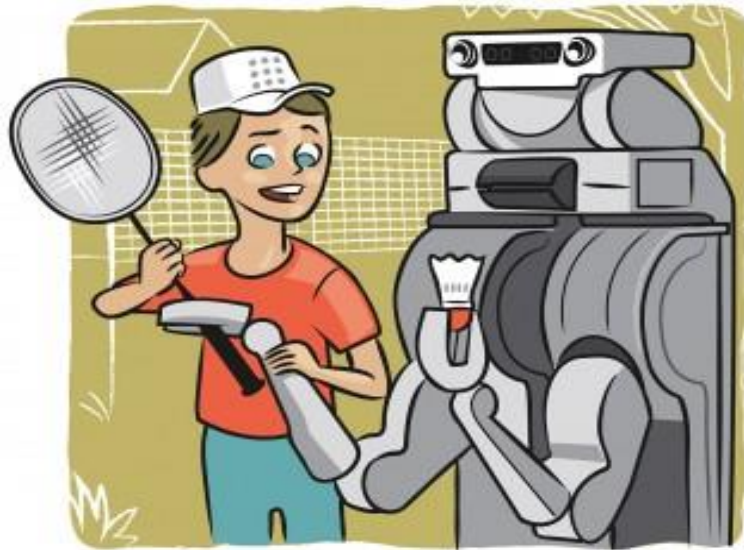


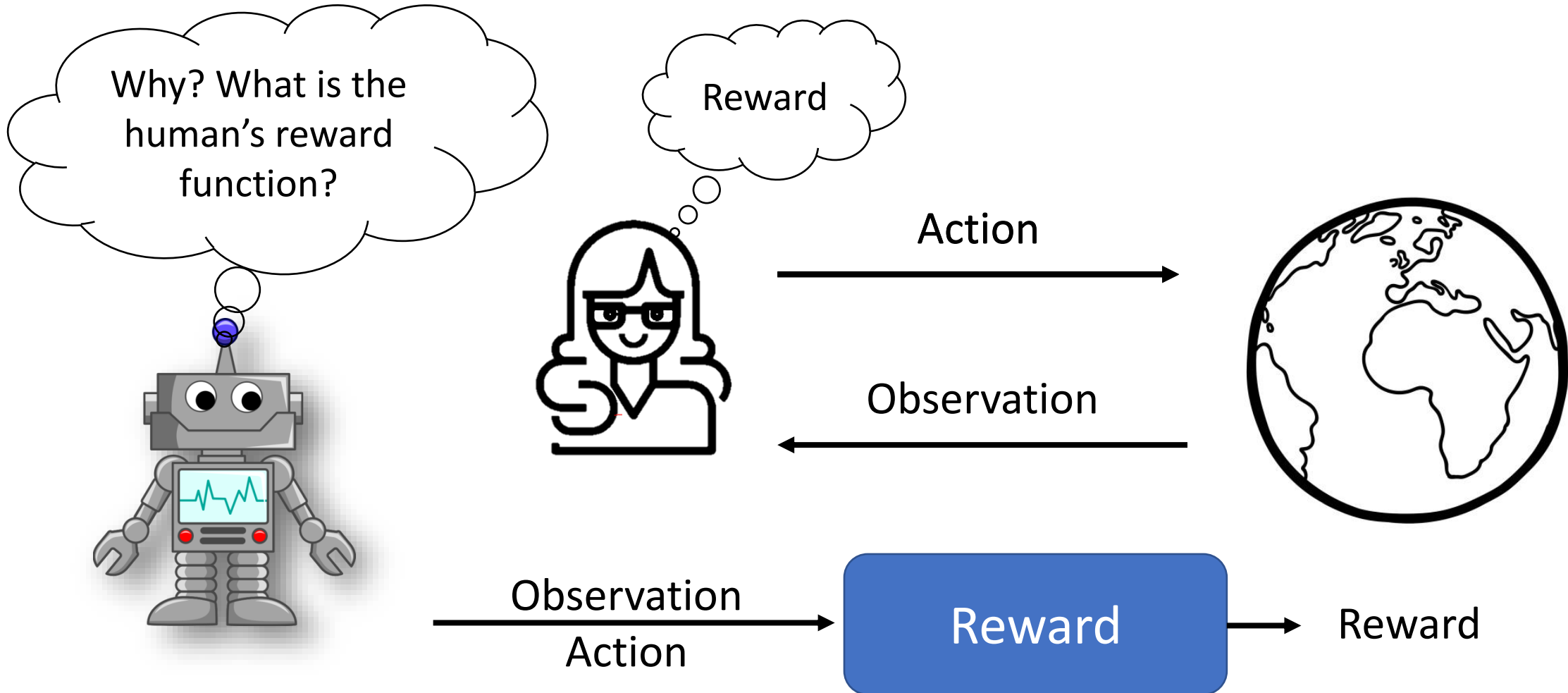
Inverse RL and Reward Learning



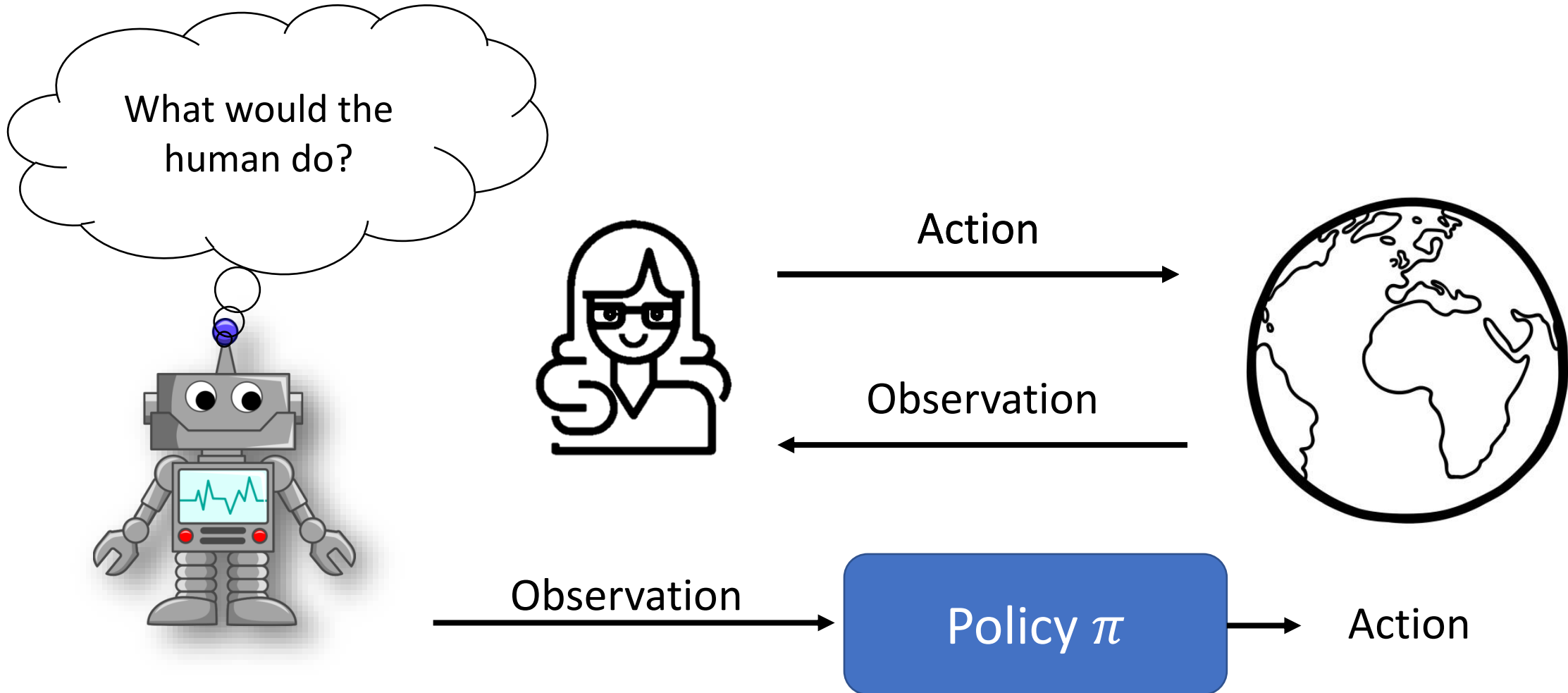
Instructor: Daniel Brown

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

Reward Learning (Inverse Reinforcement Learning)



Why not just imitate behavior? (Behavioral Cloning)





Human Intent Inference

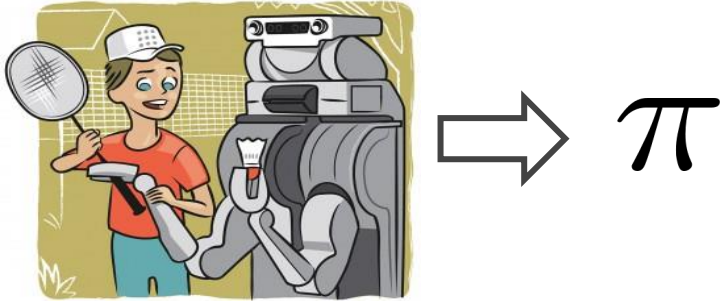


Inverse Reinforcement Learning

- Given
 - MDP without a reward function
 - Demonstrations from an optimal policy π^*
- Recover the reward function R that makes π^* optimal

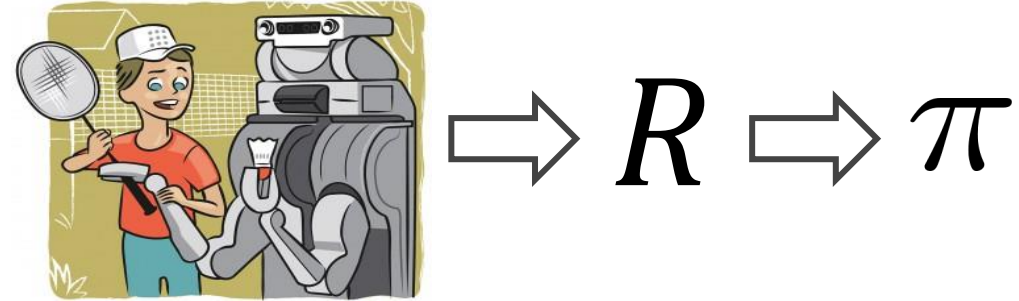
Imitation Learning

Behavioral Cloning



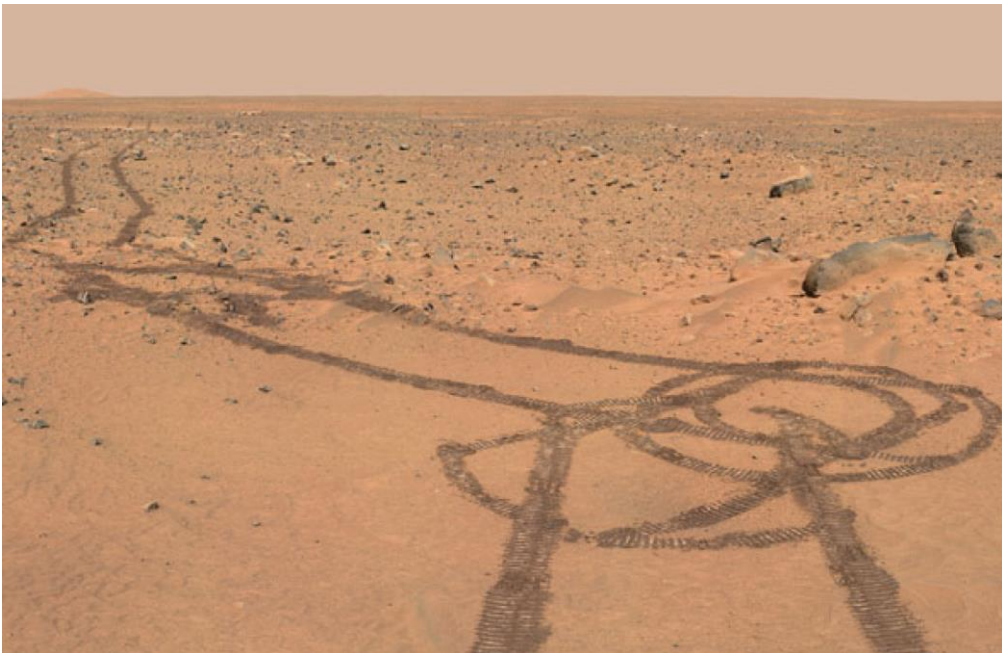
- Answers the “How?” question
- Mimic the demonstrator
- Learn mapping from states to actions
- Computationally efficient
- Compounding errors

Inverse Reinforcement Learning

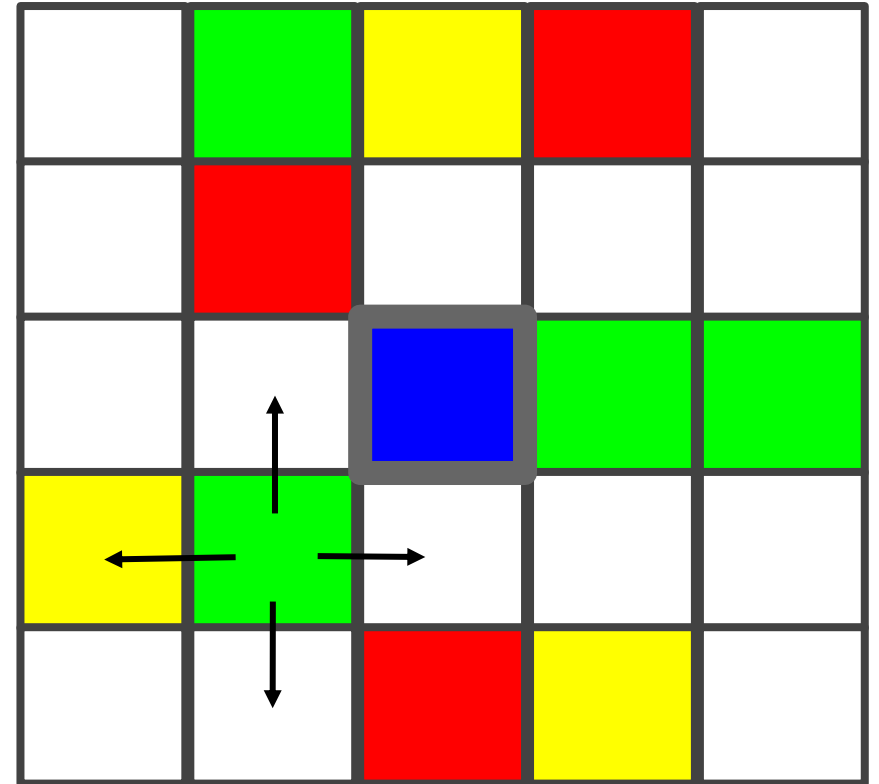


- Answers the “Why?” question
- Explain the demonstrator’s behavior
- Learn a reward function capturing the demonstrator’s intent
- Can require lots of data and compute
- Better generalization. Can recover from arbitrary states

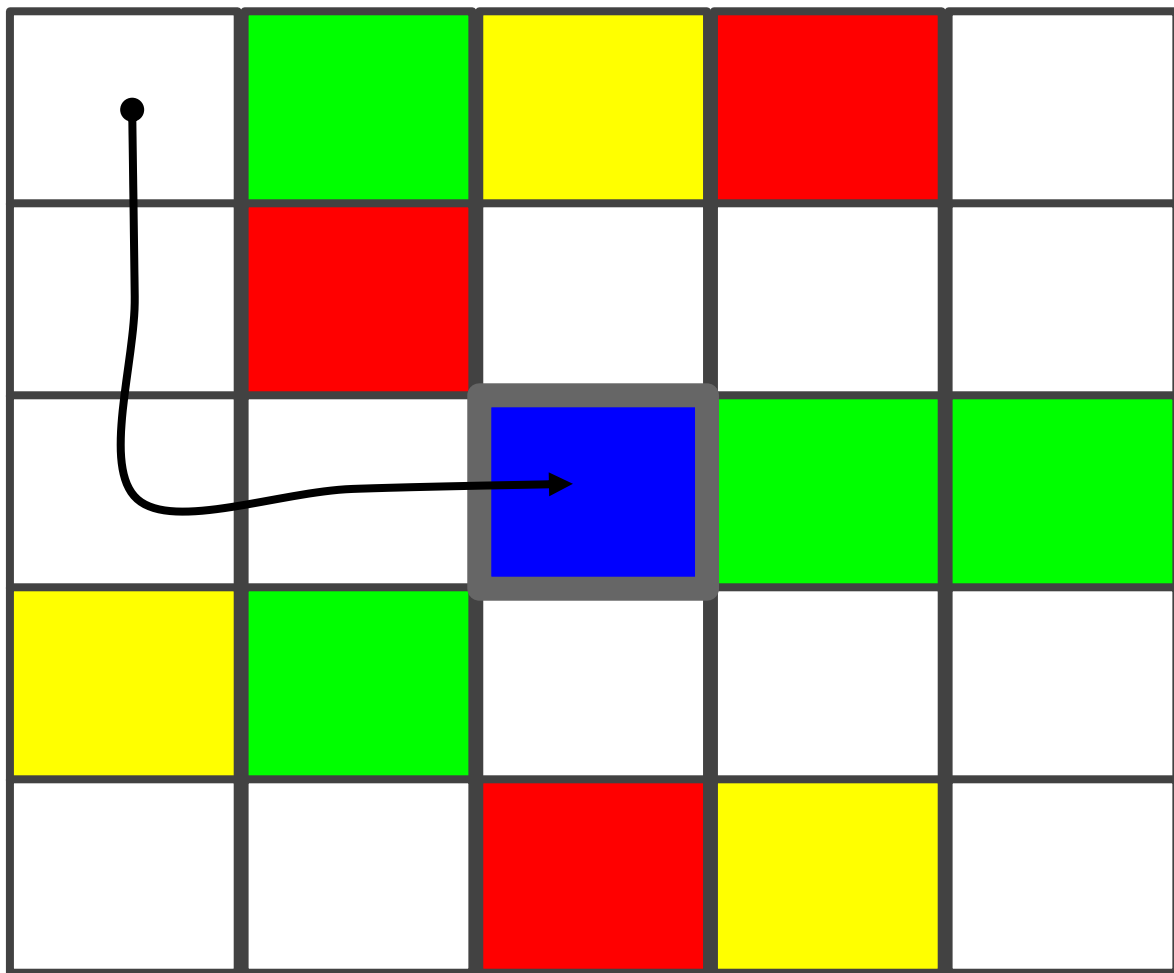
IRL Example: Teaching a robot to navigate through demonstrations



Toy version



What is the reward?



$$R(\text{Blue}) = ?$$

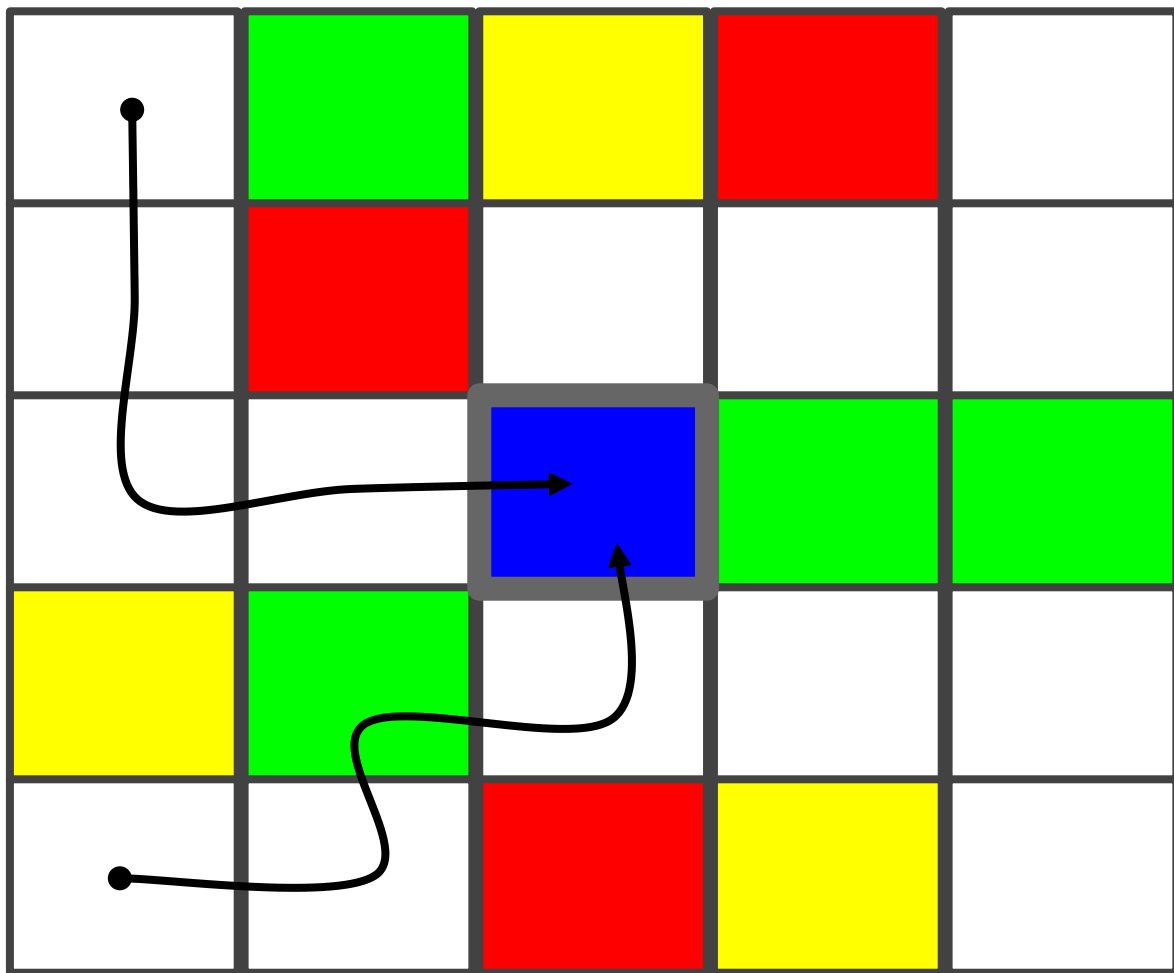
$$R(\text{Yellow}) = ?$$

$$R(\text{Green}) = ?$$

$$R(\text{Red}) = ?$$

$$R(\text{White}) = ?$$

What is the reward?



