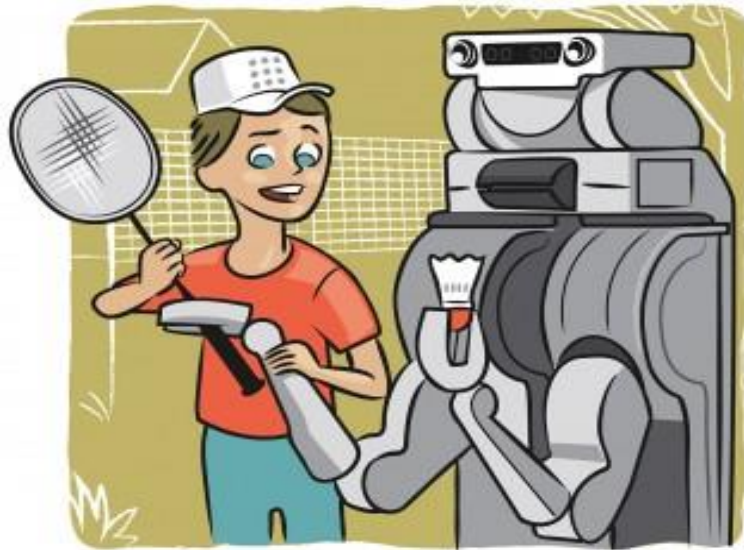
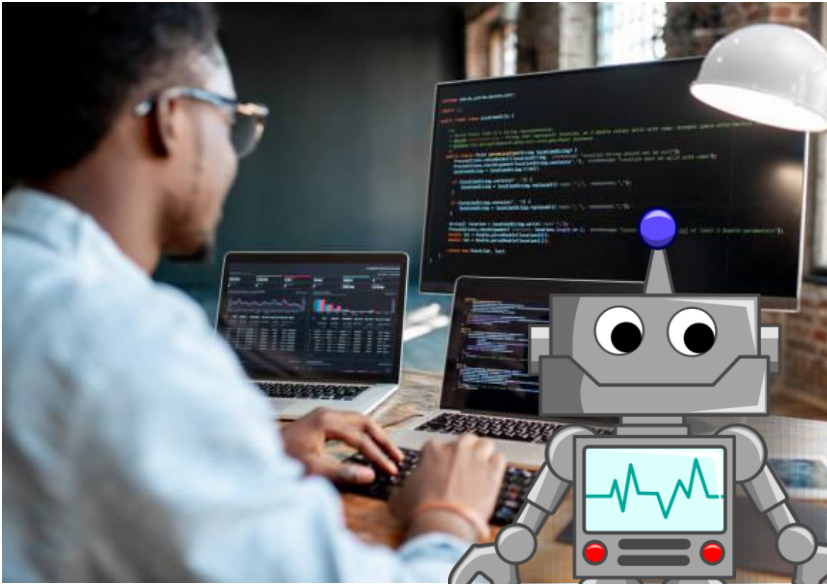


# Behavioral Cloning and Interactive Imitation Learning



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

ML  
Data mining  
This class



$$D = \{ (x, y), \dots \}$$
$$f(x) = y$$



# 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.

Discrete # of labels

Continuous output

$$x \in \mathbb{R}^2$$
$$y = w_0 + w_1 x_1 + w_2 x_2$$

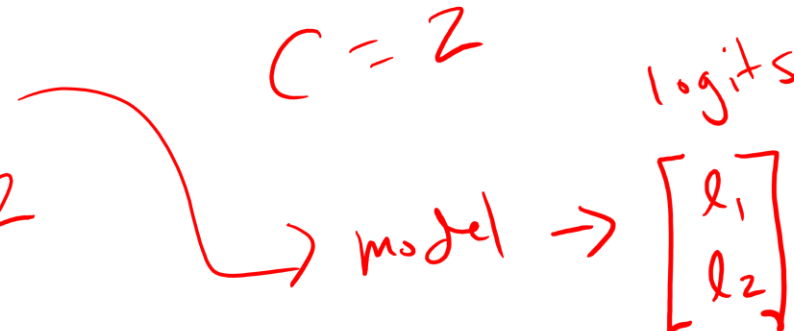
# Classification

$C = \# \text{ classes}$

$y = \text{cat} = 1$   
 $x = \boxed{\text{cat}}$

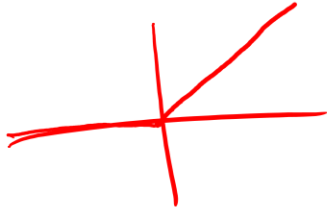
$C = 2$

$y = \text{dog} = 2$   
 $x = \boxed{\text{dog}}$



minimize this  
Cross-Entropy Loss =  $-\sum_{i=1}^C y_i \log(\hat{y}_i)$  True label  $\downarrow$

$$\hat{y}_1 = \text{Pr}(y = \text{cat}) = \frac{\exp(l_1)}{\exp(l_1) + \exp(l_2)}$$
$$= -y_1 \log(\hat{y}_1) + y_2 \log(\hat{y}_2)$$
$$= -1 \cdot \log(\hat{y}_1) + 0$$



# PyTorch Example

```
import torch.nn as nn
import torch.optim as optim
```

```
class ClassificationNetwork(nn.Module):
    def __init__(self, input_dim, num_classes):
        super(ClassificationNetwork, self).__init__()
        self.fc1 = nn.Linear(input_dim, num_classes) = 3
```

```
    def forward(self, x):
        return self.fc2(self.relu(self.fc1(x)))
```

```
model = ClassificationNetwork(input_dim, num_classes)
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=0.001)
```

```
for epoch in range(num_epochs):
    for inputs, labels in dataloader:
        optimizer.zero_grad()
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
```

3# is describing house  
West Coast, Middle, East Coast  
Input      output



$$w_i \leftarrow w_i - \frac{\partial L}{\partial w_i} \cdot \alpha$$

```
self.fc2 = nn.Linear(x, y)
self.relu = nn.ReLU()
```

# Regression



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

pred  
true

# PyTorch Example

```
import torch.nn as nn
import torch.optim as optim
```


```
class ClassificationNetworkkey(nn.Module):
    def __init__(self, input_dim, num_classes):
        super(ClassificationNetwork, self).__init__()
        self.fc = nn.Linear(input_dim, num_classes)
        1

    def forward(self, x):
        return self.fc(x)
```

```
model = ClassificationNetwork(input_dim, num_classes)
criterion = nn.CrossEntropyLossMSE()
optimizer = optim.Adam(model.parameters(), lr=0.001)
```

```
for epoch in range(num_epochs):
    for inputs, labels in dataloader:
        optimizer.zero_grad()
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
```

ReLU

$$\text{output}_1 = (x^T w_1 + b_2)$$
$$\hat{y} = \text{output}_1^T w_2 + b_2$$




# Unsupervised Learning

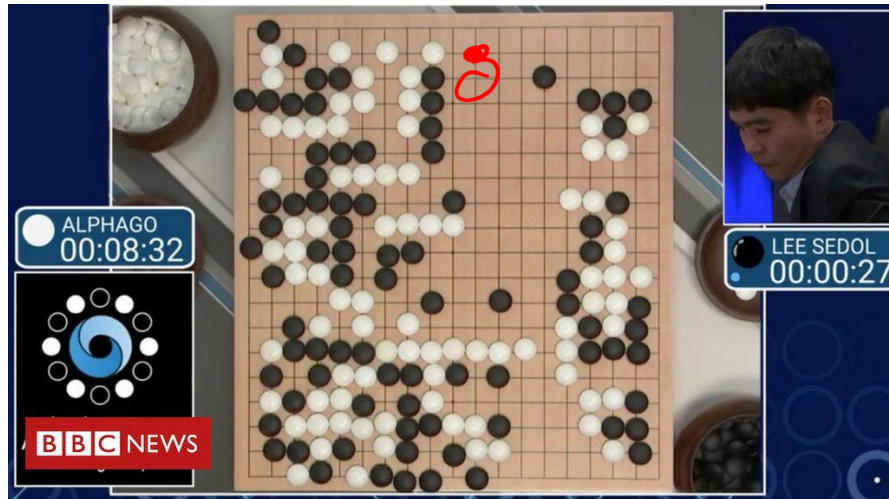
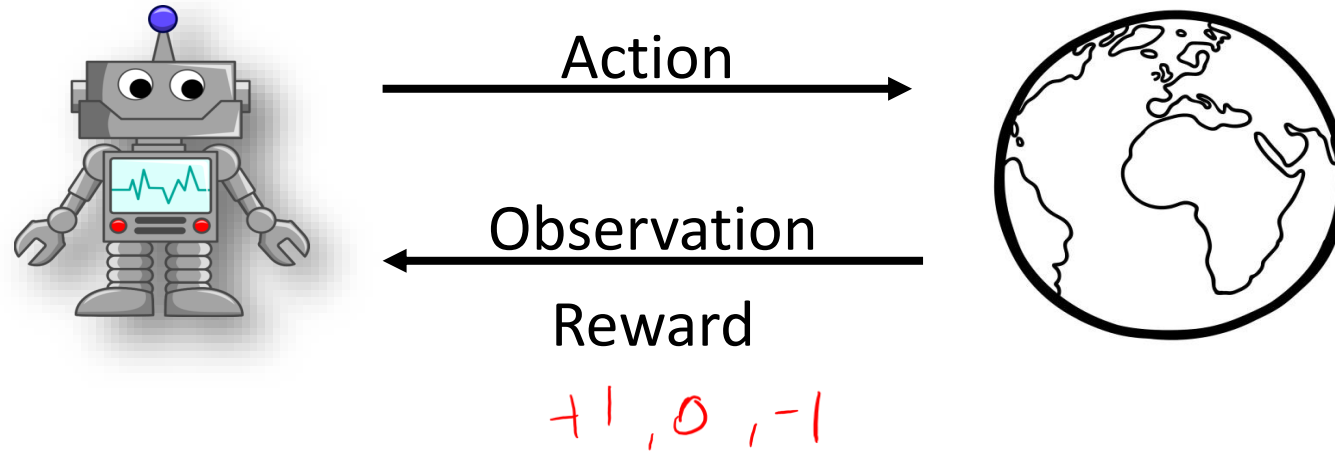
- **Setting/Assumptions:** In unsupervised learning, the model is trained on data without labeled outputs. It seeks to find patterns, structures, or relationships in the data. No “ground truth” labels.
- **Goal:** To explore the data and identify meaningful clusters, associations, or representations.
- **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.

# Reinforcement Learning

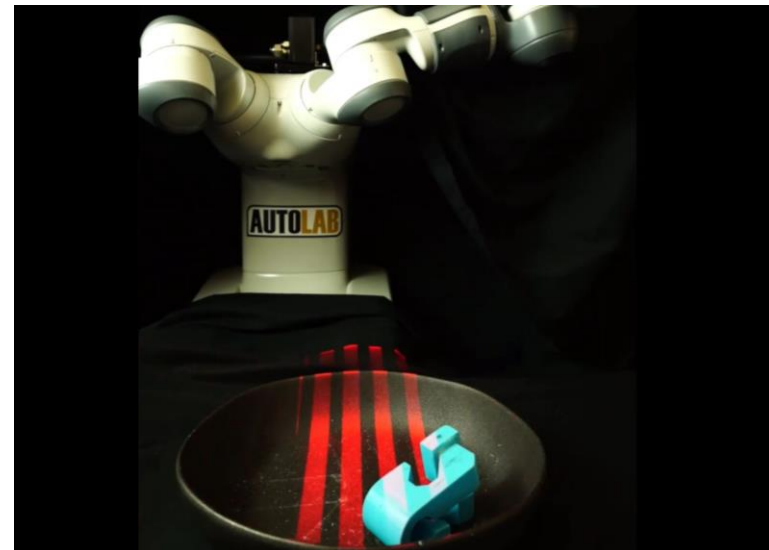
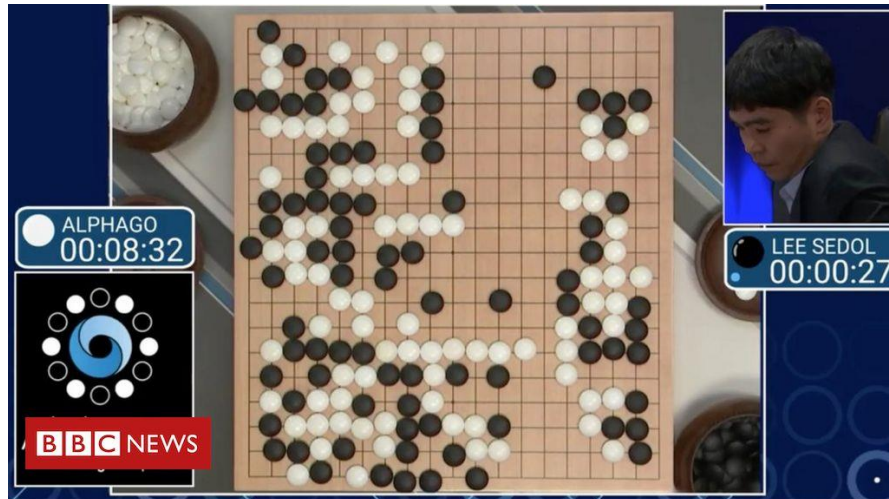
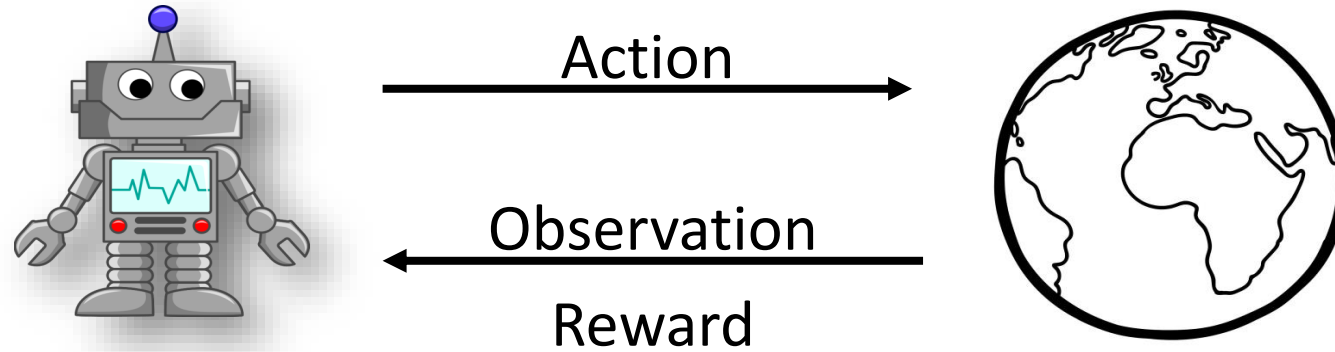
- **Setting/Assumptions:** Reinforcement learning (RL) involves training an agent to make decisions by interacting with an environment. The agent learns through trial and error (receiving rewards and penalties), optimizing its behavior to maximize cumulative rewards.
- **Goal:** To learn a policy that 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).

LLMs

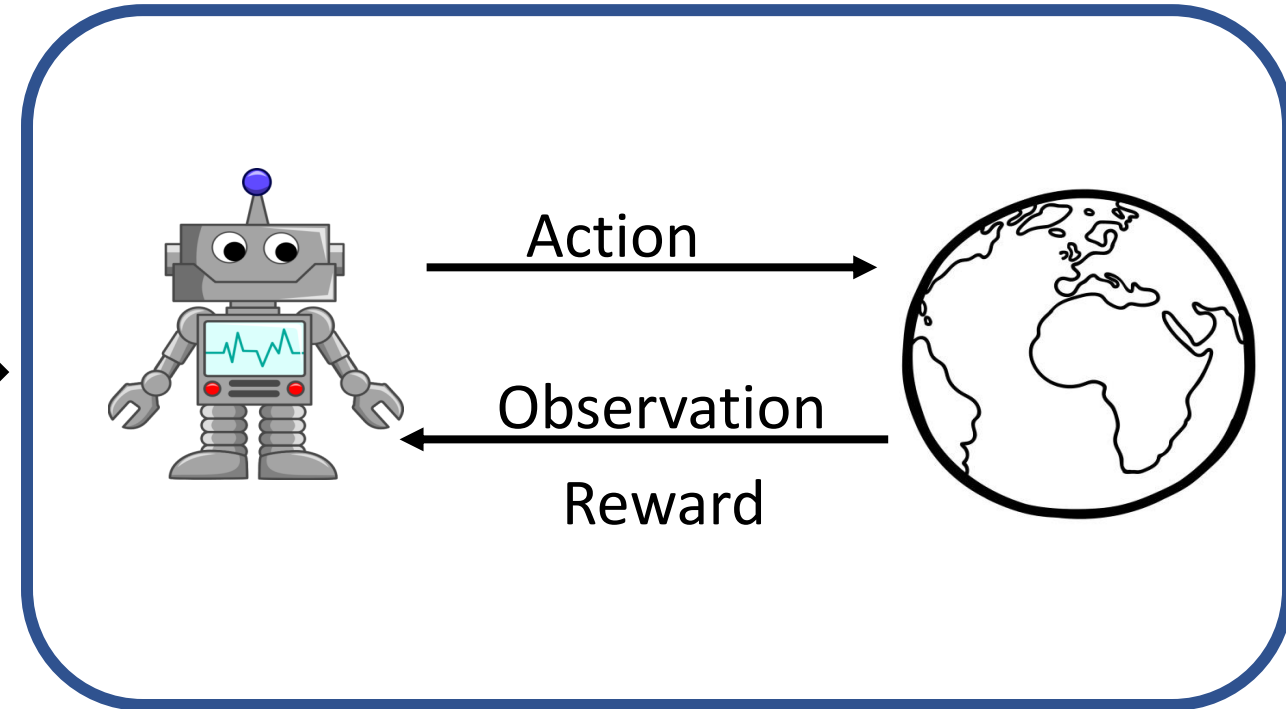
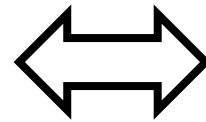
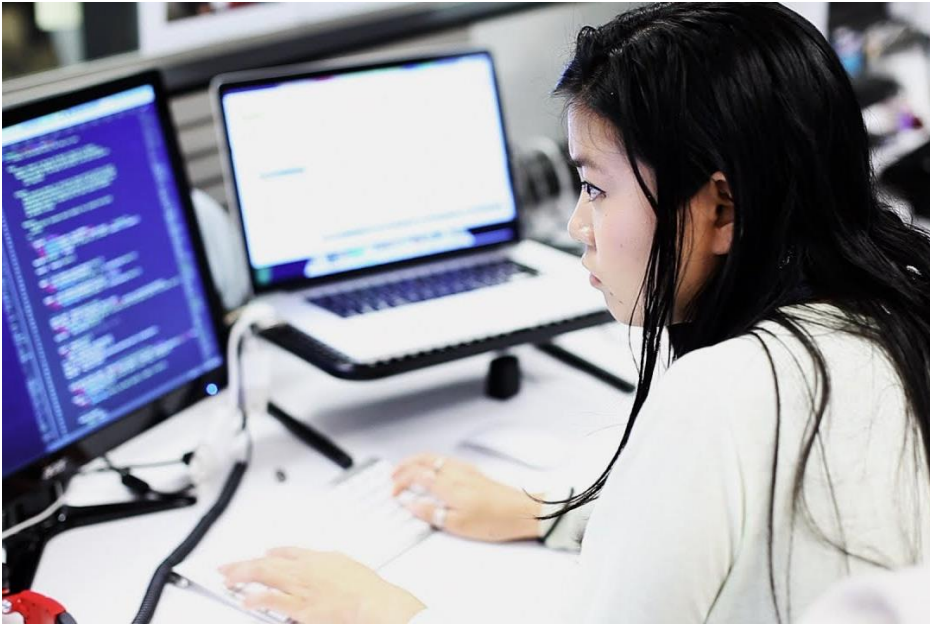
# Reinforcement Learning



