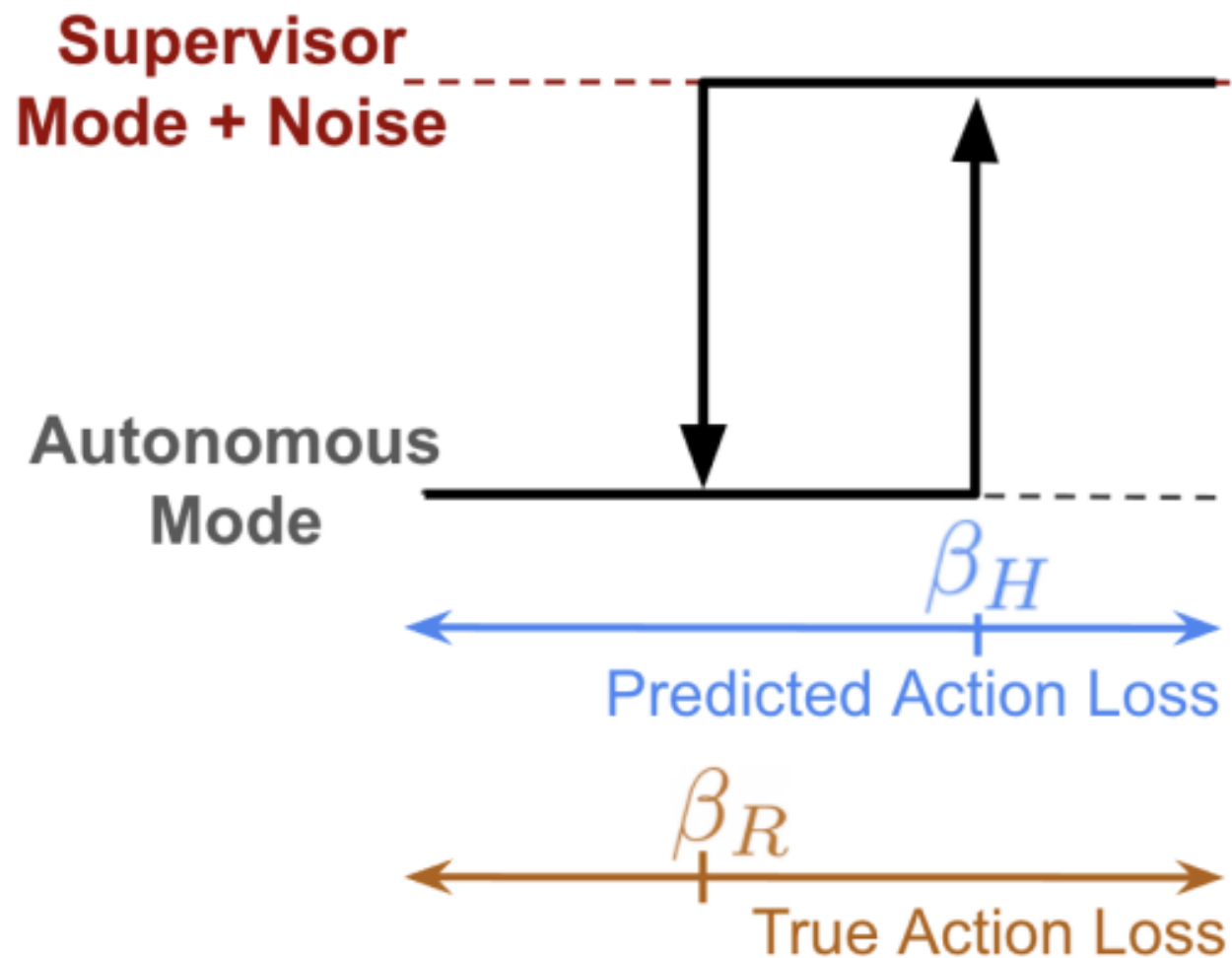
Trained using held-out set of data from human.