$$R(\text{Blue}) = ?$$

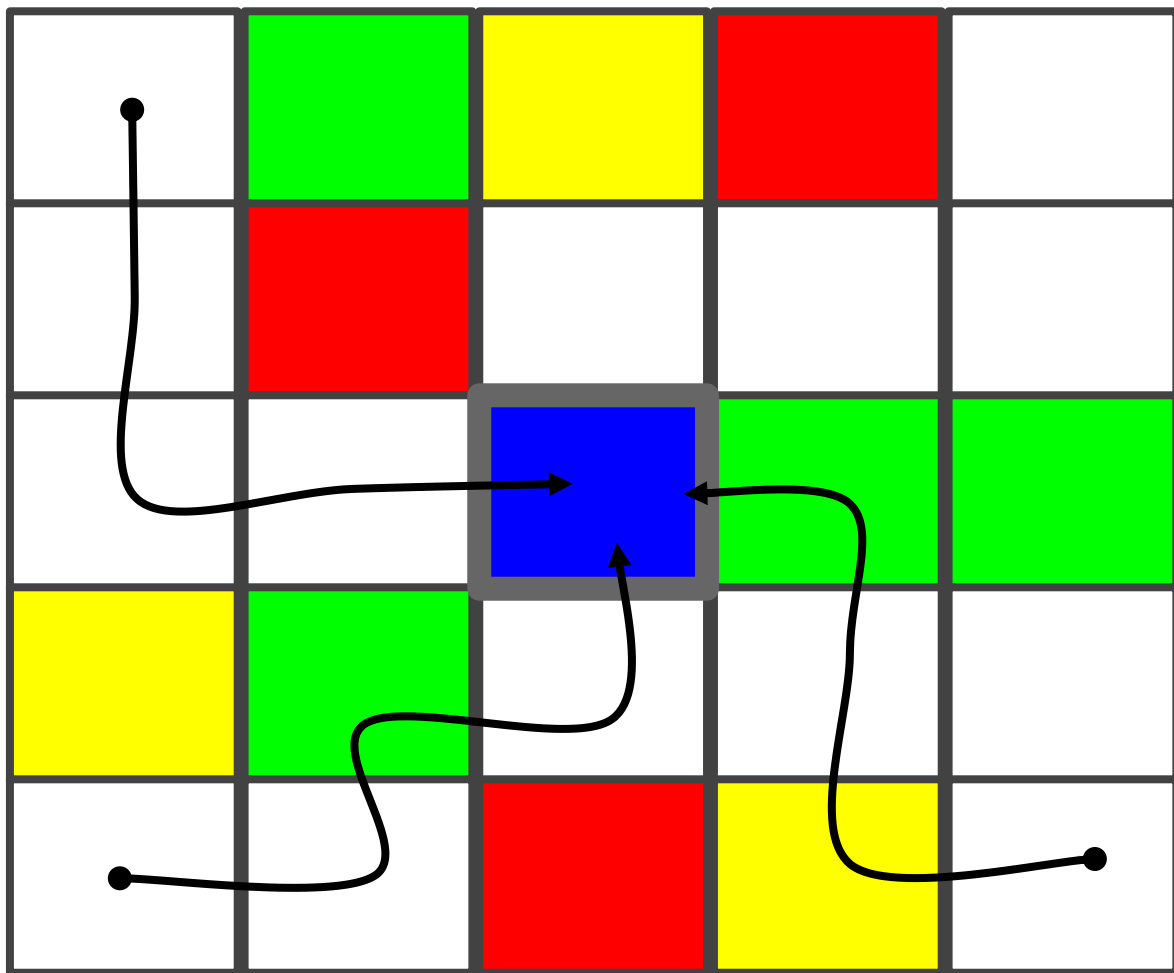
$$R(\text{Yellow}) = ?$$

$$R(\text{Green}) = ?$$

$$R(\text{Red}) = ?$$

$$R(\text{White}) = ?$$

What is the reward?



$$R(\text{Blue}) = ?$$

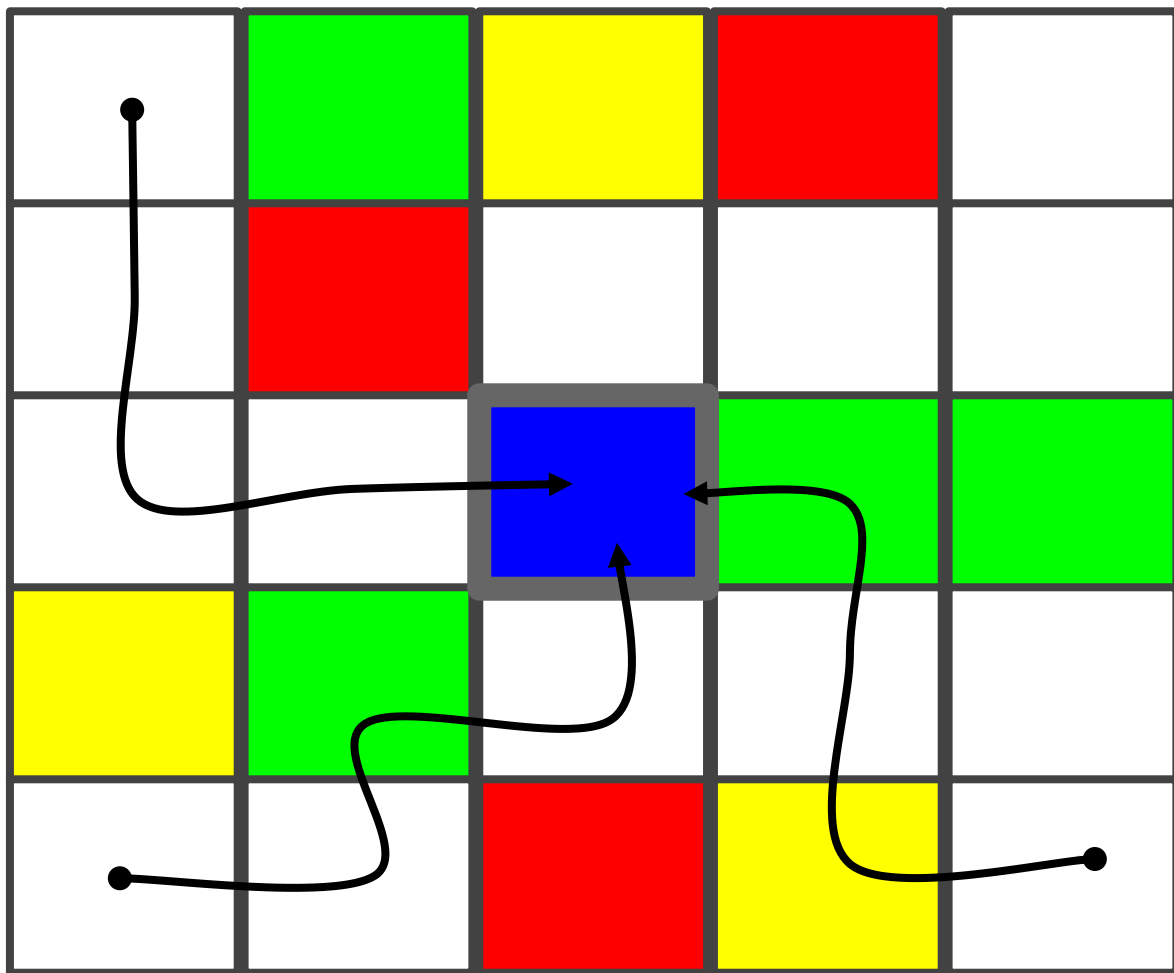
$$R(\text{Yellow}) = ?$$

$$R(\text{Green}) = ?$$

$$R(\text{Red}) = ?$$

$$R(\text{White}) = ?$$

What is the reward?



$$R(\text{Blue Square}) = +1$$

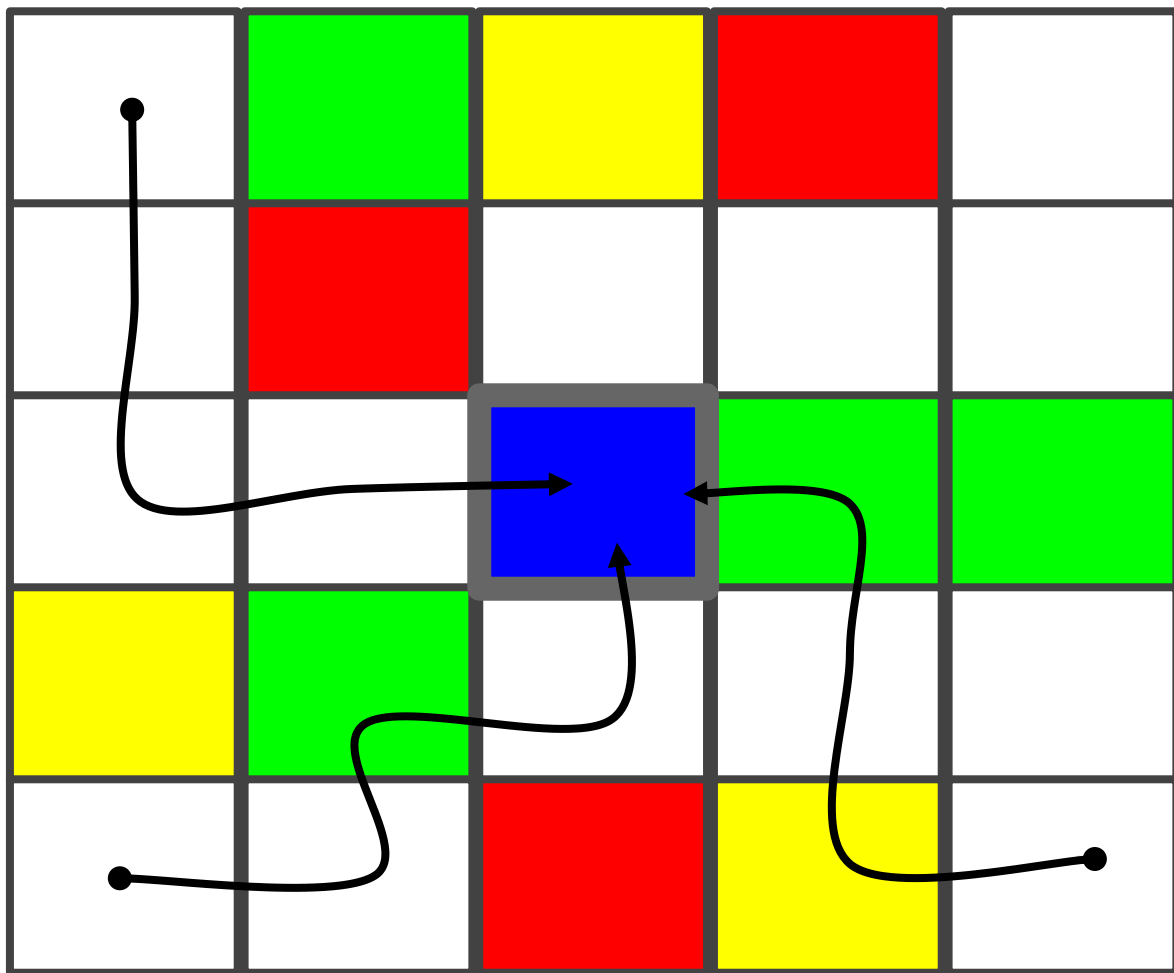
$$R(\text{Yellow Square}) = 0$$

$$R(\text{Green Square}) = 0$$

$$R(\text{Red Square}) = -1$$

$$R(\text{White Square}) = 0$$

What is the reward?



$$R(\text{Blue}) = +10$$

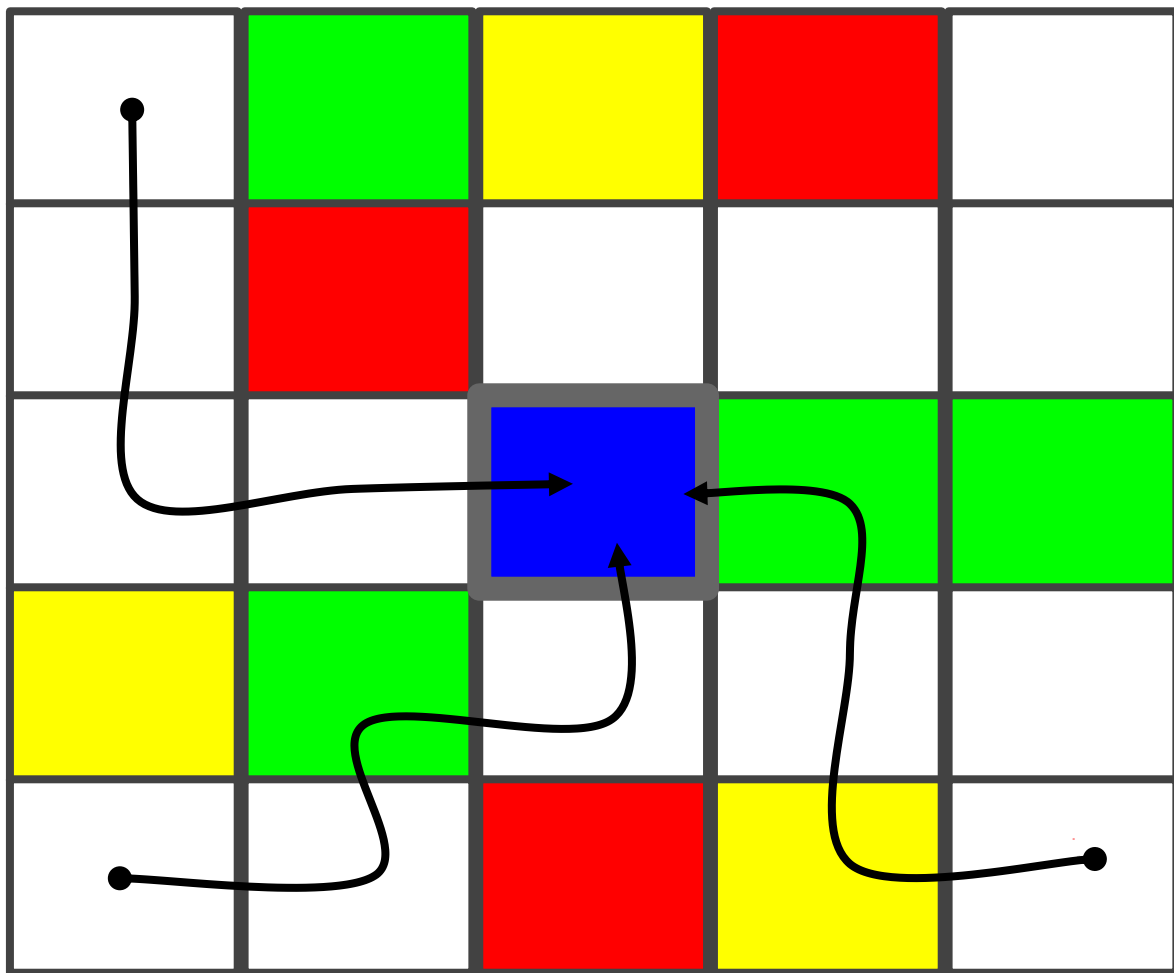
$$R(\text{Yellow}) = 0$$

$$R(\text{Green}) = 0$$

$$R(\text{Red}) = -10$$

$$R(\text{White}) = 0$$

What is the reward?



$$R(\text{Blue}) = +10$$

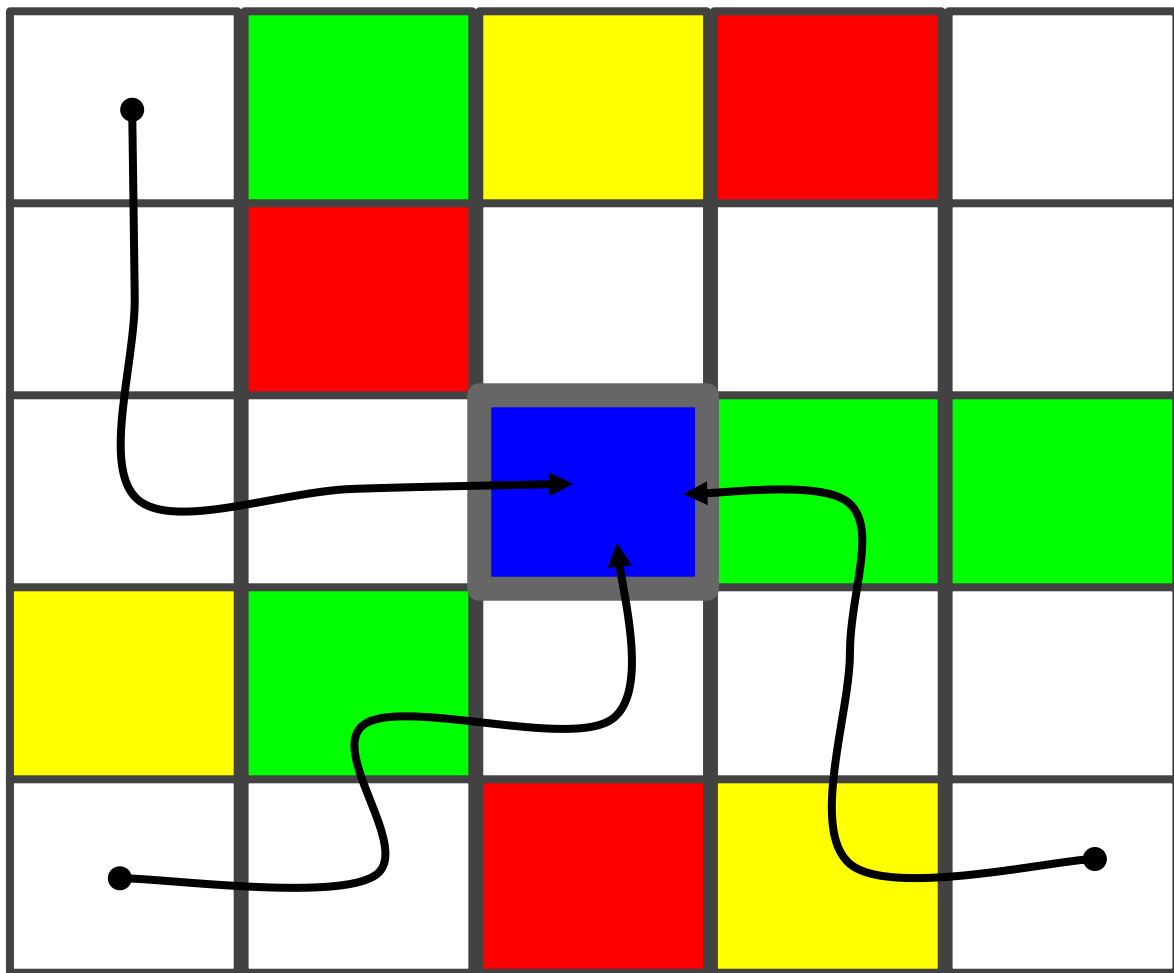
$$R(\text{Yellow}) = -1$$

$$R(\text{Green}) = -1$$

$$R(\text{Red}) = -10$$

$$R(\text{White}) = -1$$

What is the reward?



$$R(\text{Blue}) = 0$$

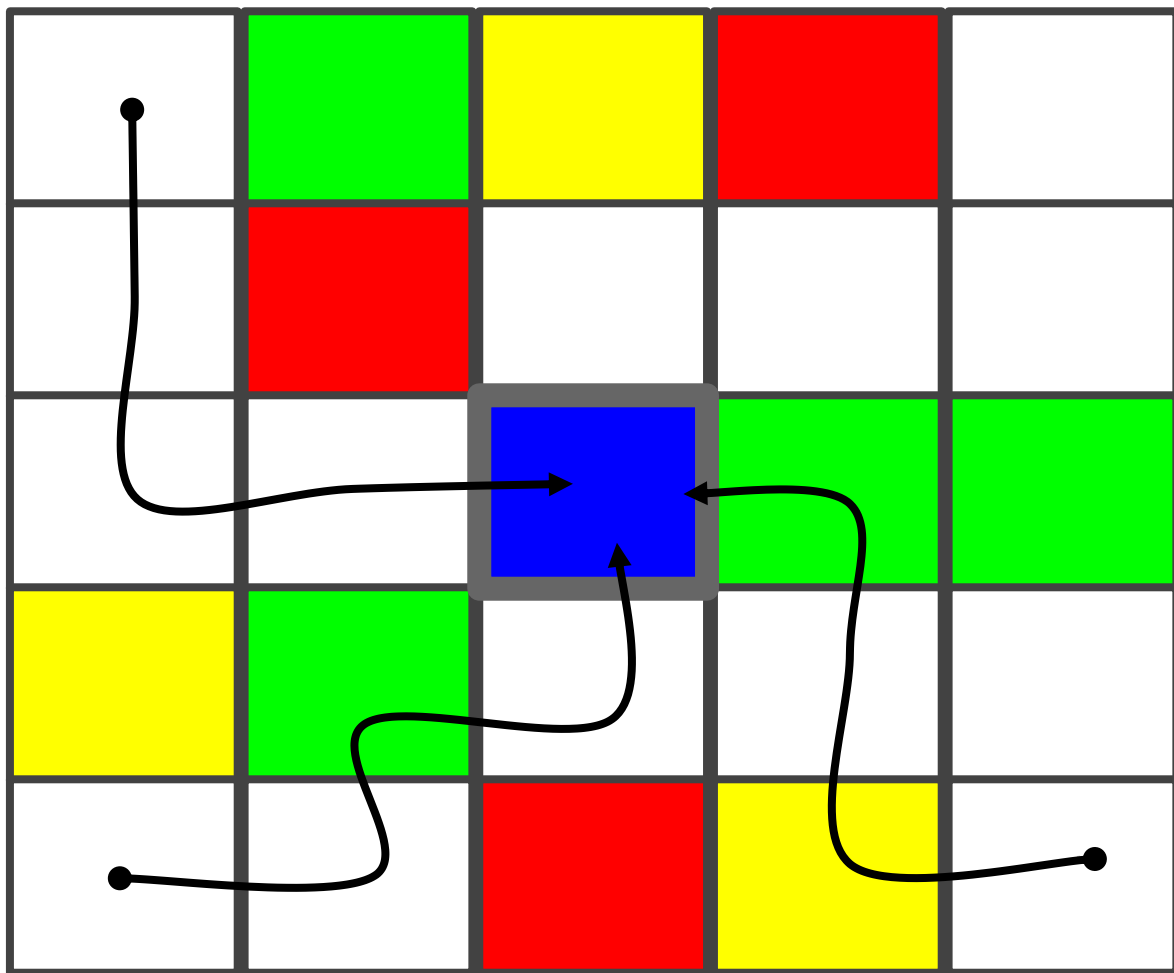
$$R(\text{Yellow}) = 0$$

$$R(\text{Green}) = 0$$

$$R(\text{Red}) = 0$$

$$R(\text{White}) = 0$$

What is the reward?



$$R(\text{Blue}) = C$$

$$R(\text{Yellow}) = C$$

$$R(\text{Green}) = C$$

$$R(\text{Red}) = C$$

$$R(\text{White}) = C$$

Inverse Reinforcement Learning Formalism

- Given
 - MDP without a reward function
 - Demonstrations from an optimal policy π^*
- Recover a reward function R that makes π^* optimal
- **Ill-Posed Problem**
 - Infinite number of reward functions that can make π^* optimal
 - Trivial all zero reward
 - Constant reward
 - $aR + c$ (positive scaling $a > 0$, and affine shifts)

Simpler problem: What if you know the policy?

How would you do this more generally?