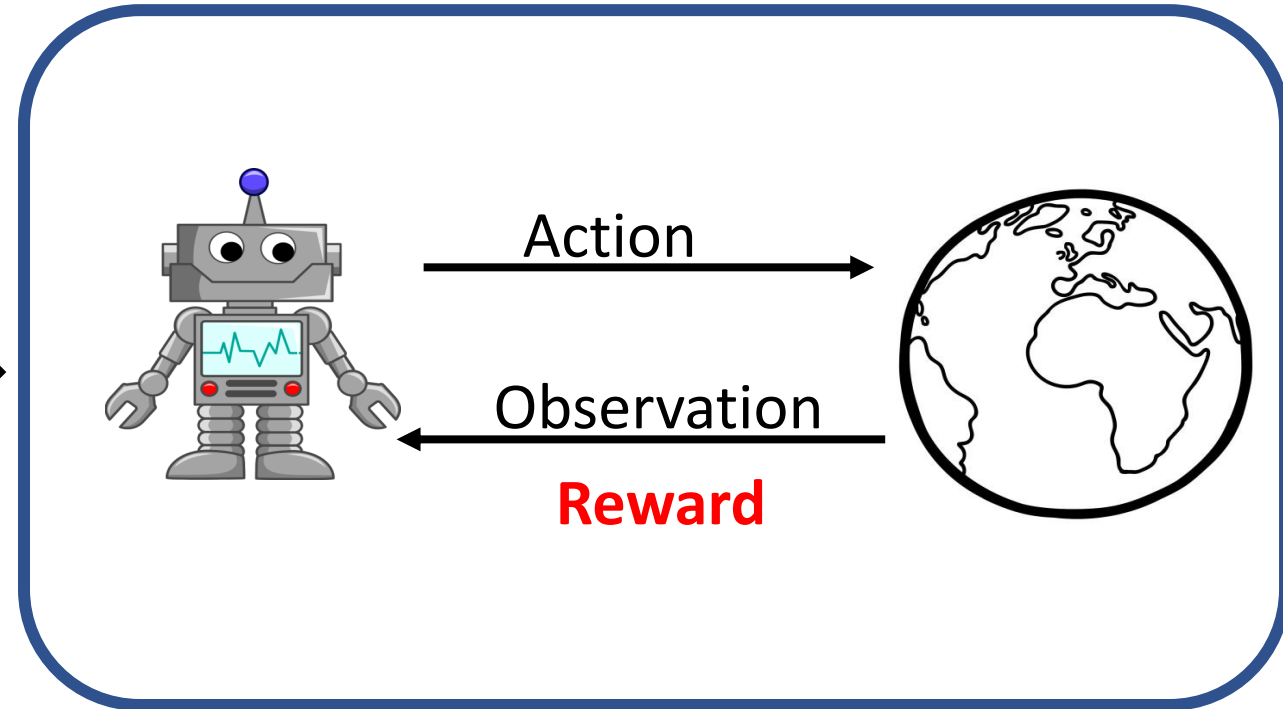
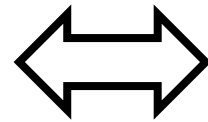
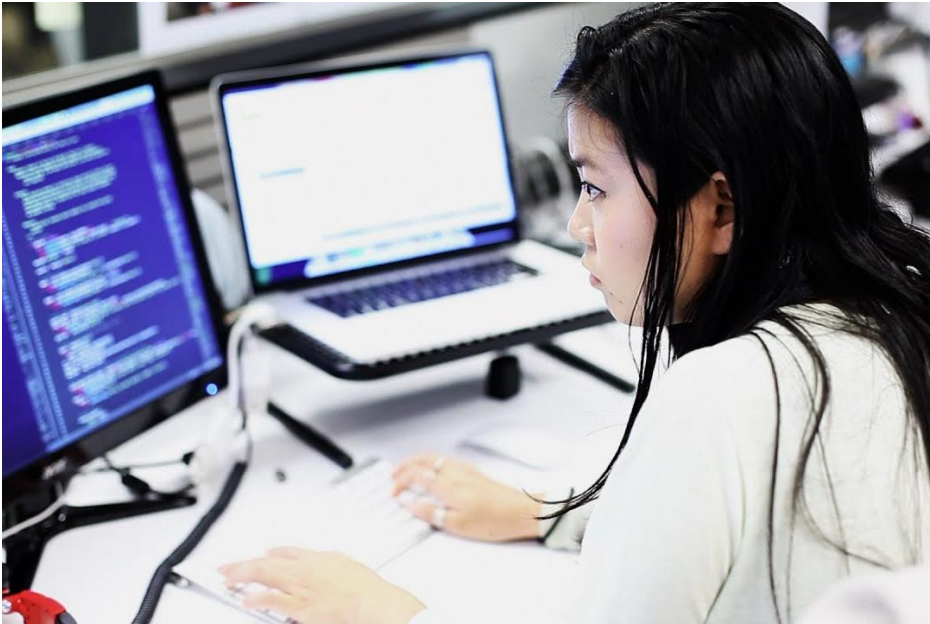
# Reinforcement Learning



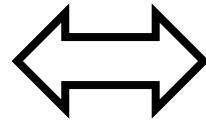
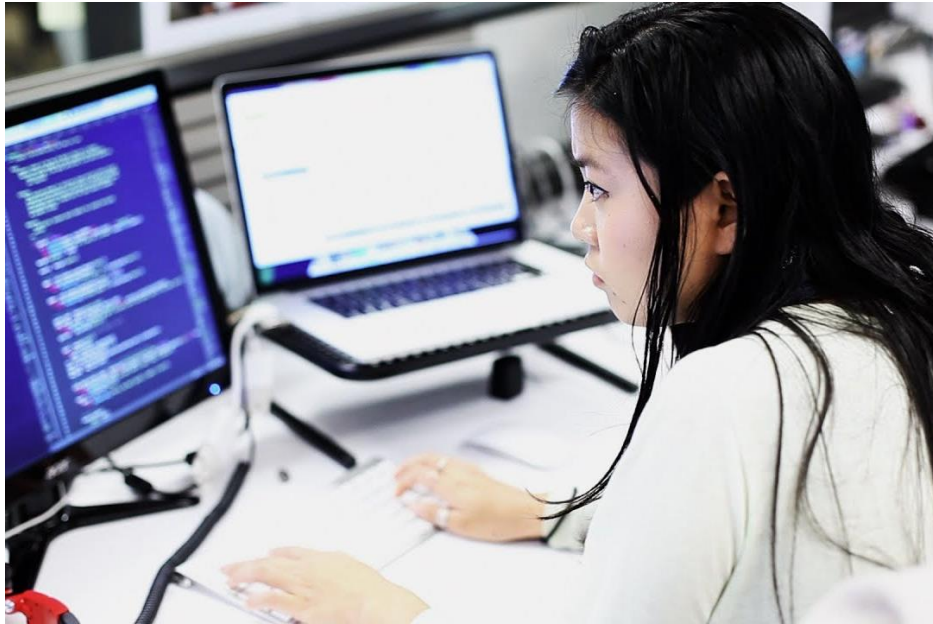
# Reward engineering is hard!



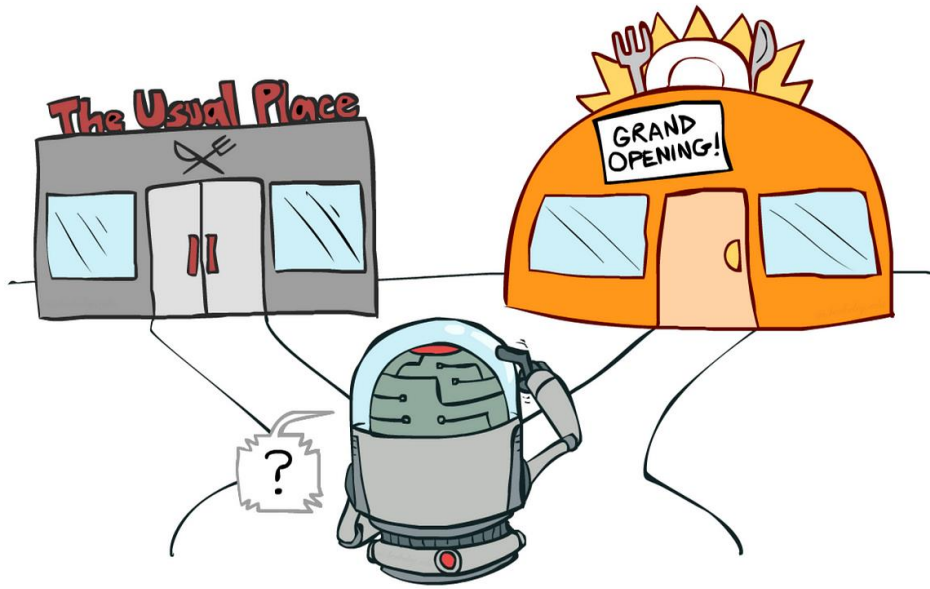
# Reward engineering is hard!



# Reward engineering is hard!



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





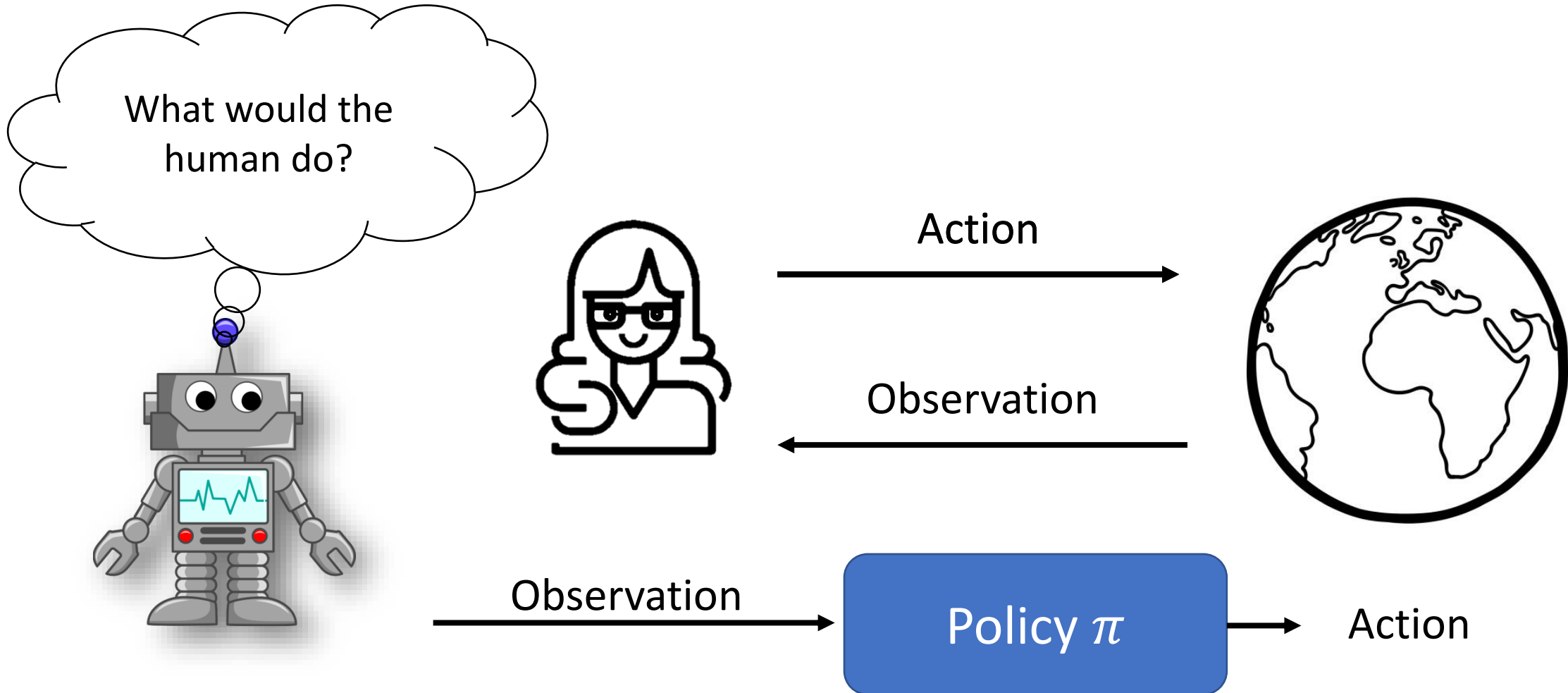
# Imitation Learning (Learning from Demonstrations):

Learn a policy from examples of good behavior.



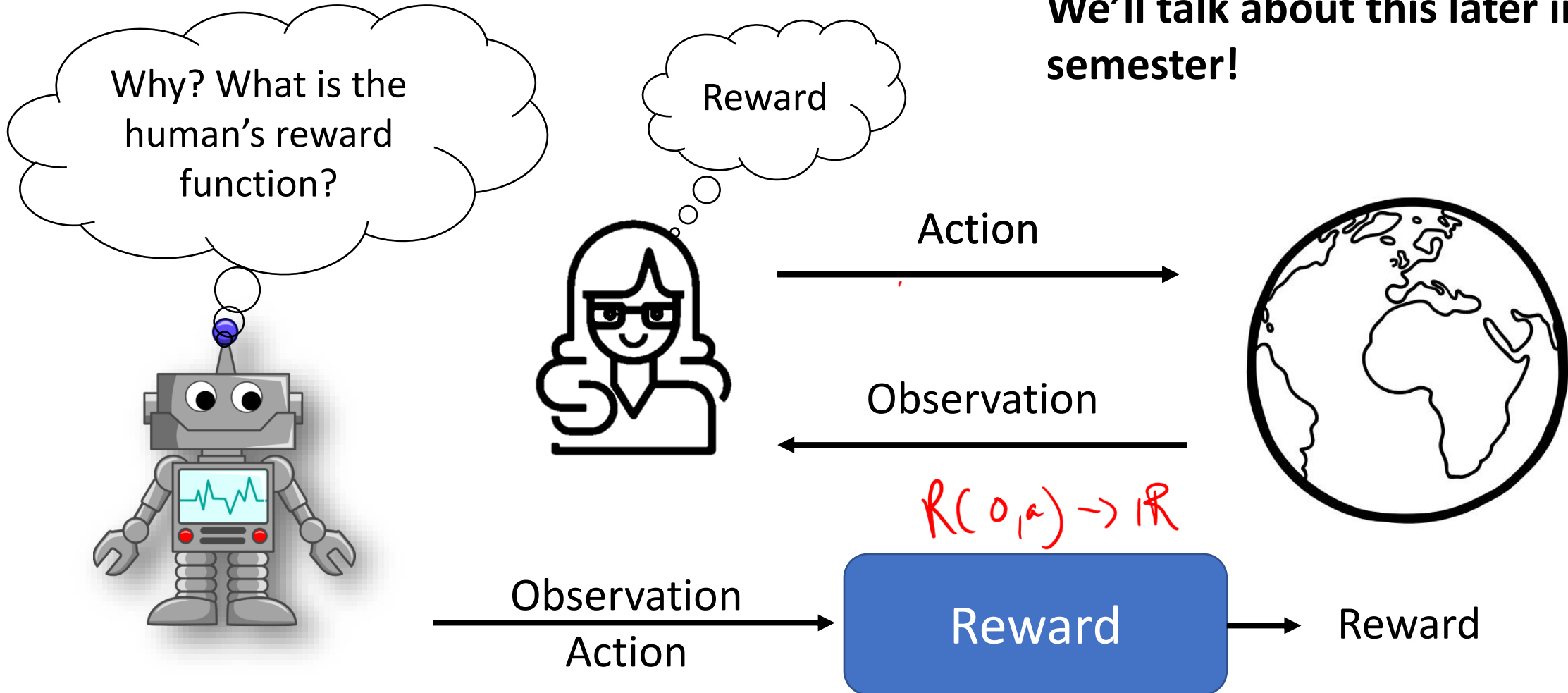
- Often showing is easier than telling.
- Alleviates problem of exploration.