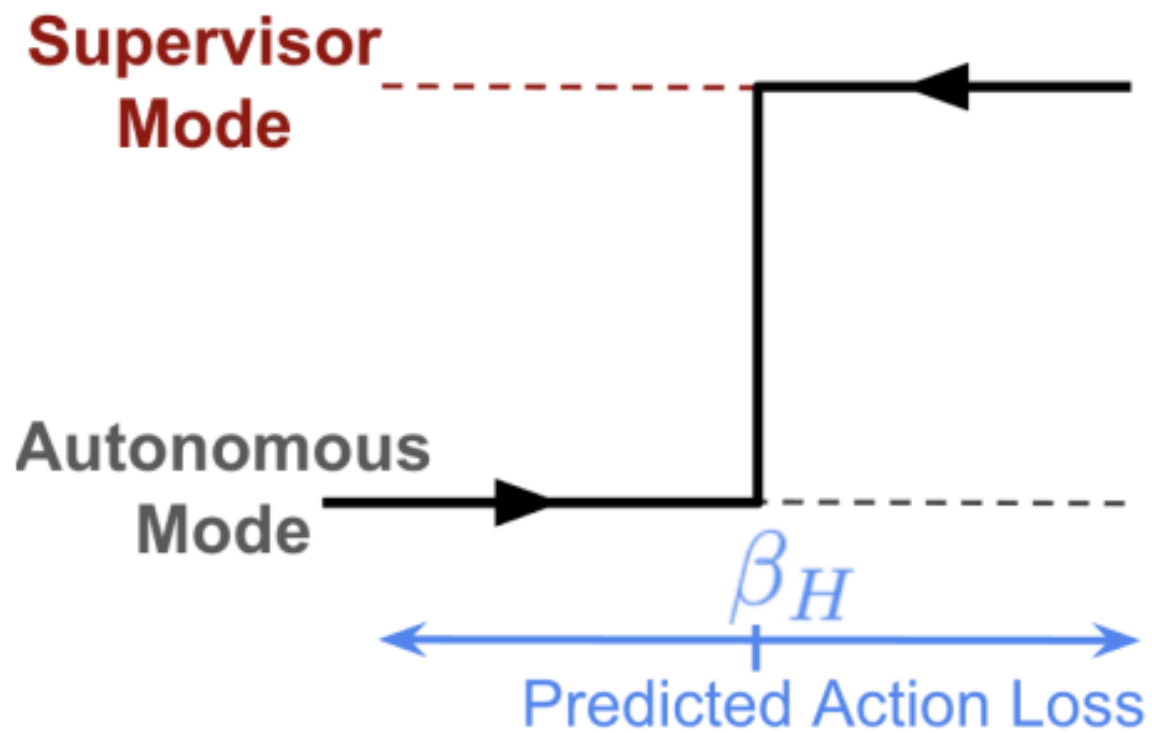
SafeDAgger



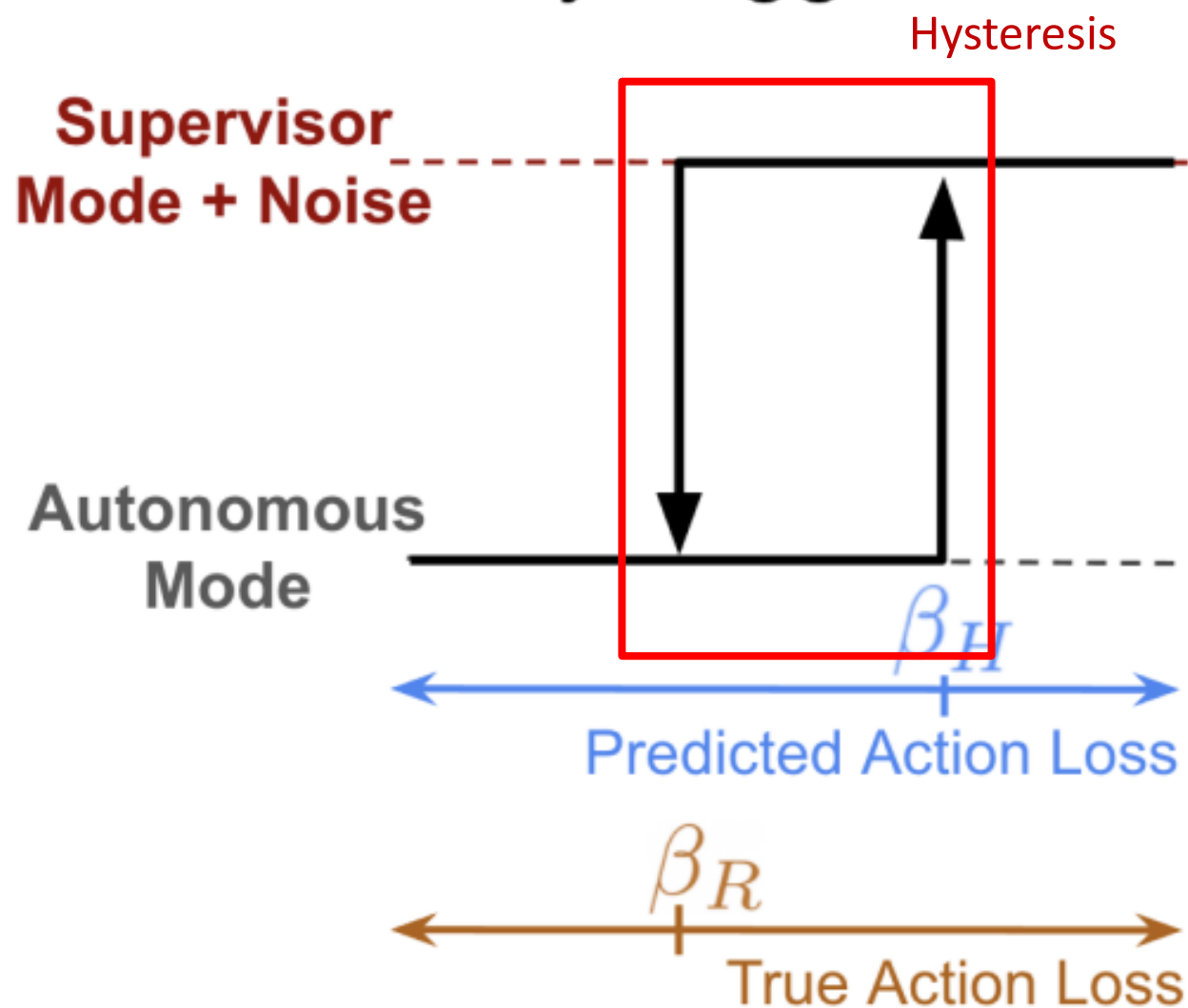
LazyDAgger



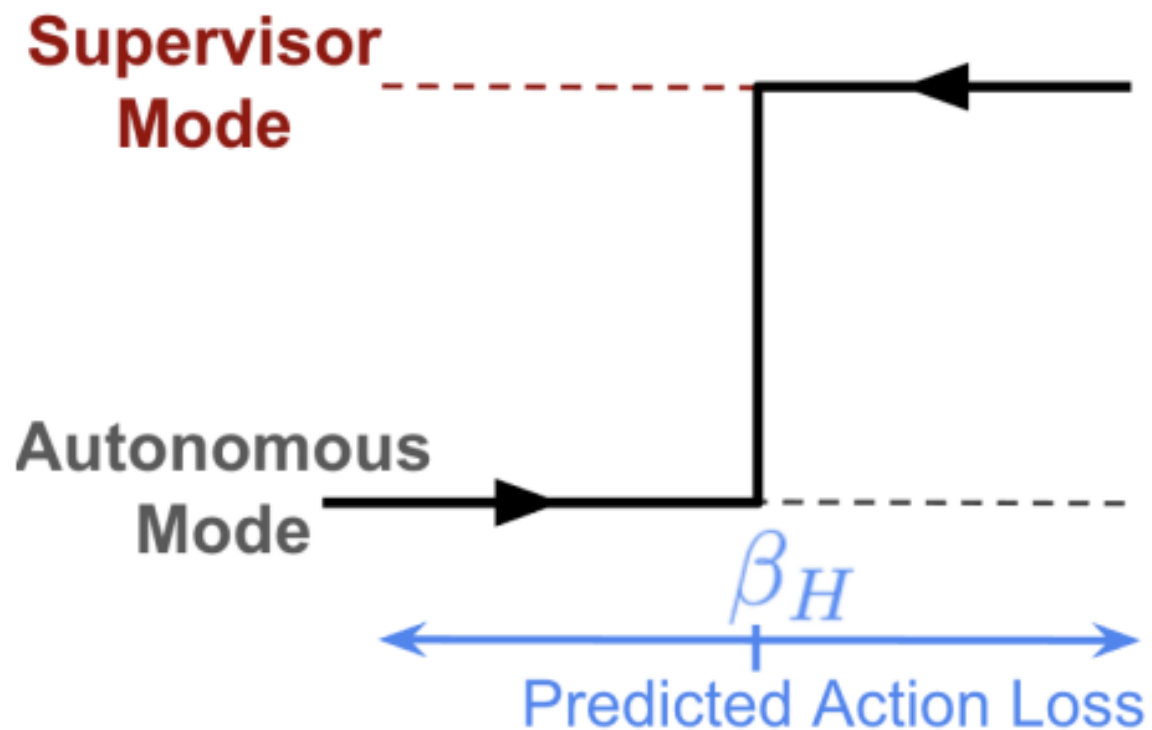
SafeDAgger



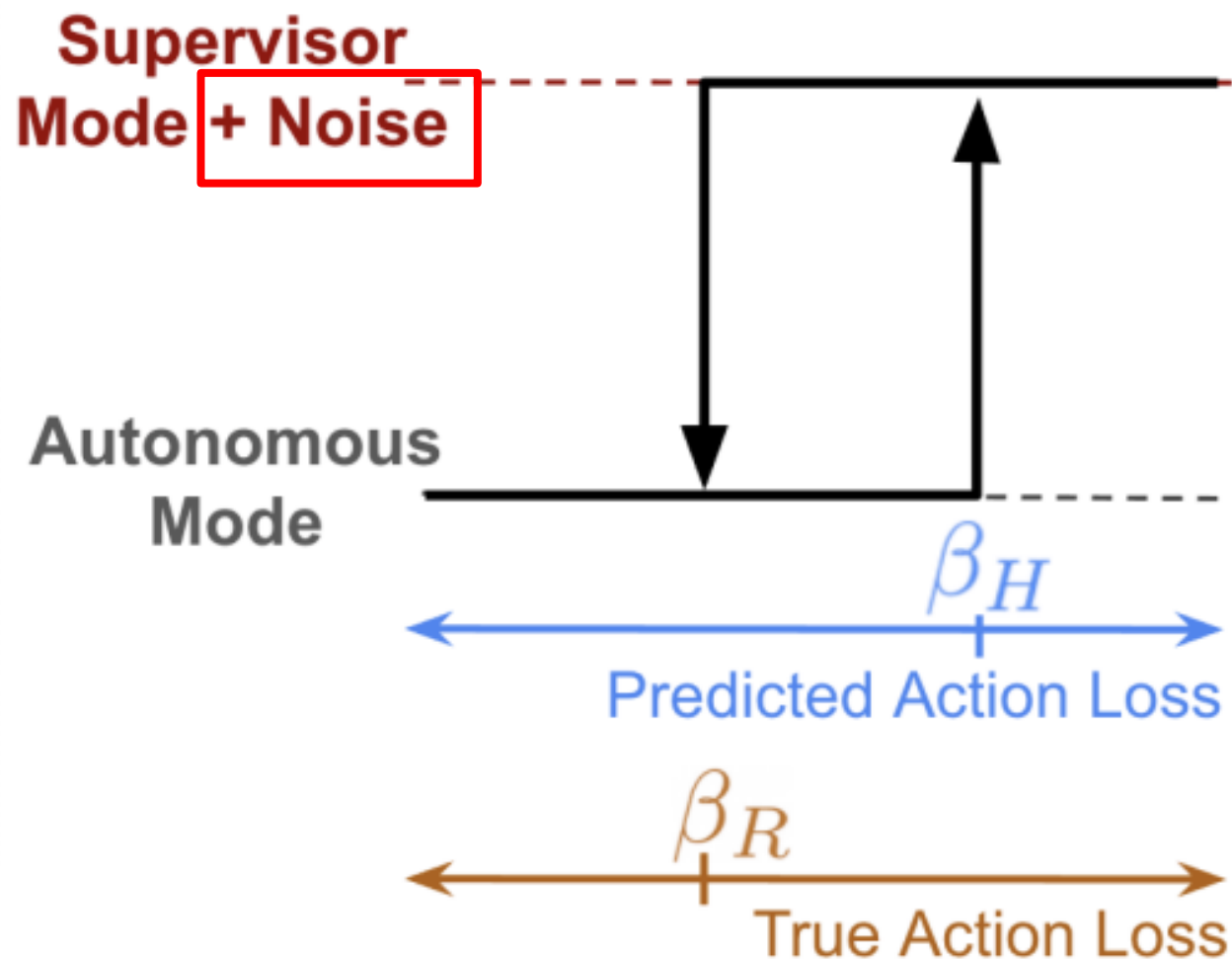
LazyDAgger



SafeDAgger

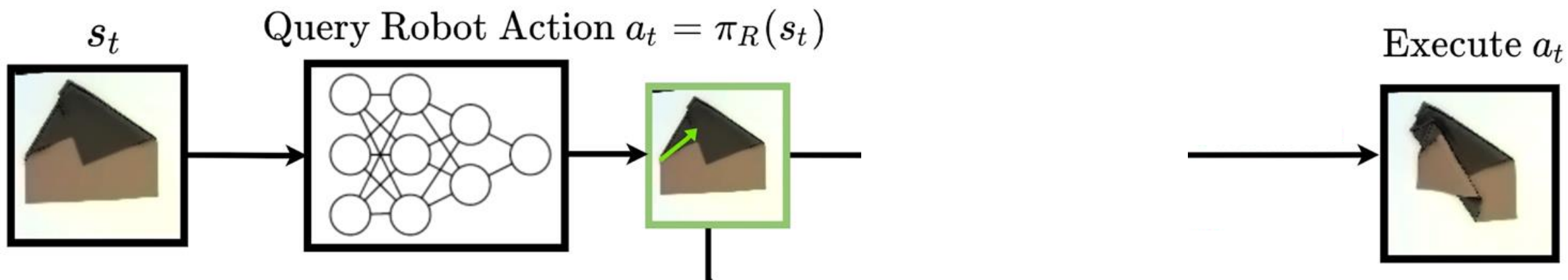


LazyDAgger

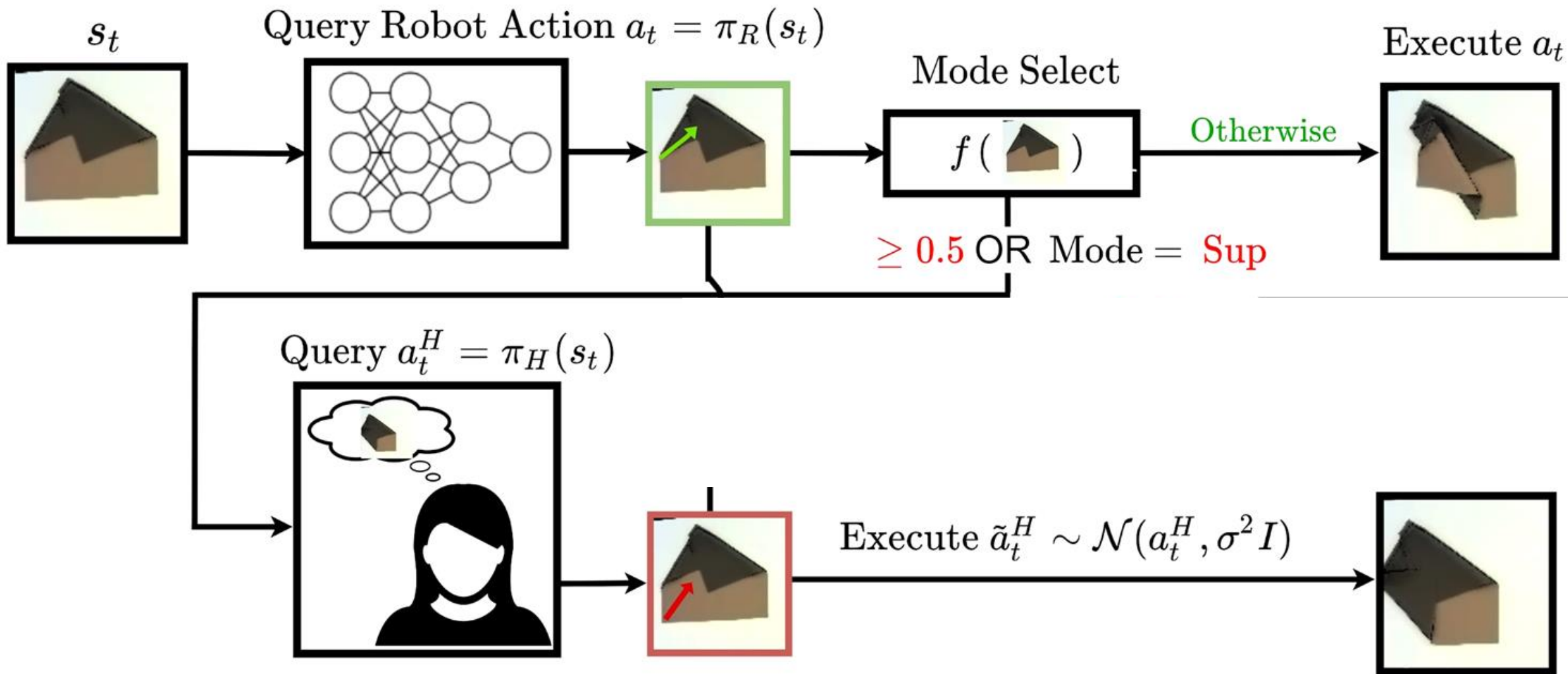


[5] M. Laskey, J. Lee, R. Fox, A. Dragan, K. Goldberg. DART: Noise Injection for Robust Imitation Learning. CoRL 2017.

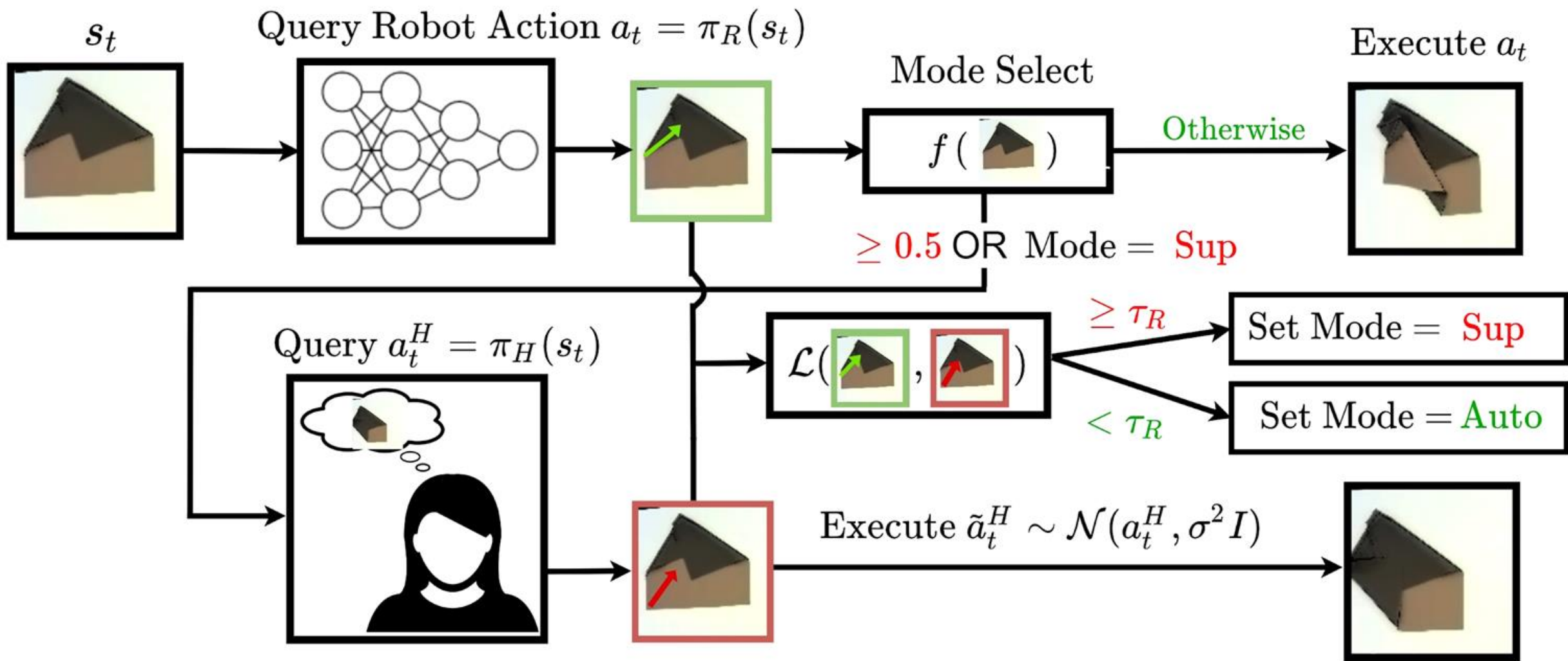
LazyDAgger



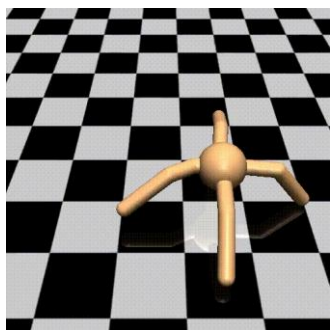
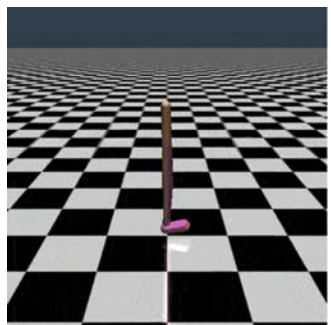
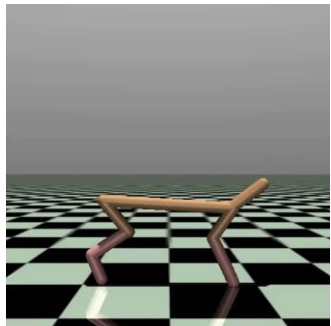
LazyDAgger



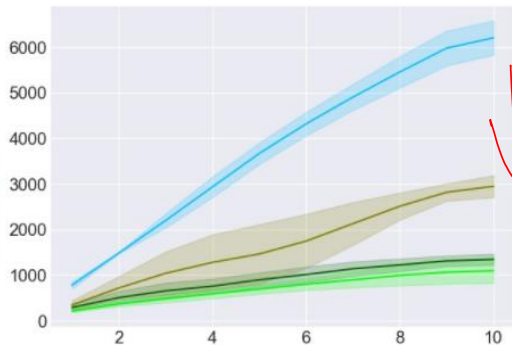
LazyDAgger



Simulation Experiments

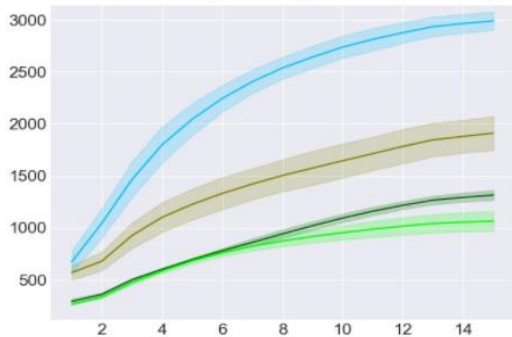


Number of Context Switches

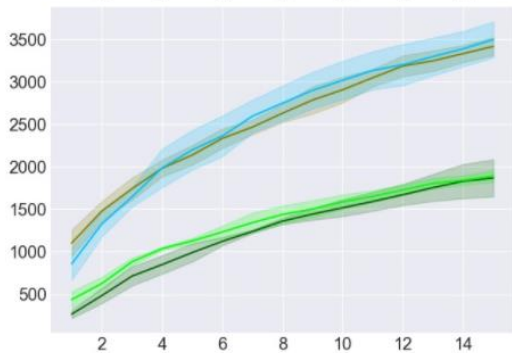


Context Switching Reduction

79%



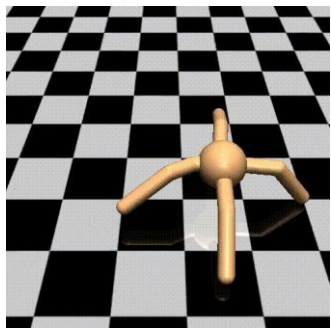
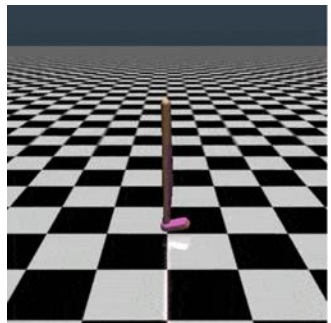
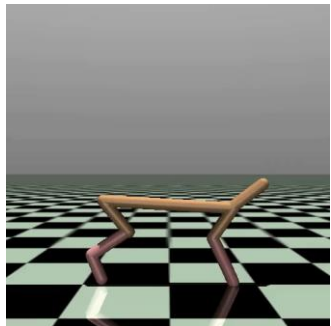
56%