Basic IRL Algorithm

- Start with demonstrations, D
- Guess initial reward function R_0
- $\hat{R} = R_0$
- Loop:
 - Solve for optimal policy $\pi_{\hat{R}}^*$
 - Compare D and $\pi_{\hat{R}}^*$
 - Update \hat{R} to try and make D and $\pi_{\hat{R}}^*$ more similar

Feature count matching

- Assume the reward function is a linear combination of features:

$$R(s) = \mathbf{w}^T \phi(s)$$

- Value function becomes linear combination of (discounted) feature expectations:

$$V_R^\pi = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

Feature count matching

- Assume the reward function is a linear combination of features:

$$R(s) = \mathbf{w}^T \phi(s)$$

- Value function becomes linear combination of (discounted) feature expectations:

$$V_R^\pi = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{w}^t \phi(s_t) \right]$$

Feature count matching

- Assume the reward function is a linear combination of features:

$$R(s) = \mathbf{w}^T \phi(s)$$

- Value function becomes linear combination of (discounted) feature expectations:

$$V_R^\pi = \mathbf{w}^T \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_t) \right] = \mathbf{w}^T \mu_\pi$$

Inverse reinforcement learning: feature matching

(Abbeel and Ng 2004, Syed and Schapire 2007)

- If $\|\mathbf{w}\|_1 \leq 1$, then

$$|x^\top y| \leq \|x\|_1 \|y\|_\infty$$

$$\begin{aligned} V_R^{\pi^*} - V_R^{\pi_{\text{robot}}} &= \mathbf{w}^\top (\mu_{\pi^*} - \mu_{\pi_{\text{robot}}}) \\ &\leq \|\mu_{\pi^*} - \mu_{\pi_{\text{robot}}}\|_\infty \end{aligned}$$

- If feature expectations match, then expected returns are identical.
- Idea: Can we update the reward guess \hat{R} so the feature counts get closer?

Problem: Many different policies can lead to same expected feature counts

Maximum Entropy IRL (Ziebart et al. 2008)

$$P(\tau) = \frac{e^{R_{\mathbf{w}}(\tau)}}{Z}$$

- Collect M demonstrations $D = \{\tau_1, \dots, \tau_M\}$
- Initialize reward weights \mathbf{w}

$$R(s) = \mathbf{w}^T \phi(s)$$

- **Loop**

- Solve for (soft) optimal policy $\pi(a|s)$ via Value Iteration
- Solve for expected feature counts of $\pi(a|s)$
- Compute weight update $\mathbf{w} \leftarrow \mathbf{w} + \alpha(\mu_D - \mu_\pi)$

Soft Value Iteration

$$\pi_{\Theta}(A_t|S_t) = e^{Q_{\pi_{\Theta}}^{\text{soft}}(A_t, S_t) - V_{\pi_{\Theta}}^{\text{soft}}(S_t)}$$

Policy is a softmax policy.

$$V_{\pi_{\Theta}}^{\text{soft}}(S_t) = \log \sum_{A_t \in \mathcal{A}} e^{Q_{\pi_{\Theta}}^{\text{soft}}(A_t, S_t)}$$

↙
Soft Maximum

Softmax is a Soft Maximum

- Assume $b > a$
- $\log(e^a + e^b) =$

- If $a = b$
- $\log(e^a + e^b) =$

- In general $\max\{x_1, x_2, \dots, x_n\} \leq \log \sum_i \exp(x_i) \leq \max\{x_1, \dots, x_n\} + \log n$

Soft Value Iteration

$$\pi_{\Theta}(A_t|S_t) = e^{Q_{\pi_{\Theta}}^{\text{soft}}(A_t, S_t) - V_{\pi_{\Theta}}^{\text{soft}}(S_t)}$$

Soft Maximum

$$V_{\pi_{\Theta}}^{\text{soft}}(S_t) = \log \sum_{A_t \in \mathcal{A}} e^{Q_{\pi_{\Theta}}^{\text{soft}}(A_t, S_t)}$$


$$Q_{\pi_{\Theta}}^{\text{soft}}(A_t, S_t) = R_{\Theta}(S_t, A_t) + \sum_{S' \in \mathcal{S}} P_T(S'|A_t, S_t) V_{\pi_{\Theta}}^{\text{soft}}(S')$$

- Initialize values
- Repeat until convergence:
 - Solve for Q
 - Solve for V

Watch This: Scalable Cost-Function Learning for Path Planning in Urban Environments

Markus Wulfmeier¹, Dominic Zeng Wang¹ and Ingmar Posner¹

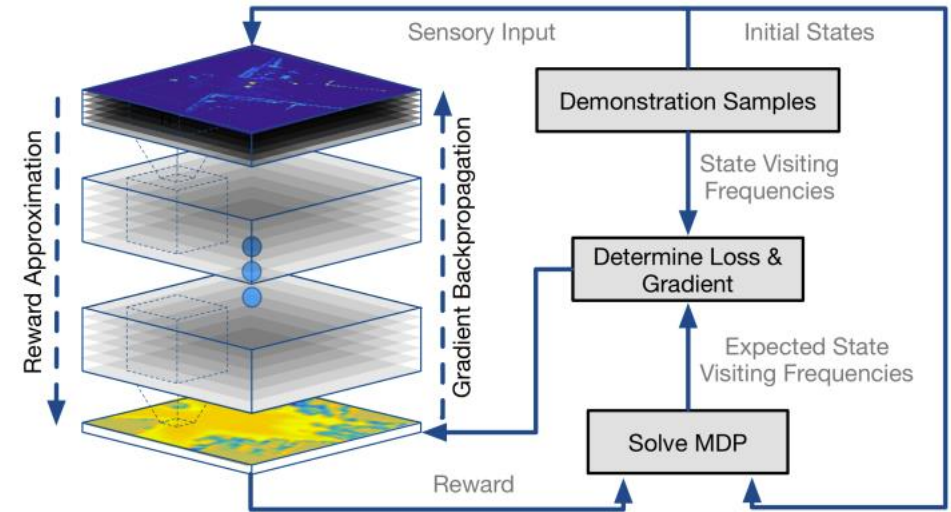
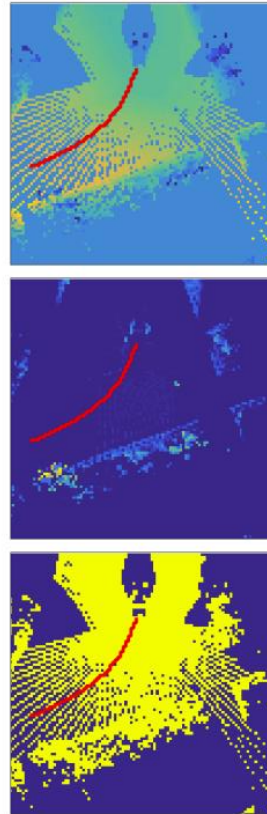


Fig. 1: Schema for training neural networks in the Maximum Entropy paradigm for IRL.

Another way to look at MaxEnt IRL

$$P(\tau) = \frac{e^{R_w(\tau)}}{Z} \quad Z = \int e^{R_w(\tau)} d\tau$$

- Maximum Likelihood Estimation
- Find reward function that maximizes the log likelihood of the demonstration trajectories:

$$\max_{\theta} \frac{1}{N} \sum_{\tau \in D} R_w(\tau) - \log Z$$

How to avoid fully solving MDP

$$\max_{\theta} \frac{1}{N} \sum_{\tau \in D} R_w(\tau) - \log Z \quad Z = \int e^{R_w(\tau)} d\tau$$

- Estimate Z with a finite set of trajectories Z_τ .
- Loop:
 - Update parameters w so demonstrations have higher reward than trajectories in Z_τ .
 - Update Z_τ

How to make this more tractable

Relative Entropy Inverse Reinforcement Learning

Abdeslam Boularias

Jens Kober

Jan Peters

Max-Planck Institute for Intelligent Systems

72076 Tübingen, Germany

{abdeslam.boularias,jens.kober,jan.peters}@tuebingen.mpg.de

Learning Objective Functions for Manipulation

Mrinal Kalakrishnan*, Peter Pastor*, Ludovic Righetti*[†], and Stefan Schaal*[†]

kalakris@usc.edu, pastorsa@usc.edu, ludovic.righetti@a3.epfl.ch, sschaal@usc.edu

*CLMC Lab, University of Southern California, Los Angeles CA 90089

[†]Max Planck Institute for Intelligent Systems, Tübingen, Germany 72076

Guided Cost Learning: Deep Inverse Optimal Control via Policy Optimization

Chelsea Finn
Sergey Levine
Pieter Abbeel

CBFINN@EECS.BERKELEY.EDU
SVLEVINE@EECS.BERKELEY.EDU
PABBEEL@EECS.BERKELEY.EDU

University of California, Berkeley, Berkeley, CA 94709 USA

$$P(\tau) = \frac{e^{R_w(\tau)}}{Z}$$

Uniform sampling to approximate Z.

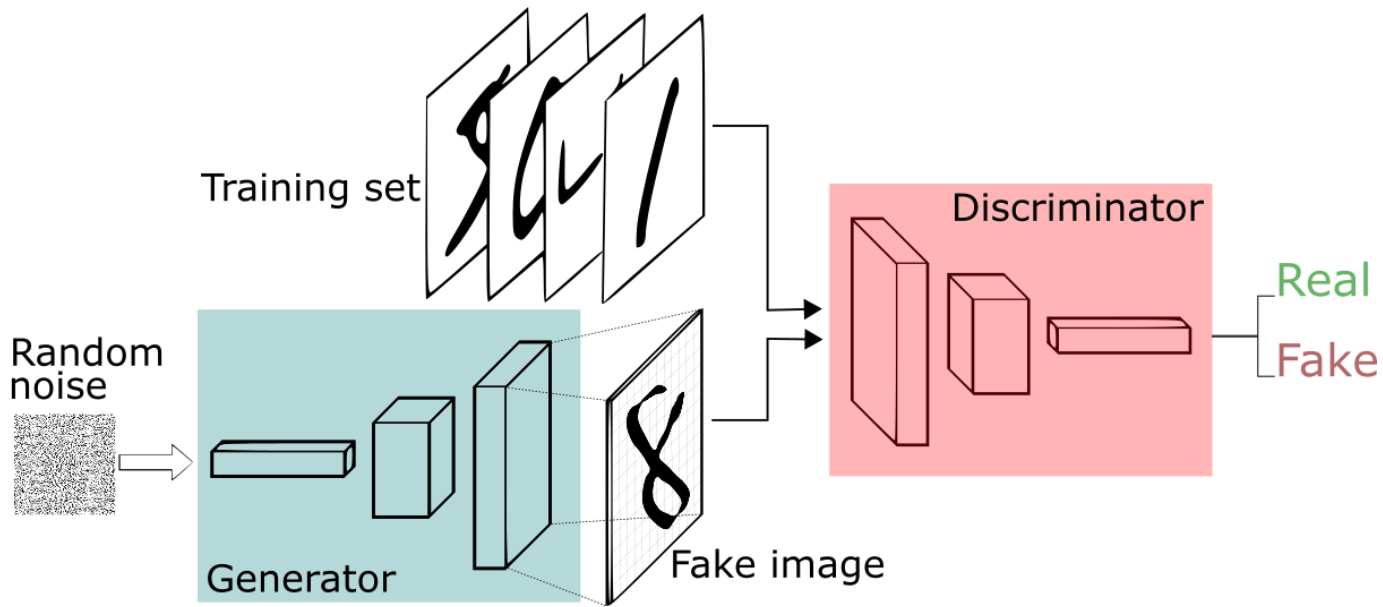
Noisy perturbations of demonstrations to approximate Z

Use current policy to approximate Z.
Alternate between a few steps of reward updates and a few steps of policy updates.

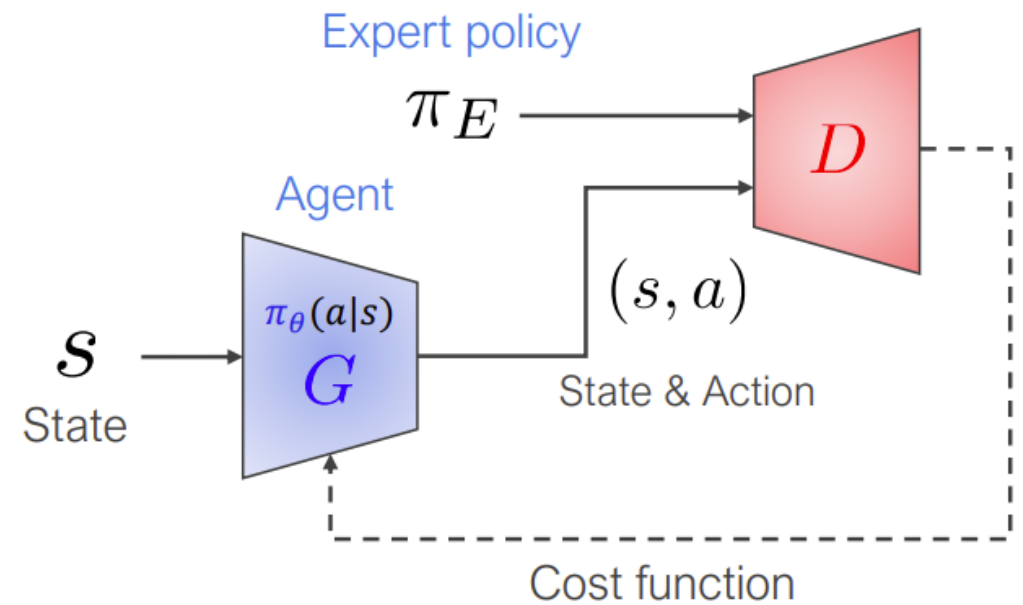
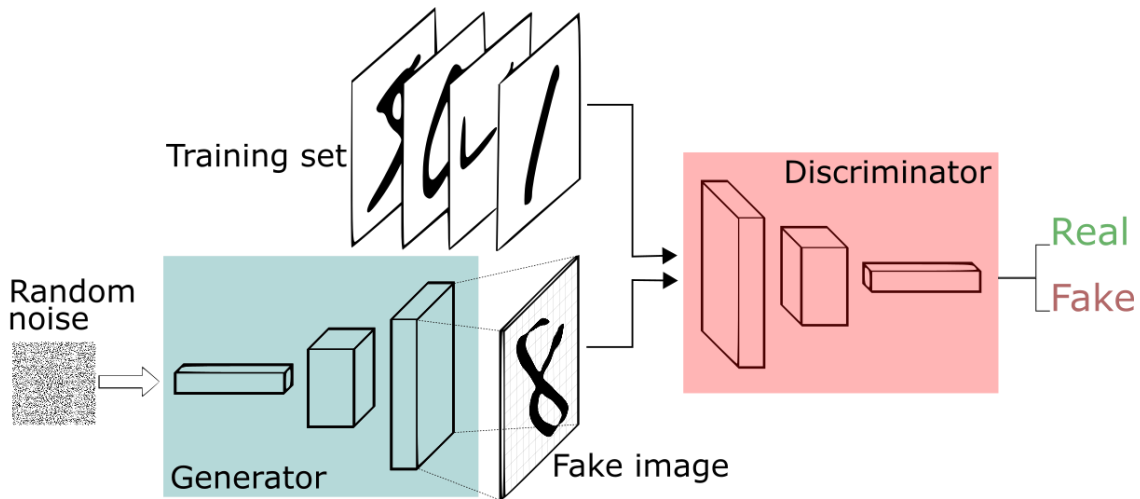
Finn et al. "Guided Cost Learning." 2016



GANs (Generative Adversarial Networks)



GAIL (Generative Adversarial Imitation Learning)



What if we don't want just a single reward estimate?

- Can we get a samples from the full Bayesian posterior?

$$P(R|D) \propto P(D|R)P(R)$$

Bayesian Inverse Reinforcement Learning

(Ramachandran and Amir 2007)

- Assume demonstrator is Boltzmann rational
 - Demonstrator follows a softmax policy with inverse temperature c

$$P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q^*(s,a,R)}}{\sum_{b \in A} e^{\beta Q^*(s,b,R)}}$$

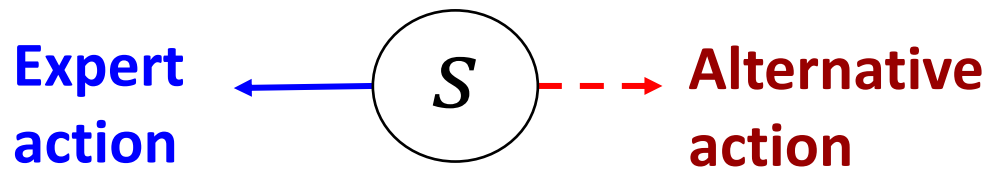
$Q^*(s, a, R) =$ How much reward will I expect to see if I take action a in state s and act optimally thereafter.

Bayesian Inverse Reinforcement Learning

(Ramachandran and Amir 2007)

- Assume demonstrator is Boltzmann rational
 - Demonstrator follows a softmax policy with inverse temperature β

$$P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q^*(s,a,R)}}{\sum_{b \in A} e^{\beta Q^*(s,b,R)}}$$



$$P((s, \leftarrow) | R) = \frac{e^{Q^*(s, \leftarrow, R)}}{e^{Q^*(s, \leftarrow, R)} + e^{Q^*(s, \dashrightarrow, R)}}$$

Bayesian Inverse Reinforcement Learning

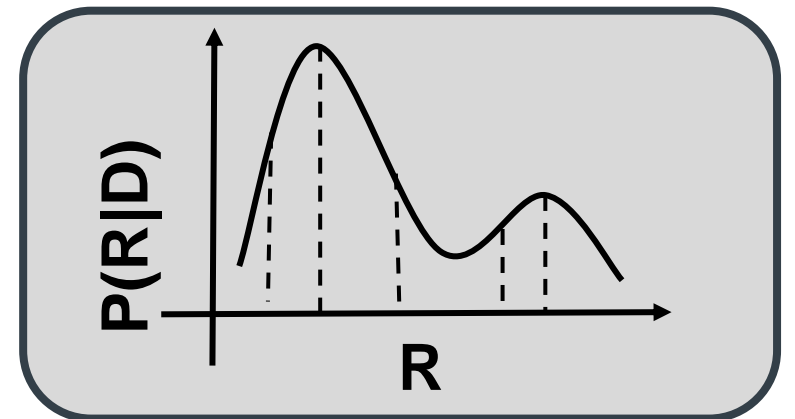
(Ramachandran and Amir 2007)

- Assume demonstrator is Boltzmann rational
 - Demonstrator follows a softmax policy with inverse temperature β

$$P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q^*(s,a,R)}}{\sum_{b \in A} e^{c Q^*(s,b,R)}}$$

- Perform Bayesian inference (MCMC) to sample from posterior distribution

$$P(R|D) \propto P(D|R)P(R)$$



Applications of Bayesian IRL

- Active Learning
- Uncertainty Estimation
- Demonstration Sufficiency



Autonomous Assessment of Demonstration Sufficiency via Bayesian Inverse Reinforcement Learning

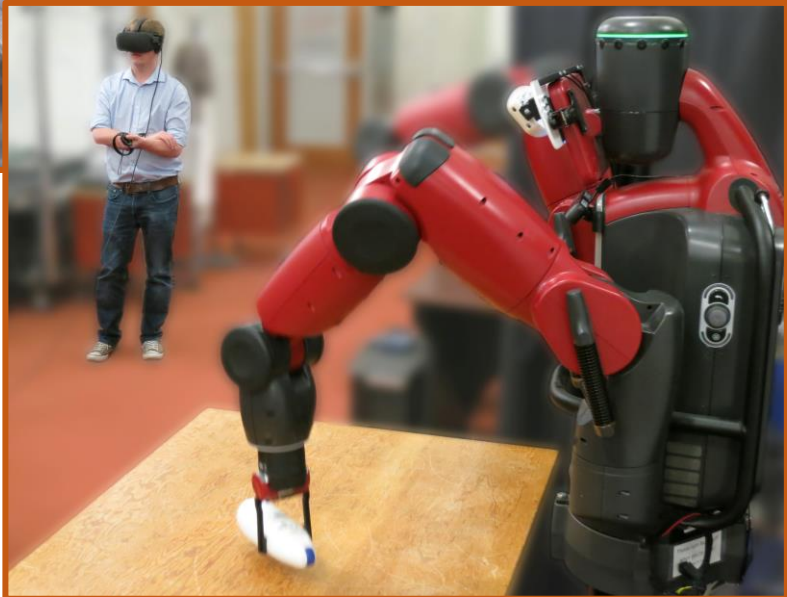
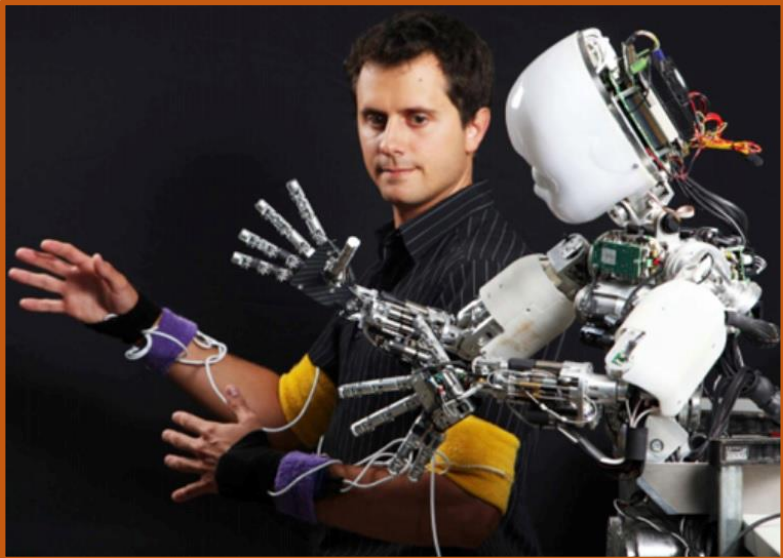
Tu (Alina) Trinh
University of
California, Berkeley