# Behavioral Cloning

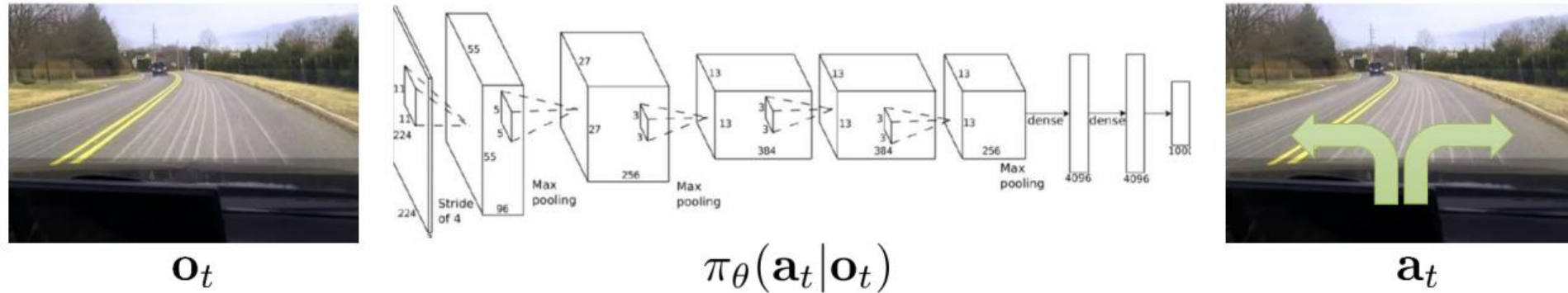


# Inverse Reinforcement Learning

**We'll talk about this later in the semester!**



# Imitation Learning via Behavioral Cloning



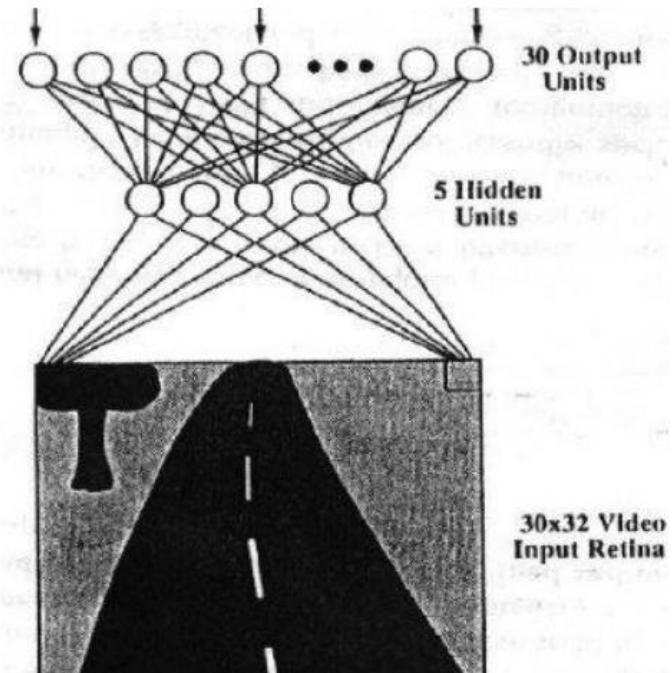
# Live demo

```
python test_gym.py
```

```
python mountain_car_bc.py --num_demos 1
```

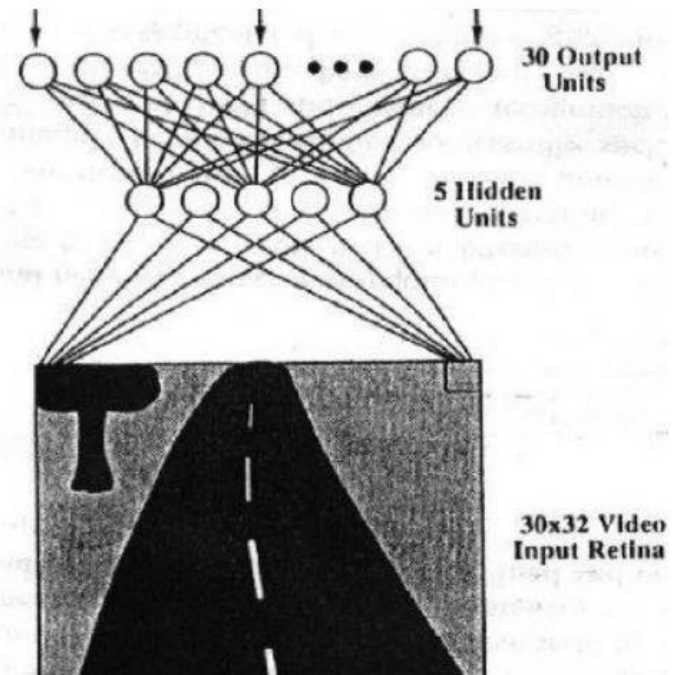
# ALVINN: One of the first imitation learning systems

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?

Chrome File Edit View History Bookmarks Profiles Tab Window Help

youtube.com/watch?v=joBmbh0AGSQ

YouTube

tire change

Incognito

Sign in

AHSOKA Now streaming Disney+ STREAM NOW

Ad - www.disneyplus.com Watch

How to change a tire | Dad, how do I? Dad, how do I? 886K views · 3 years ago 13:24

How to Replace your Flat Tire Pushing Pistons 117K views · 2 years ago 0:57

Steer Tire Change The Tire Doctor 571K views · 1 year ago 1:01

Winter Tire Swap How To Change Your Tires Yourself - Winter/Summer Tire... Your Home Garage 48K views · 2 years ago 12:47

How to Plug a Flat Tire (easily) ChrisFix 922K views · 9 months ago 1:00

POV Tire Change with your Dad #shorts Charlie Berens 7.4M views · 6 months ago

How to Change a Tire | Change a flat car tire step by step

Howdini 714K subscribers

Subscribe

36K

Share

Save

3.5M views 15 years ago

Nothing takes the joy out of a road trip like a flat tire. Do you know how to change it? We didn't, but we've learned from Allan Stanley of AAA. Download this video to your mobile phone just in case.



# Behavioral Cloning from Observation (Torabi et al. 2018)

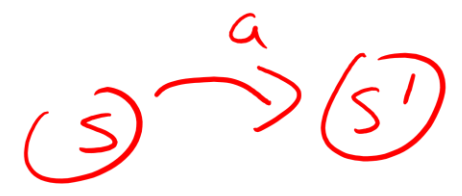
Demo =  $(s_1, s_2, \dots, s_{N-1}, s_N)$   
 $\hat{a}_1 = g(s_1, s_2)$        $\hat{a}_{N-1}$

Train action predictor  
Inverse dynamics model

$$g(s, s') \rightarrow a$$

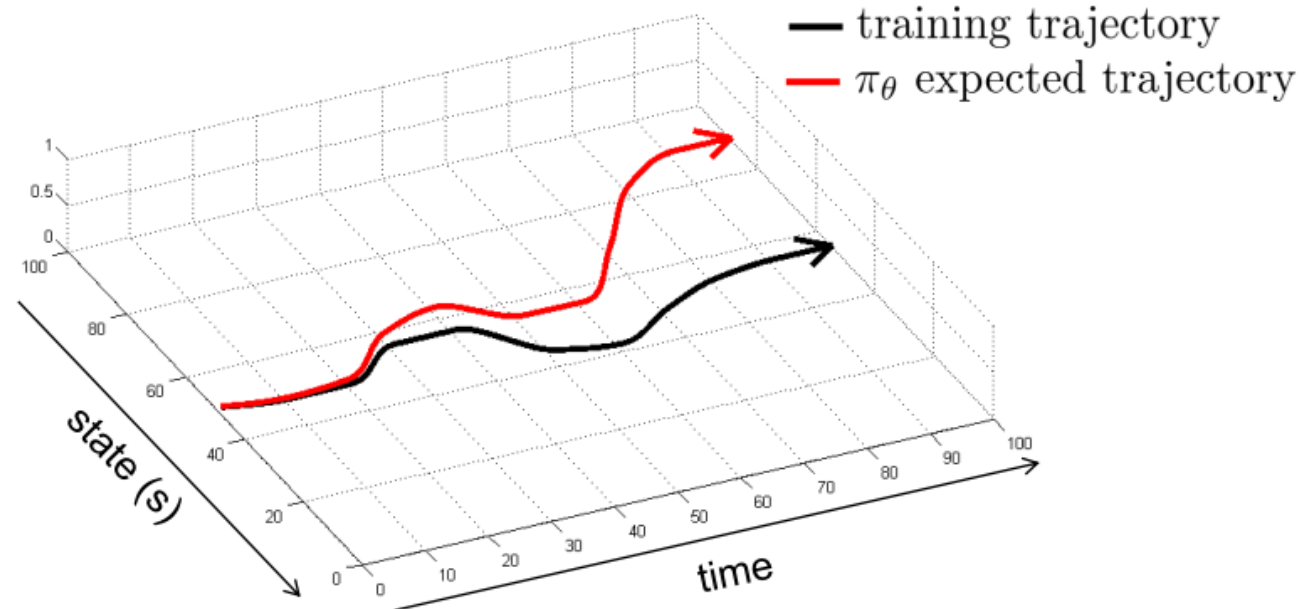
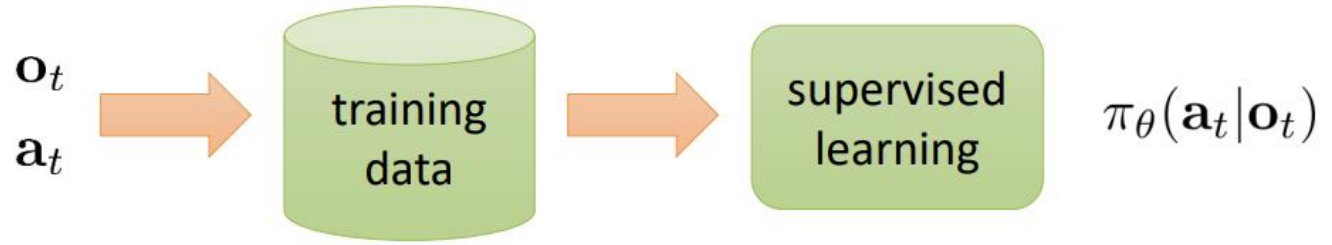
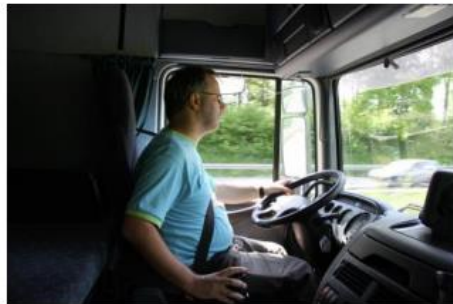
Explore  $\rightarrow \underline{s, a, s', a', s''}$

Dynamics model  
 $f(s, a) \rightarrow s'$



BC(0) Explore, learn  $g$ , relabel, BC

# What could go wrong?

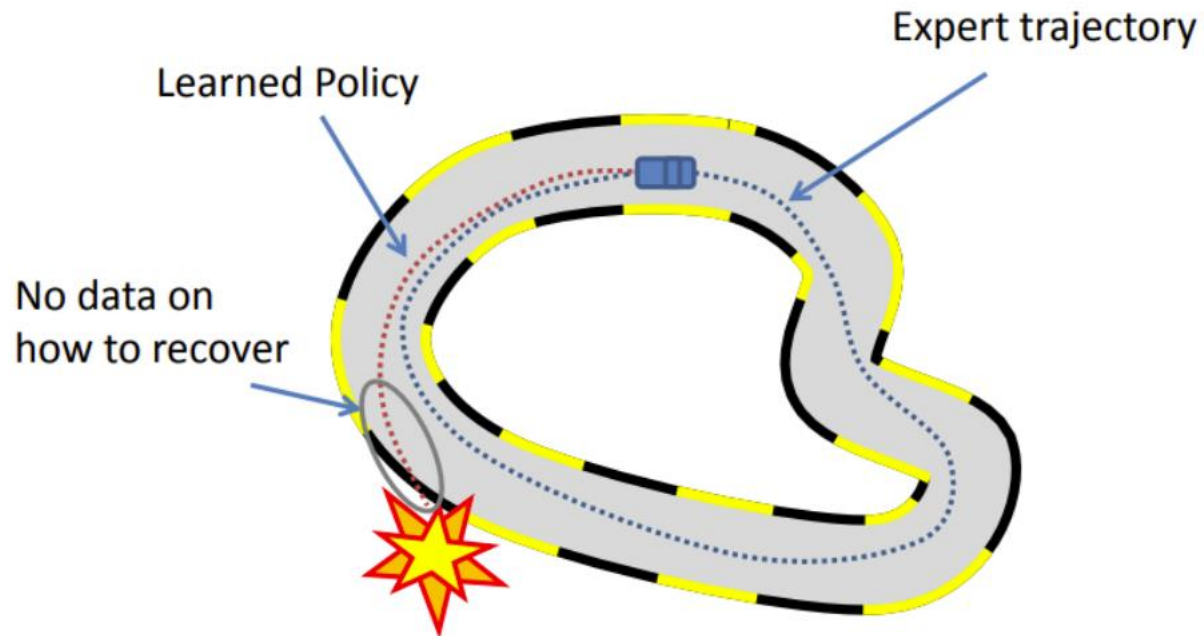


# Distribution Shift

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

$$\pi(s) \rightarrow a$$

$\pi^*$  expert

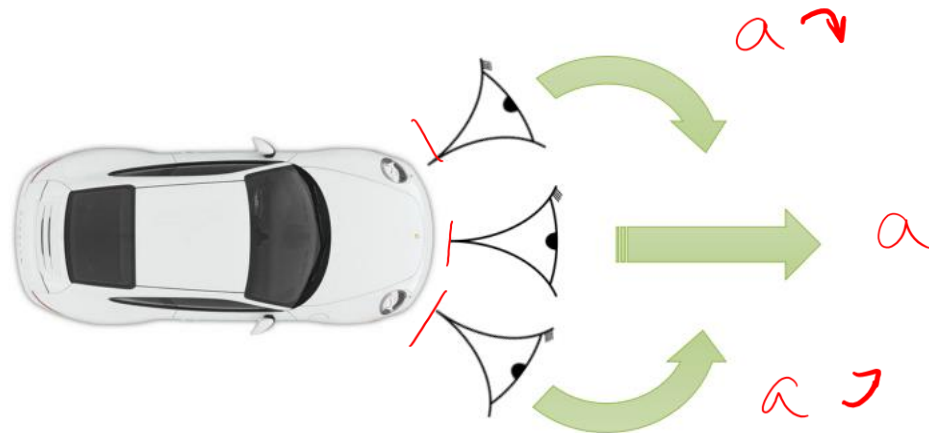
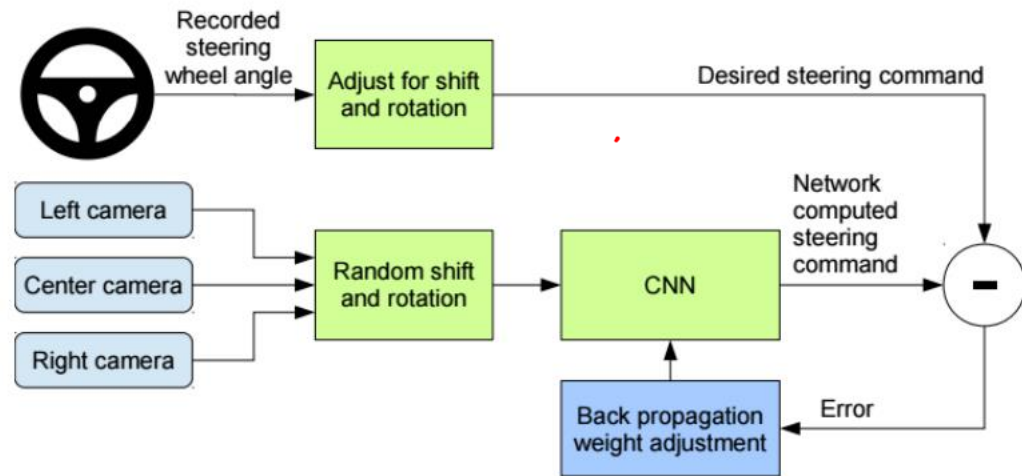


	Supervised Learning	Supervised Learning + Control
Train	$(x, y) \sim D$	$s \sim P(\cdot   s, \pi^*(s))$
Test	$(x, y) \sim D$	$s \sim P(\cdot   s, \pi(s))$

But it still can work in practice...