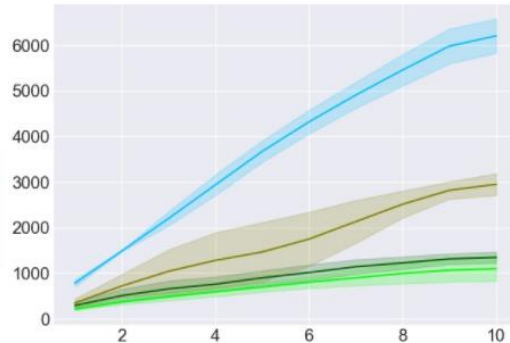
46%

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

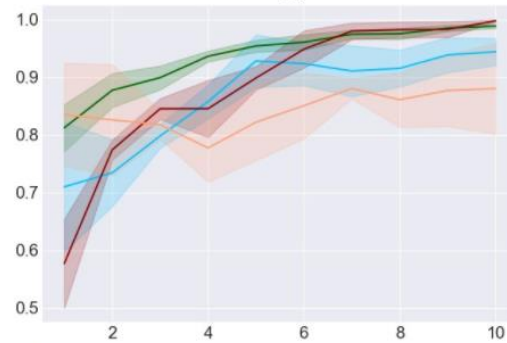
Simulation Experiments



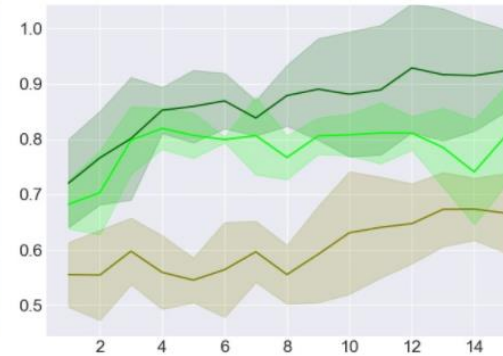
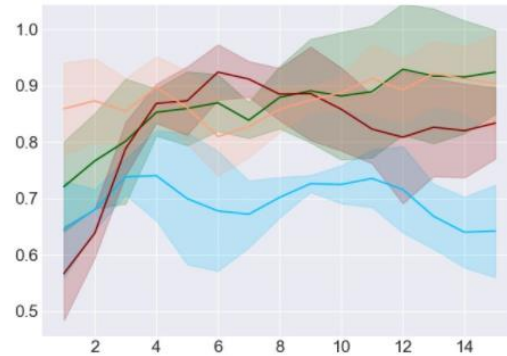
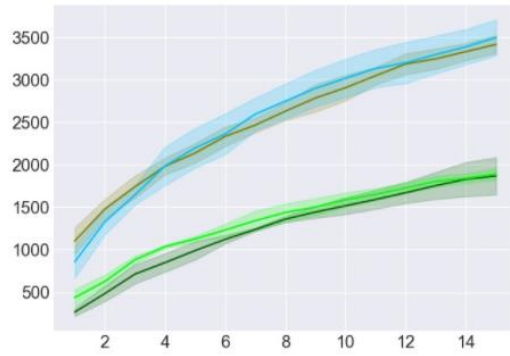
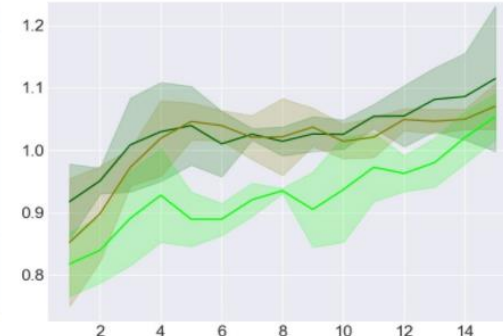
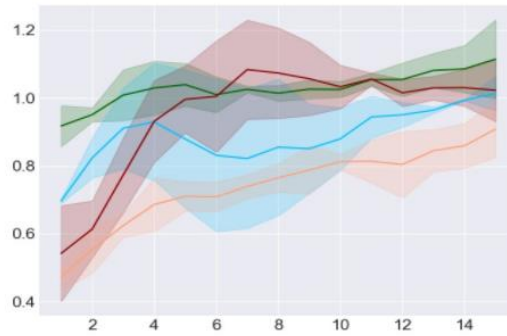
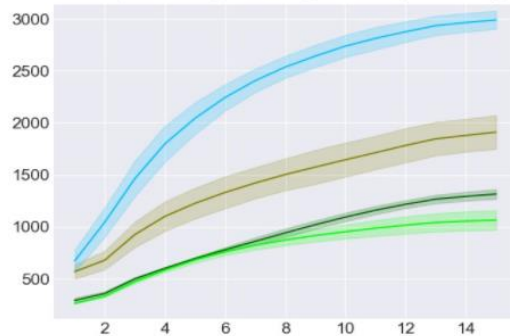
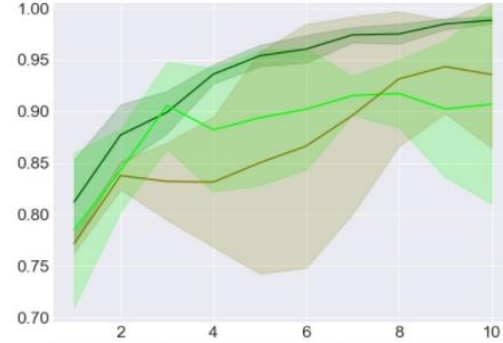
Number of Context Switches



Learning Curves

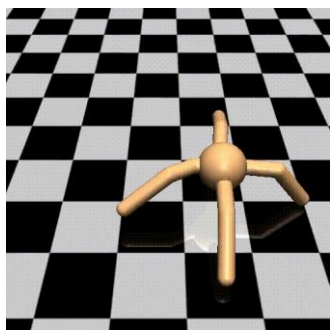
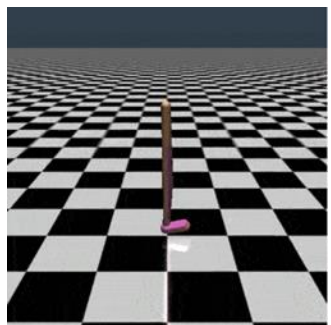
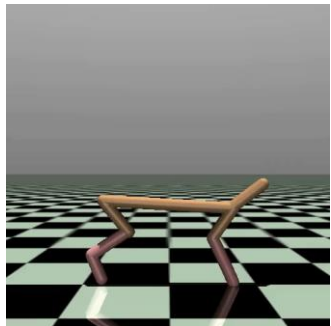


Learning Curves: Ablations

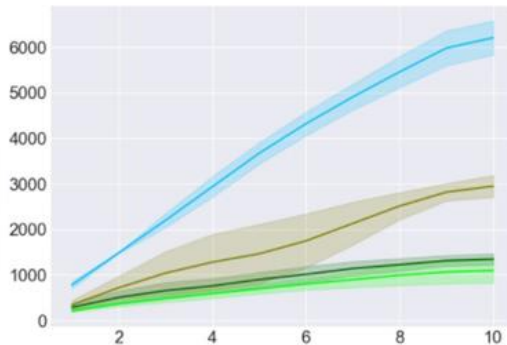


Behavior Cloning DAgger SafeDAgger Ours Ours (- Noise) Ours (- Switch to Auto)

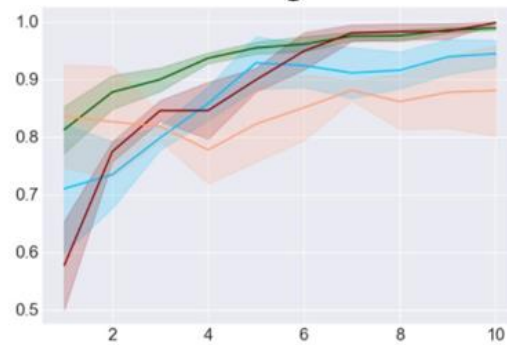
Simulation Experiments



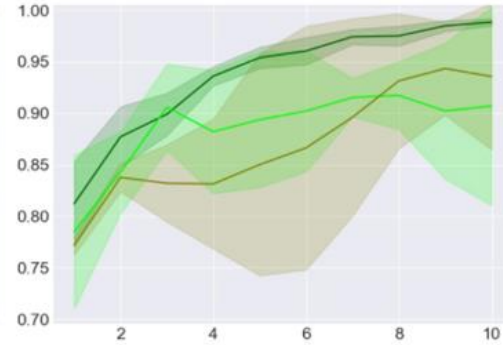
Number of Context Switches



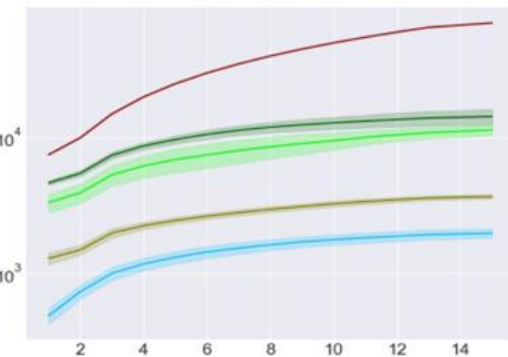
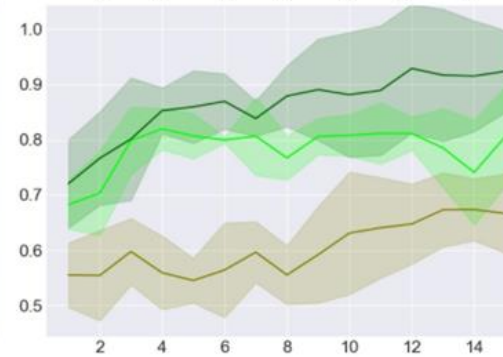
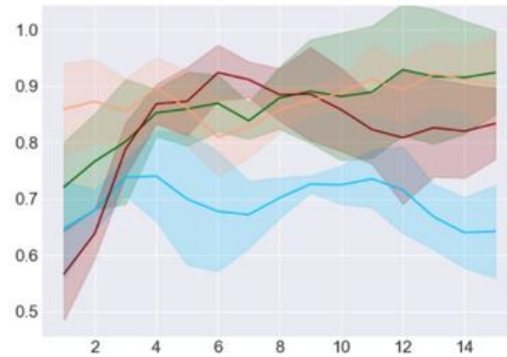
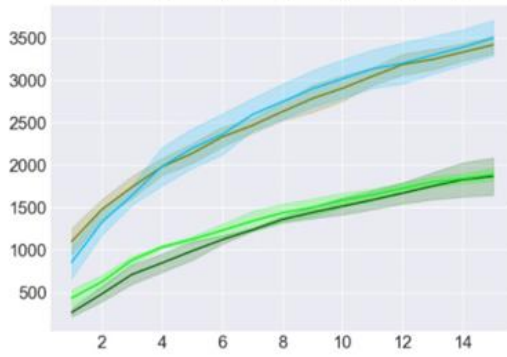
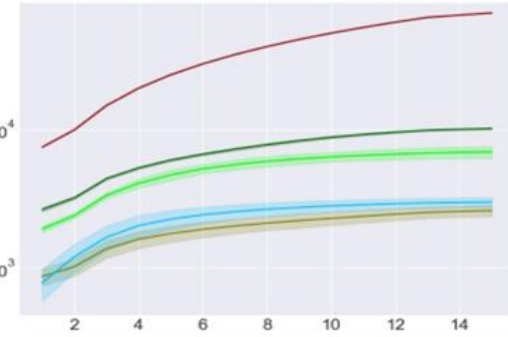
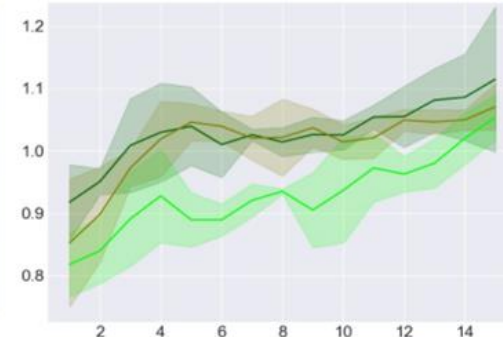
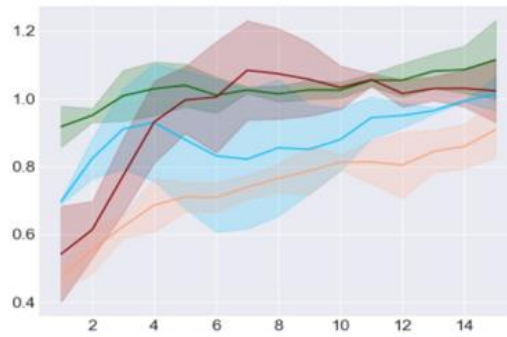
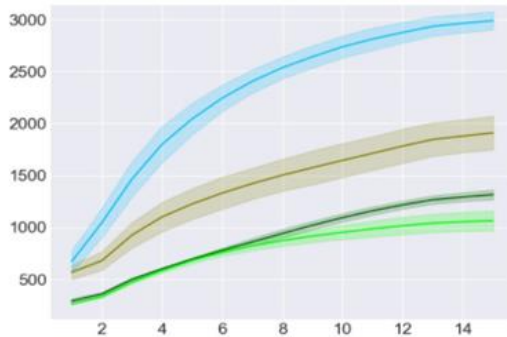
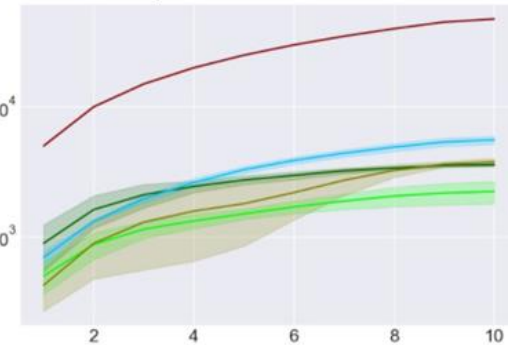
Learning Curves



Learning Curves: Ablations



Number of Supervisor Actions



Behavior Cloning DAGger SafeDagger Ours Ours (- Noise) Ours (- Switch to Auto)

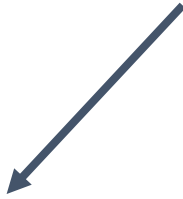
$$L = \frac{\text{(Time to perform one context switch)}}{\text{(Time to perform one action)}}$$

$$L = \frac{\text{(Time to perform one context switch)}}{\text{(Time to perform one action)}}$$

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

$$L = \frac{\text{(Time to perform one context switch)}}{\text{(Time to perform one action)}}$$