Haoyu Chen
University of Utah

Daniel S. Brown
University of Utah

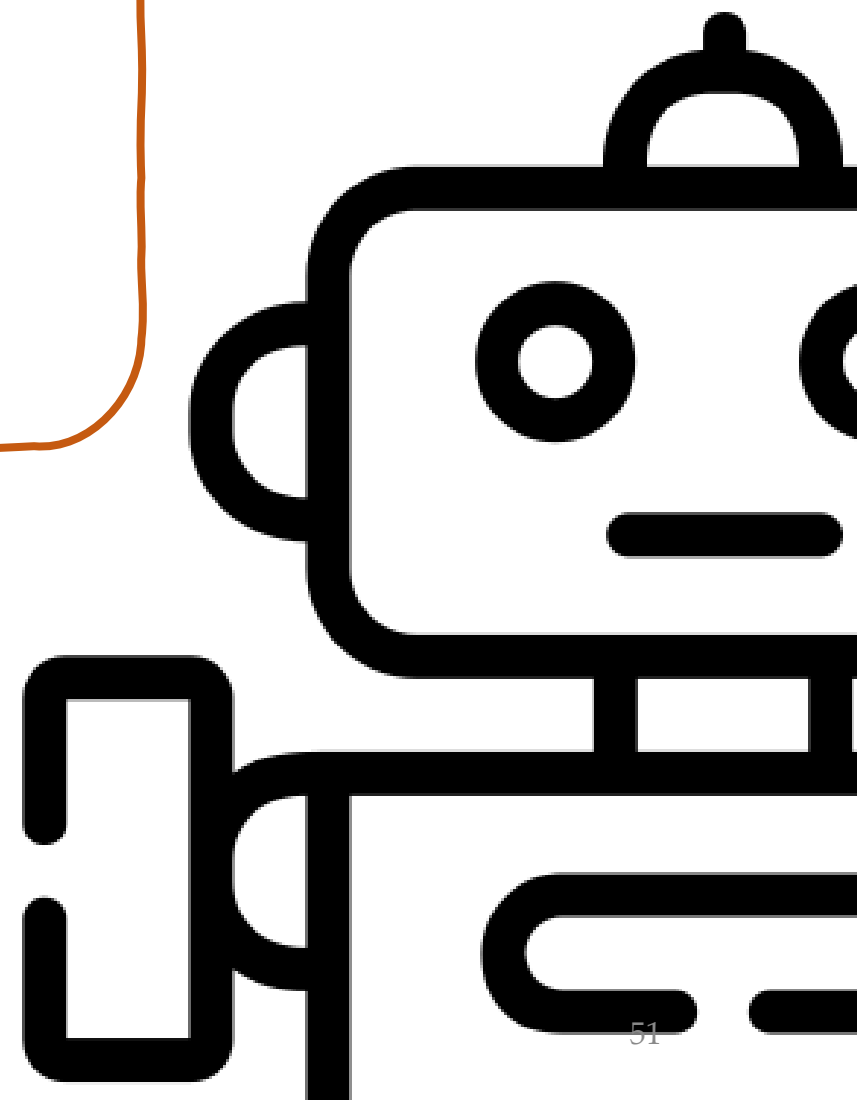


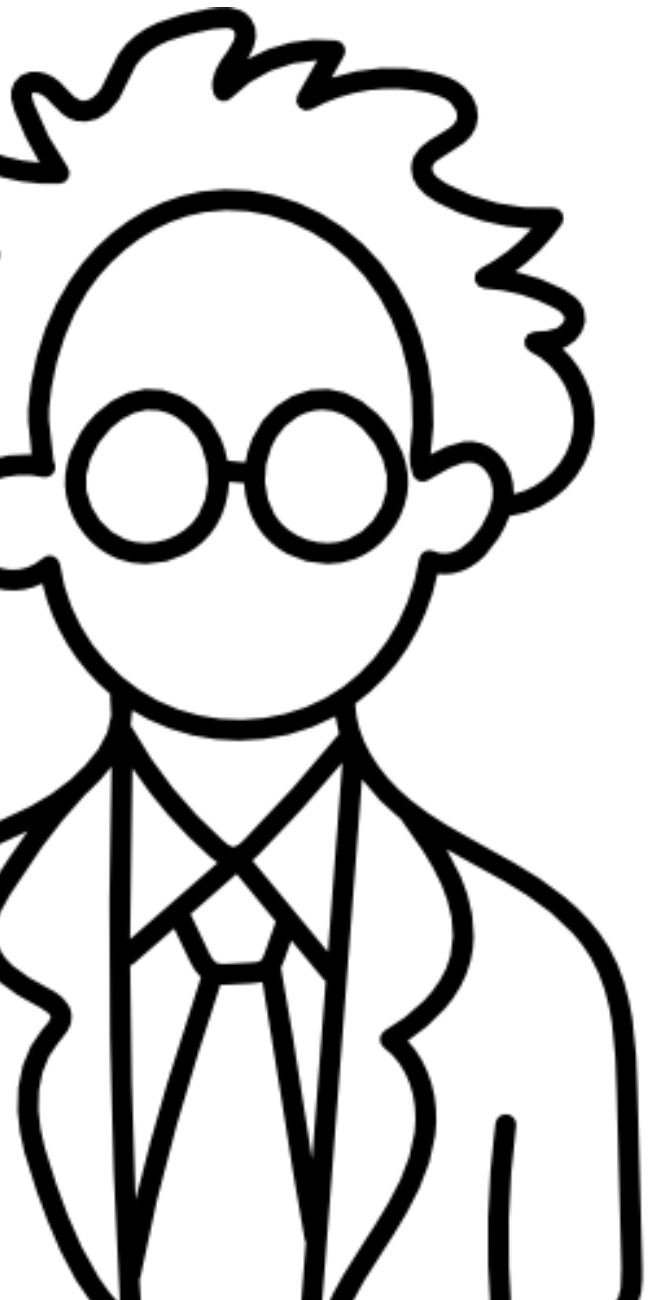
Learning From Demonstration (LfD)



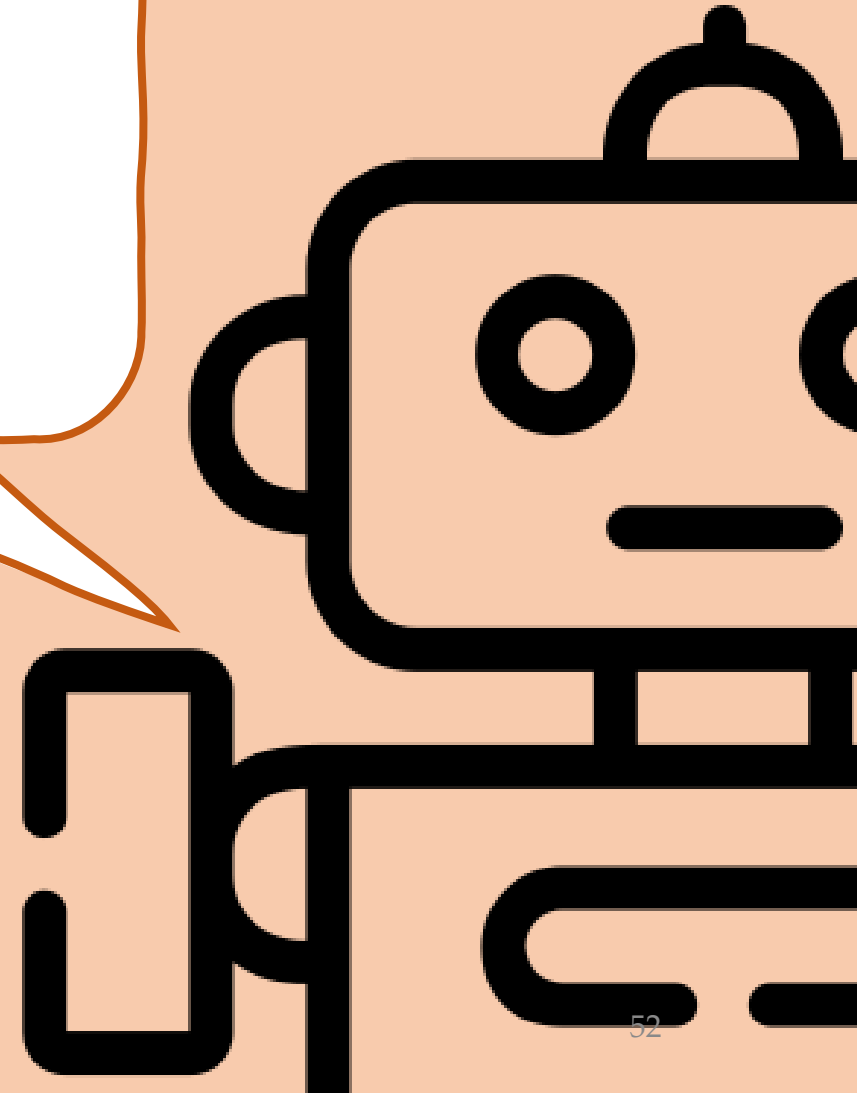


- Have I provided enough demonstrations?
- Are my demonstrations informative enough?
- Should I just supervise the robot?





- Have I received enough demonstrations?
- Are these demonstrations informative enough?

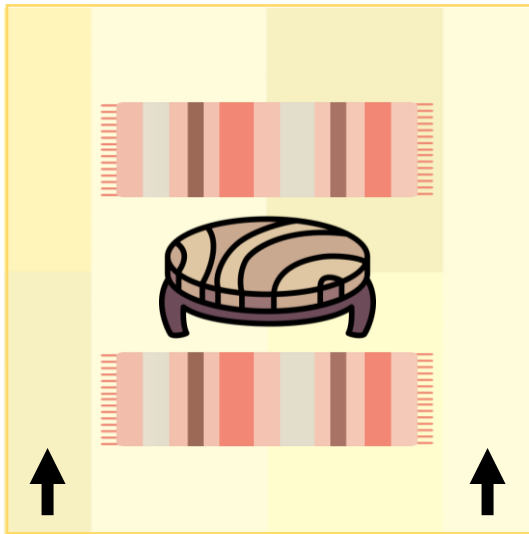


Determining Demonstration Sufficiency

Demonstration Insufficiency

Determining Demonstration Sufficiency

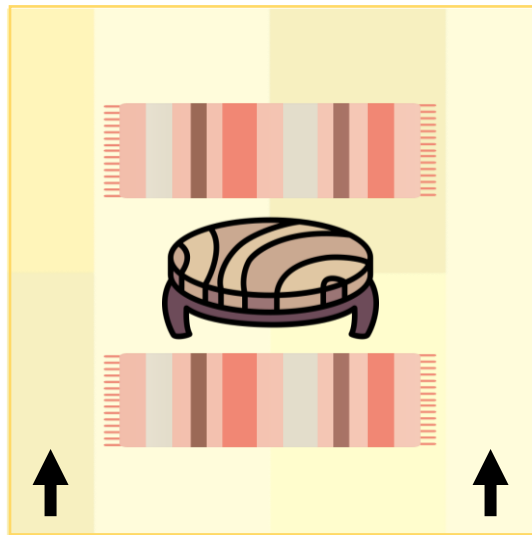
Demonstration Insufficiency



Uninformative Demos

Determining Demonstration Sufficiency

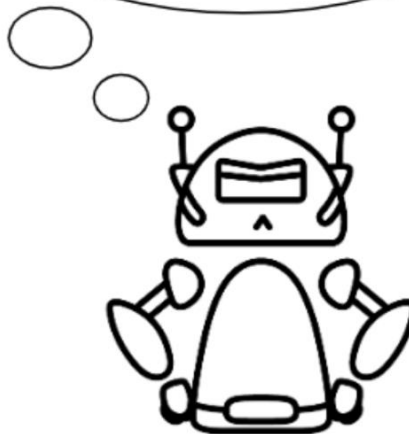
Demonstration Insufficiency



Uninformative Demos

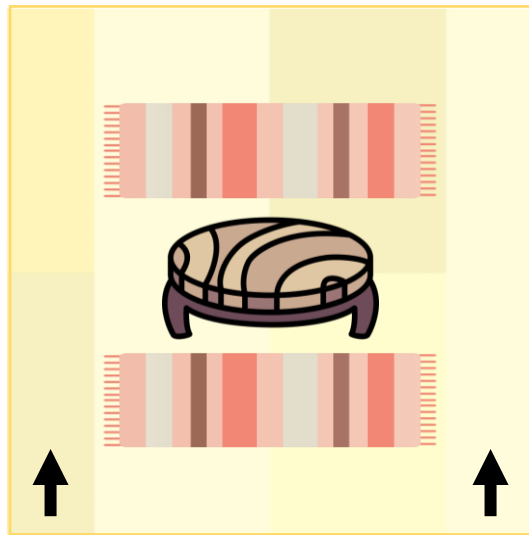


Am I able to learn a policy that performs within 5% of the expert with high confidence?
NO.

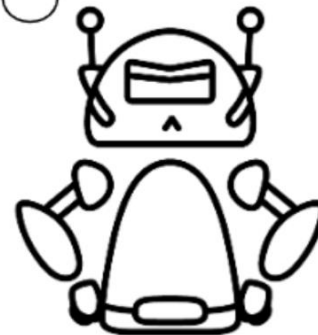
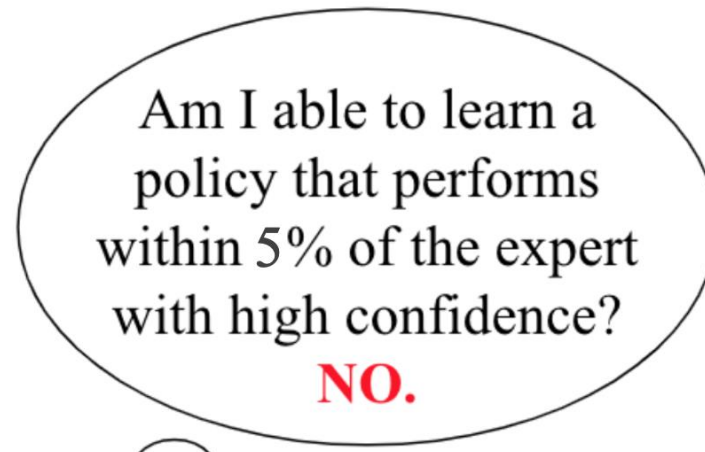


Determining Demonstration Sufficiency

Demonstration Insufficiency



Uninformative Demos



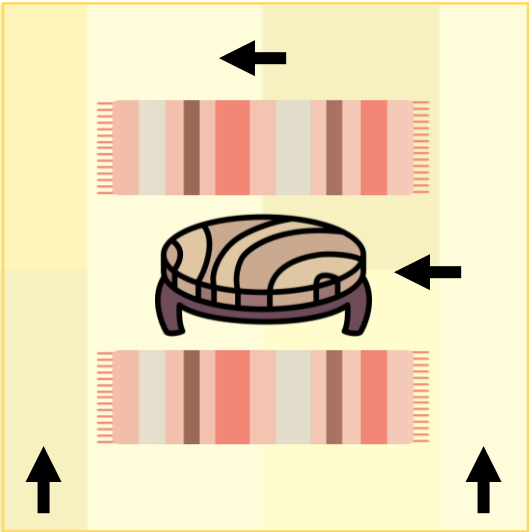
More demos
needed

Determining Demonstration Sufficiency

Demonstration Sufficiency

Determining Demonstration Sufficiency

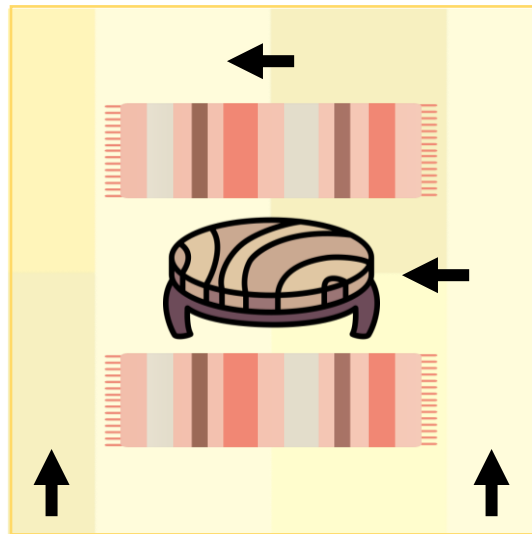
Demonstration Sufficiency



Demos

Determining Demonstration Sufficiency

Demonstration Sufficiency

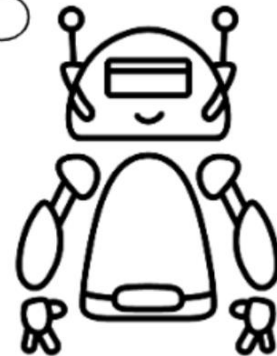


Demos



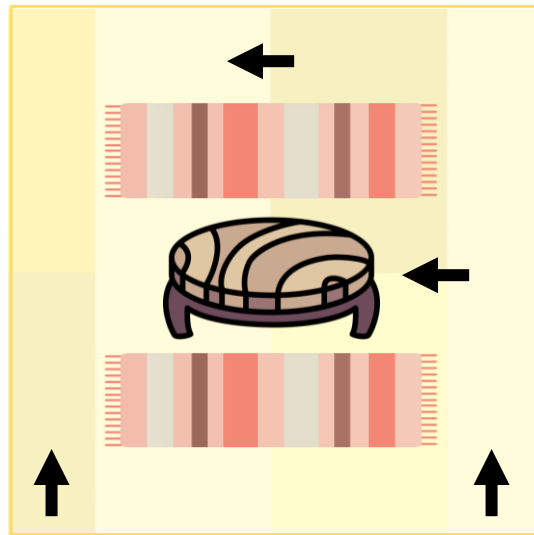
Am I able to learn a policy that performs within 5% of the expert with high confidence?

YES.

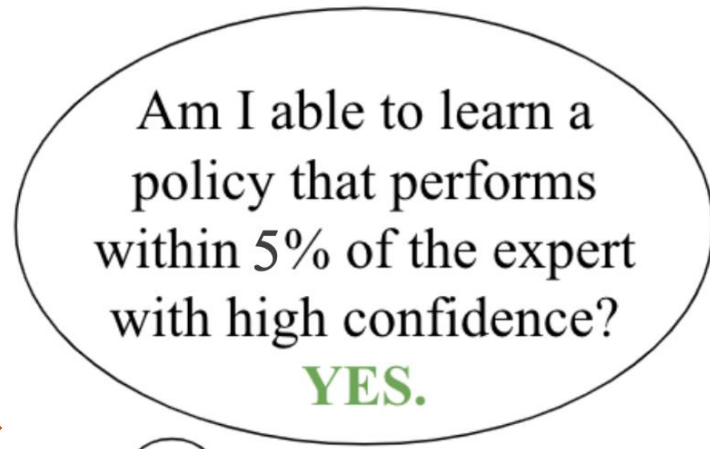
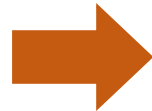


Determining Demonstration Sufficiency

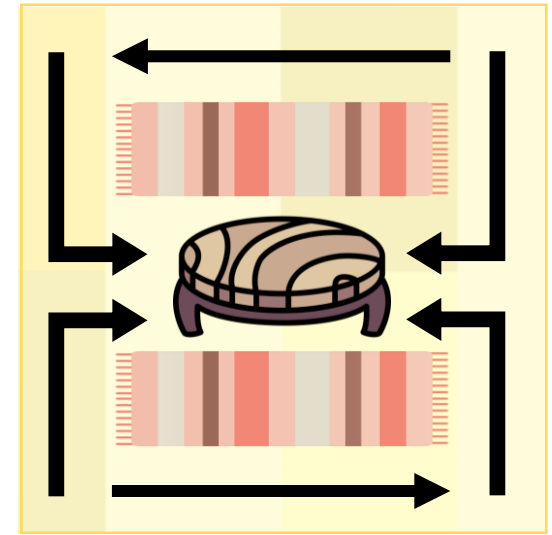
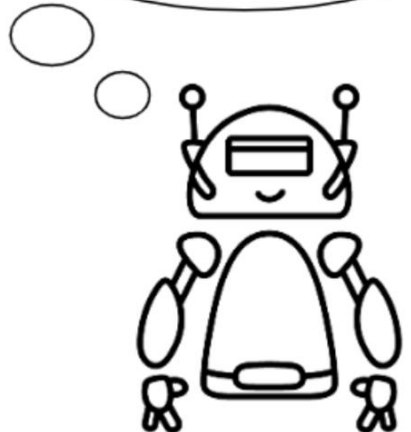
Demonstration Sufficiency



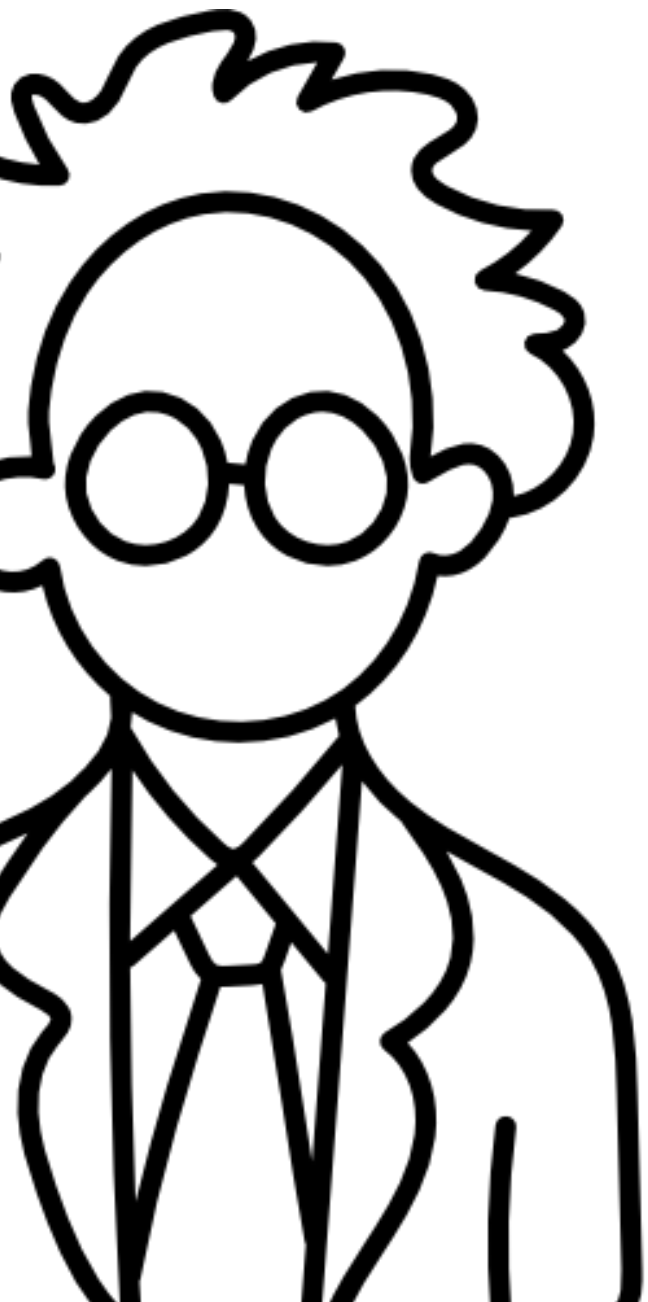
Demos



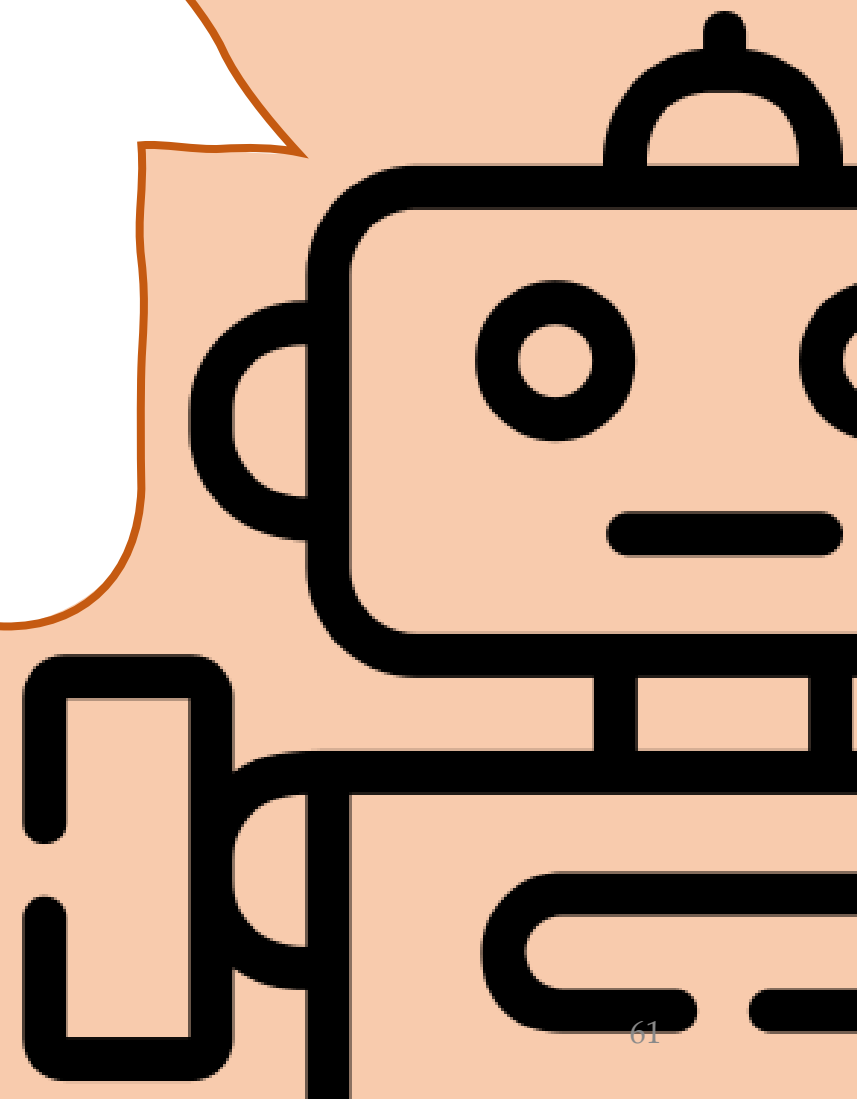
YES.