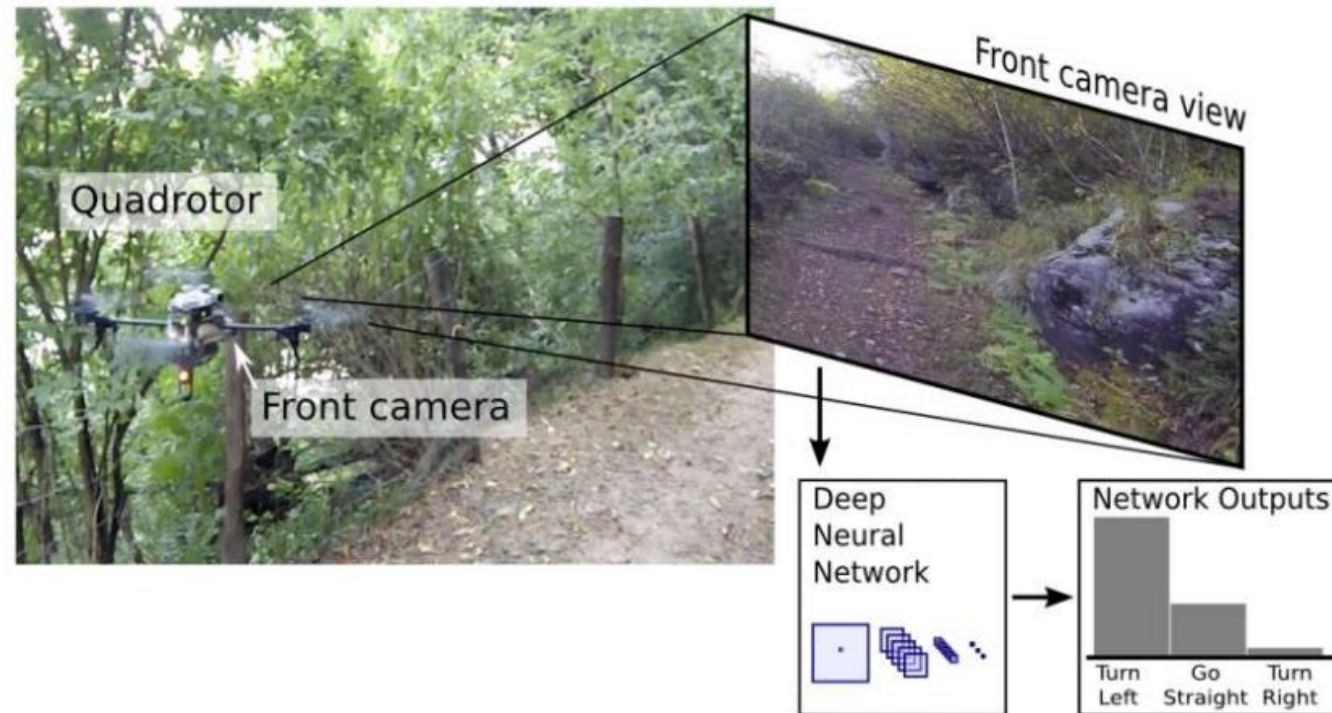


# How?



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

Alessandro Giusti<sup>1</sup>, Jérôme Guzzi<sup>1</sup>, Dan C. Cireşan<sup>1</sup>, Fang-Lin He<sup>1</sup>, Juan P. Rodríguez<sup>1</sup>  
Flavio Fontana<sup>2</sup>, Matthias Faessler<sup>2</sup>, Christian Forster<sup>2</sup>  
Jürgen Schmidhuber<sup>1</sup>, Gianni Di Caro<sup>1</sup>, Davide Scaramuzza<sup>2</sup>, Luca M. Gambardella<sup>1</sup>



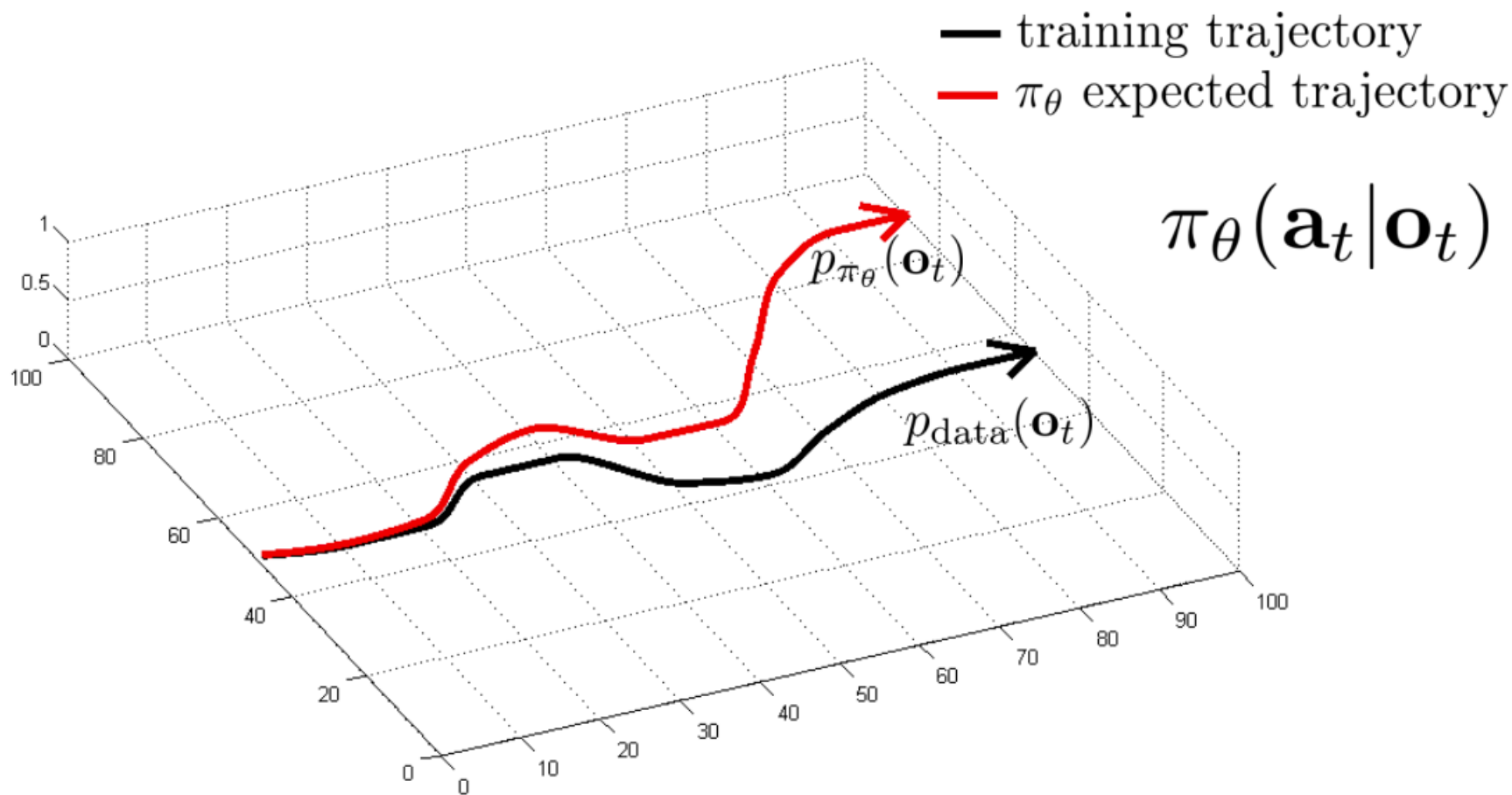
### Deep Neural Network



### Control Signal

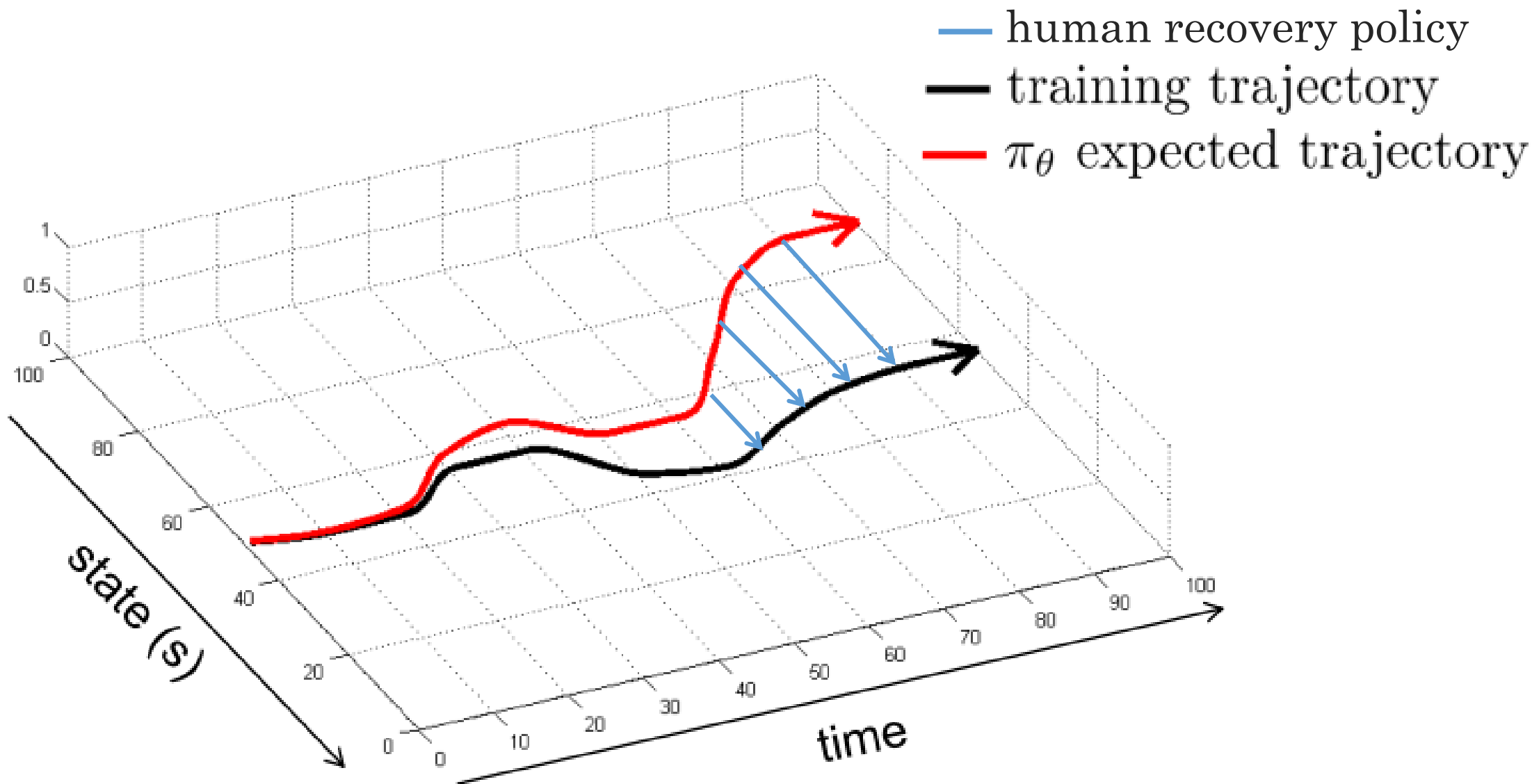


# Can we make it work more often?



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





# Dataset Aggregation

## 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  $\mathcal{D}_\pi$  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?

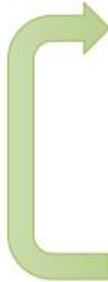
## **D**agger: Dataset Aggregation

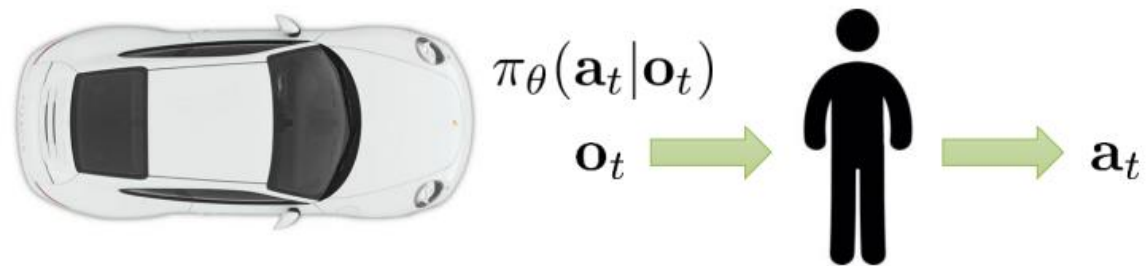
goal: collect training data from  $p_{\pi_{\theta}}(\mathbf{o}_t)$  instead of  $p_{\text{data}}(\mathbf{o}_t)$

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  $\mathcal{D}_{\pi}$  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  $\mathcal{D}_{\pi}$  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  $\mathbf{a}_t$ !

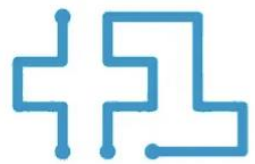
- 
1. train  $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$  from human data  $\mathcal{D} = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_N, \mathbf{a}_N\}$
  2. run  $\pi_{\theta}(\mathbf{a}_t|\mathbf{o}_t)$  to get dataset  $\mathcal{D}_{\pi} = \{\mathbf{o}_1, \dots, \mathbf{o}_M\}$
  3. Ask human <sup>oracle</sup> to label  $\mathcal{D}_{\pi}$  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