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



C = 20 switches

I = 10 actions

B = 20L + 100

$$L = \frac{\text{(Time to perform one context switch)}}{\text{(Time to perform one action)}}$$

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



$C = 20$ switches
 $I = 2$ actions
 $B = 20L + 40$

$C = 4$ switches
 $D = 20$ actions
 $B = 4L + 80$

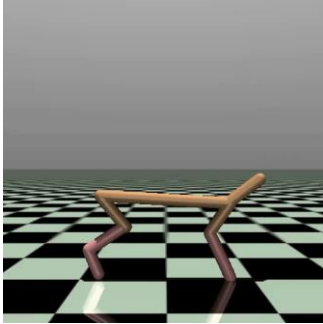
$$L = \frac{\text{(Time to perform one context switch)}}{\text{(Time to perform one action)}}$$

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

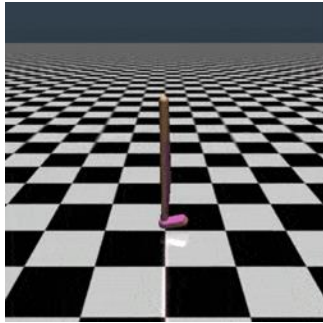
Define “cut-off latency” $L^* \geq 0$, as the minimum value such that

$$B(\text{SafeDAgger}) > B(\text{LazyDAgger}) \text{ for all } L \geq L^*$$

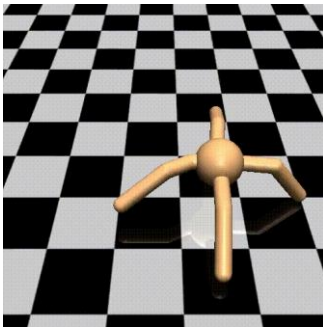
Simulation Experiments



$$L^* = 0.0$$



$$L^* = 4.3$$

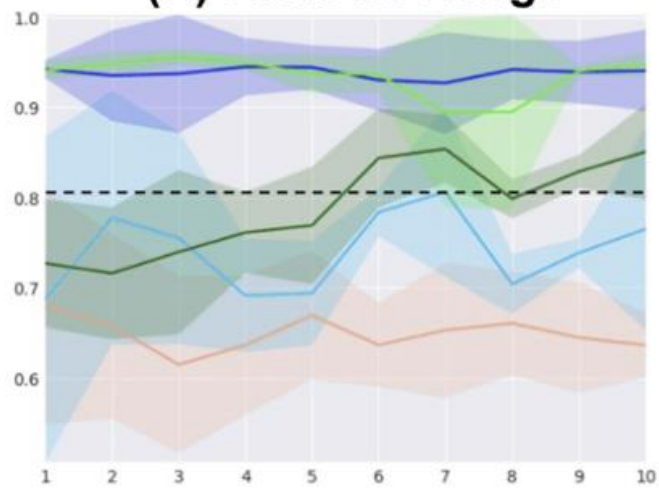


$$L^* = 7.6$$

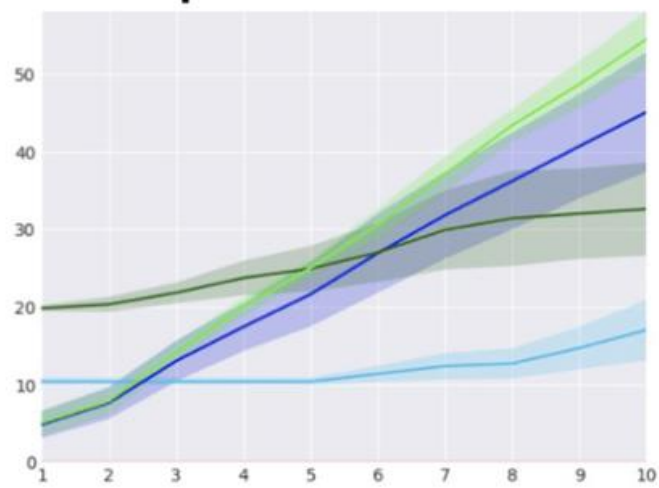
Gym-Cloth



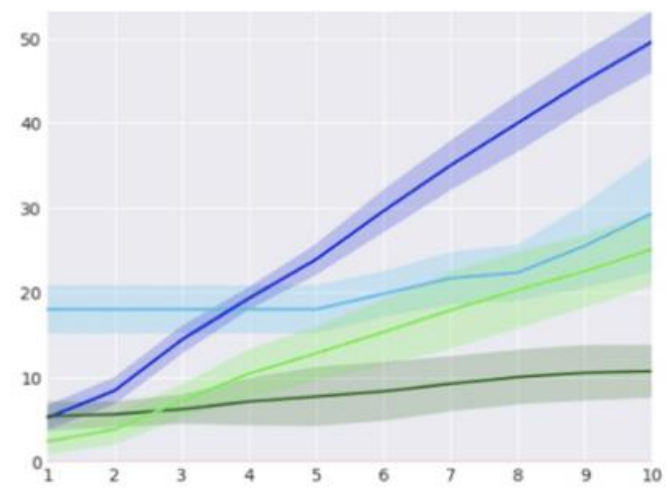
(A) Final Coverage



(B) Number of Supervisor Actions



(C) Number of Context Switches

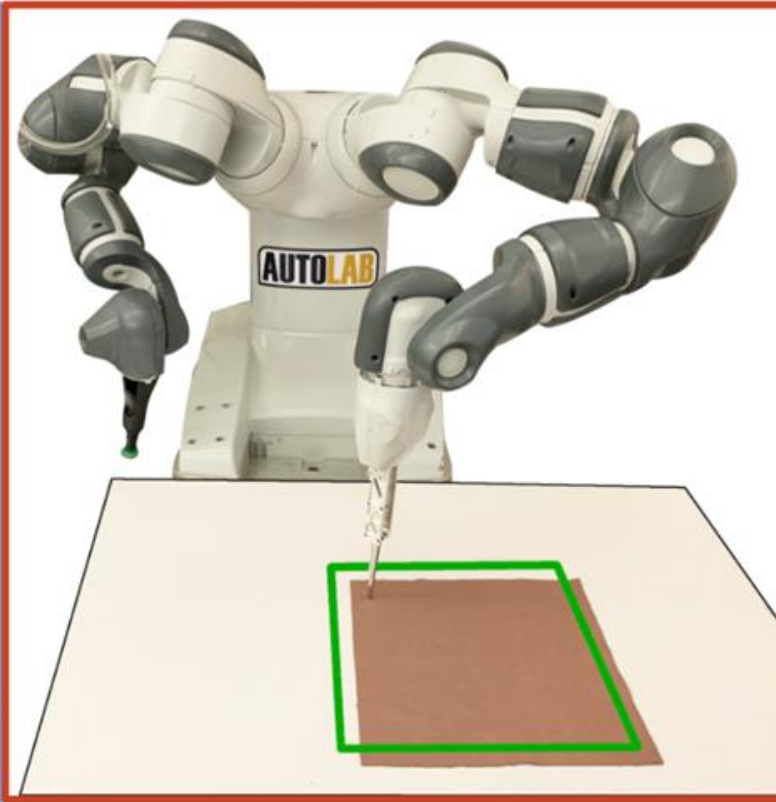


Behavior Cloning SafeDagger SafeDagger Execution LazyDagger LazyDagger Execution

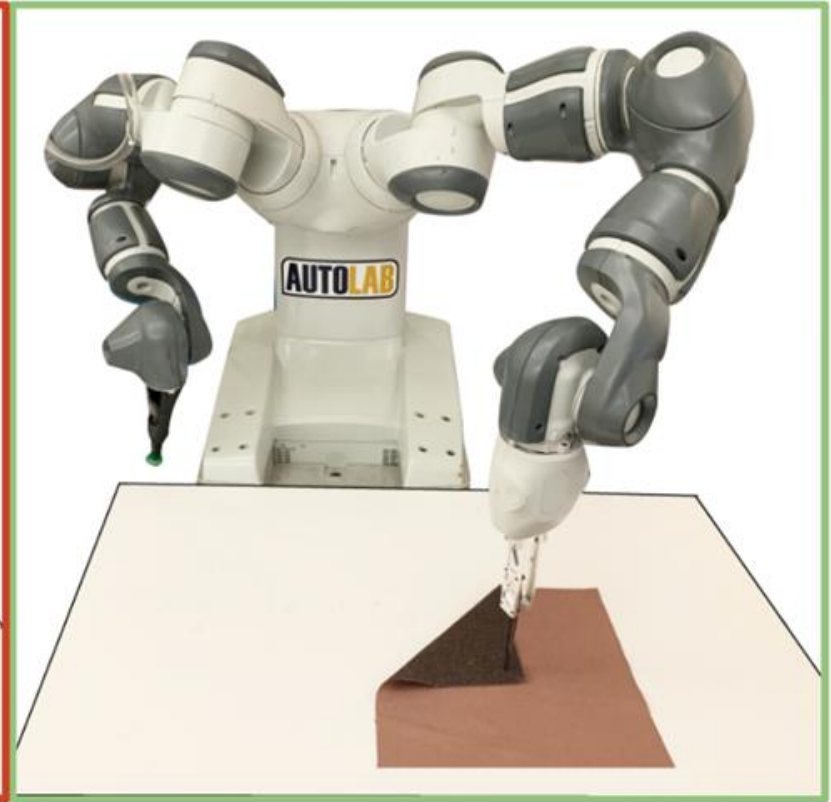
(1) SMOOTH



(2) ALIGN

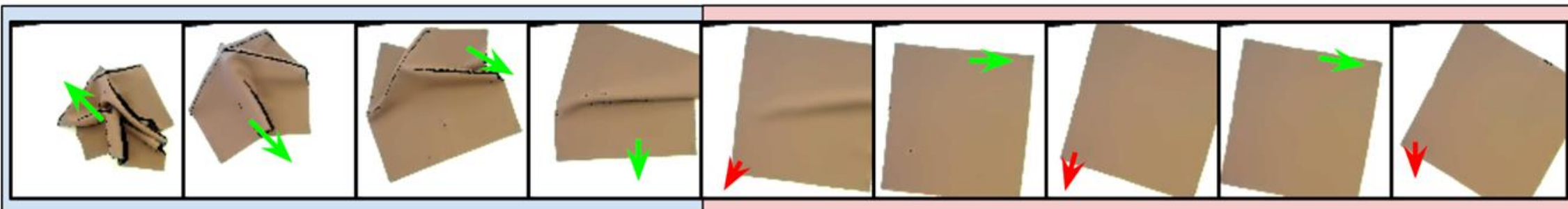


(3) FOLD

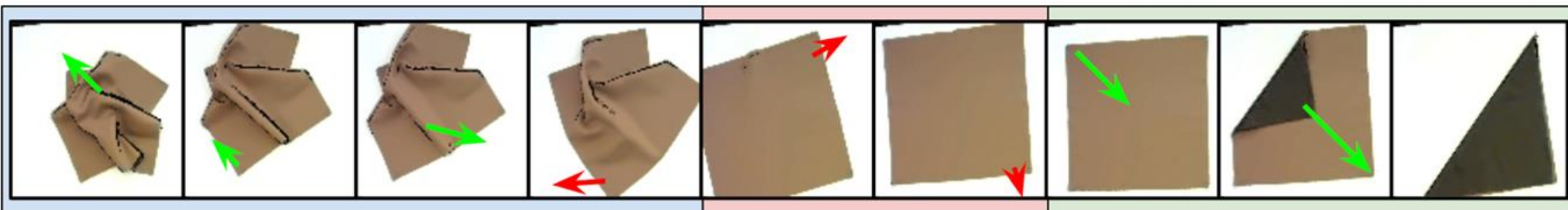




**Safe
Dagger**



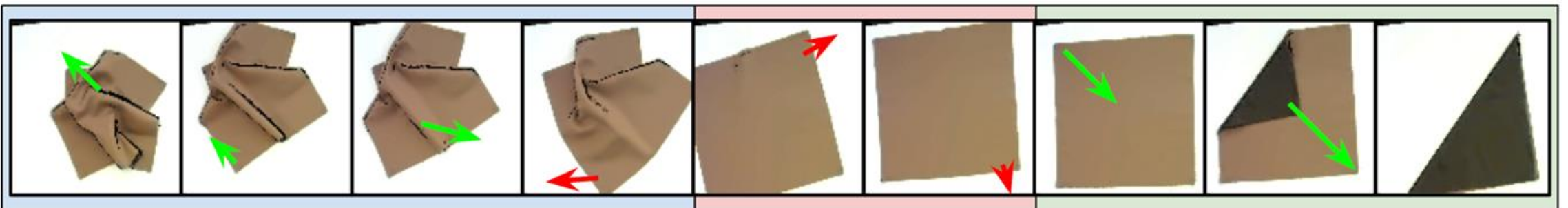
**Lazy
Dagger**



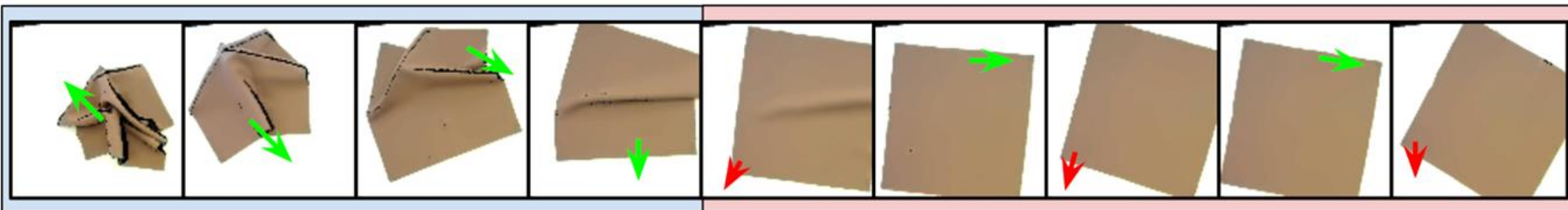
**Safe
Dagger**



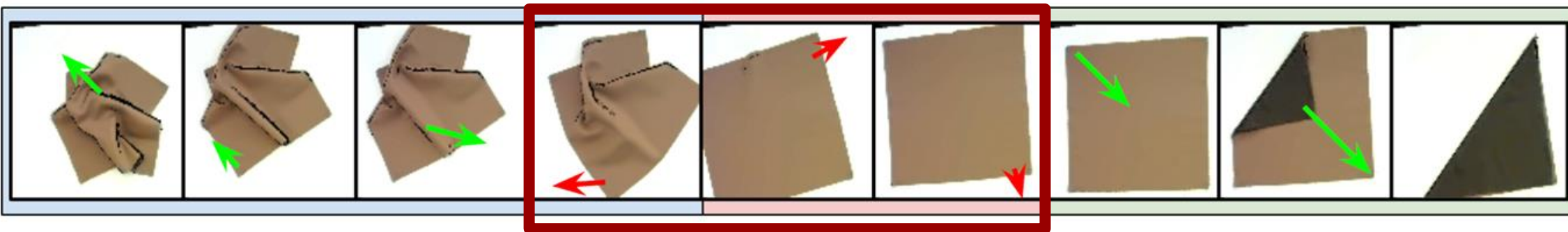
**Lazy
Dagger**



**Safe
Dagger**

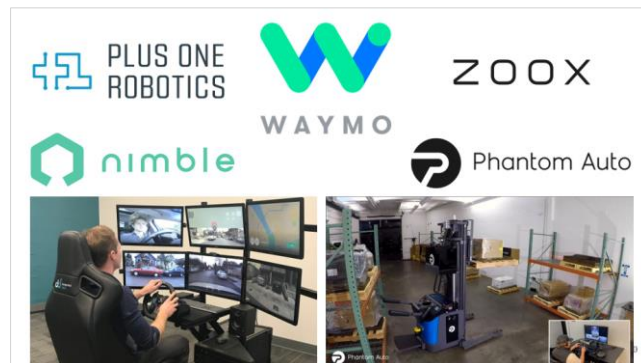
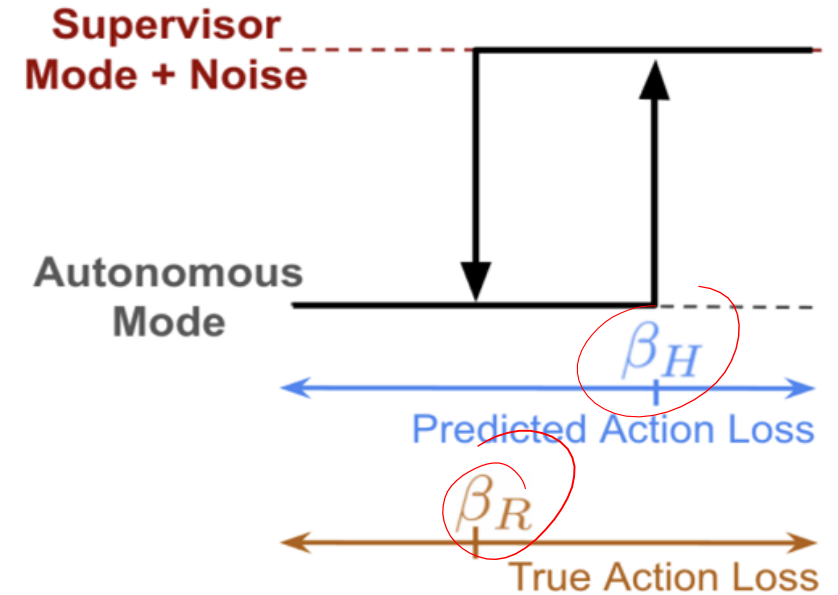


**Lazy
Dagger**



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

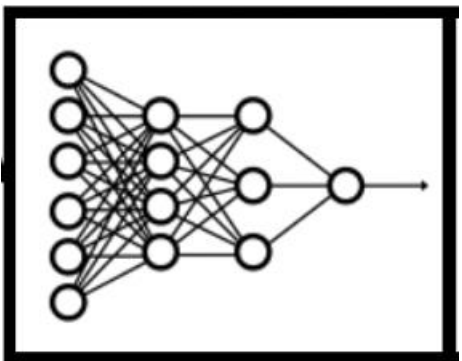
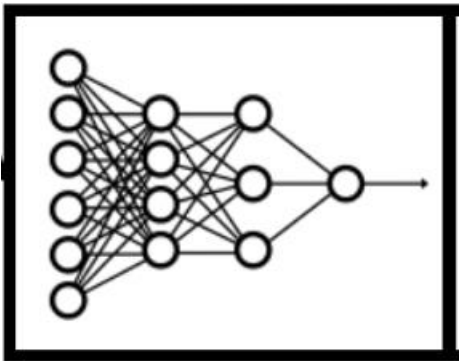
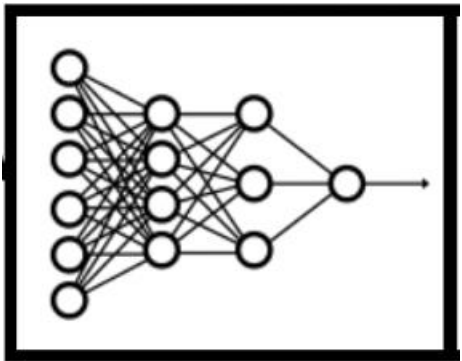
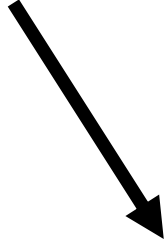
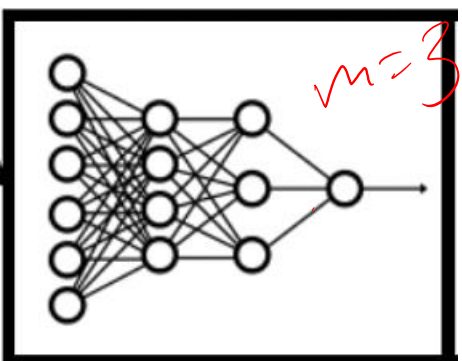
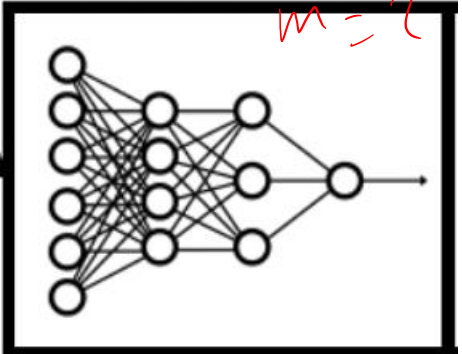
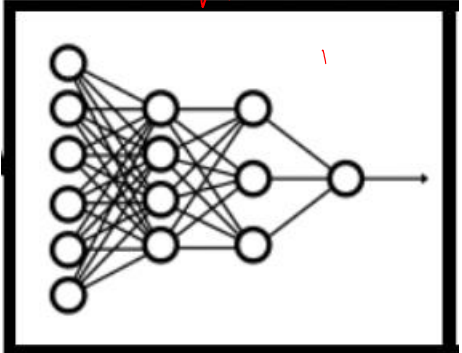
in distribution

Ensembling

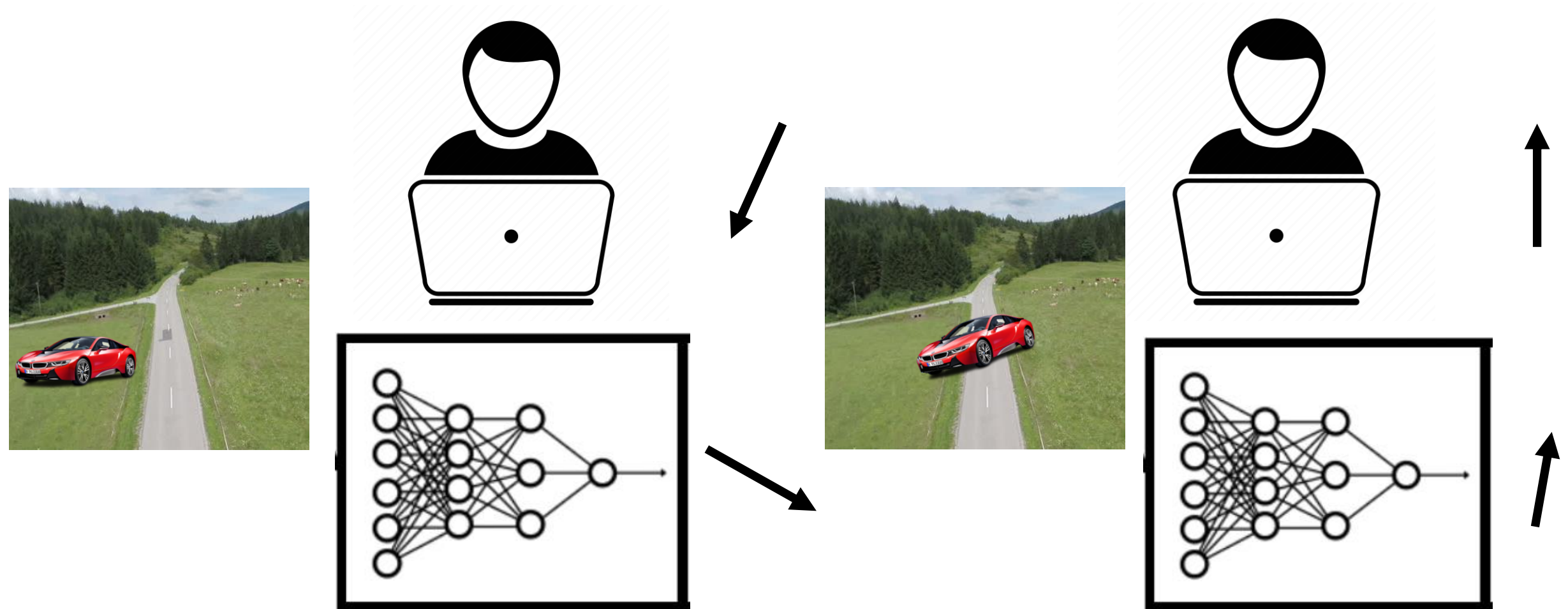
m=1

m=2

m=3



Novelty Estimation: Supervisor Mode



Assume
Goal identifier

$R = 1$ if at goal, 0 otherwise

Risk Estimation

$$Q_G^{\pi_r}(s_t, a_t) = \mathbb{E}_{\pi_r} \left[\sum_{t'=t}^{\infty} \gamma^{t'-t} \mathbb{1}_G(s_{t'}) \mid s_t, a_t \right]$$

Expected probs of
reaching goal / task
success

Risk Estimation