Learned policy



- Have I received enough demonstrations?
- Are these demonstrations informative enough?
- **What is the reward function?**
- **How do I measure policy “goodness”?**



Estimating the Reward

Bayesian
IRL

Estimating the Reward

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL

Estimating the Reward

$$P(R|D) \propto P(D|R)$$

$$= \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL

Estimating the Reward

$$P(R|D) \propto P(D|R)$$

$$= \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL



R_{robot}
estimated
reward function

Measuring Policy Goodness

- Normalized expected value difference (**nEVD**)

$$nEVD(\pi_{\text{robot}}, R^*) = \frac{V_{R^*}^* - V_{R^*}^{\pi_{\text{robot}}}}{V_{R^*}^* - V_{R^*}^{\pi_{\text{rand}}}}$$

- Puts policy regret in interpretable percentage form

Measuring Policy Goodness

- Normalized expected value difference (**nEVD**)

$$nEVD(\pi_{\text{robot}}, R^*) = \frac{V_{R^*}^* - V_{R^*}^{\pi_{\text{robot}}}}{V_{R^*}^* - V_{R^*}^{\pi_{\text{rand}}}}$$

- Puts policy regret in interpretable percentage form
- We only have an estimate of R^* , so...

Measuring Policy Goodness

$$P(R|D) \propto P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL



R_{robot}
estimated
reward function



R_1, \dots, R_n



π_{robot}

Measuring Policy Goodness

$$P(R|D) \propto P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL



R_{robot}
estimated
reward function



R_1, \dots, R_n



π_{robot}



$\nu_{\alpha} (nEVD(\pi_{\text{robot}}, R))$



α worst-case bound on nEVD

Measuring Policy Goodness

$$P(R|D) \propto P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

$(s, a) \in \mathcal{D}$
demonstrations



Bayesian
IRL



R_{robot}
estimated
reward function



R_1, \dots, R_n



π_{robot}



**demonstration
sufficiency!**

$$\varepsilon > \nu_\alpha (n \text{EVD}(\pi_{\text{robot}}, R))$$

α worst-case bound on nEVD

Comparing With Theoretical Bounds

How many demonstrations is enough for a simple gridworld?

nEVD Threshold	Ours	Abbeel and Ng '04	Syed and Schapire '07
0.1	17	1,600,000	3,700,000
0.3	16	180,000	410,000
0.5	15	65,000	150,000

based on Chernoff-Hoeffding bound

[1] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. ICML 2004.

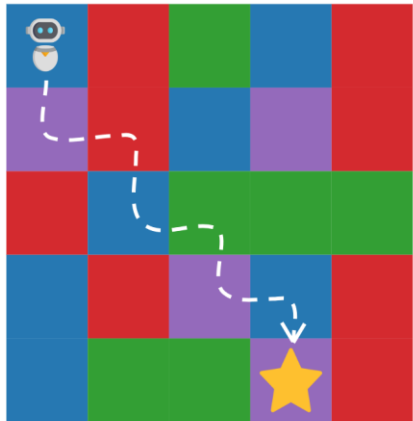
[2] Umar Syed and Robert E Schapire. A game-theoretic approach to apprenticeship learning. NeurIPS 2007.

Baselines and Environments

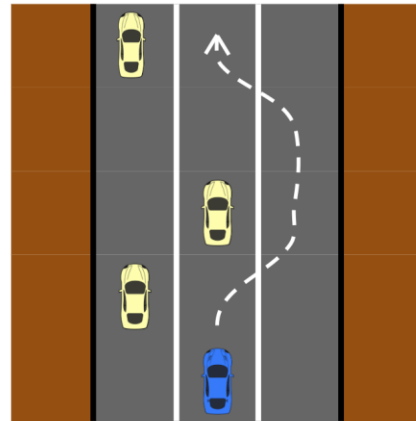
- Convergence heuristic
- Validation set heuristic

Baselines and Environments

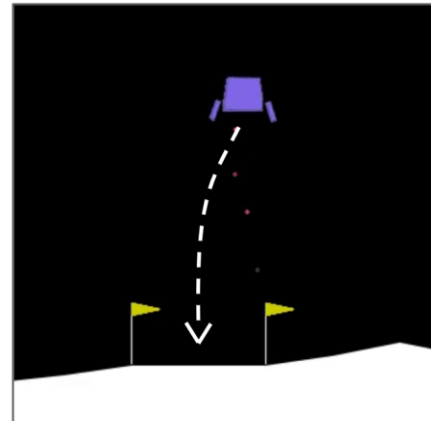
- Convergence heuristic
- Validation set heuristic



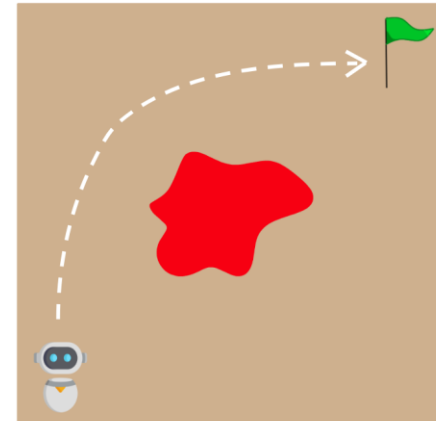
Gridworld



Driving

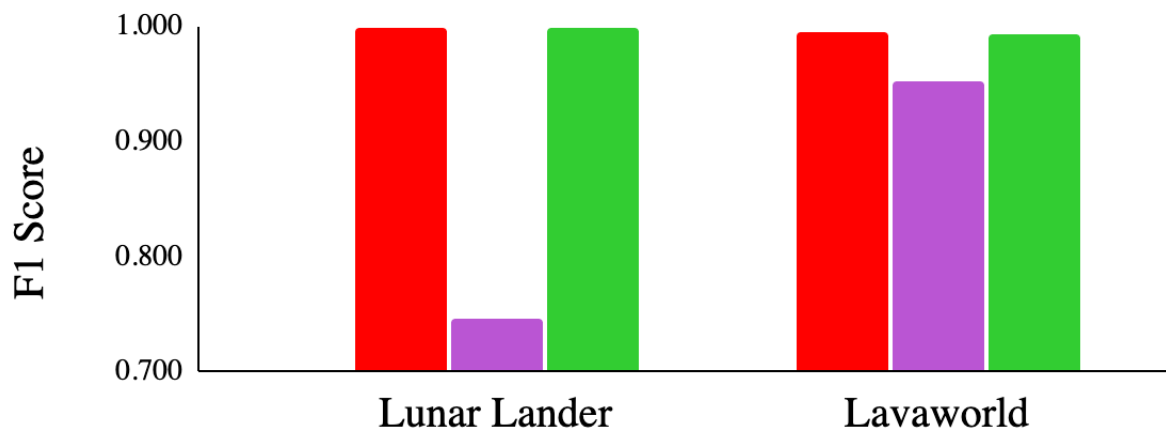
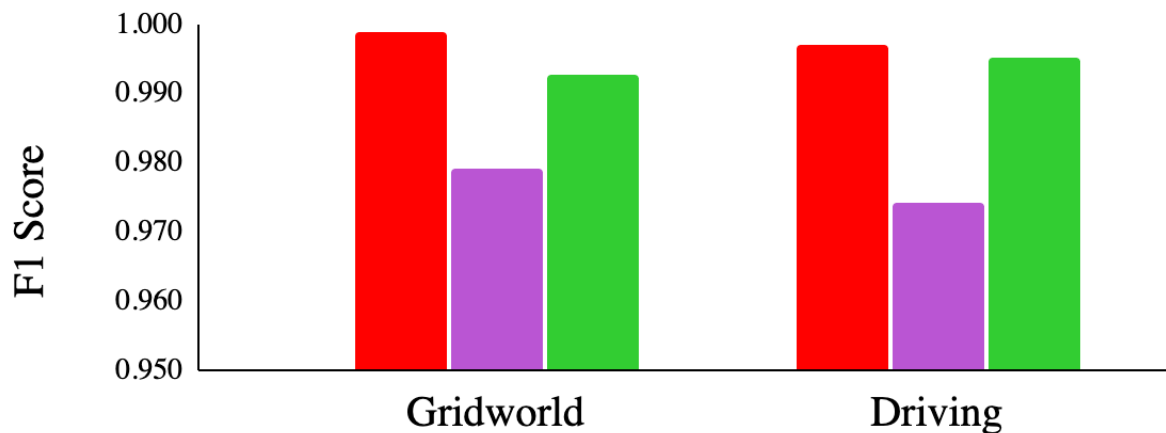


Lunar lander

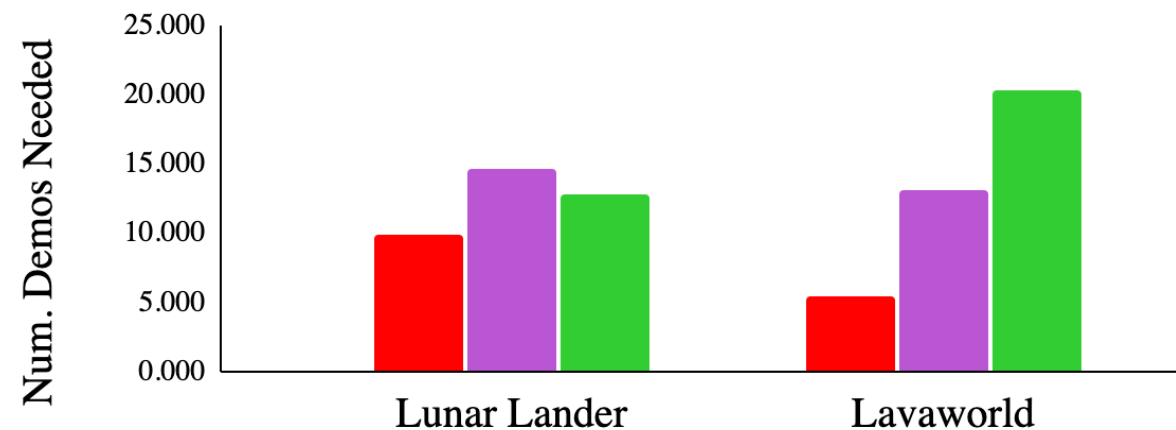
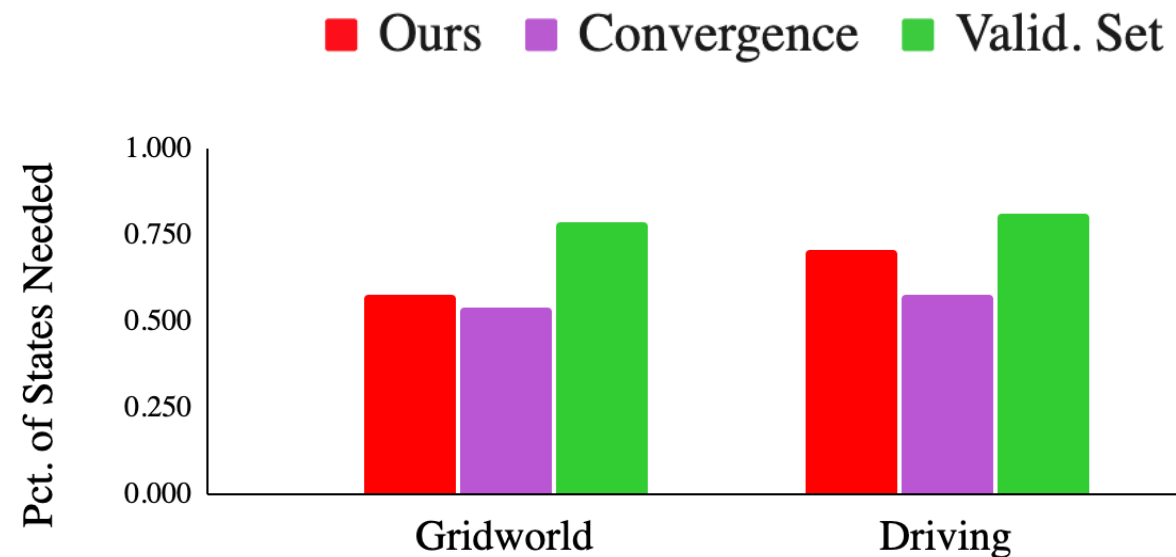