WAYMO

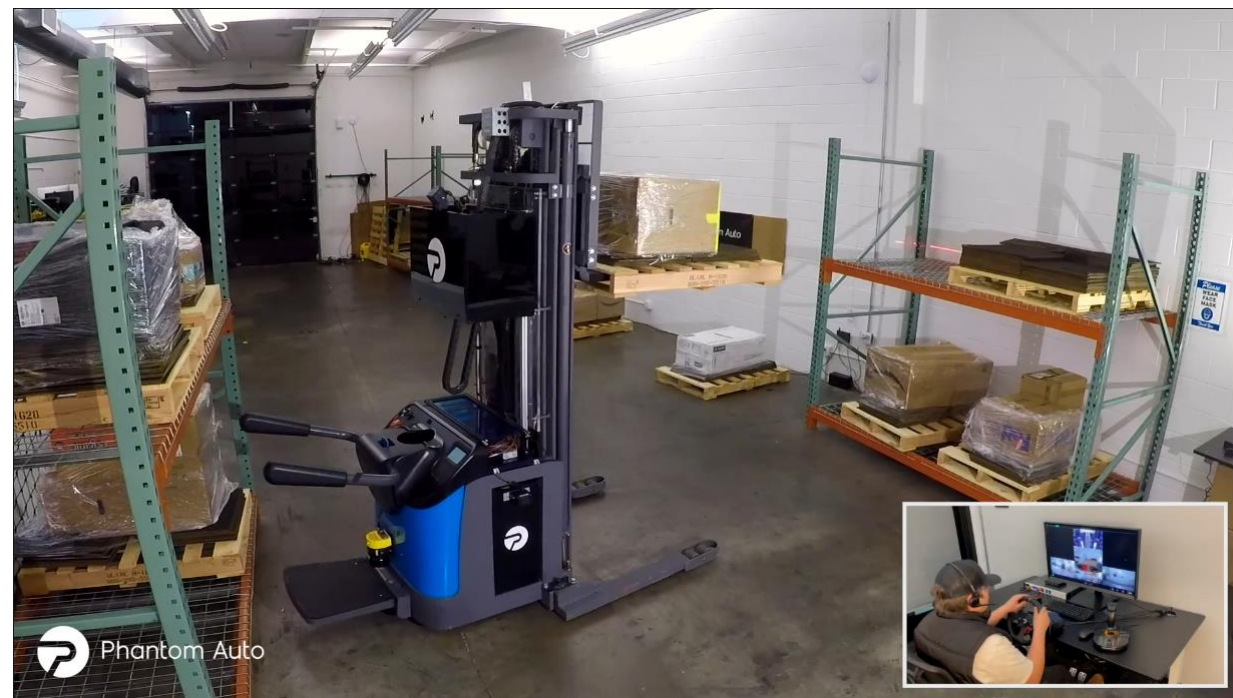
ZOOX



nimble

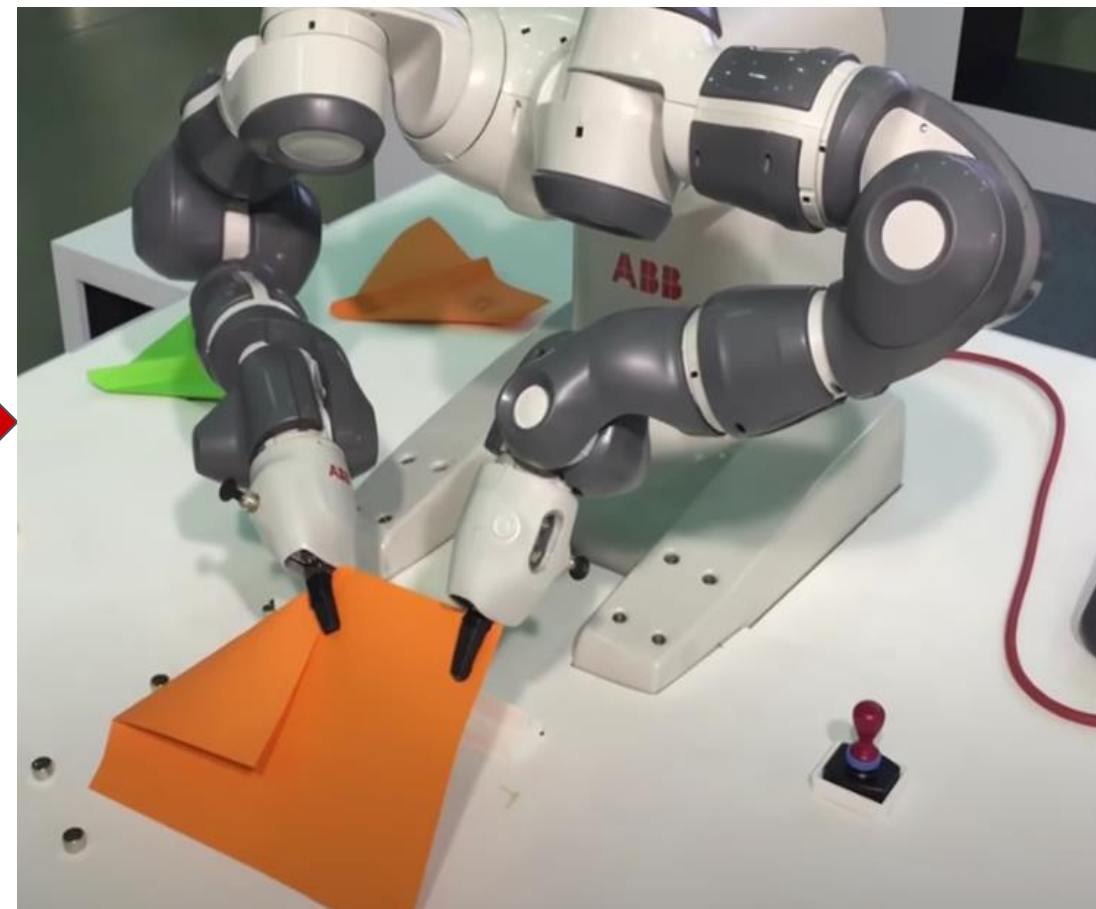


Phantom Auto





# Interactive IL



# Interactive IL

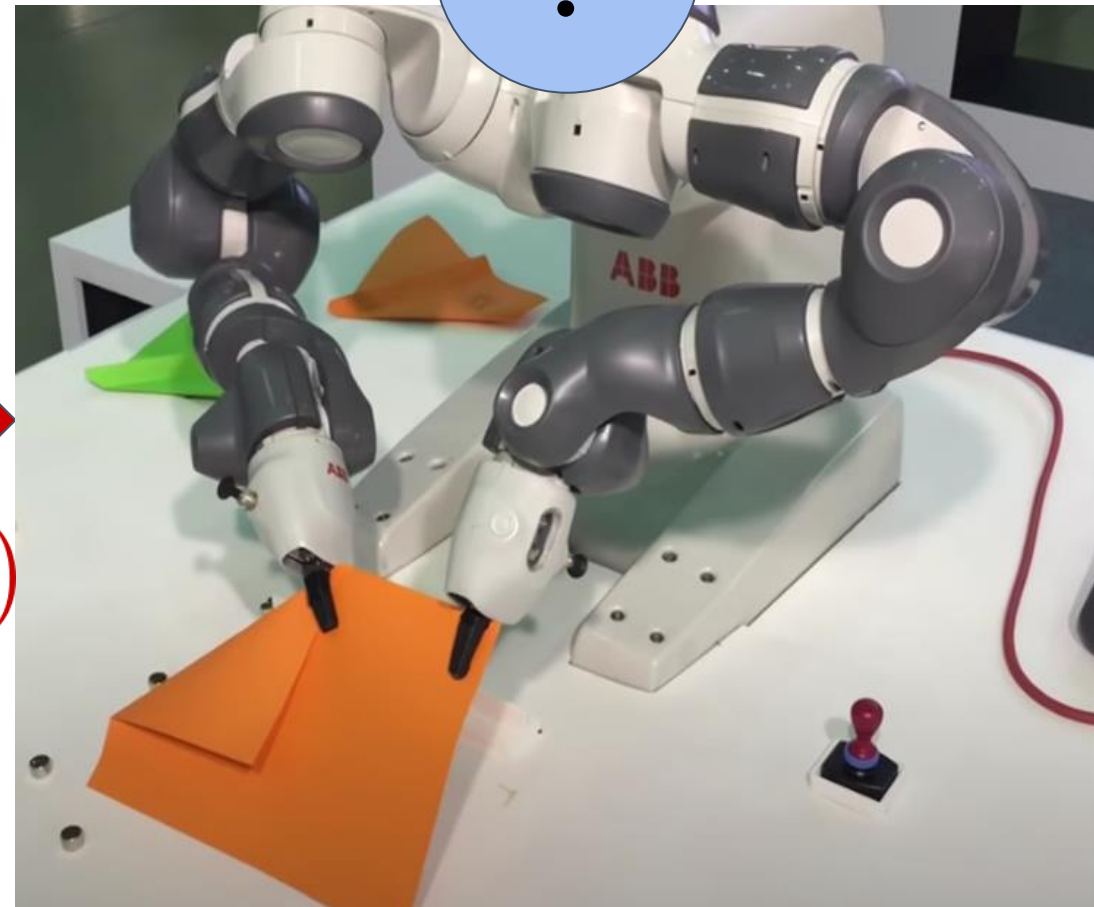


$\pi_H(s)$



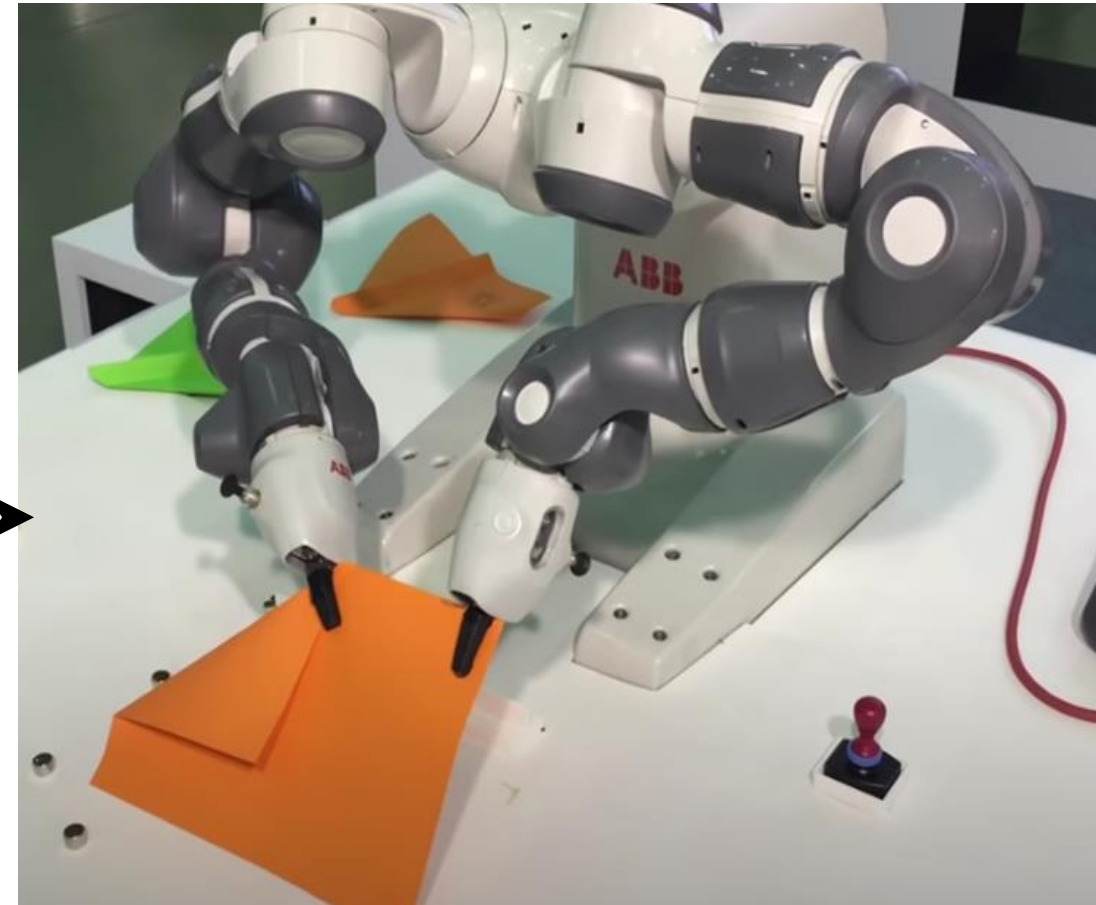
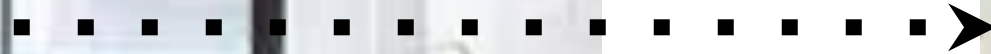
$\pi_{\text{meta}}(s)$

???



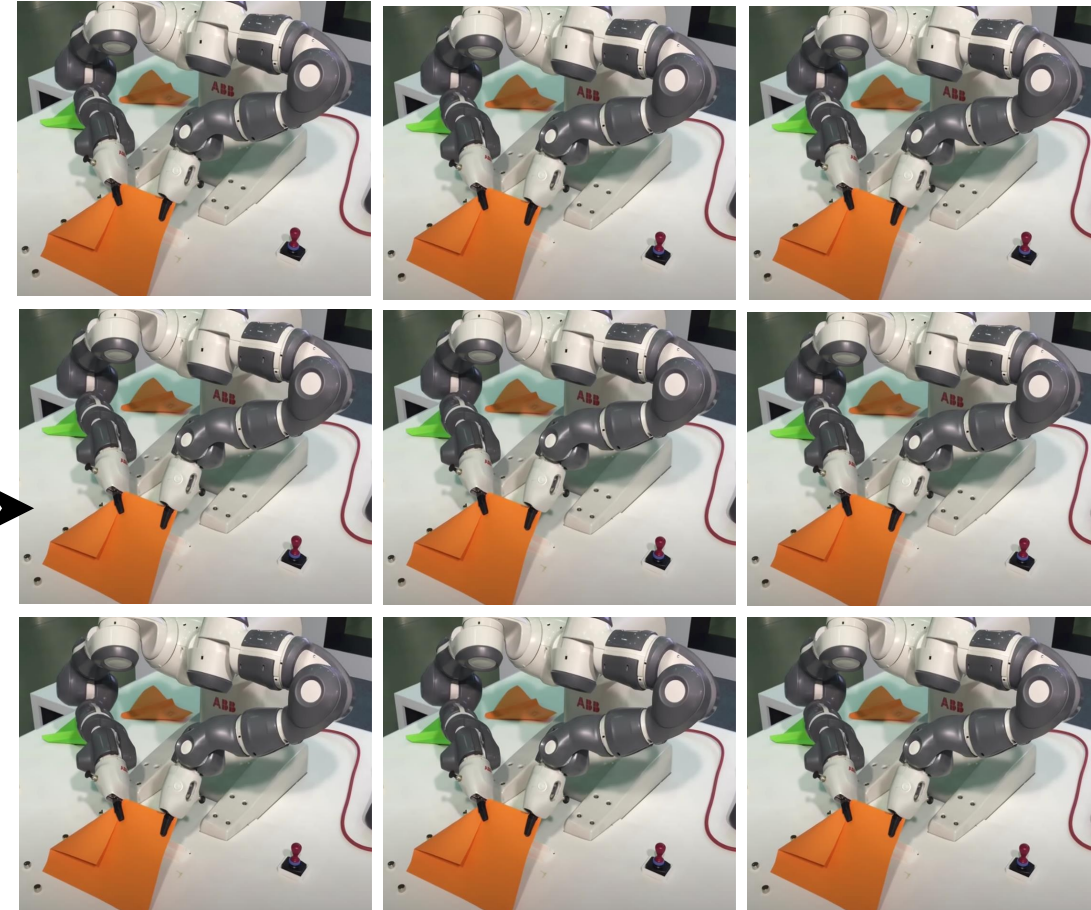
$\pi_R(s)$

# Human-Gated Interactive IL



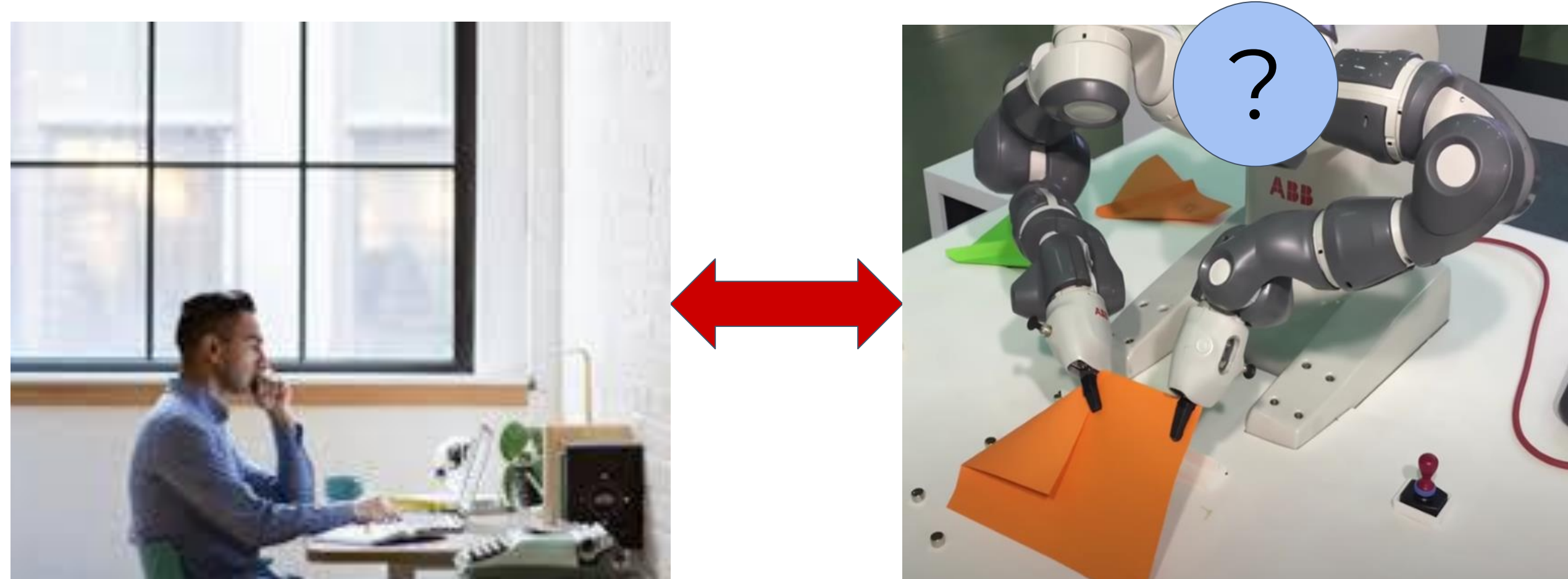
[3] M. Kelly, C. Sidrane, K. Driggs-Campbell, and M. J. Kochenderfer. HG-Dagger: Interactive Imitation Learning with Human Experts. ICRA 2019.

# 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.

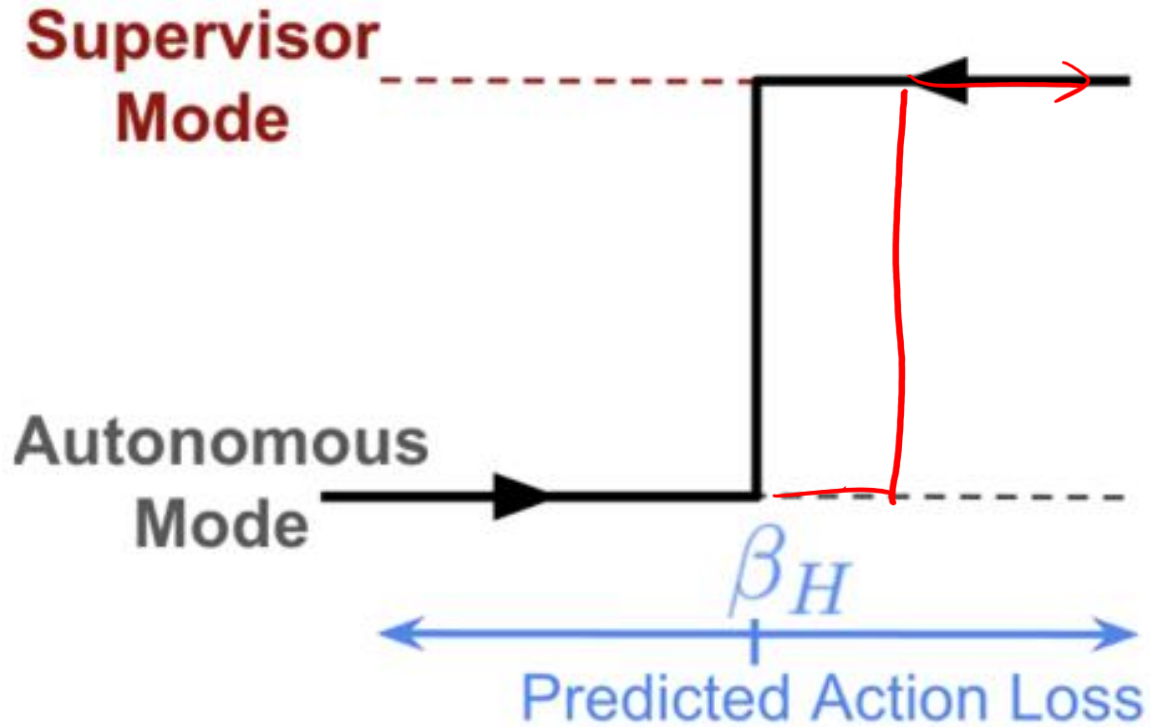
# Robot-Gated Interactive IL



[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.

# SafeDAgger



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

Trained using held-out set of data from human.