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

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

Risk Estimation

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

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

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

Putting it all together...

**AUTONOMOUS
MODE**

$$\begin{aligned} & \text{Novelty}(s_t) > \delta_h \\ & \text{OR} \\ & \text{Risk}^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h \end{aligned}$$



Switch to
**SUPERVISOR
MODE**

Putting it all together...

**AUTONOMOUS
MODE**

$$\begin{aligned} & \text{Novelty}(s_t) > \delta_h \\ & \text{OR} \\ & \text{Risk}^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h \end{aligned}$$



Switch to
**SUPERVISOR
MODE**

**SUPERVISOR
MODE**

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



Switch to
**AUTONOMOUS
MODE**

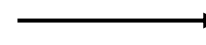
not novel

not risky

Putting it all together...

**AUTONOMOUS
MODE**

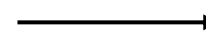
$$\begin{aligned} & \text{Novelty}(s_t) > \delta_h \\ & \text{OR} \\ & \text{Risk}^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h \end{aligned}$$



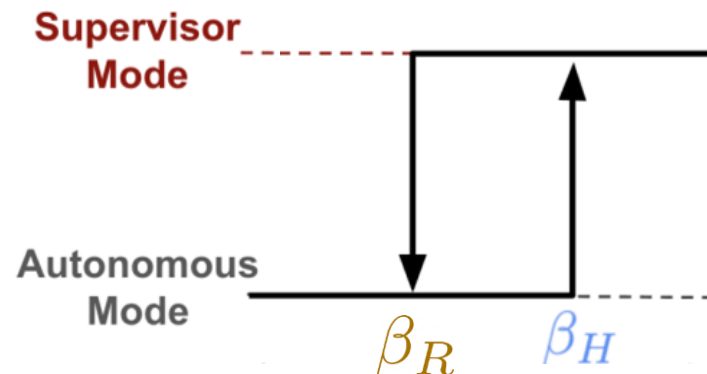
Switch to
**SUPERVISOR
MODE**

**SUPERVISOR
MODE**

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



Switch to
**AUTONOMOUS
MODE**



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

Putting it all together...

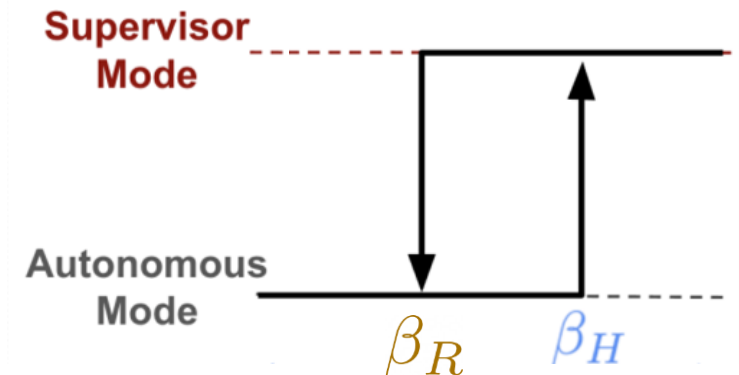
**AUTONOMOUS
MODE**

Novelty(s_t) > δ_h — Set to $1-\alpha$ quantiles of empirical data
 OR
 Risk $^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h$ — \longrightarrow Switch to SUPERVISOR MODE

**SUPERVISOR
MODE**

$\|\pi_r(s_t) - \pi_h(s_t)\|_2^2 < \delta_r$
 AND
 Risk $^{\pi_r}(s_t, \pi_r(s_t)) < \beta_r$ — \longrightarrow Switch to AUTONOMOUS MODE

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



Putting it all together...

**AUTONOMOUS
MODE**

Novelty(s_t) > δ_h OR
Risk $^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h$ \longrightarrow Switch to SUPERVISOR MODE

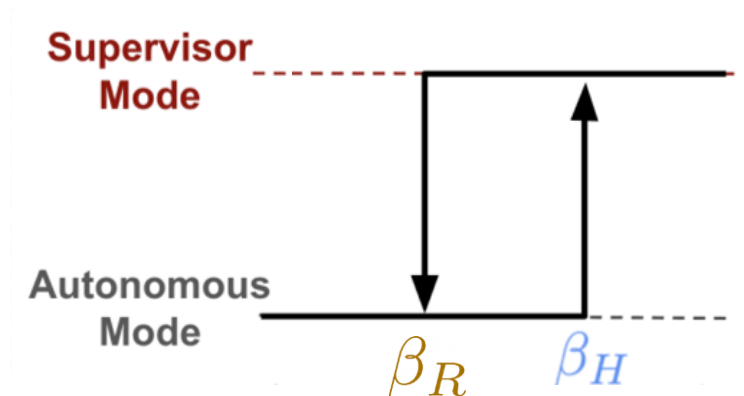
Set to 1- α quantiles of empirical data

**SUPERVISOR
MODE**

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

Set to medians of empirical data

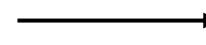
$$\alpha = \frac{\text{desired \# interventions}}{\text{\# robot actions}}$$



Putting it all together...