Lavaworld

Simulation Results

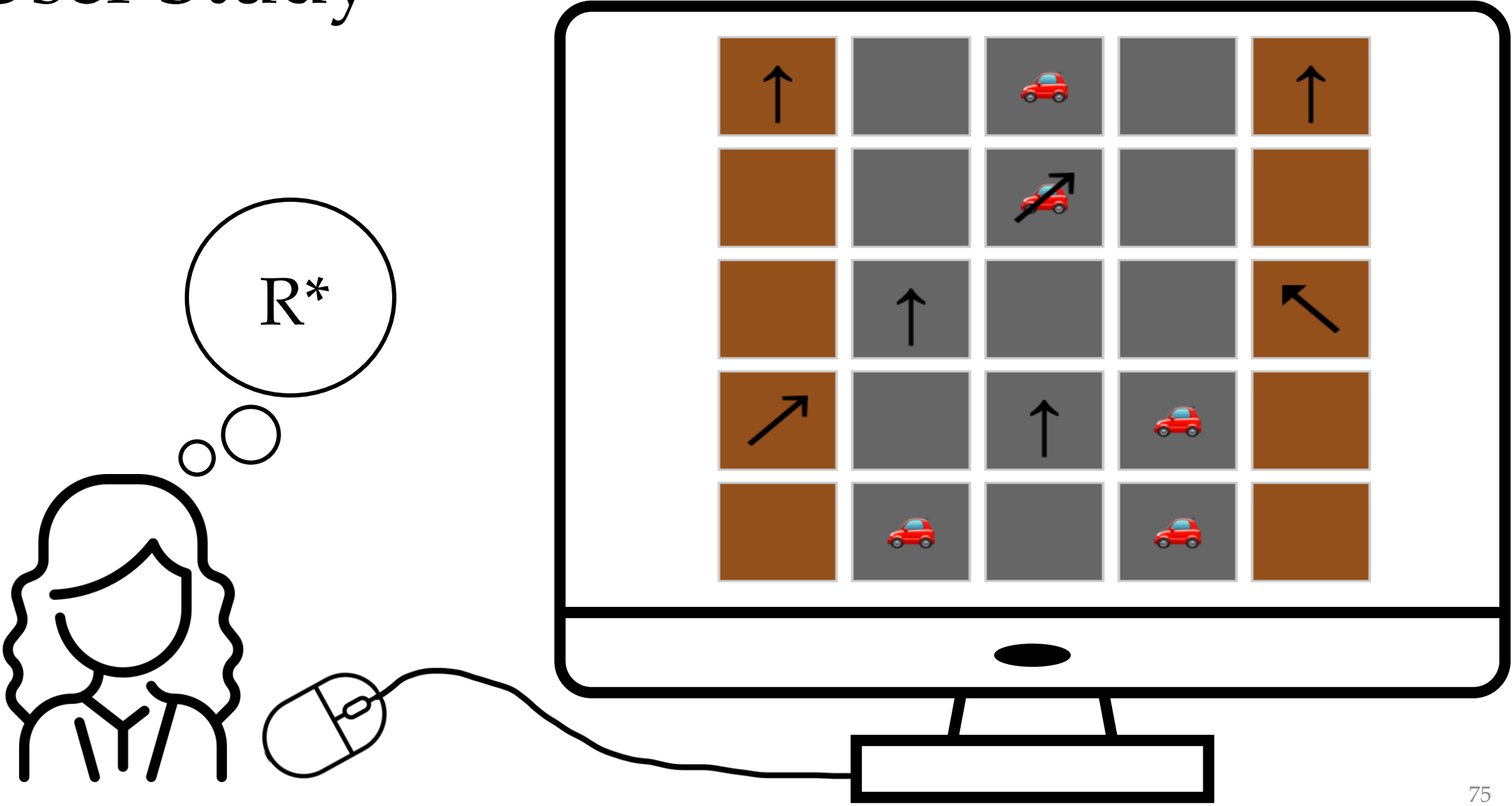


on average, 5% more accurate than baselines

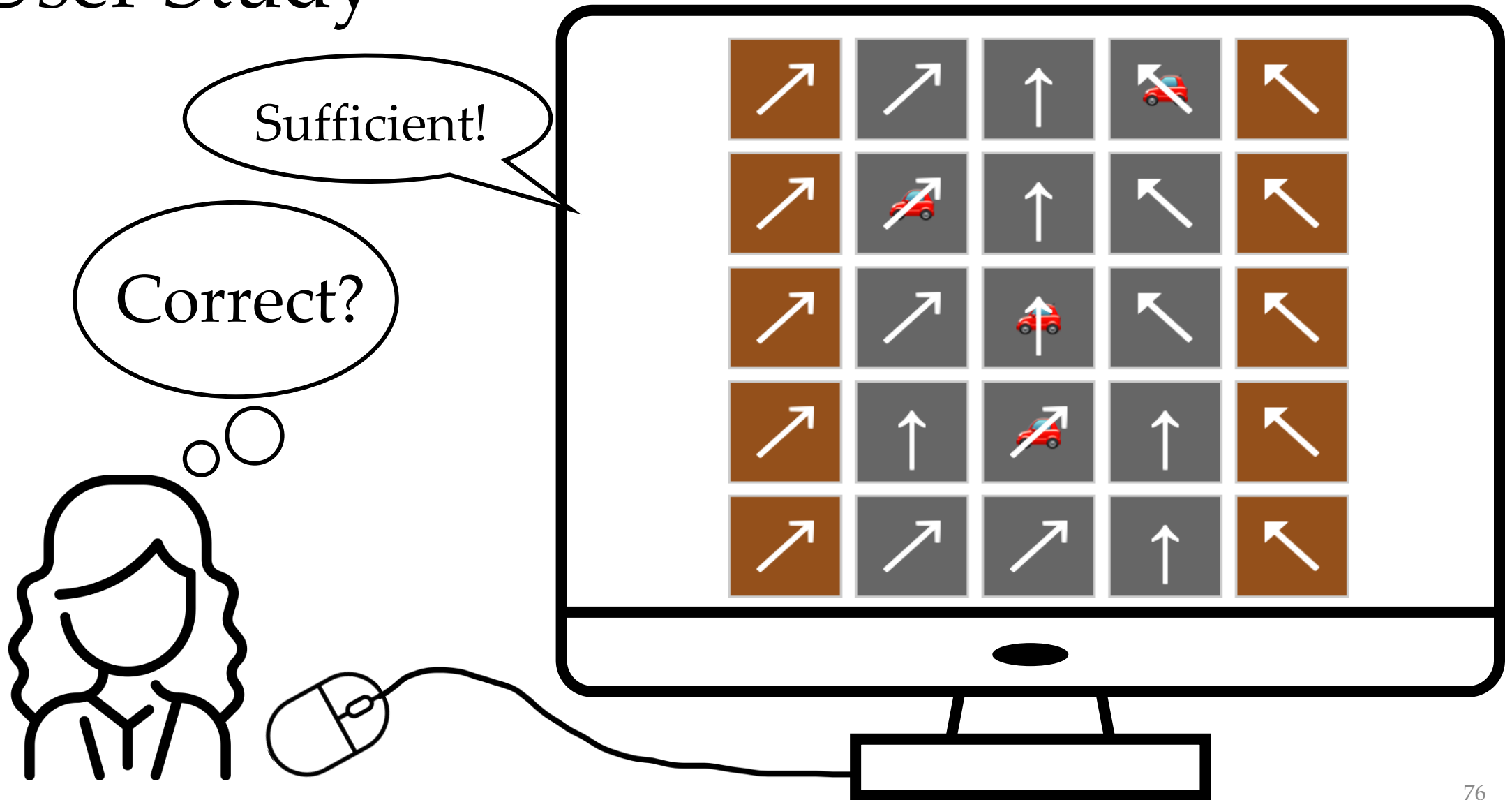


on average, 25% more sample efficient than baselines

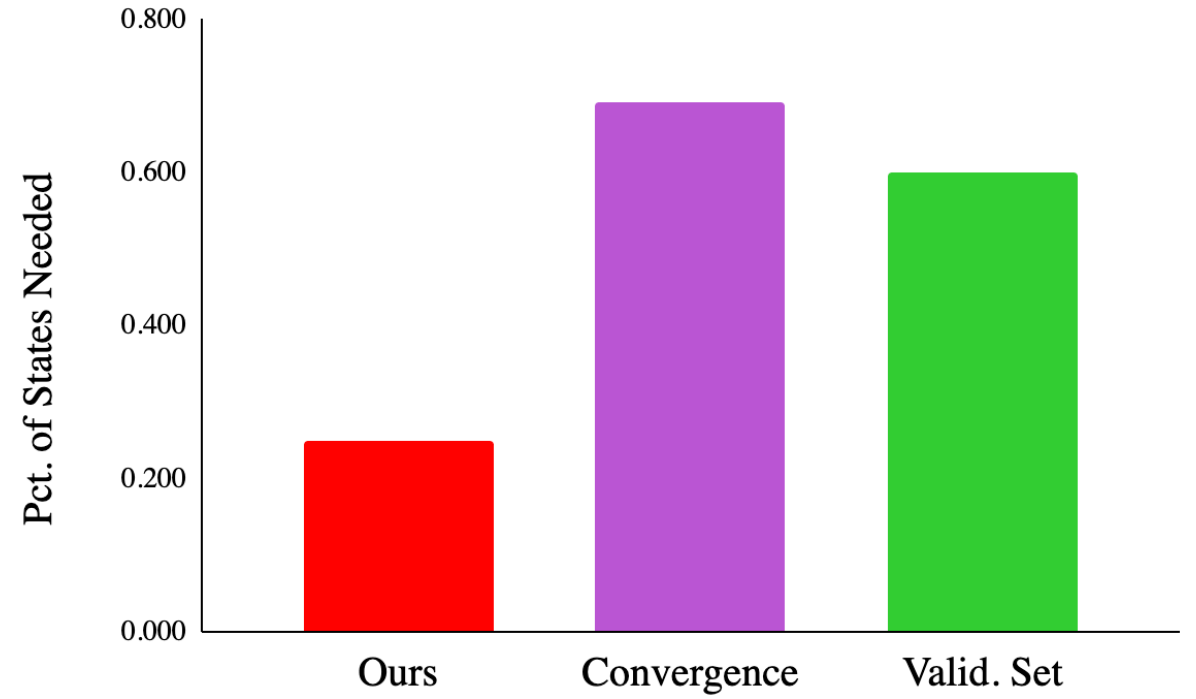
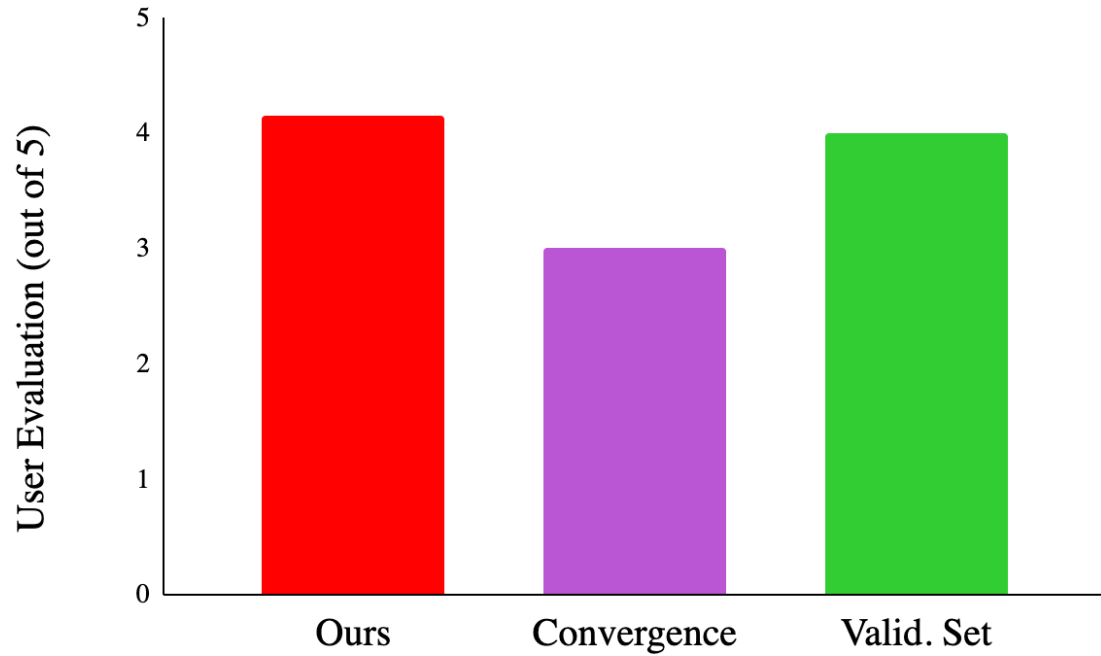
User Study



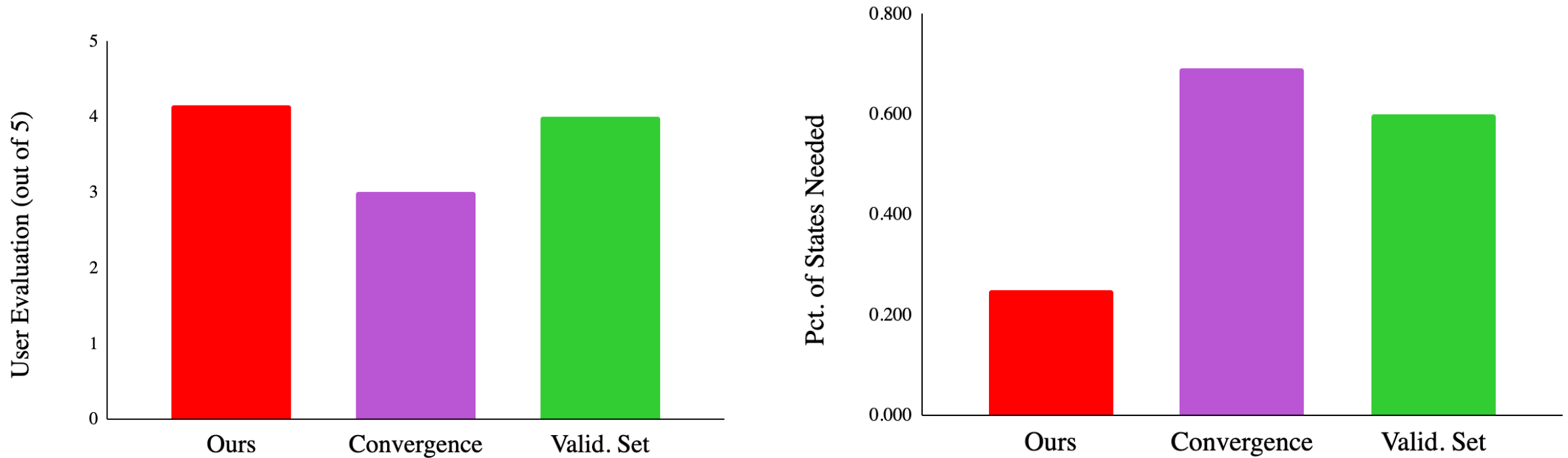
User Study



User Study Results



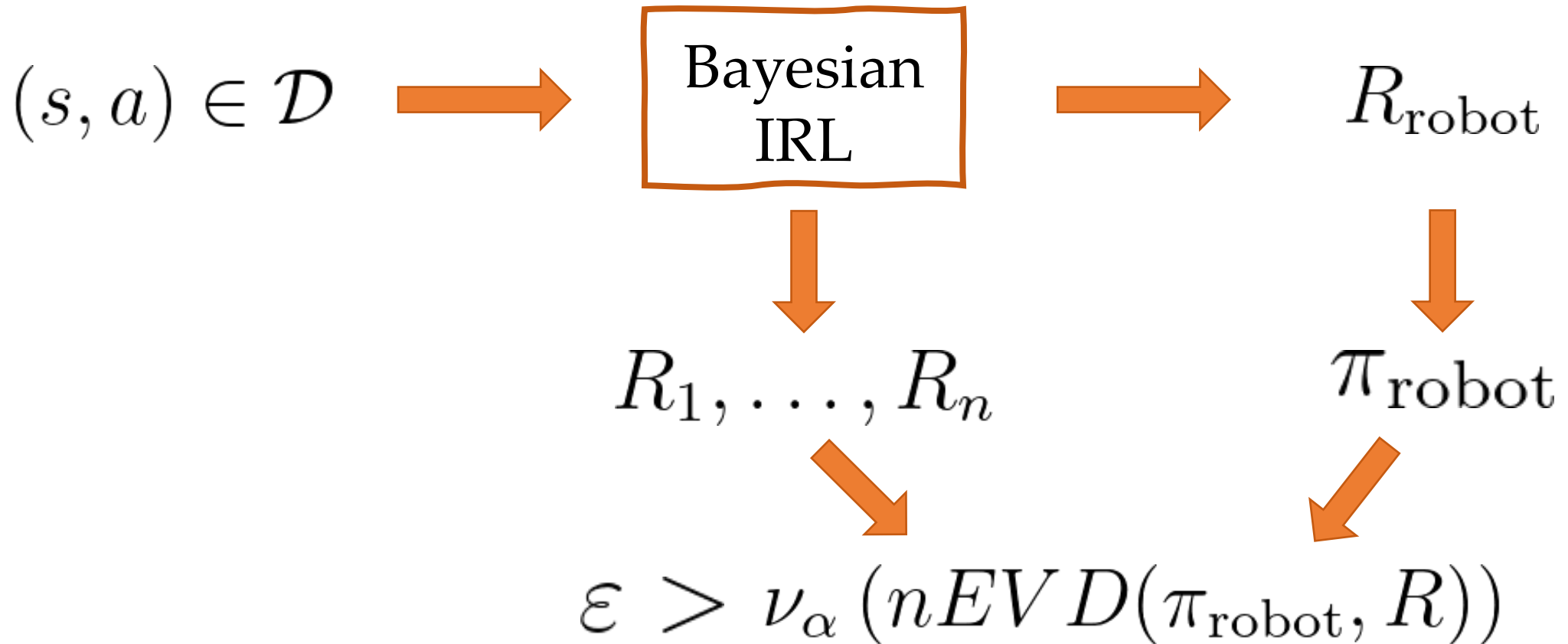
User Study Results



robust against suboptimal demonstrations: ~12%

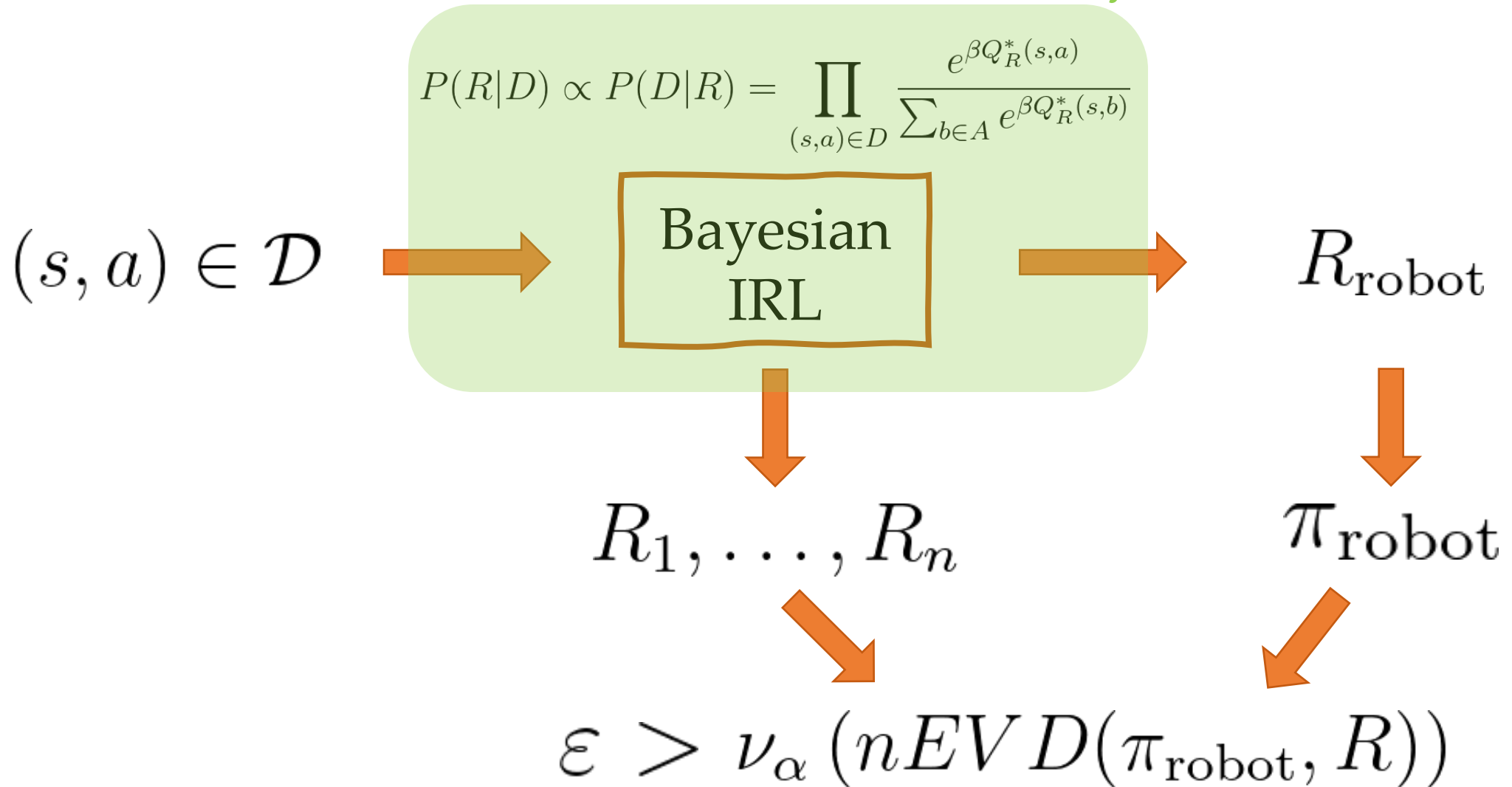
Future Work

$$P(R|D) \propto P(D|R) = \prod_{(s,a) \in D} \frac{e^{\beta Q_R^*(s,a)}}{\sum_{b \in A} e^{\beta Q_R^*(s,b)}}$$

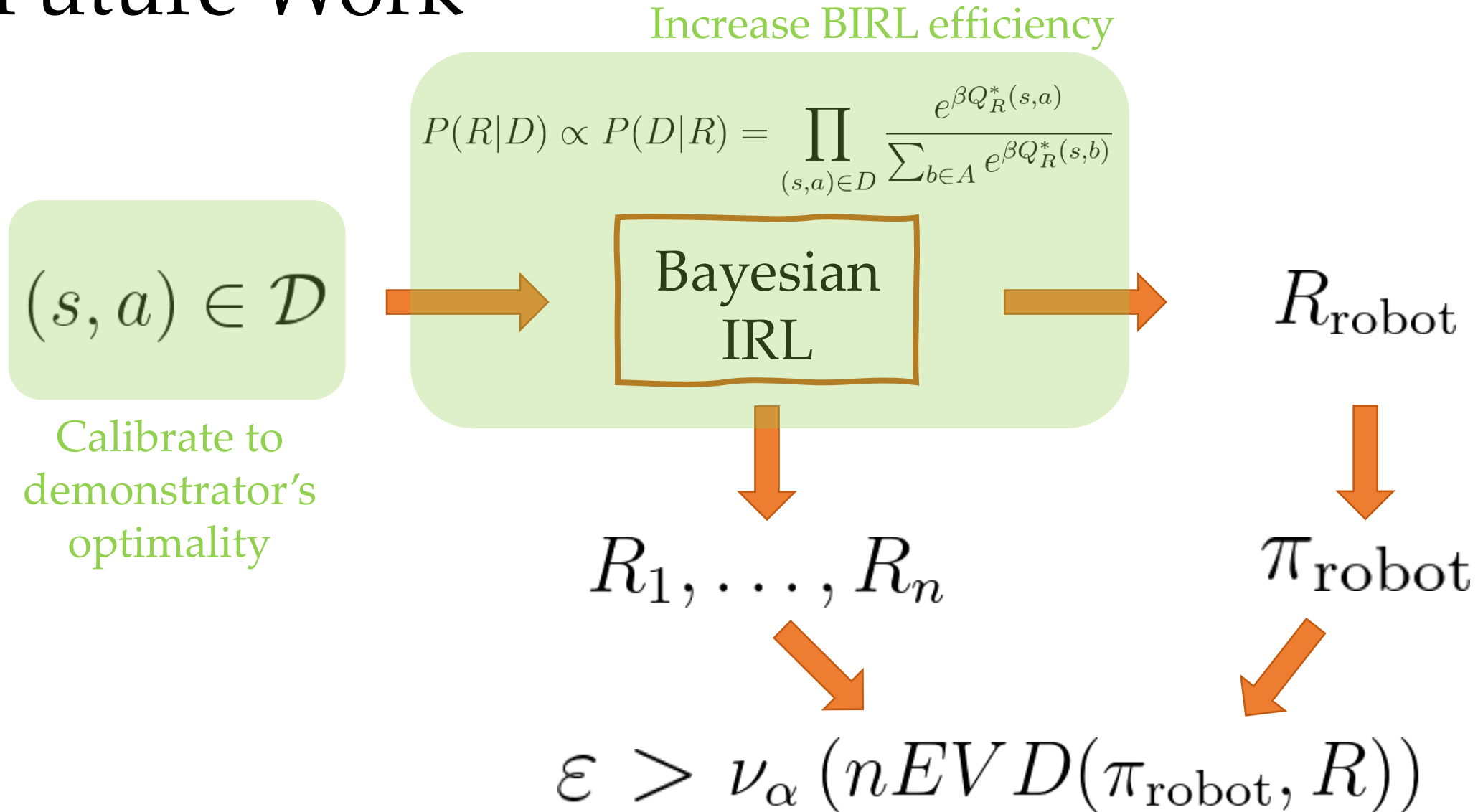


Future Work

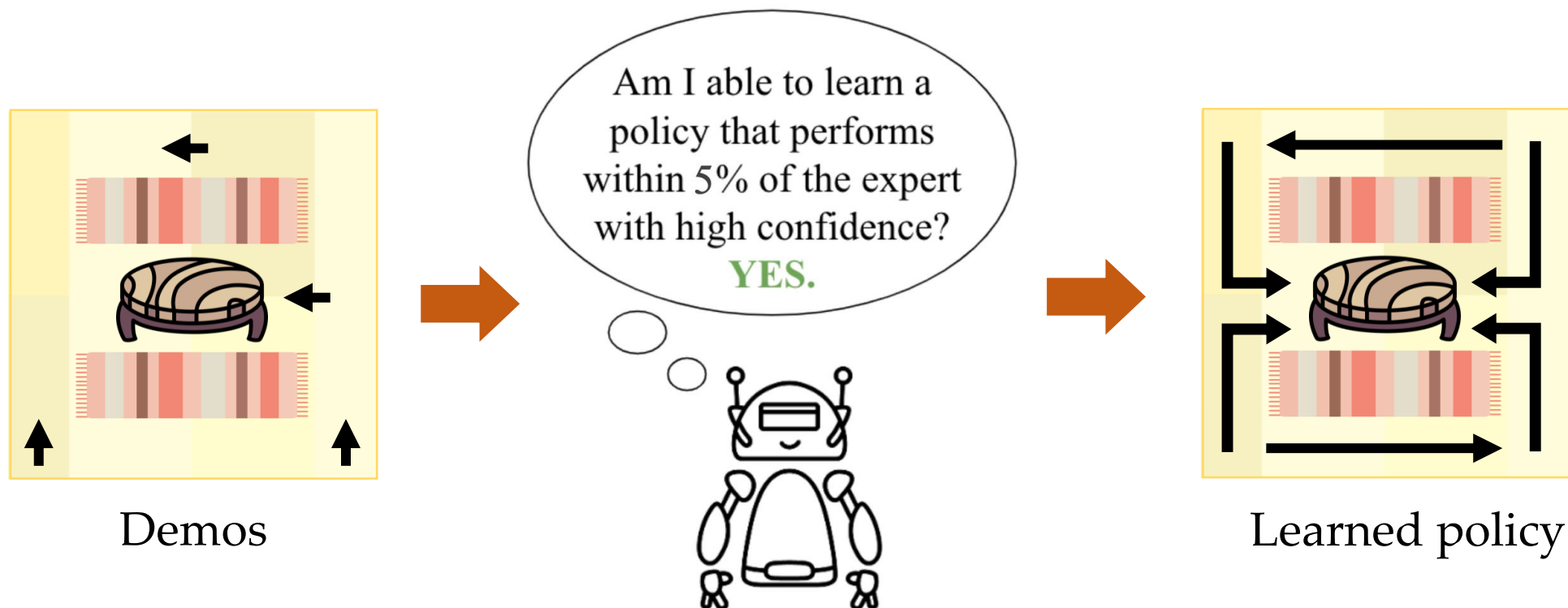
Increase BIRL efficiency



Future Work



Autonomous Assessment of Demonstration Sufficiency via Bayesian Inverse Reinforcement Learning



What if I want to learn rewards from more than demonstrations?

Reward-rational (implicit) choice: A unifying formalism for reward learning

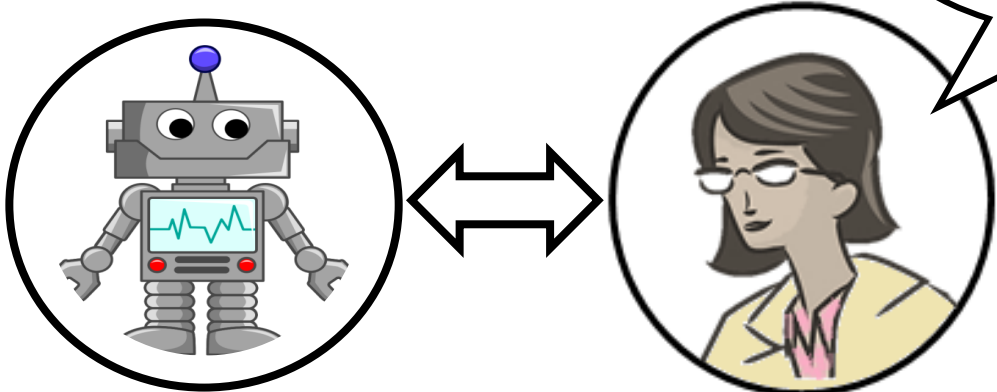
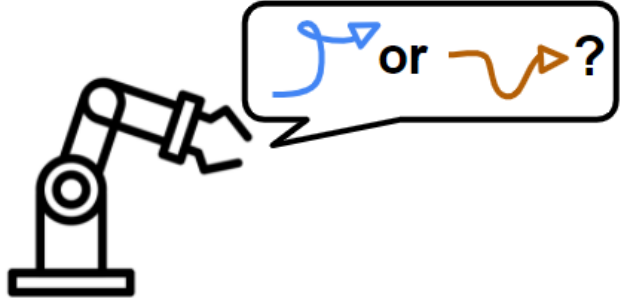
Hong Jun Jeon^{*1}, Smitha Milli^{*2}, Anca Dragan²
hjjeon@stanford.edu, smilli@berkeley.edu, anca@berkeley.edu

^{*}Equal contribution,

¹Stanford University,

²University of California, Berkeley

How do we learn from diverse types of feedback of unknown quality?