# When should a robot ask for help?



Novel (and risky)

# When should a robot ask for help?



Novel (and risky)

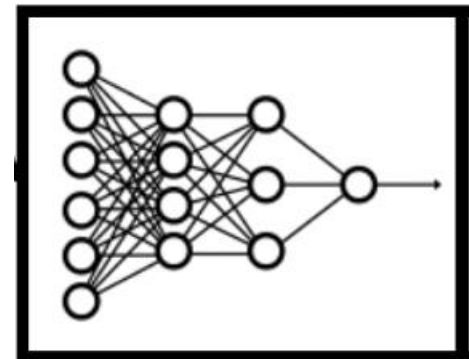
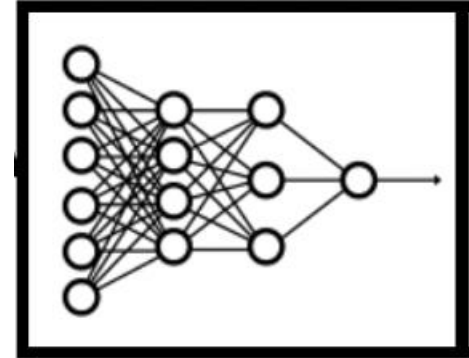
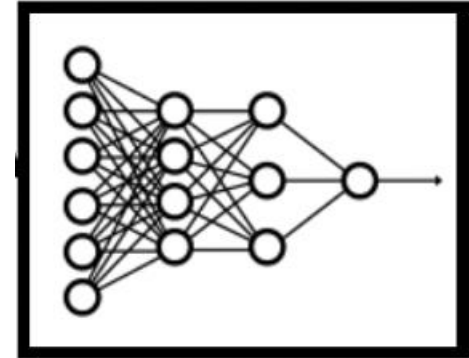
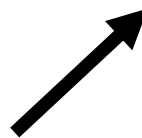
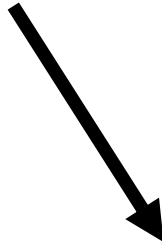
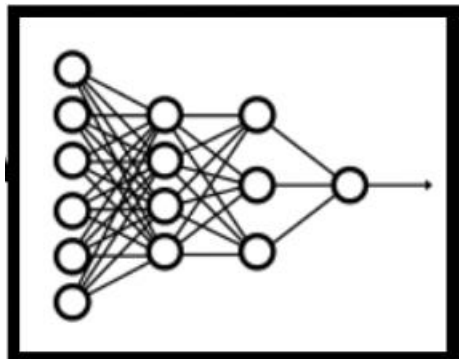
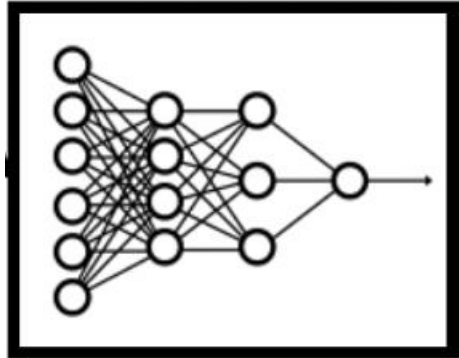
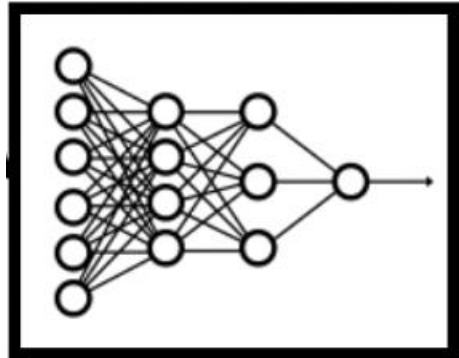


Risky (but not novel)

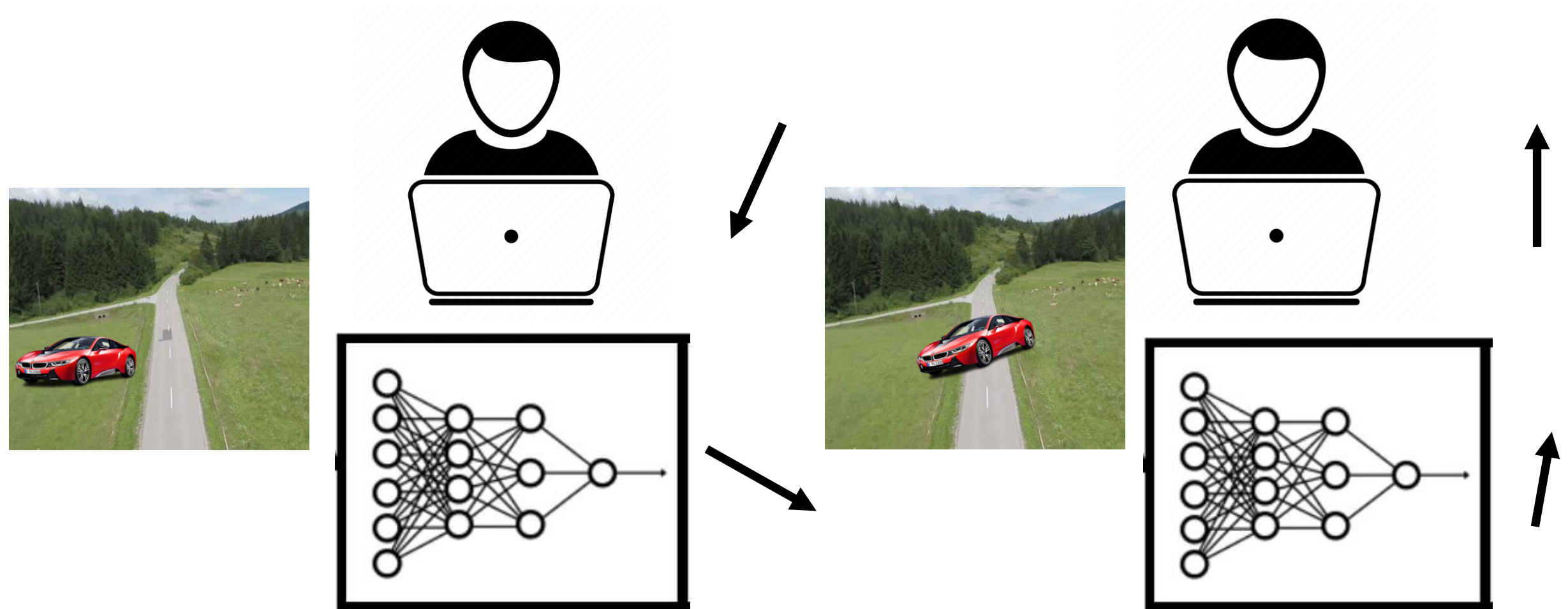


# Novelty Estimation

*Ensemble*



# Novelty Estimation: Supervisor Mode



# Risk Estimation

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

*goal classifier*

# 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

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

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

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

# 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**

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

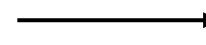


Switch to  
**AUTONOMOUS  
MODE**

# Putting it all together...

**AUTONOMOUS  
MODE**

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



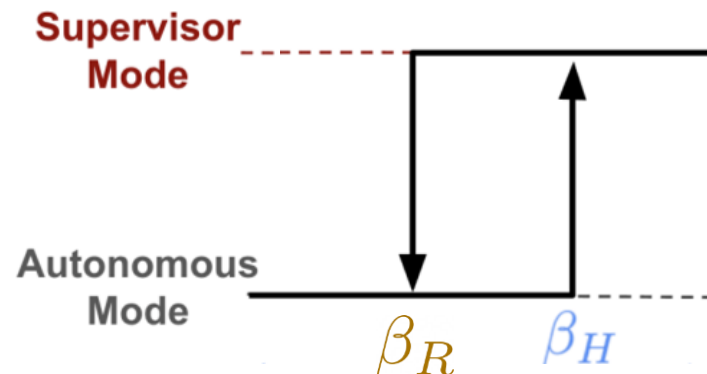
Switch to  
**SUPERVISOR  
MODE**

**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**



How do we deal with all the hyperparameters?



# Putting it all together...

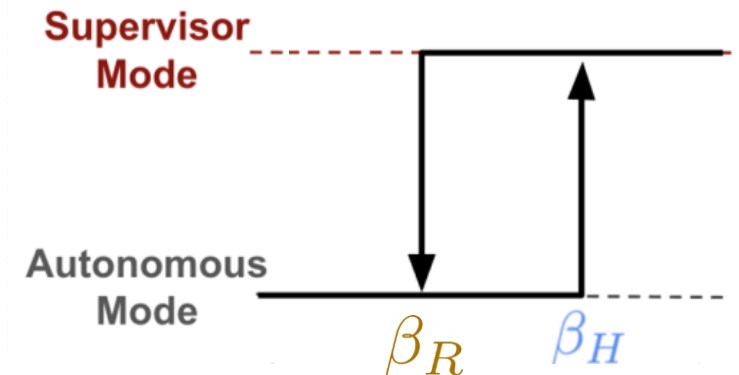
**AUTONOMOUS  
MODE**

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

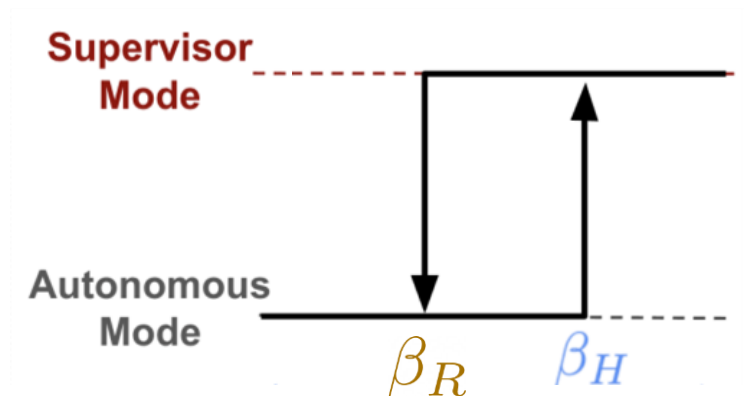
Set to 1- $\alpha$  quantiles of empirical data

**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

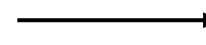
$$\alpha = \frac{\text{desired \# interventions}}{\text{\# robot actions}}$$



# 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**

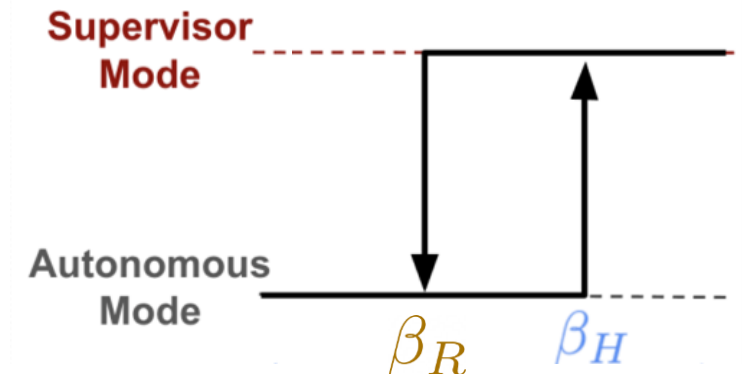
**SUPERVISOR  
MODE**

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



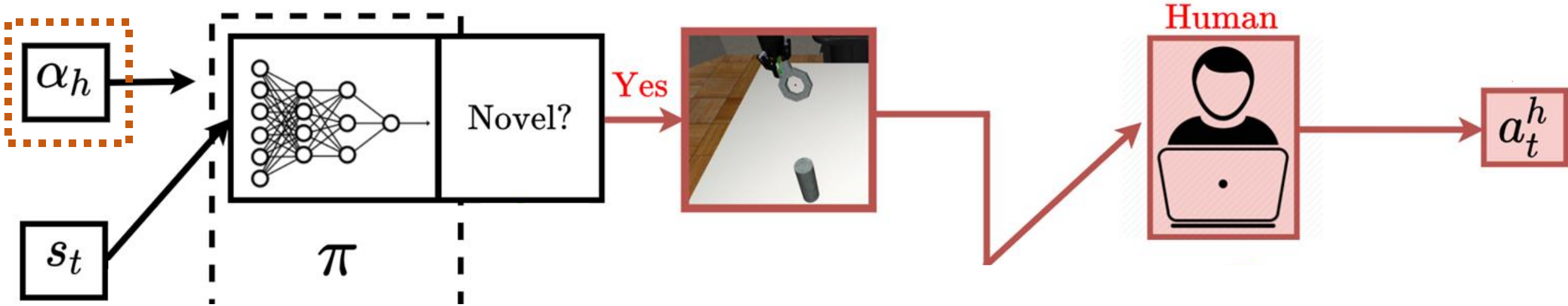
Switch to  
**AUTONOMOUS  
MODE**

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

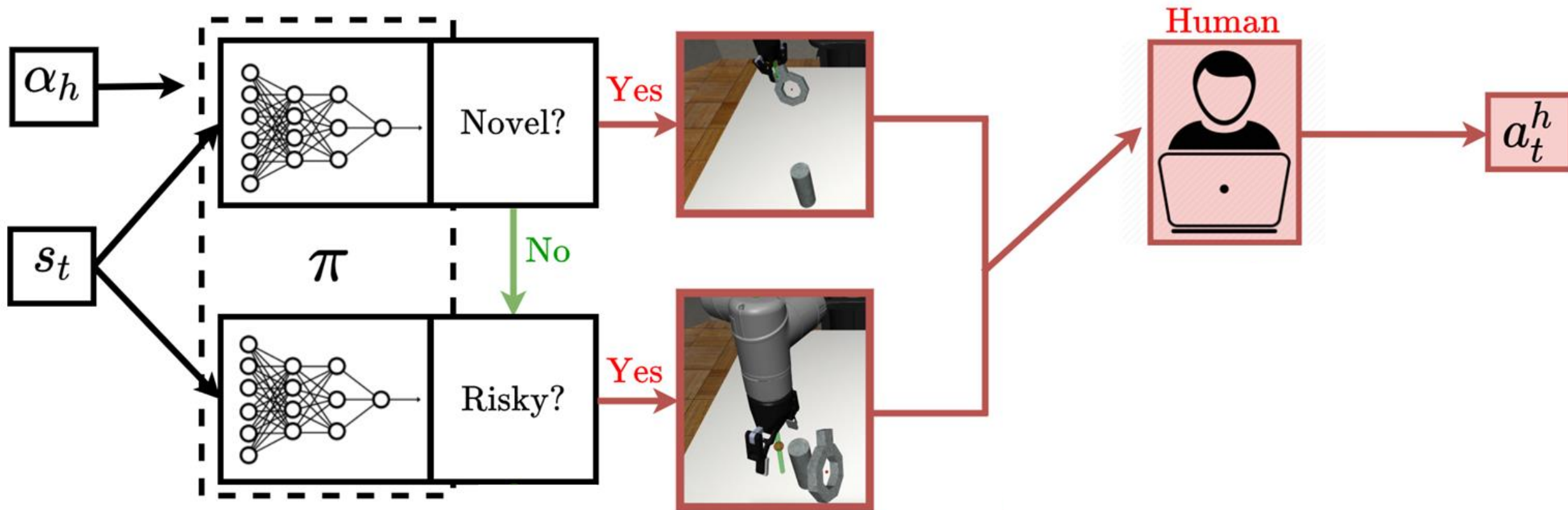


Target percent of time human wants to give interventions.