**AUTONOMOUS
MODE**

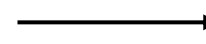
$$\begin{aligned} & \text{Novelty}(s_t) > \delta_h \\ & \text{OR} \\ & \text{Risk}^{\pi_r}(s_t, \pi_r(s_t)) > \beta_h \end{aligned}$$



Switch to
**SUPERVISOR
MODE**

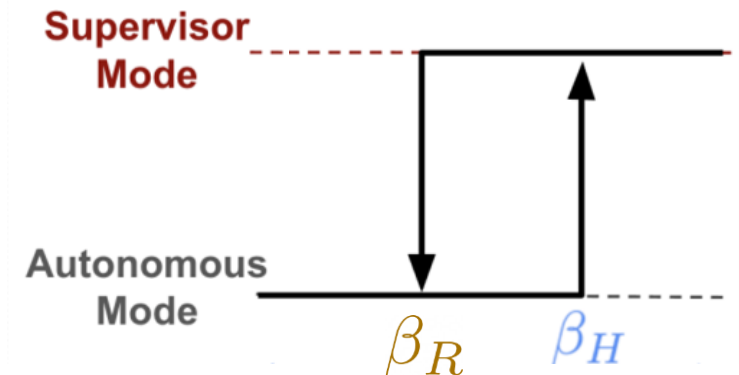
**SUPERVISOR
MODE**

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



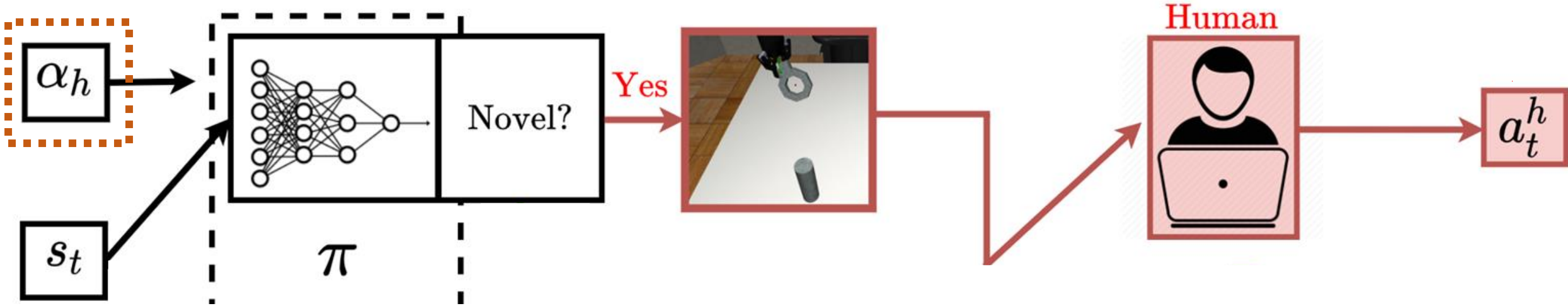
Switch to
**AUTONOMOUS
MODE**

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

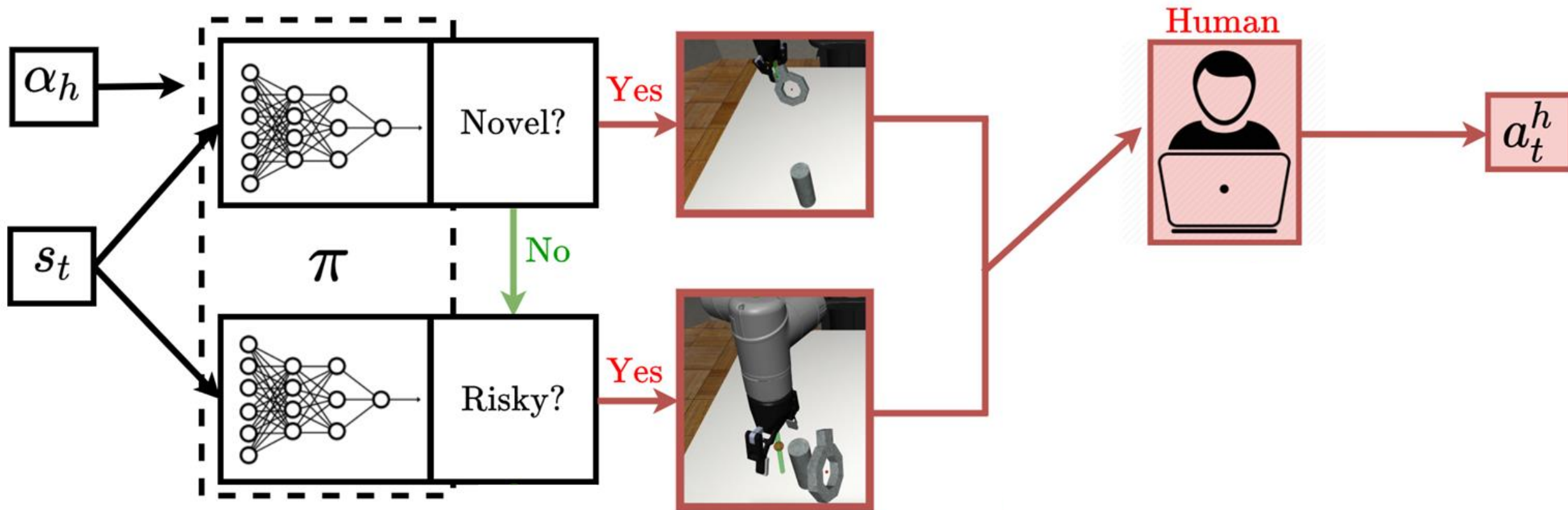


Target percent of time human wants to give interventions.