Unifying Human Feedback Types

Boltzmann Rational Choice Model

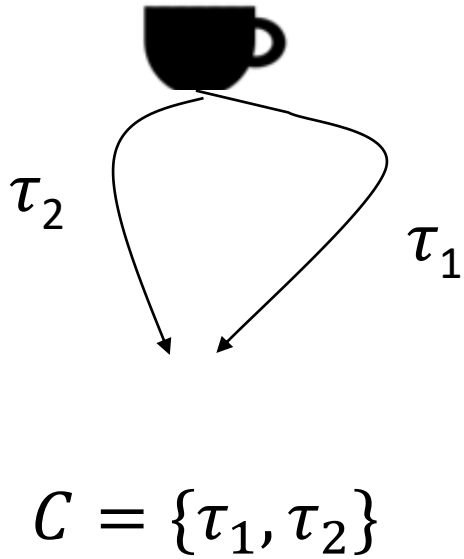
$$\mathcal{C} = \{c_1, \dots, c_n\} \quad r: \mathcal{C} \rightarrow \mathbb{R}$$

$$\mathbb{P}(c_i|r) = \frac{\exp(\beta r(c_i))}{\sum_{c'} \exp(\beta r(c'))}$$

$\beta \rightarrow 0$ random (non-rational) choices

$\beta \rightarrow \infty$ deterministic (perfectly-rational) choices

Trajectory Comparisons (Pairwise Prefs)



Boltzmann Rational Choice Model

$$C = \{c_1, \dots, c_n\} \quad r: C \rightarrow \mathbb{R}$$

$$\mathbb{P}(c_i|r) = \frac{\exp(\beta r(c_i))}{\sum_C \exp(\beta r(c'))}$$

$$\mathbb{P}(\tau_i|r) = \frac{\exp(\beta r(\tau_i))}{\exp(\beta r(\tau_1)) + \exp(\beta r(\tau_2))}$$

Demonstrations



$C = \{\text{All Possible Trajectories}\}$

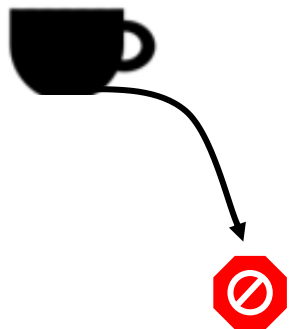
$$\mathbb{P}(\tau|r) = \frac{\exp(\beta r(\tau))}{\sum_{\tau} \exp(\beta r(\tau'))}$$

Boltzmann Rational Choice Model

$$C = \{c_1, \dots, c_n\} \quad r: C \rightarrow \mathbb{R}$$

$$\mathbb{P}(c_i|r) = \frac{\exp(\beta r(c_i))}{\sum_c \exp(\beta r(c'))}$$

E-Stops



$$C = \{\tau_{:1}, \dots, \tau_{:T}\}$$

$$\mathbb{P}(\tau_{:i}|r) = \frac{\exp(\beta r(\tau_{:i}))}{\sum_1^T \exp \beta r(\tau_{:j})}$$

Boltzmann Rational Choice Model

$$C = \{c_1, \dots, c_n\} \quad r: C \rightarrow \mathbb{R}$$

$$\mathbb{P}(c_i|r) = \frac{\exp(\beta r(c_i))}{\sum_C \exp(\beta r(c'))}$$

Reward Learning from Human Feedback

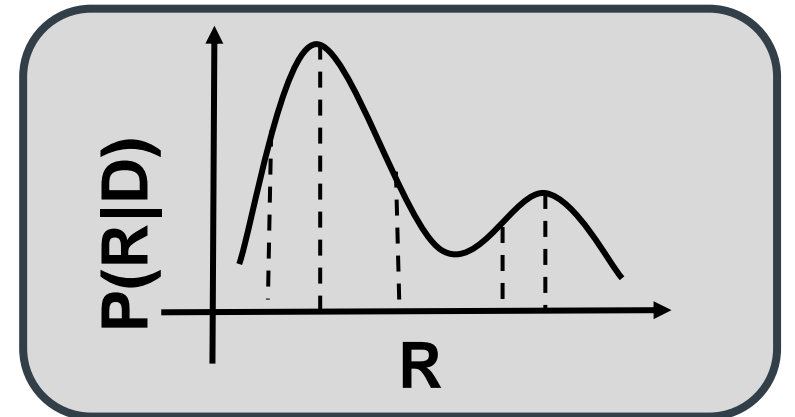
- Assume demonstrator is Boltzmann rational
 - Demonstrator follows a softmax policy with inverse temperature β

$$C = \{c_1, \dots, c_n\} \quad r: C \rightarrow \mathbb{R}$$

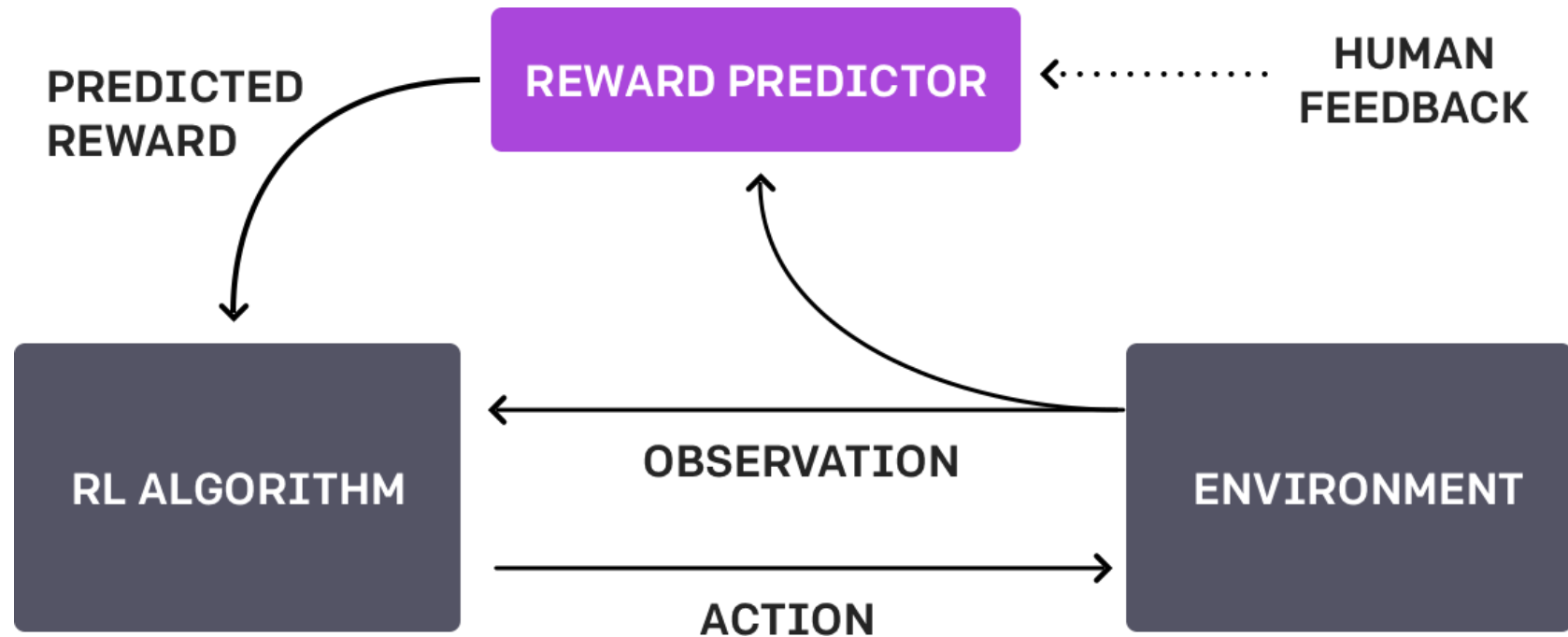
$$\mathbb{P}(c_i|r) = \frac{\exp(\beta r(c_i))}{\sum_C \exp(\beta r(c'))}$$

- Perform Bayesian inference (MCMC) to sample from posterior distribution

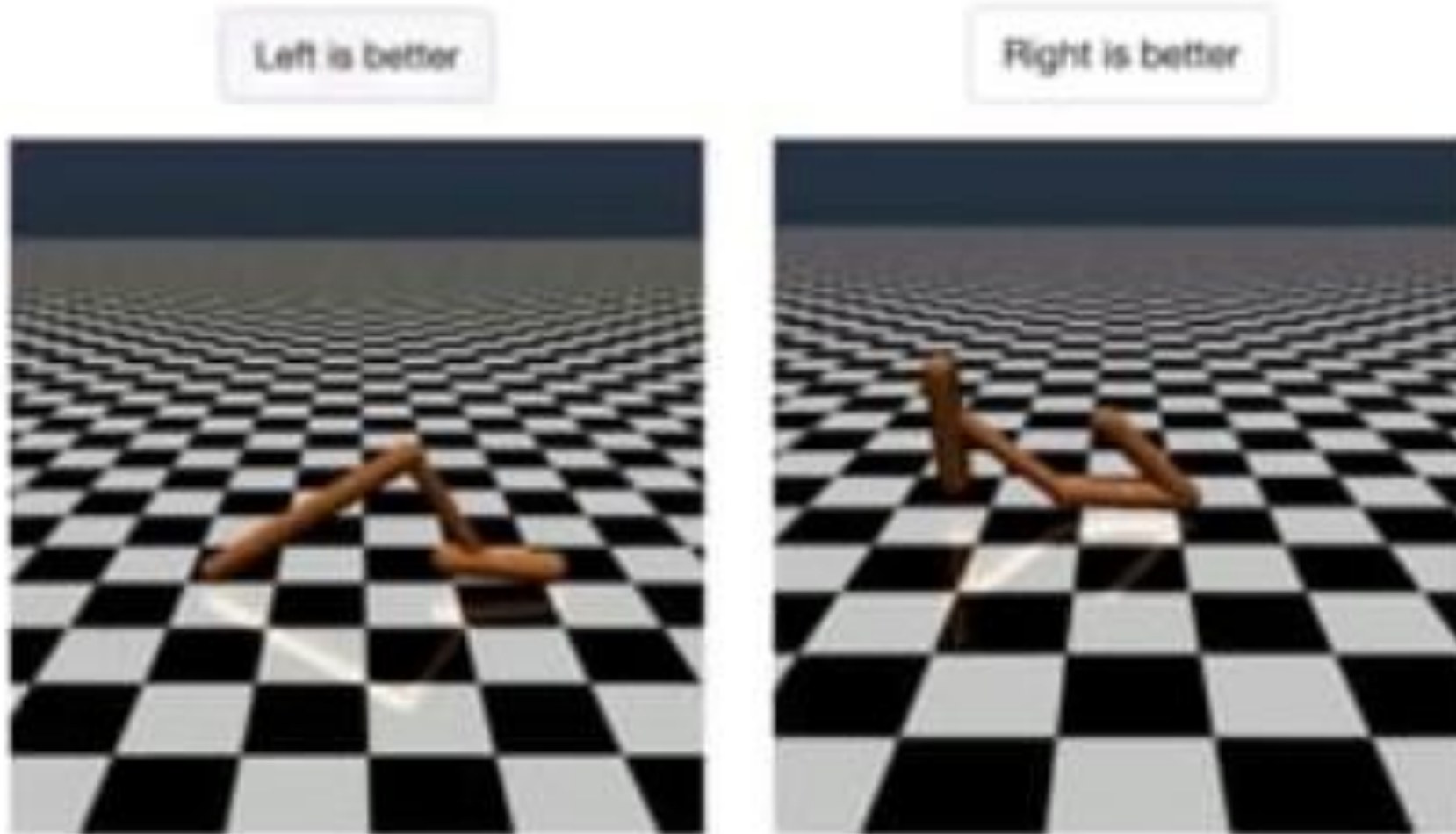
$$P(R|D) \propto P(D|R)P(R)$$



RL from Human Feedback (RLHF)



RL from Human Preferences



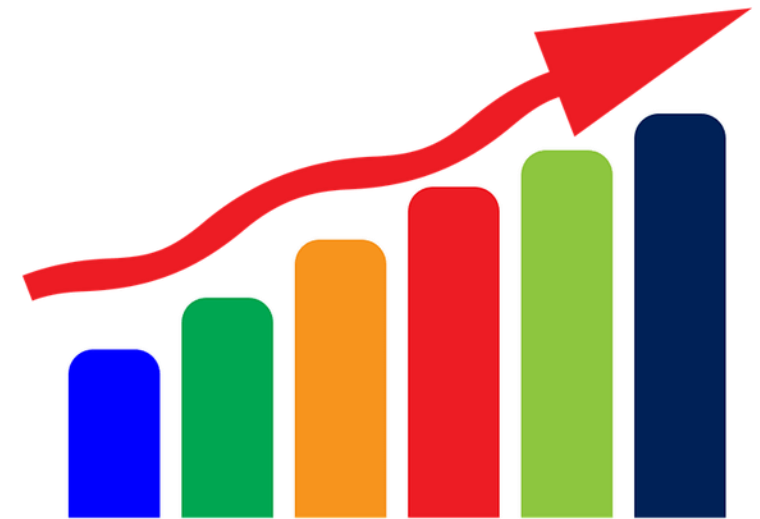
Why would you want to learn a reward from ranked examples?

Inverse Reinforcement Learning

Prior approaches ...

1. Typically couldn't do much better than the demonstrator.
2. Were hard to scale to complex problems.

Pre-Ranked
Demonstrations



Inverse Reinforcement Learning

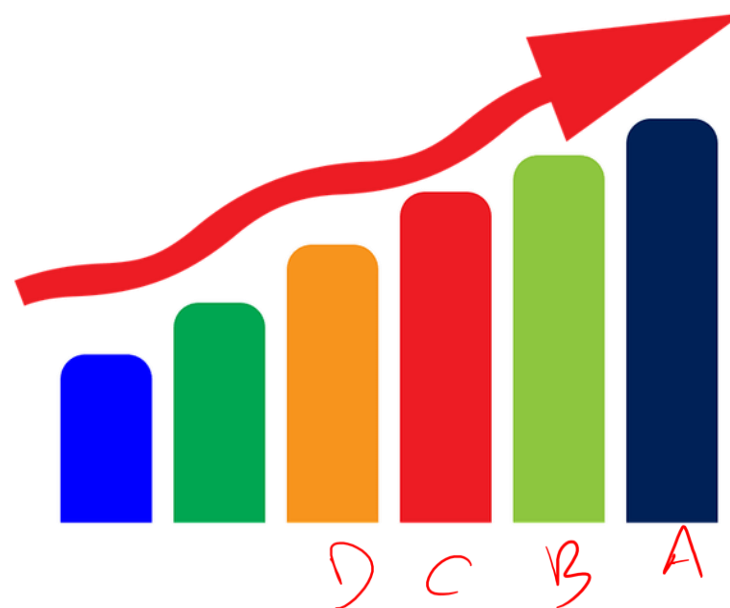
Prior approaches ...

~~1. Typically couldn't do much better than the demonstrator.~~

Find a reward function that explains the ranking, allowing for extrapolation.

2. Were hard to scale to complex problems.

Pre-Ranked
Demonstrations



Inverse Reinforcement Learning

Prior approaches ...

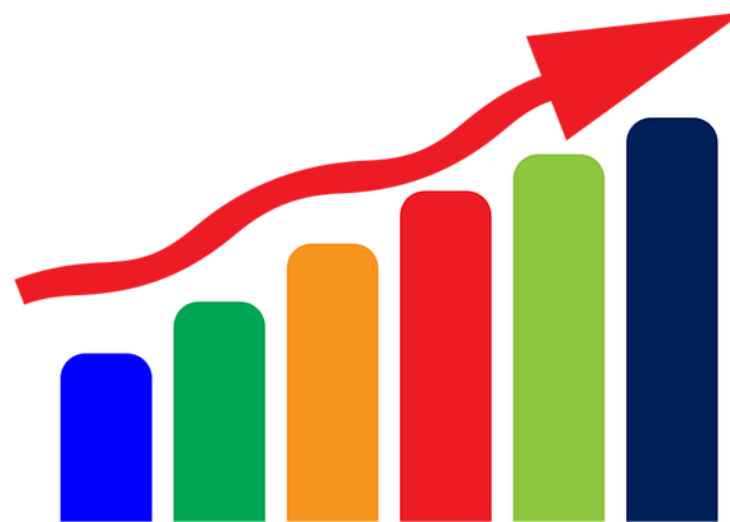
~~1. Typically couldn't do much better than the demonstrator.~~

Find a reward function that explains the ranking, allowing for extrapolation.

~~2. Were hard to scale to complex problems.~~

Reward learning becomes a supervised learning problem.

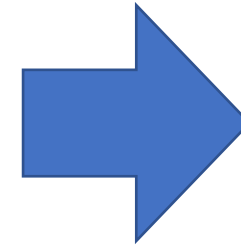
Pre-Ranked
Demonstrations



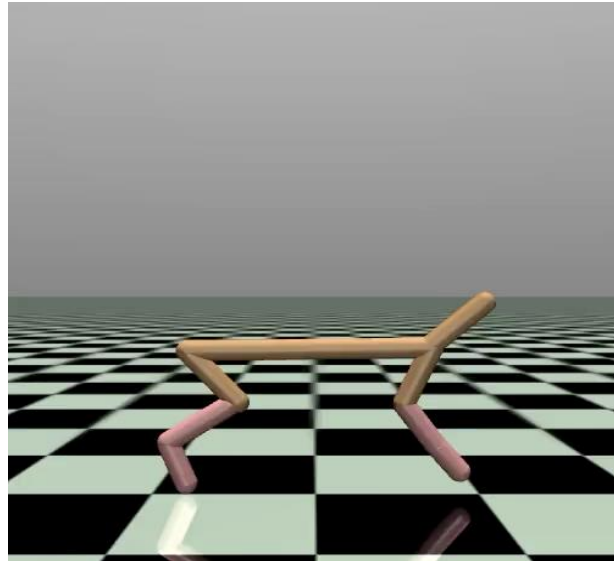
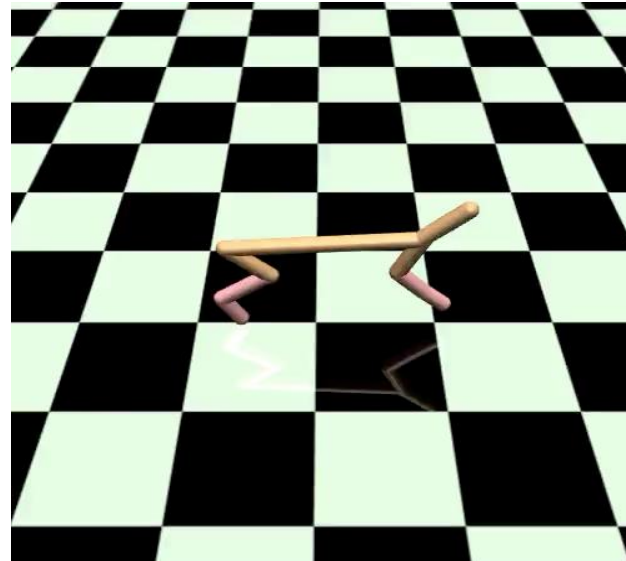
Trajectory-ranked Reward Extrapolation (T-REX)



Reward
Function

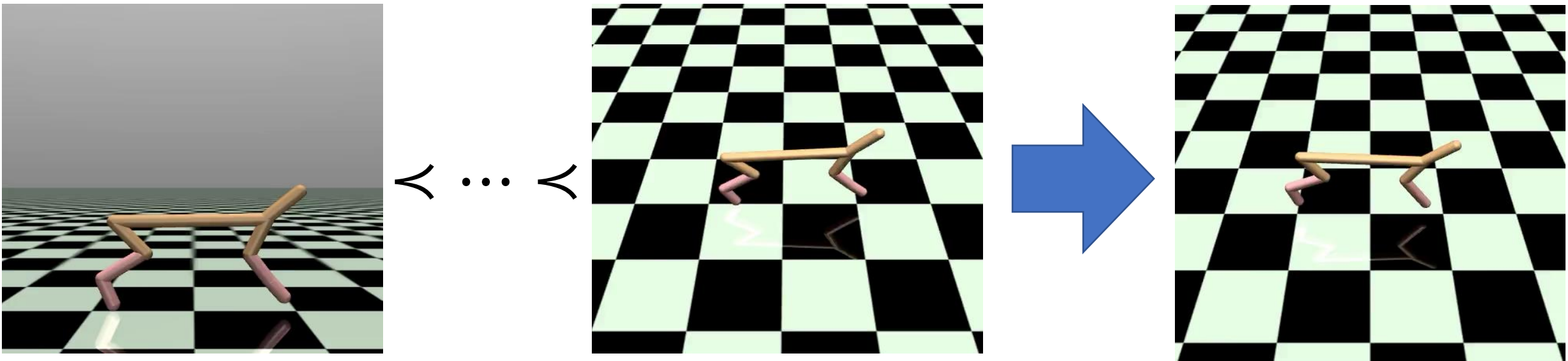


$\wedge \dots \wedge$



Pre-ranked demonstrations

Trajectory-ranked Reward Extrapolation (T-REX)



Pre-ranked demonstrations

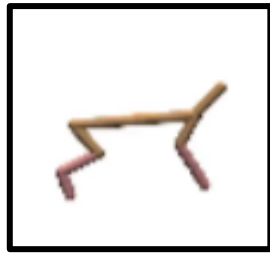
T-REX Policy

Reward Function

$$R_{\theta}: \mathcal{S} \rightarrow \mathbb{R}$$

Examples of \mathcal{S} :

Current Robot Joint
Angles and Velocities



→ 0.5



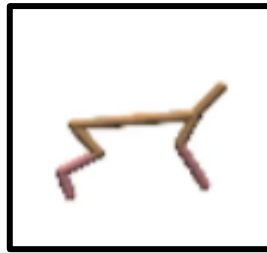
→ -0.7

Reward Function

$$R_{\theta}: \mathcal{S} \rightarrow \mathbb{R}$$

Examples of \mathcal{S} :

Current Robot Joint
Angles and Velocities



→ 0.5



→ -0.7

Short
Sequence of
Images



→ 0.9



→ -1.2

Trajectory-ranked Reward Extrapolation (T-REX)

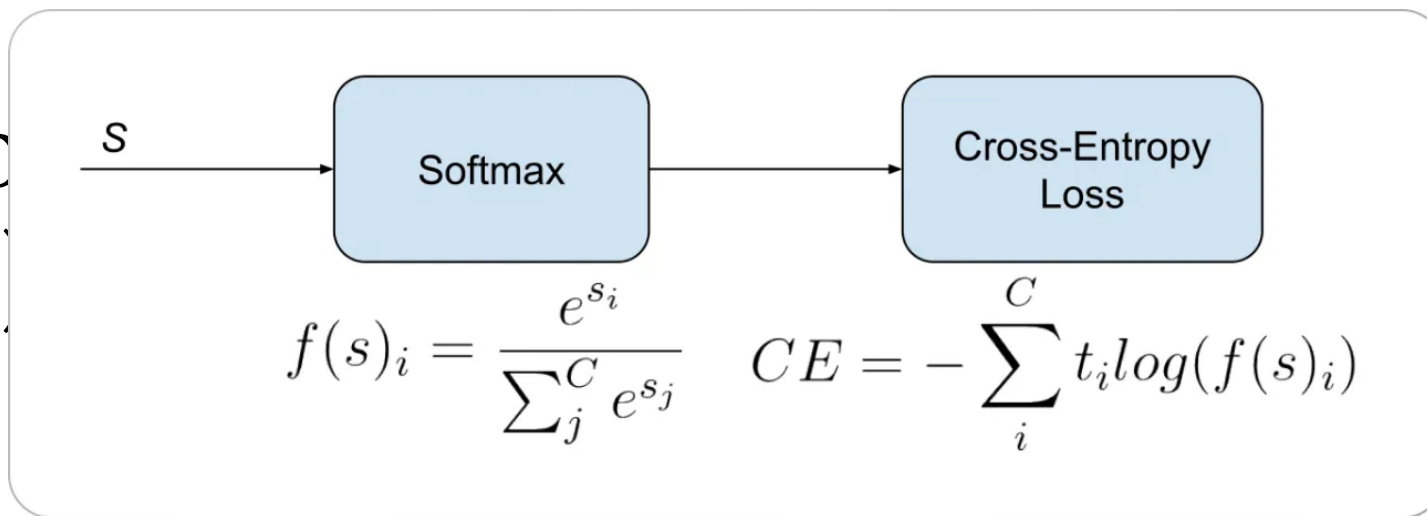
$$\tau_1 \prec \tau_2 \prec \dots \prec \tau_T$$

$$\sum_{s \in \tau_1} R_\theta(s) < \sum_{s \in \tau_2} R_\theta(s)$$

Bradley-Terry pairwise ranking loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \sum_{s \in \tau_j} R_\theta(s)}{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}$$