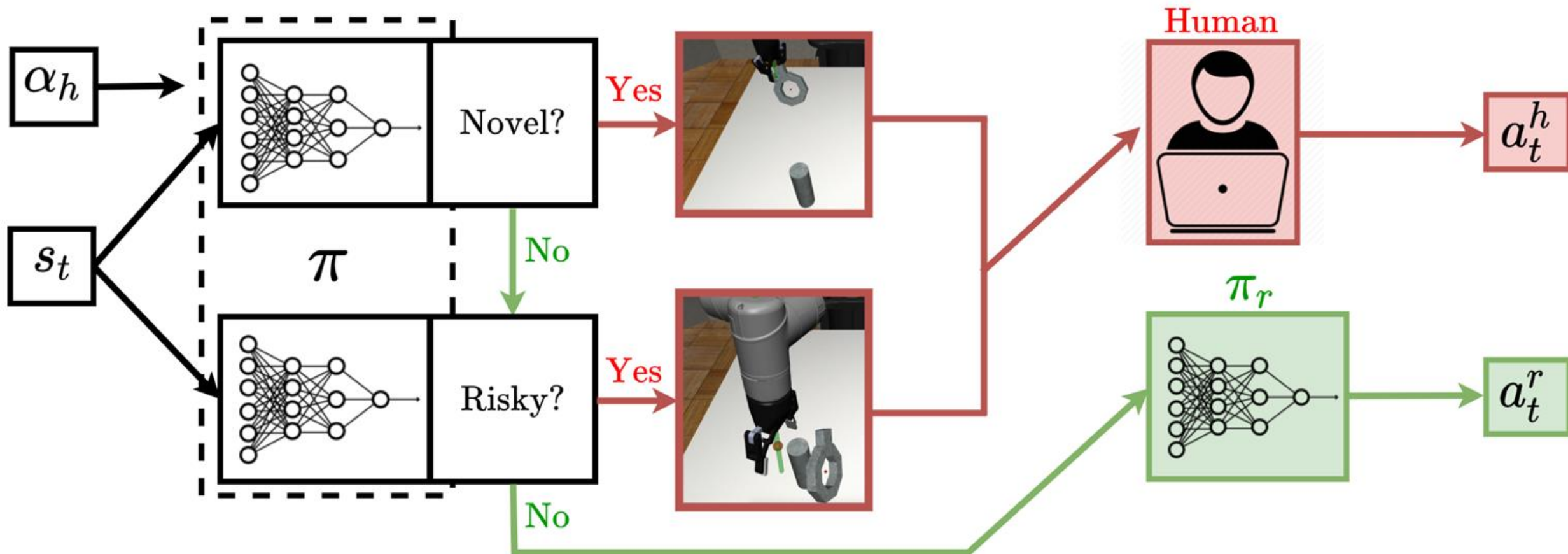
# ThriftyDAgger

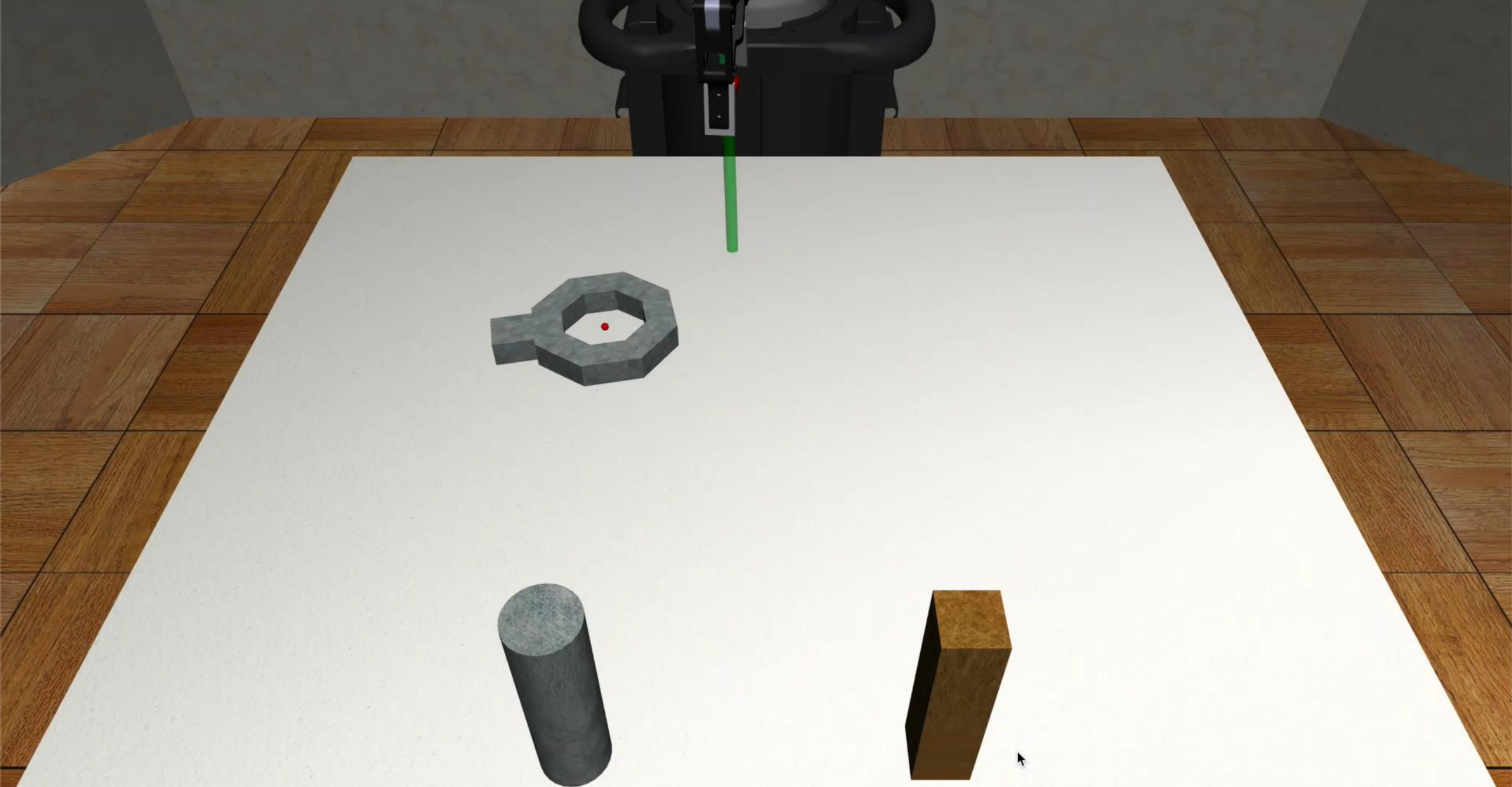


# ThriftyDAgger

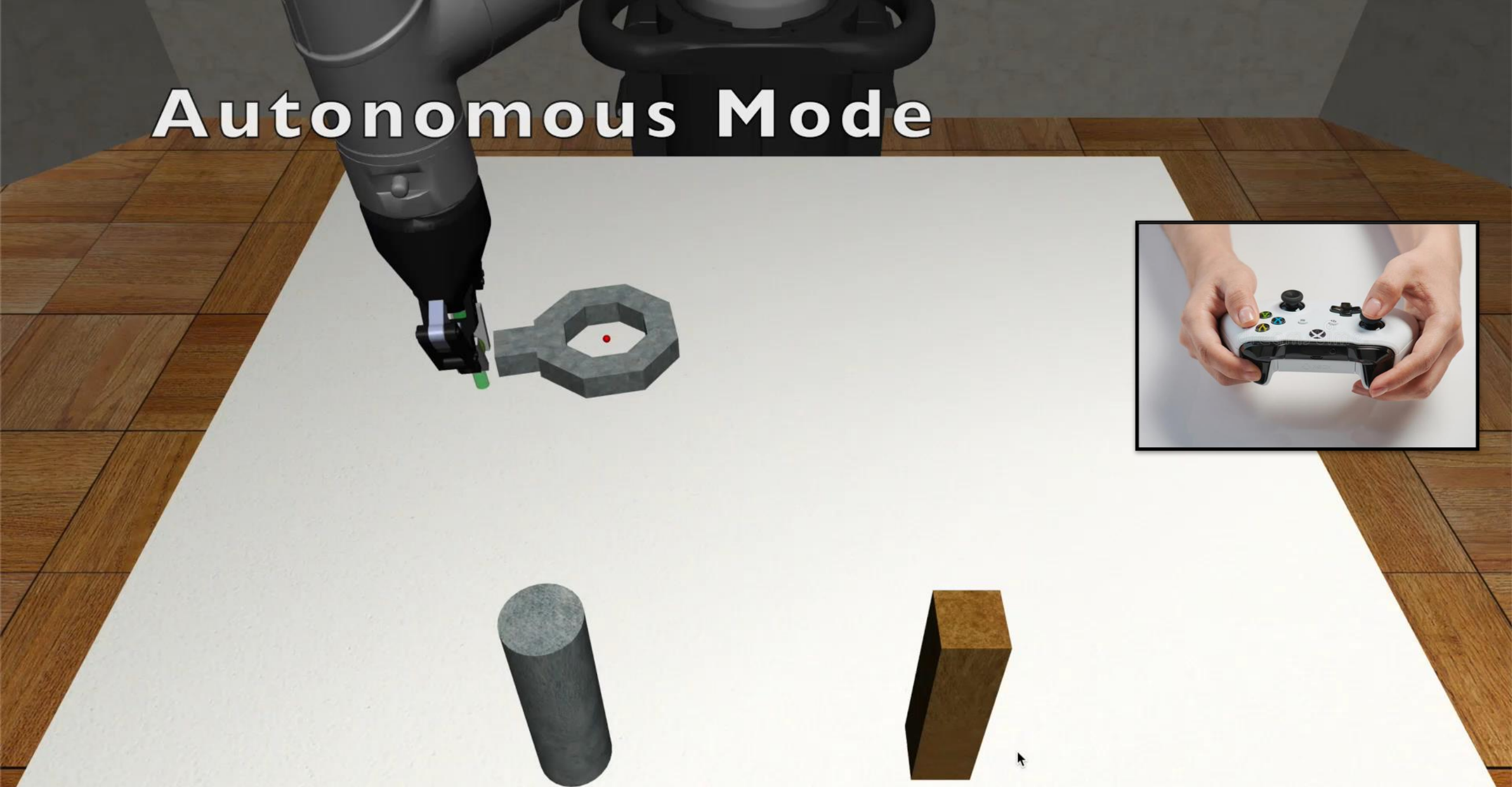


# ThriftyDAgger





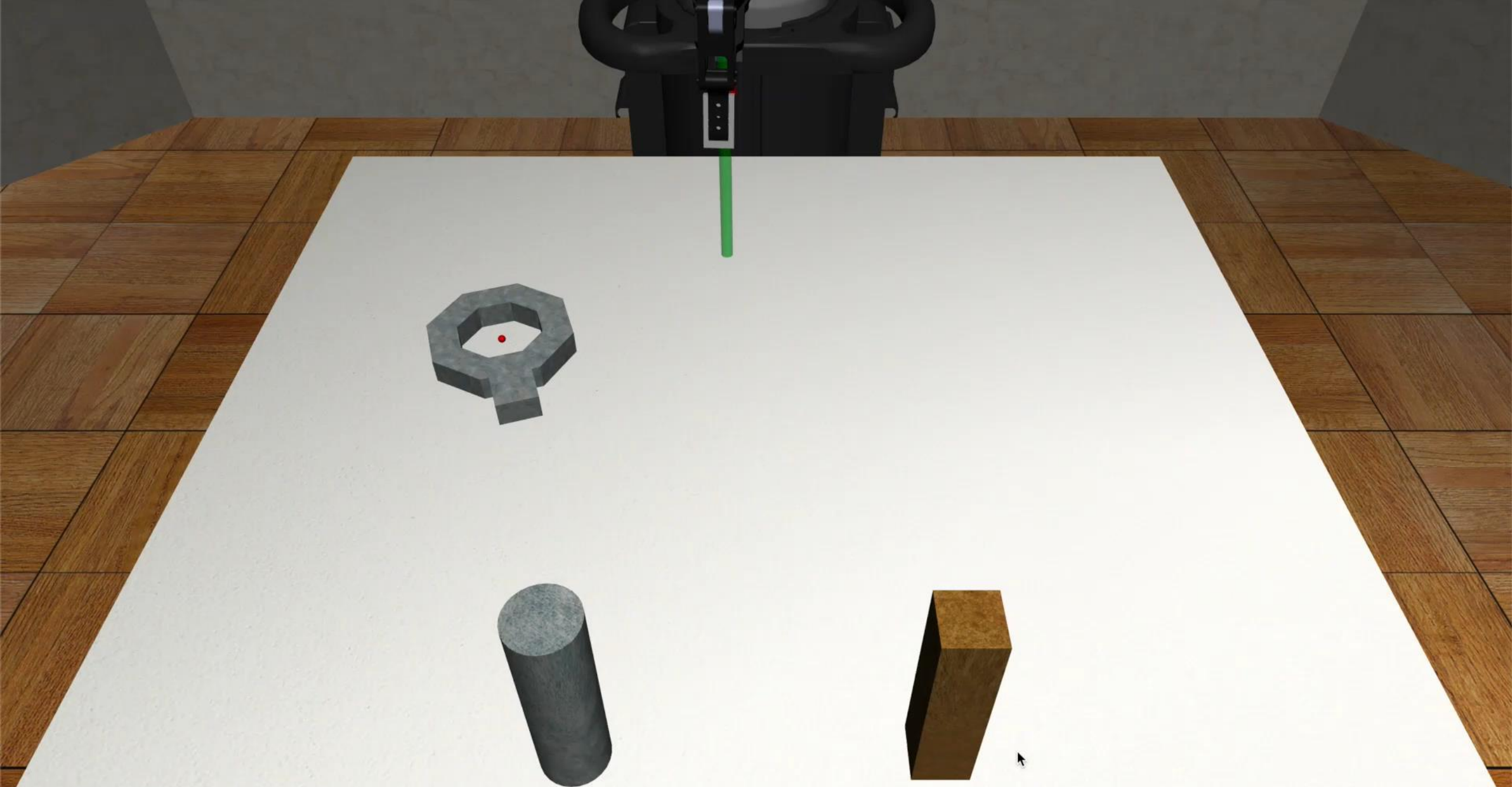
# Autonomous Mode



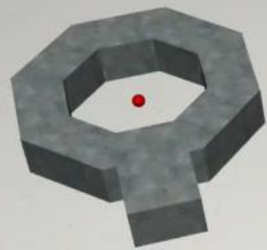


# Supervisor Mode (Novel)

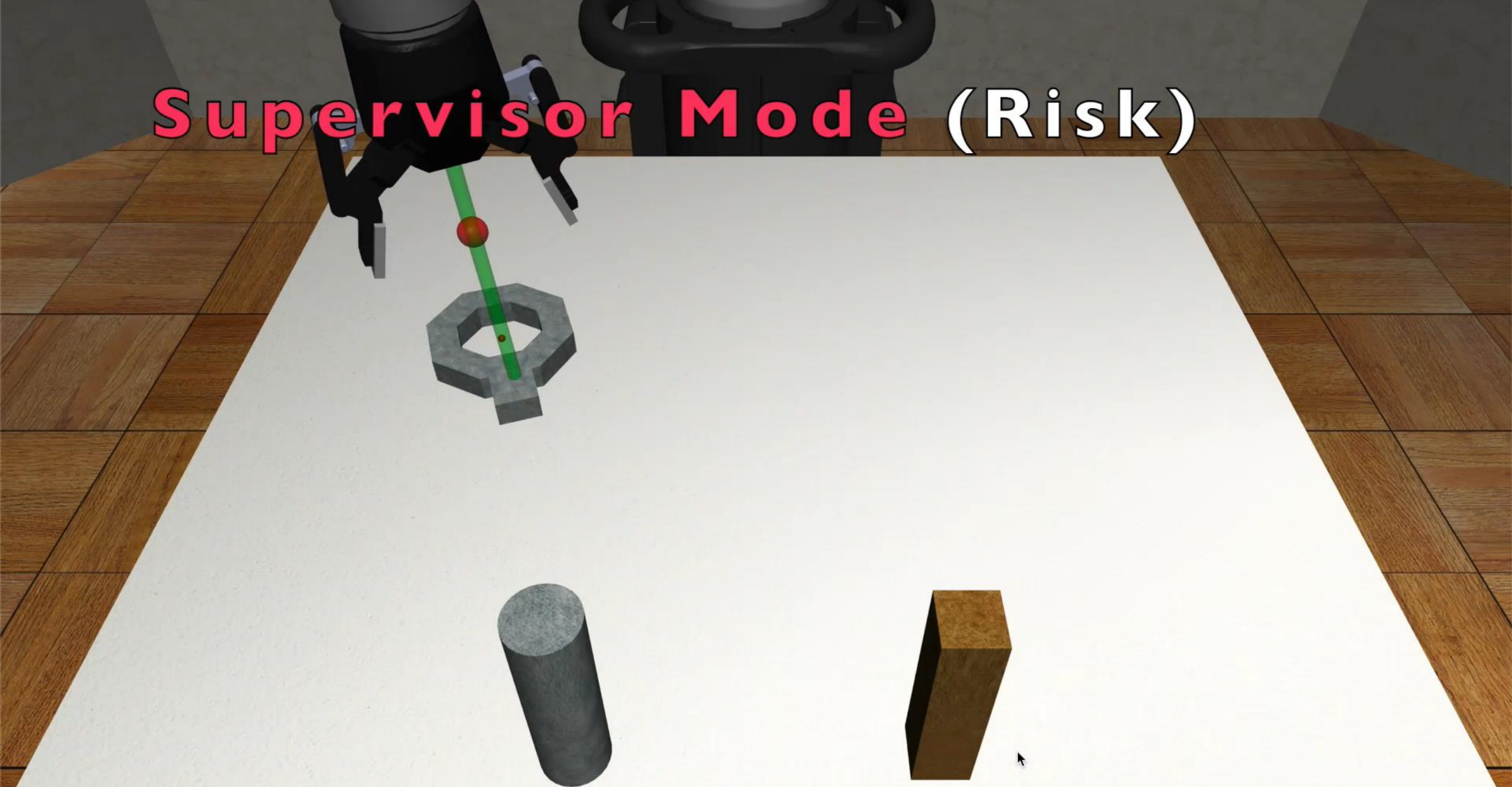




# Supervisor Mode (Risk)



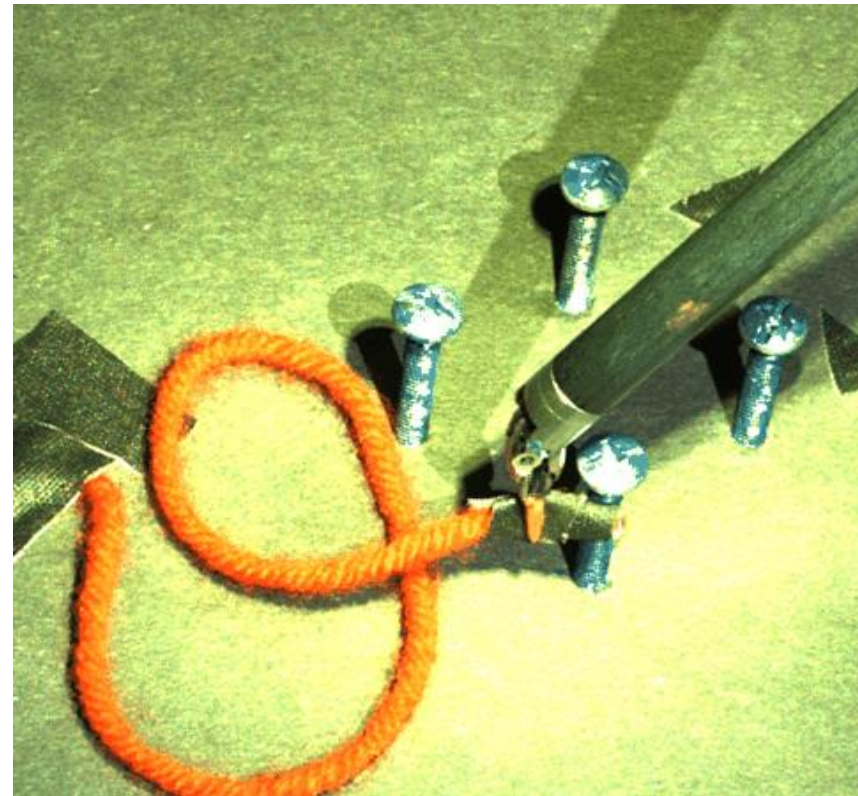
# Supervisor Mode (Risk)