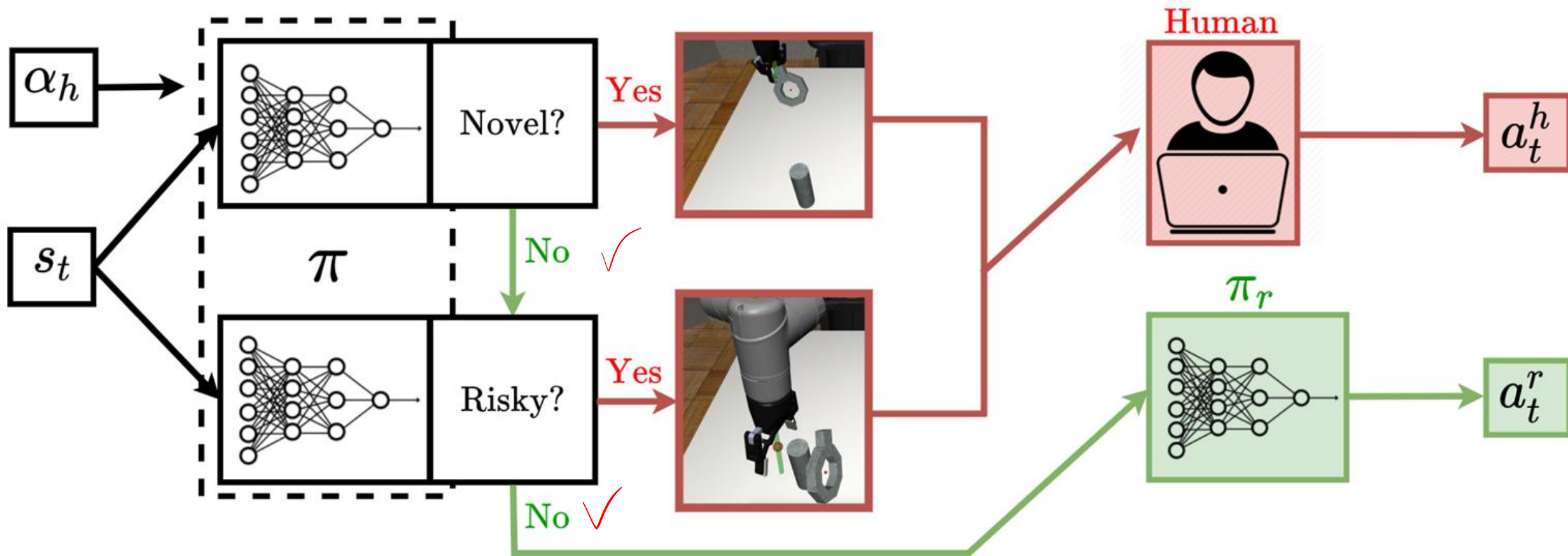
ThriftyDAgger

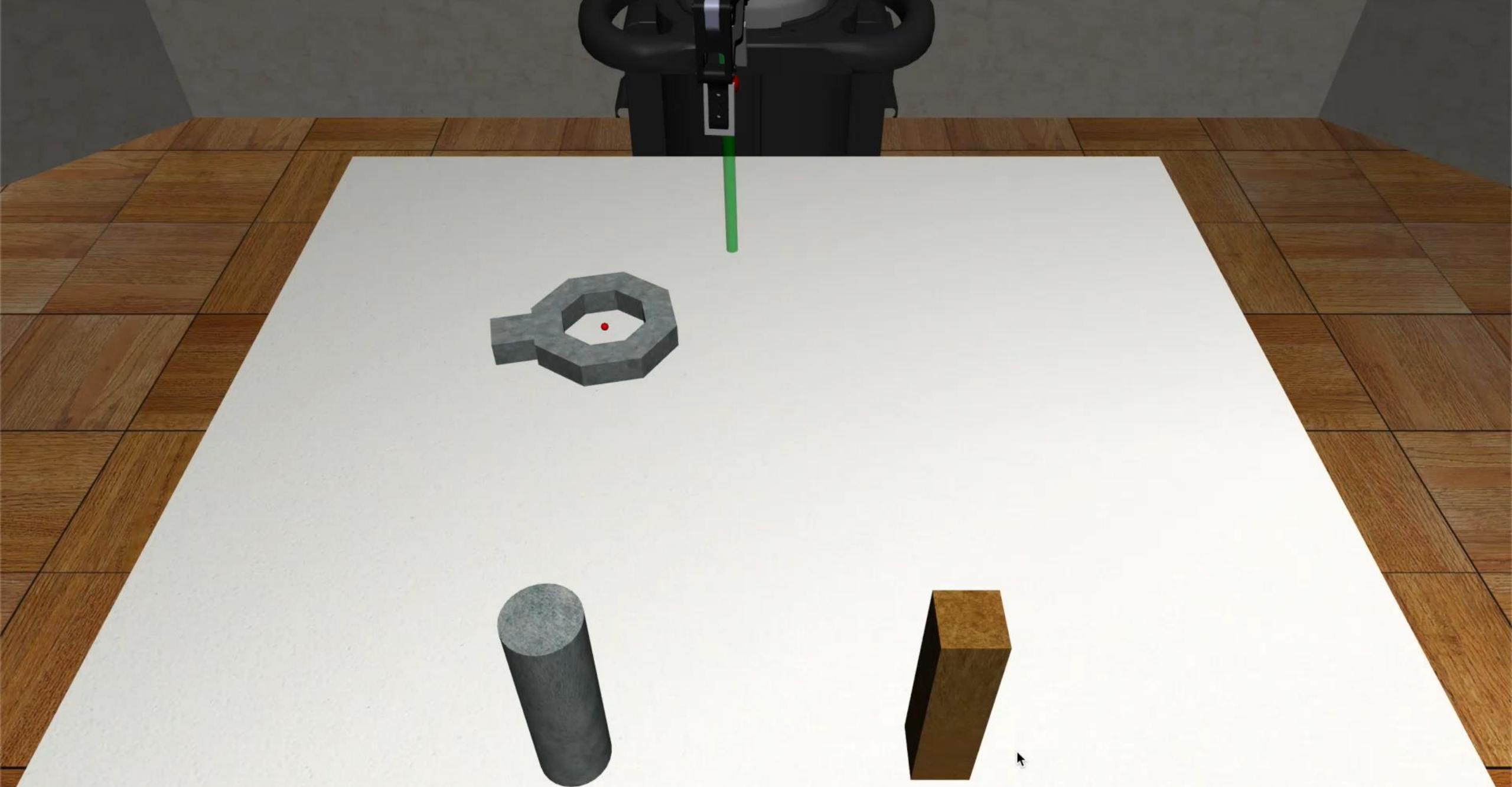


ThriftyDAgger

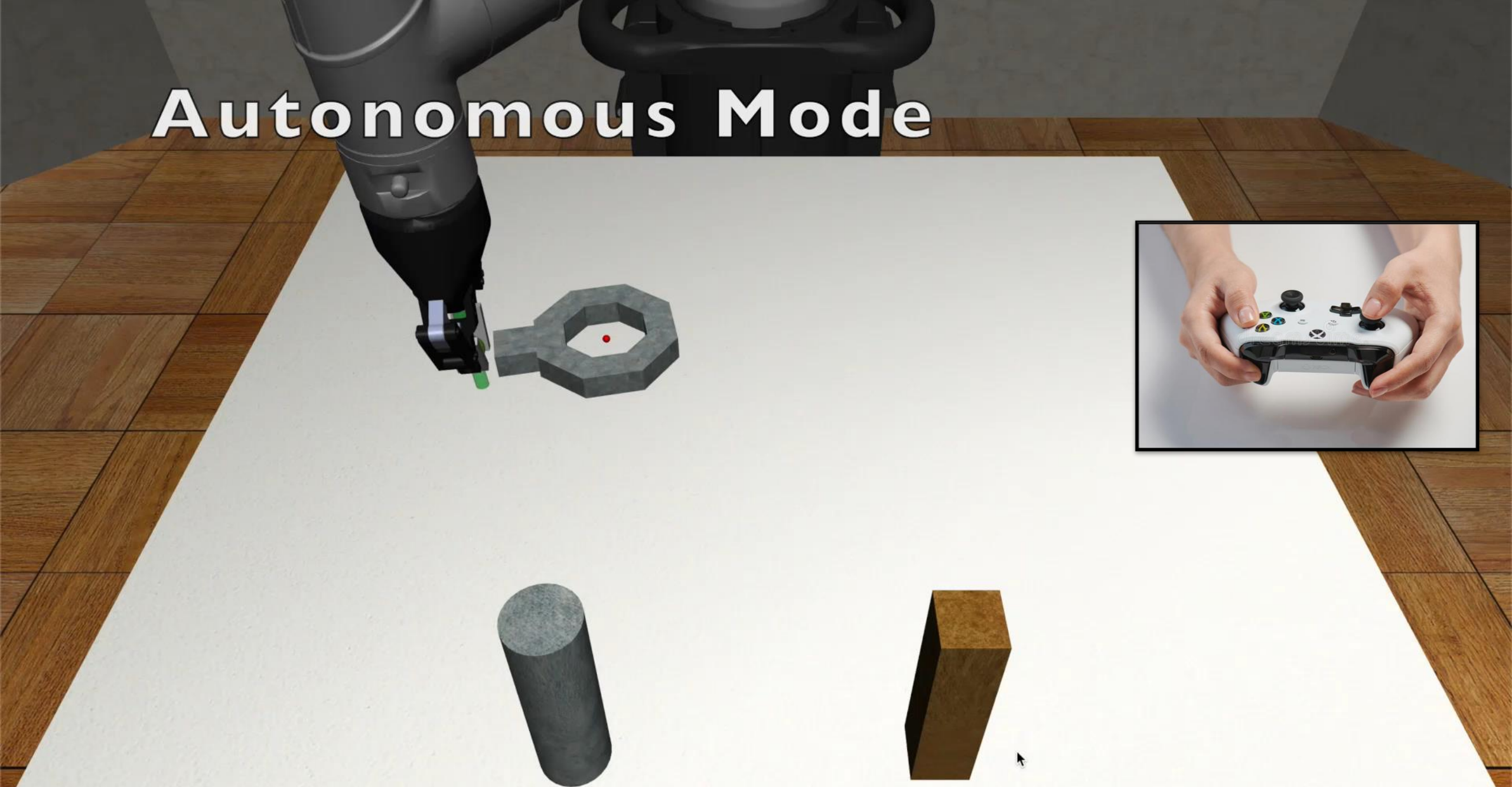


ThriftyDAgger



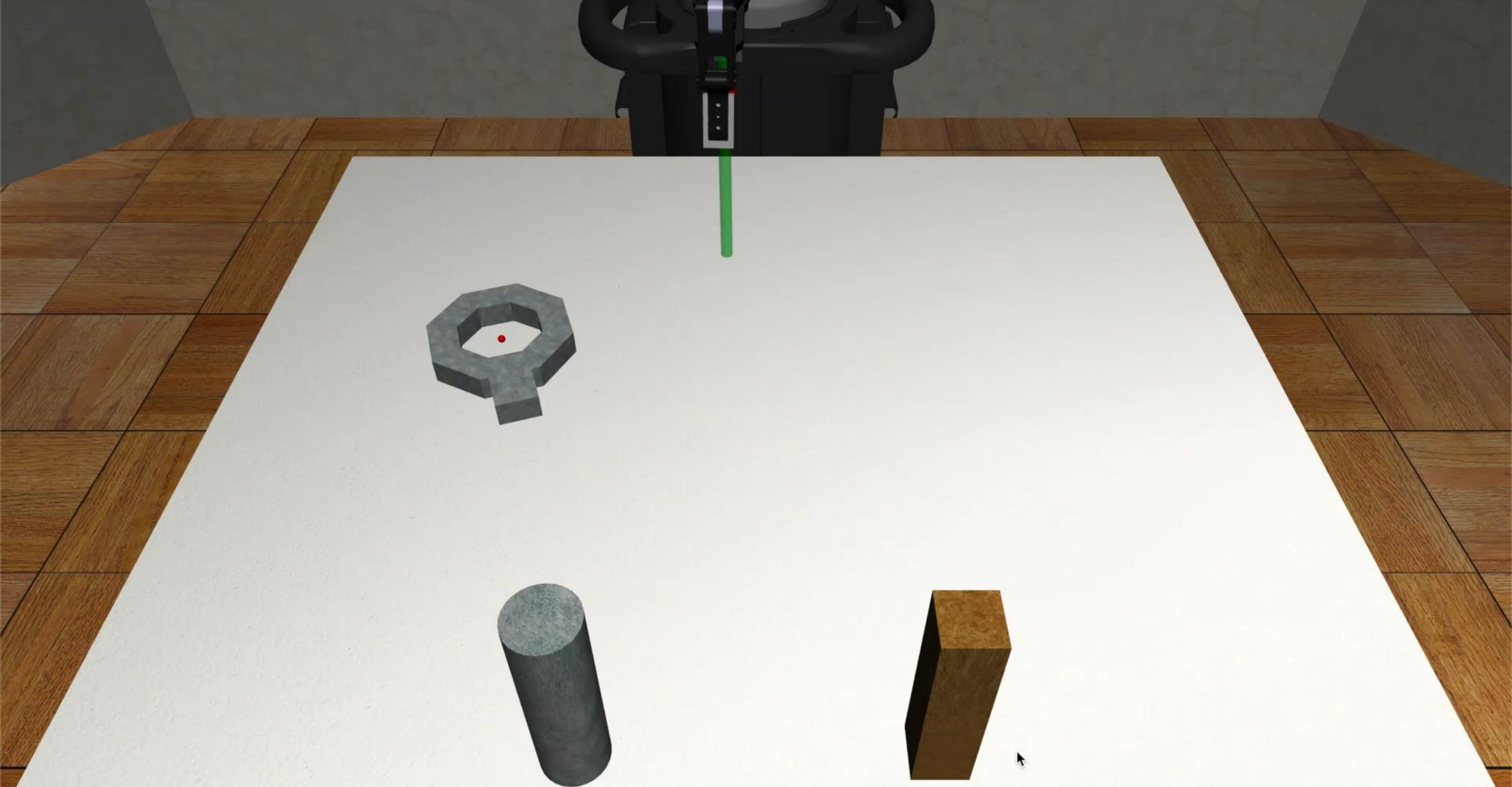


Autonomous Mode

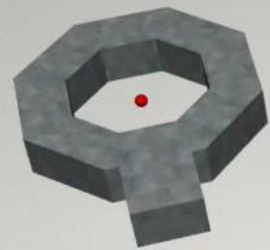


Supervisor Mode (Novel)





Supervisor Mode (Risk)



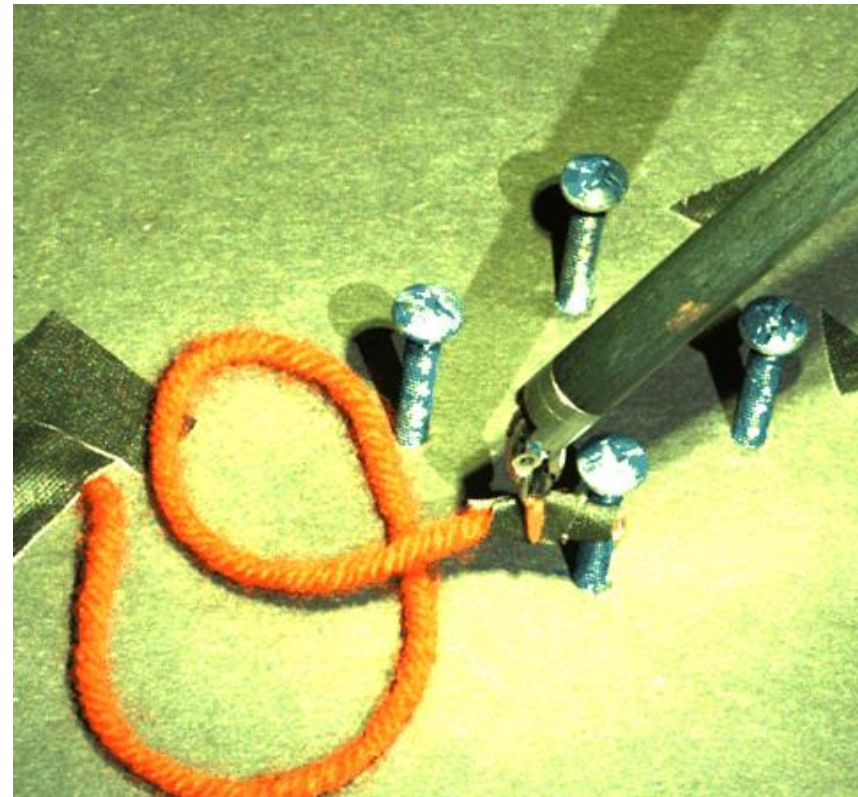
Supervisor Mode (Risk)



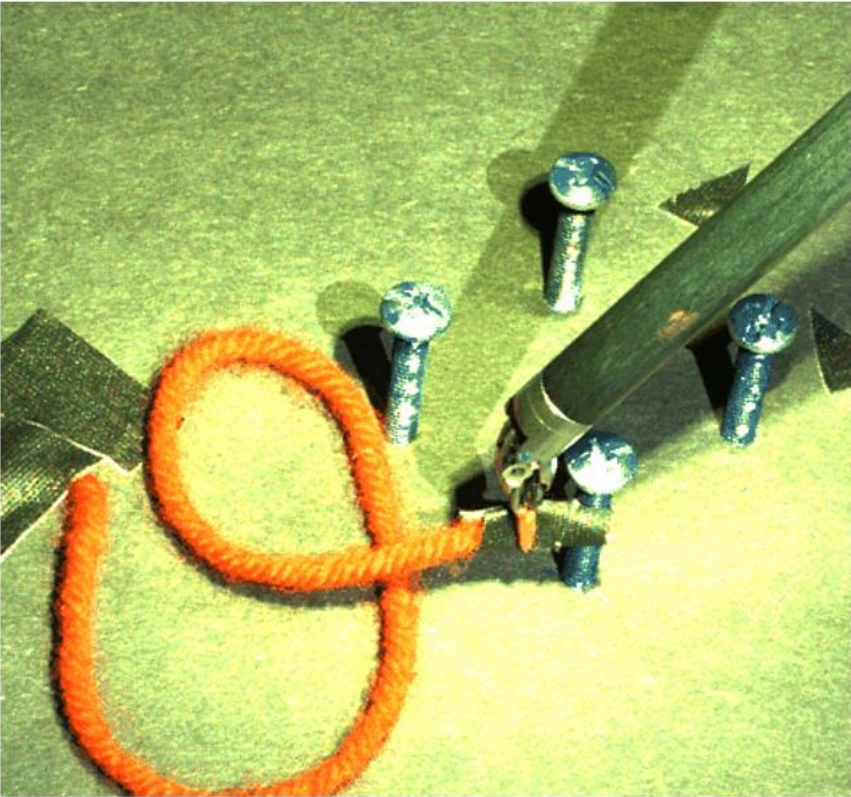
Human Demonstration



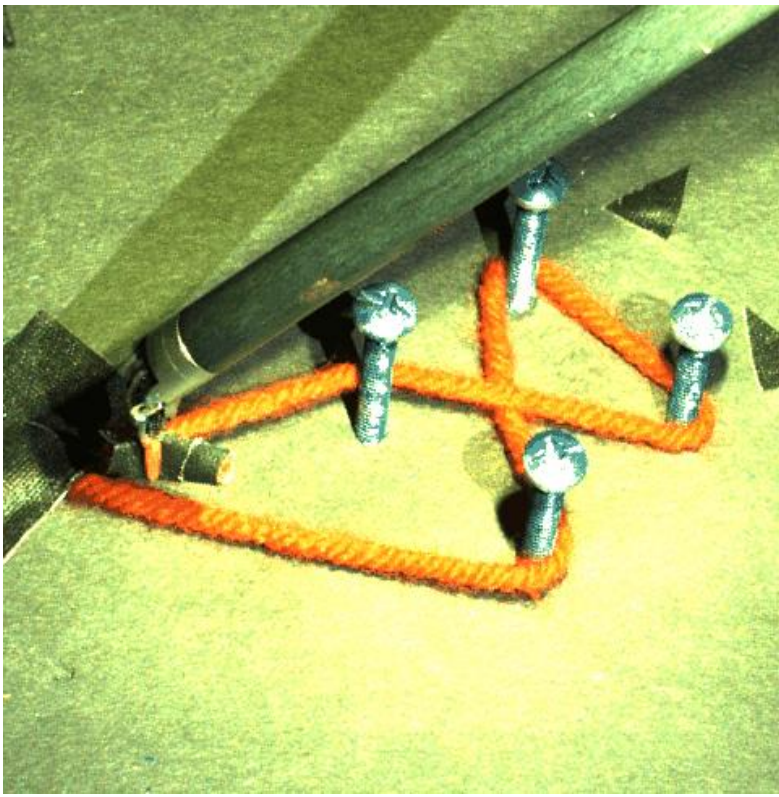
Behavior Cloning



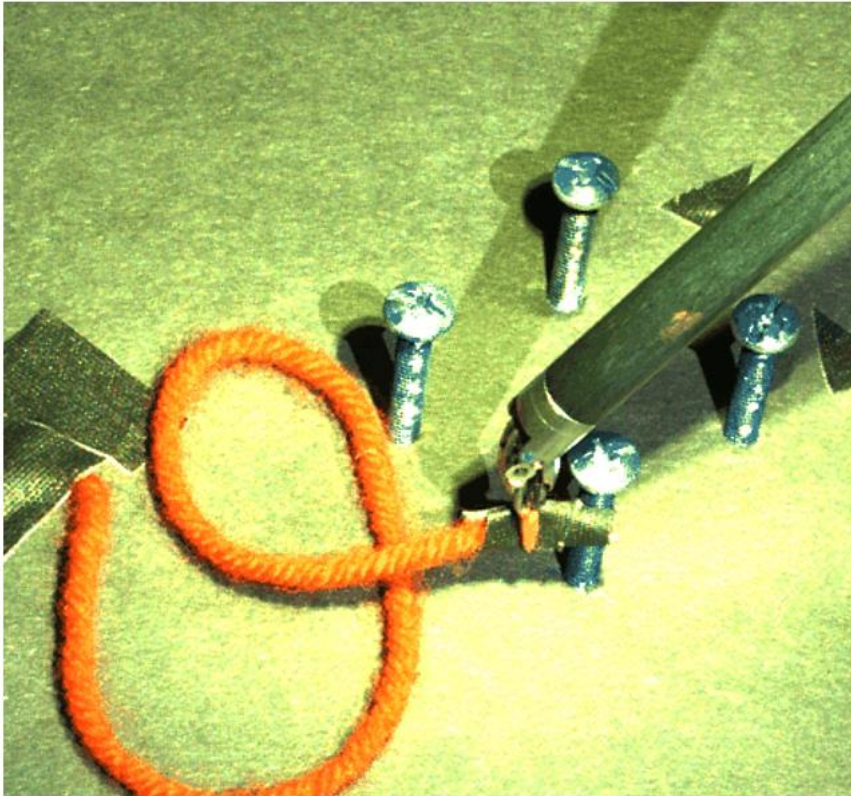
Behavior Cloning



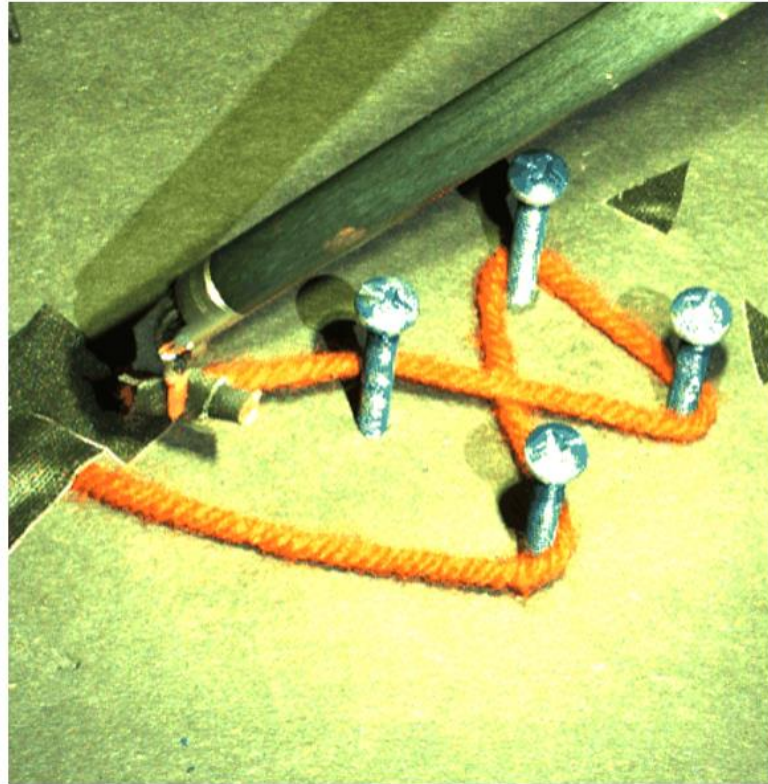
ThriftyDagger (autonomous)



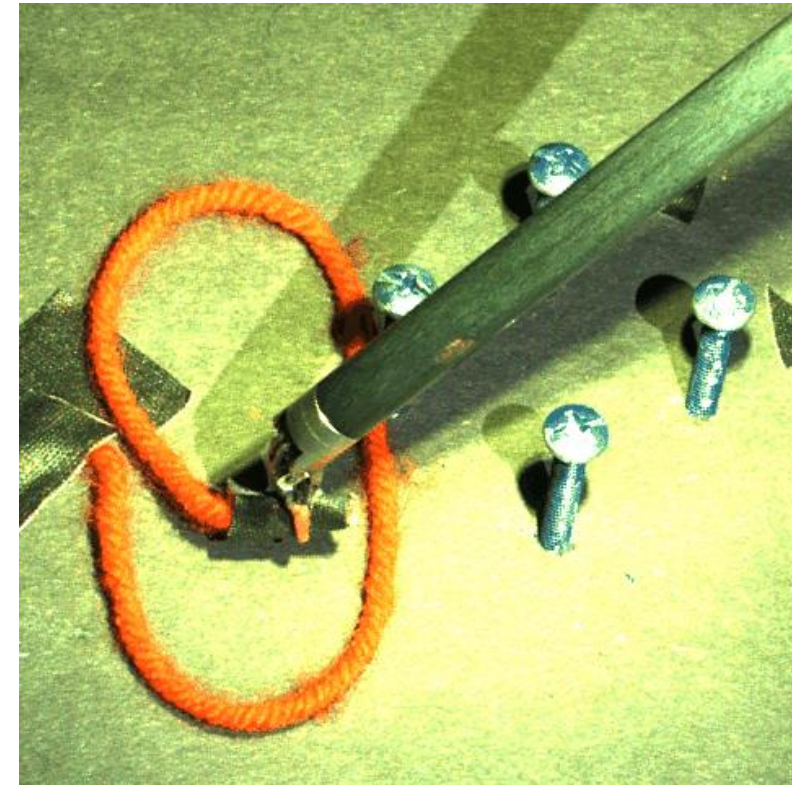
Behavior Cloning



ThriftyDAgger (autonomous)



ThriftyDAgger (+human)



User Study

N=10 subjects each control 3 robots in simulation.

humans

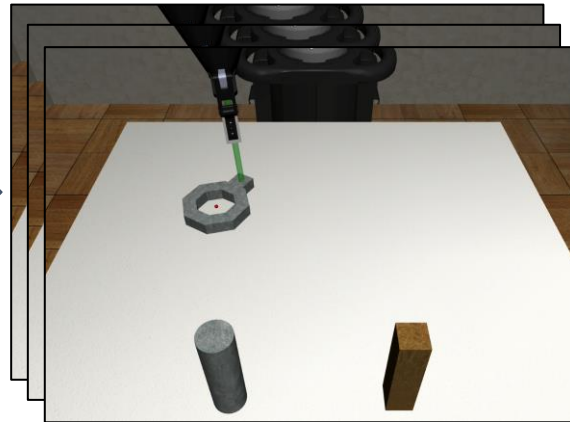
Robot-Gated

Memory: Non-Match

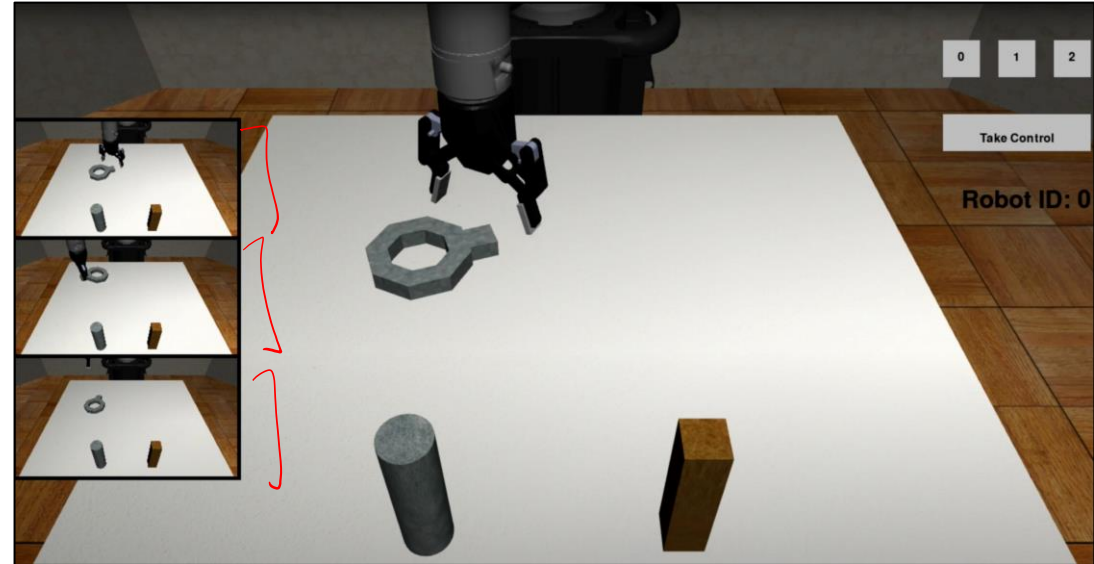
H	H	H	H	H
H	🍌	❤️	H	H
H	H	H	H	H
H	H	H	H	H

Memory: Match

H	H	H	H	H
H	H	H	🕒	H
🕒	H	H	H	H
H	H	H	H	H

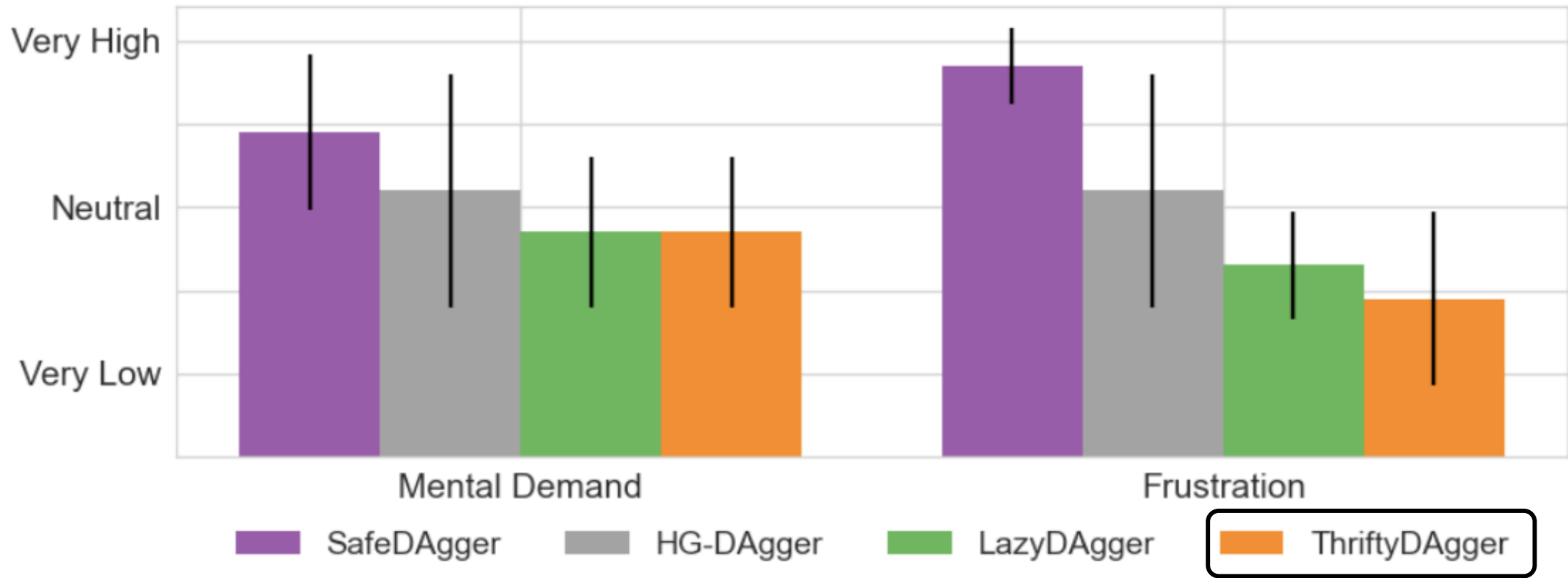


Human-Gated



ThriftyDagger Qualitative Results

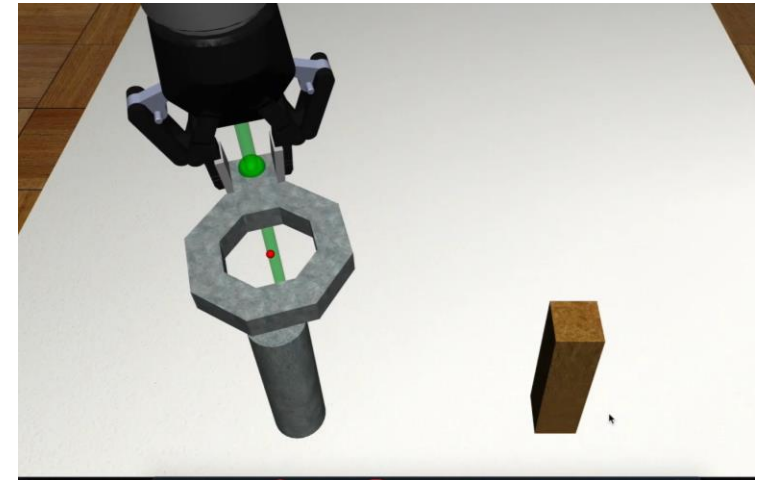
Survey Responses



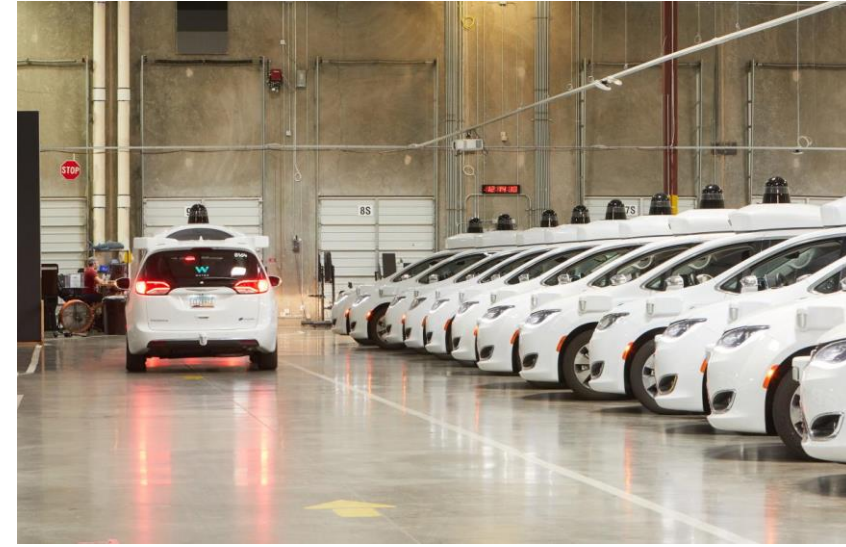
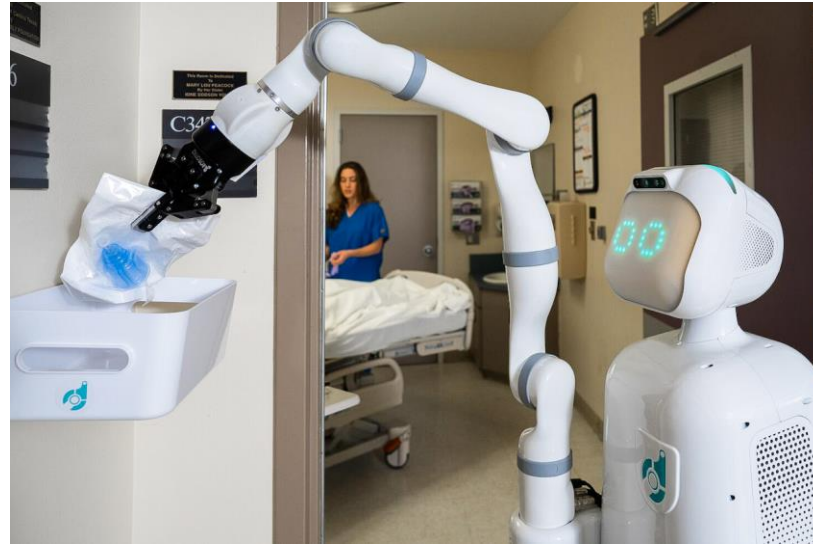
User Study Quantitative Results

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.



Next time: Inverse Reinforcement Learning!