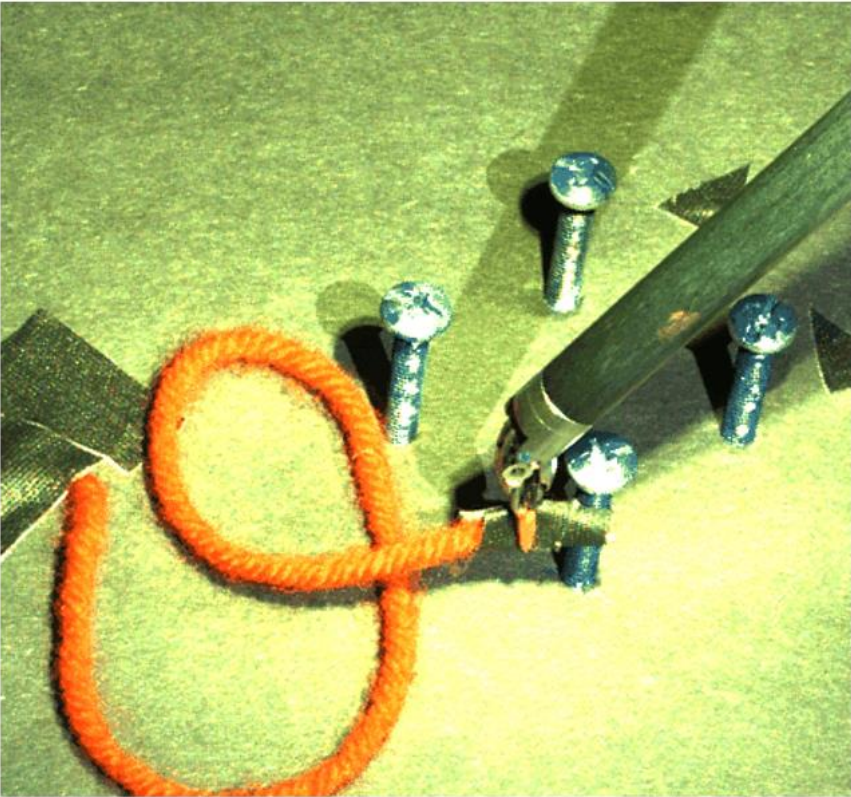
# Human Demonstration



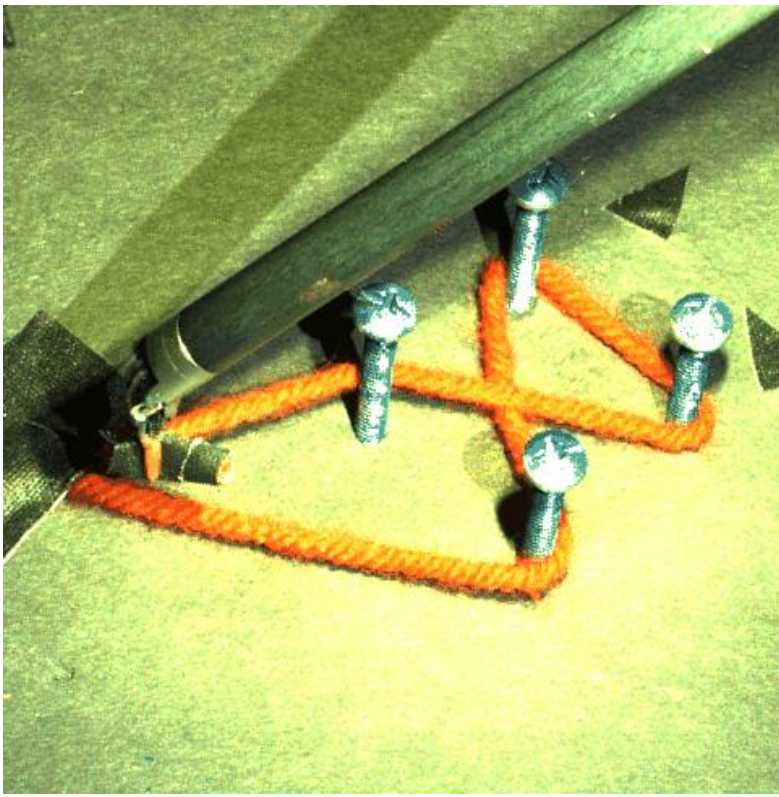
## Behavior Cloning



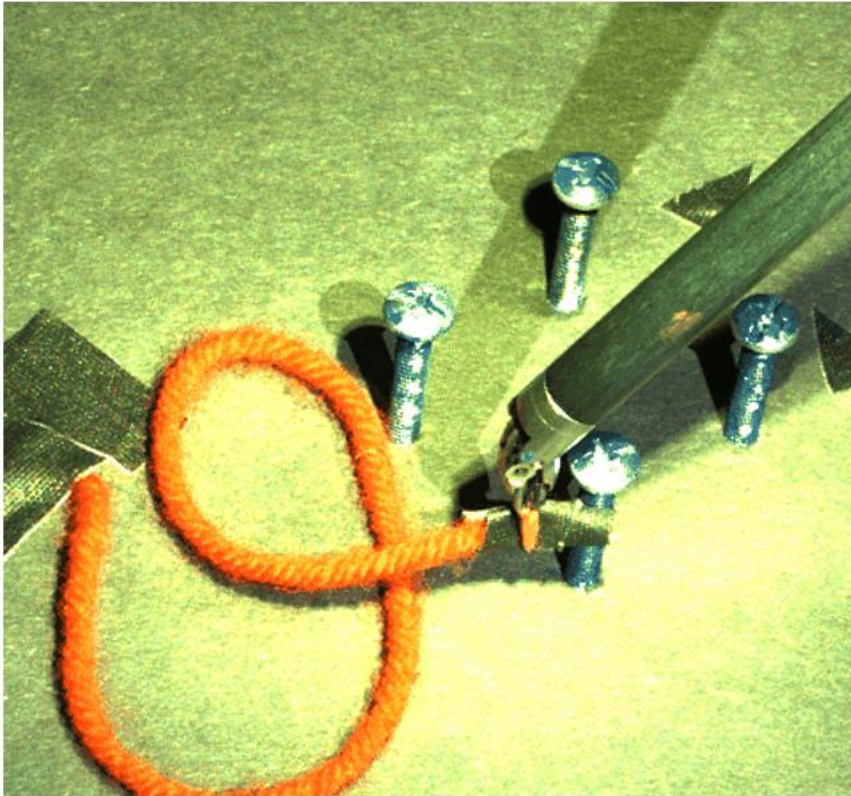
**Behavior Cloning**



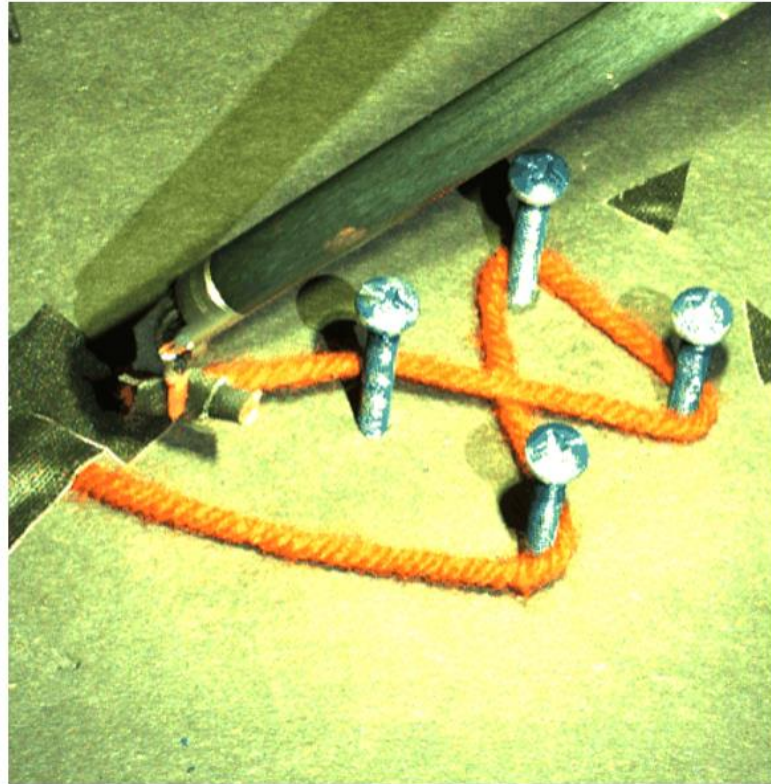
**ThriftyDagger (autonomous)**



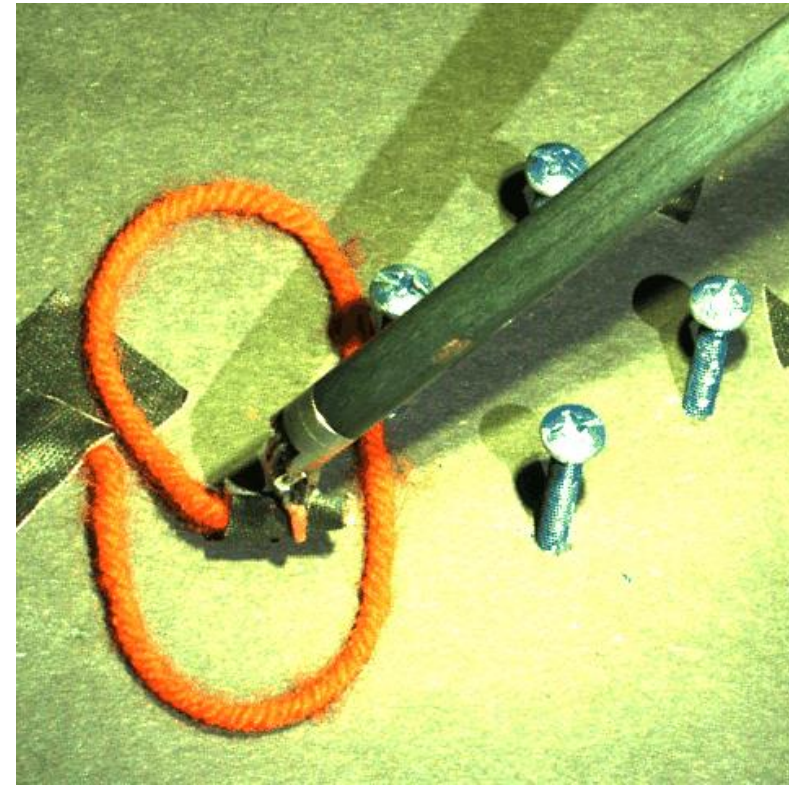
**Behavior Cloning**



**ThriftyDAgger (autonomous)**



**ThriftyDAgger (+human)**





# User Study

N=10 subjects each control 3 robots in simulation.

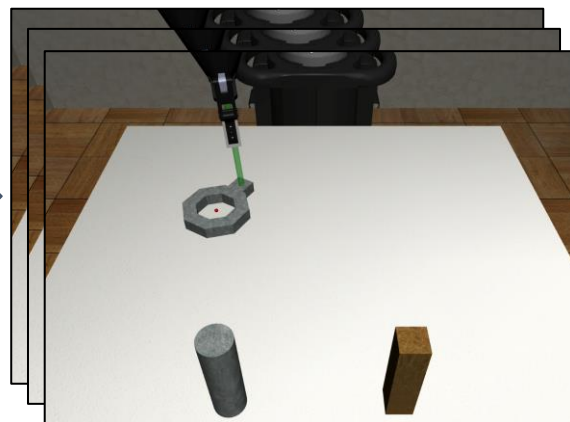
## Robot-Gated

Memory: Non-Match

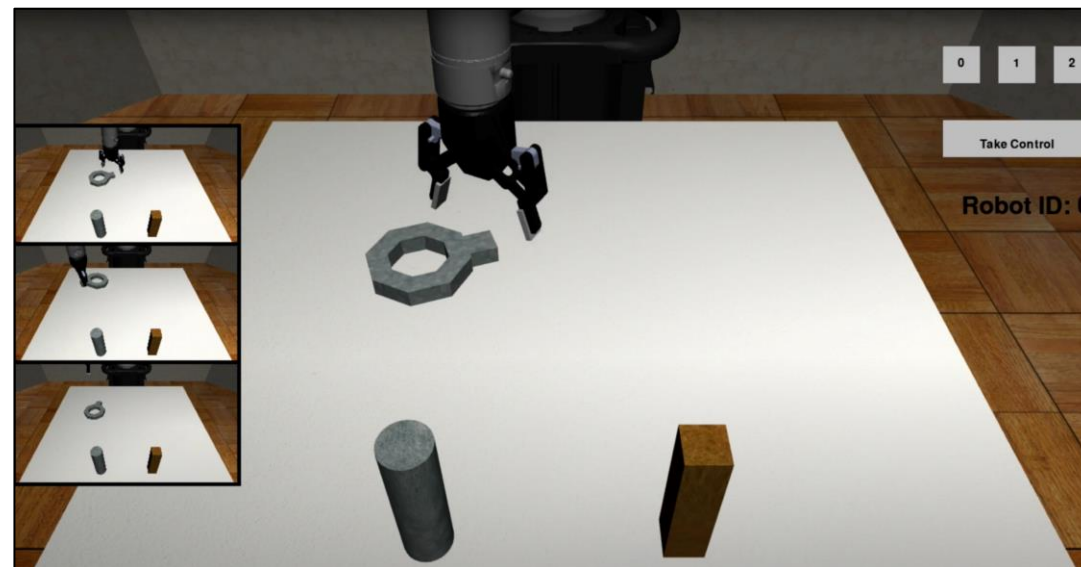
H	H	H	H	H
H	🍌	❤️	H	H
H	H	H	H	H
H	H	H	H	H

Memory: Match

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

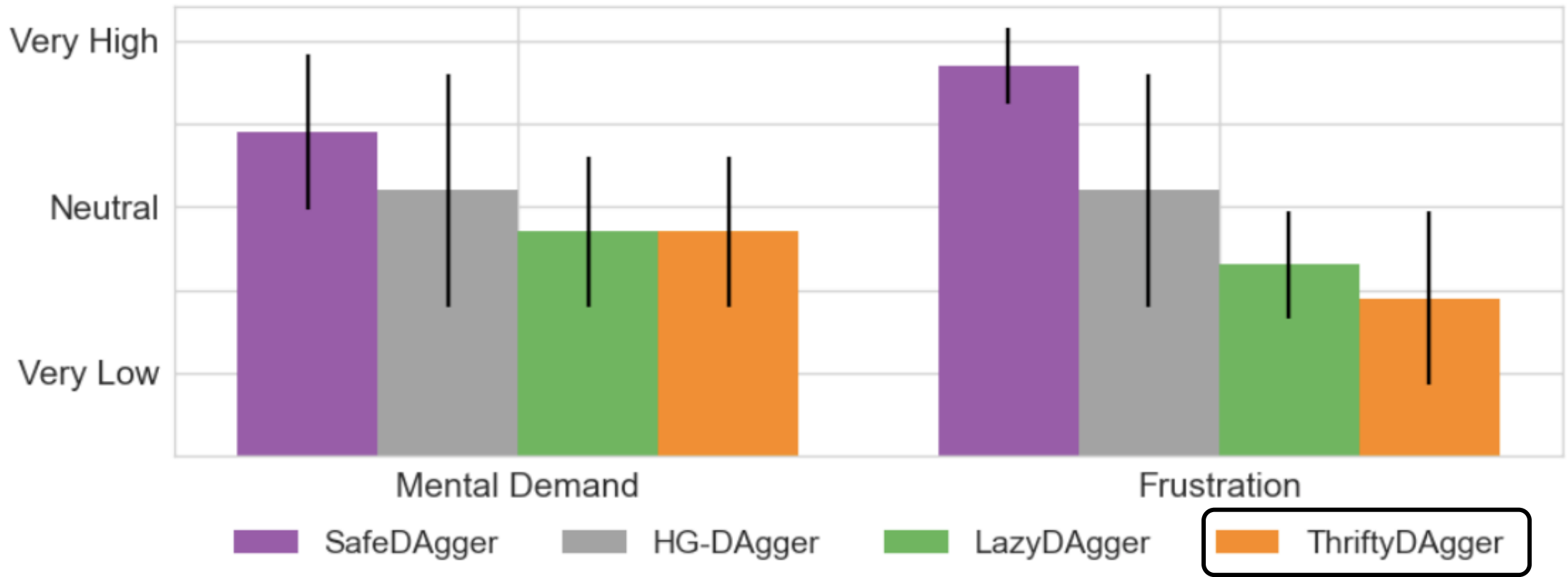


## Human-Gated



# ThriftyDagger Qualitative Results

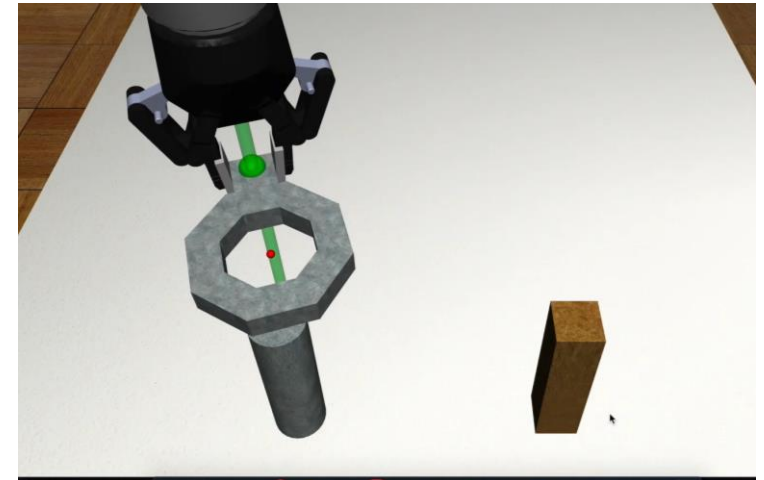
## Survey Responses



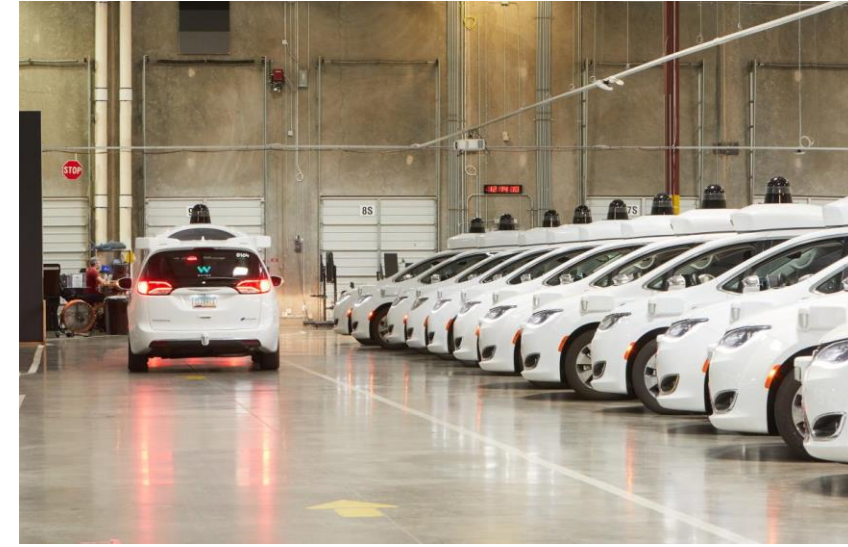
# User Study Quantitative Results

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: Survey of Recent BC methods

- Choose your own adventure reading assignment
  - Implicit Behavior Cloning
  - Action Chunking Transformer
  - Diffusion Policy
- Submit a paragraph before class summarizing at a high-level:
  - What's the problem the authors want to solve?
  - Why is important?
  - What is their proposed solution?
  - What evidence do they give that their solution is good?
  - What is one question you had about the paper?