Trajectory (T-REX)



tion

$$\sum_{s \in \tau_1} R_\theta(s) < \sum_{s \in \tau_2} R_\theta(s)$$

Logits

Minimize cross-entropy loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i < \tau_j} \frac{\exp \sum_{s \in \tau_i} R_\theta(s)}{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}$$

Trajectory-ranked Reward Extrapolation (T-REX)

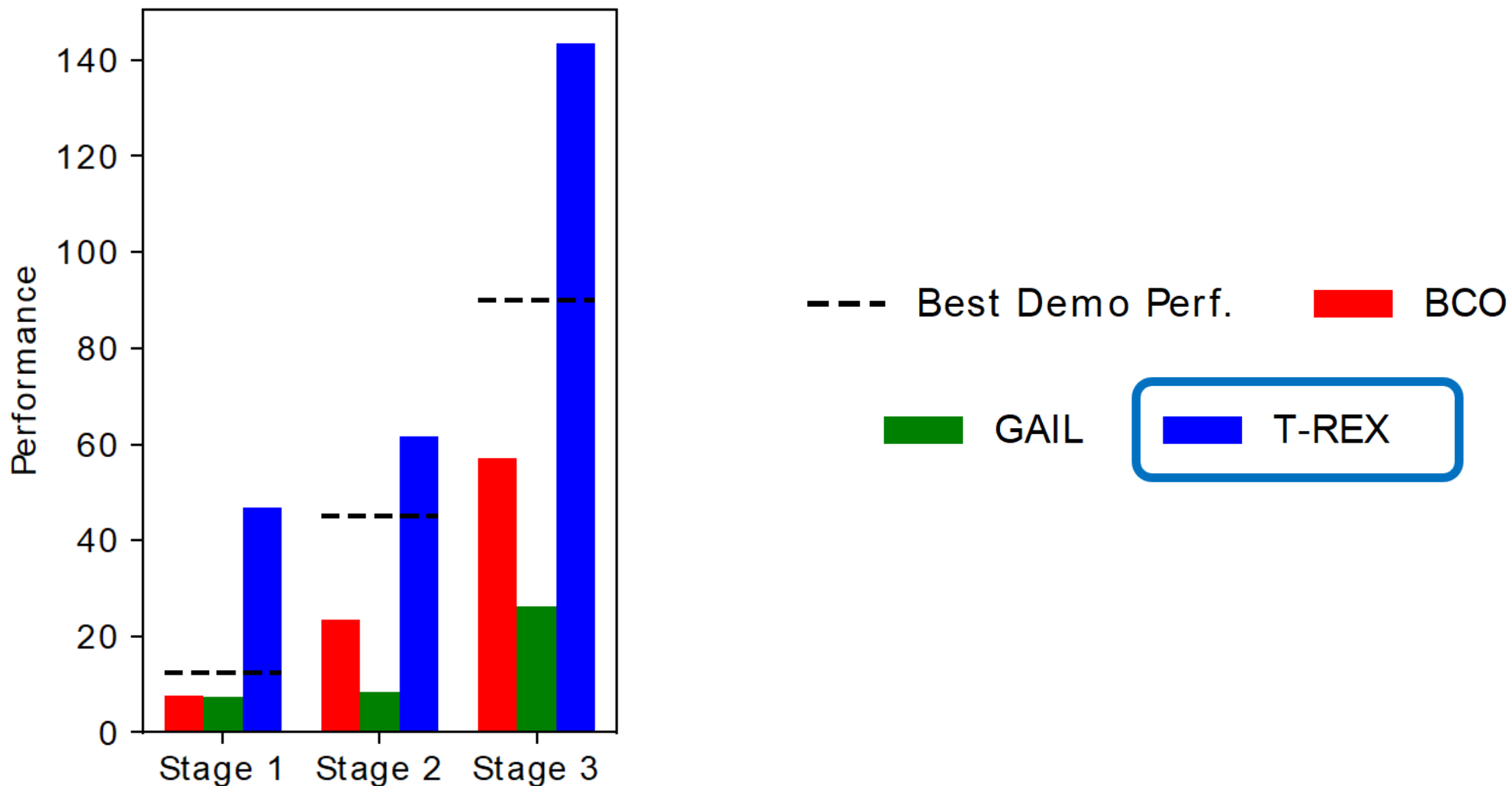
$$\tau_1 \prec \tau_2 \prec \dots \prec \tau_T$$

Given pre-ranked demos, reward learning can be formulated as a standard supervised learning task.

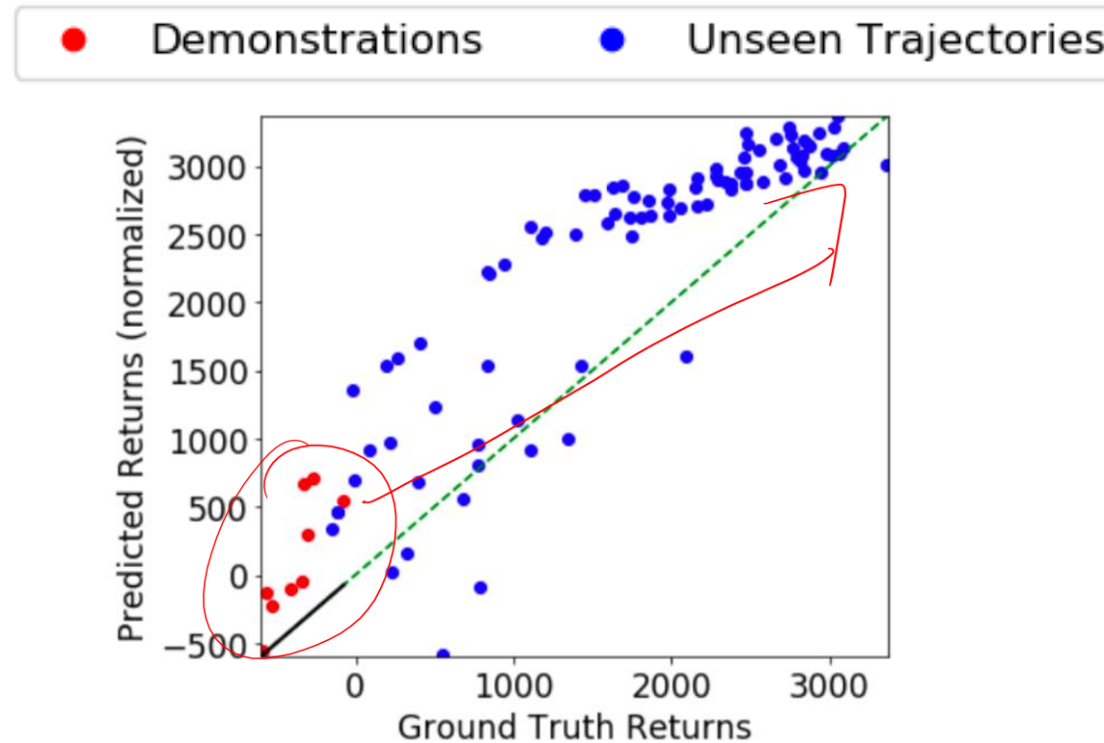
Minimize cross-entropy loss

$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \sum_{s \in \tau_j} R_\theta(s)}{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}$$

T-REX Policy Performance

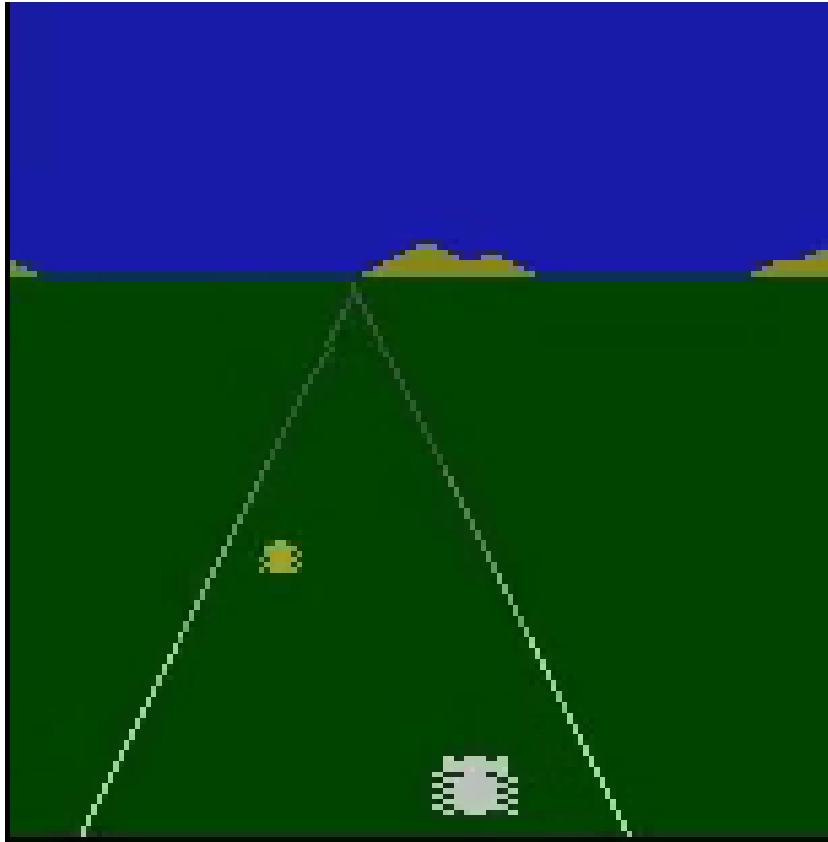


Reward Extrapolation

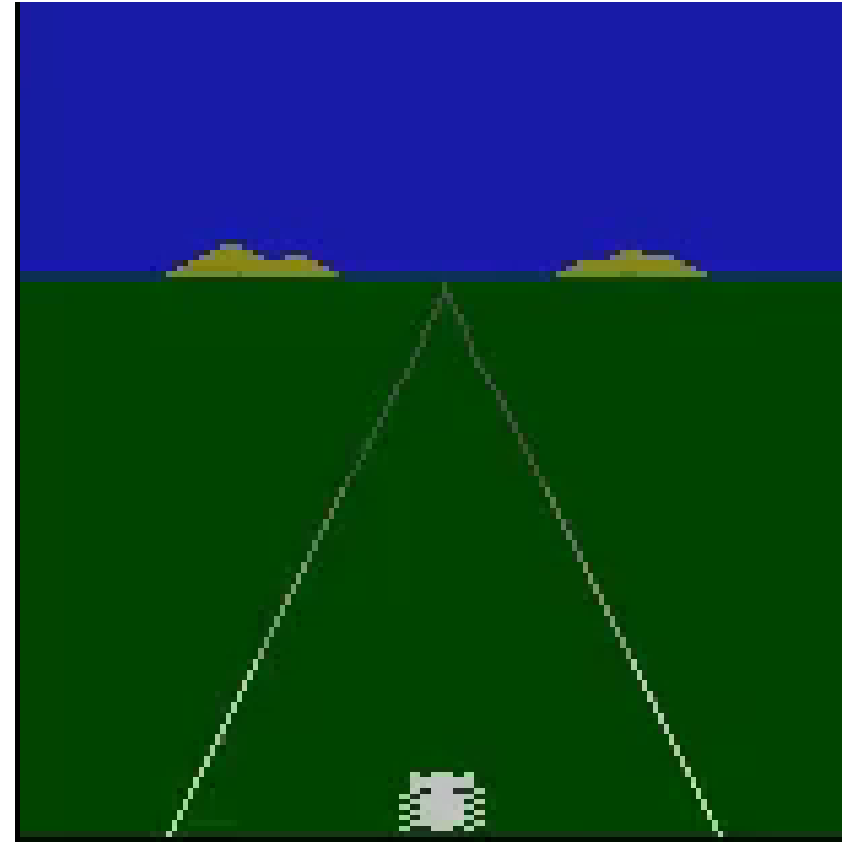


T-REX can extrapolate beyond the performance of the best demo

“Autonomous Driving” in Atari



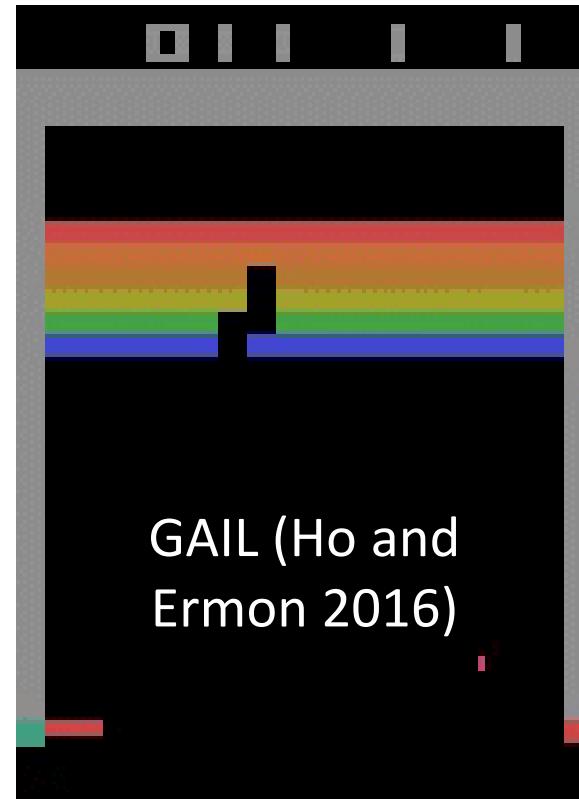
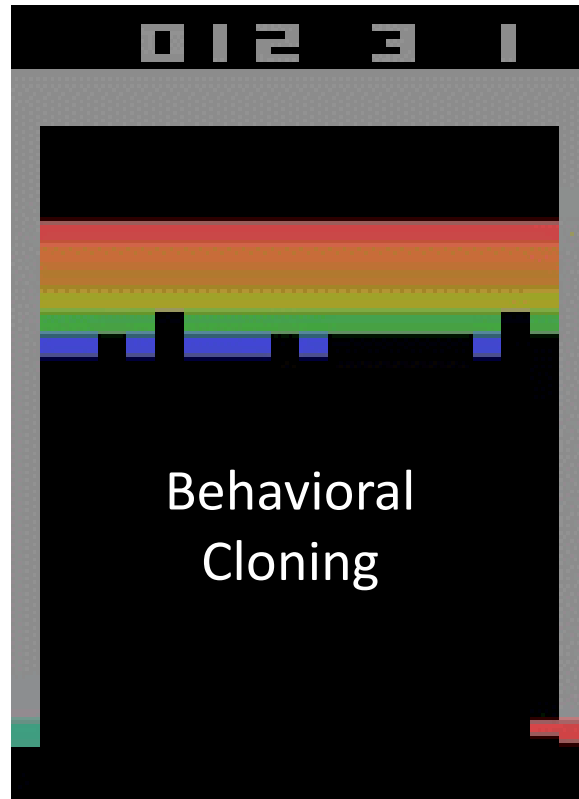
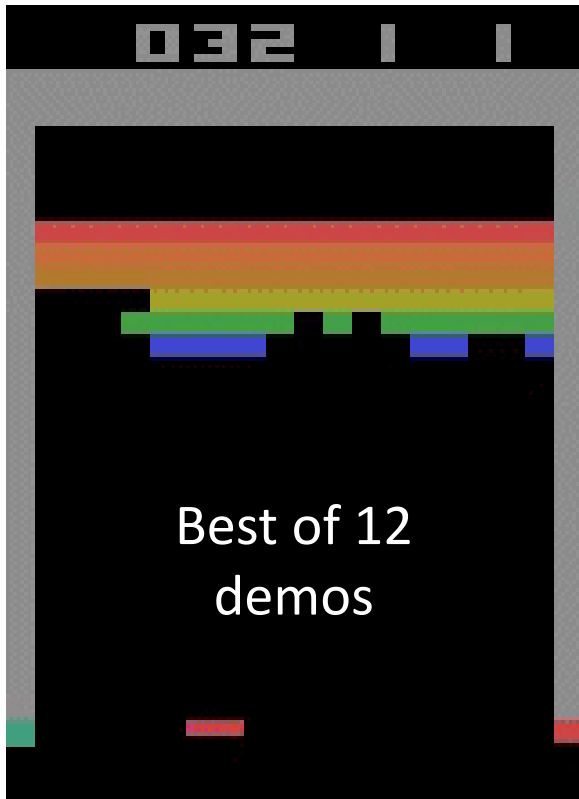
Best demo (Score = 84)



T-REX (Score = 520)

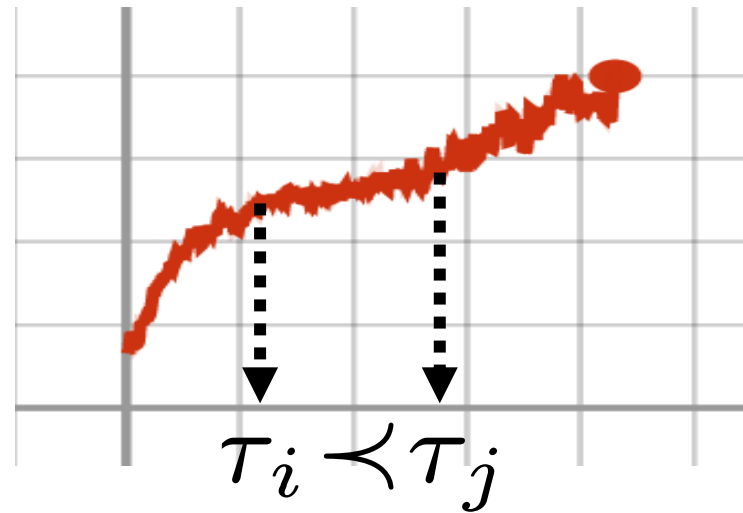
Uses only 12 ranked demonstrations

Atari Breakout

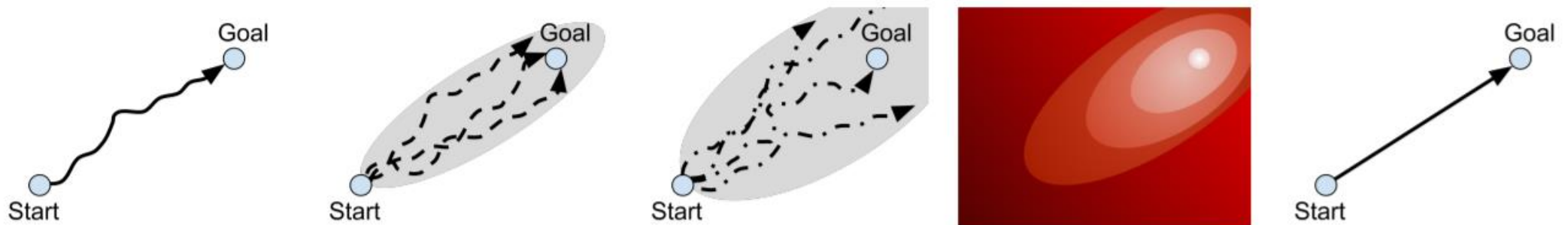


What if you don't have explicit preference labels?

Learning from a learner [ICML'19]

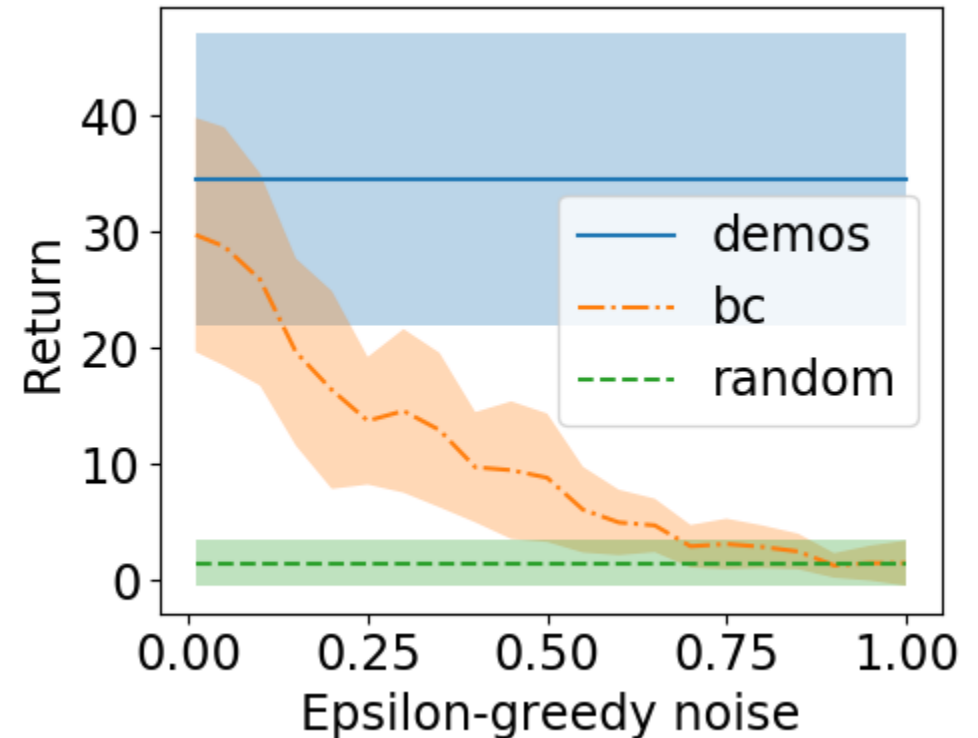


Automatic preference label generation [CoRL'20]



Automatic Rankings via Noise Injection

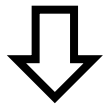
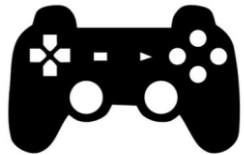
- Assumption: Demonstrator is significantly better than a purely random policy.
- Provides automatic rankings as noise increases.
- Generates a large diverse set of ranked demonstrations



Disturbance-based Reward Extrapolation (D-REX)

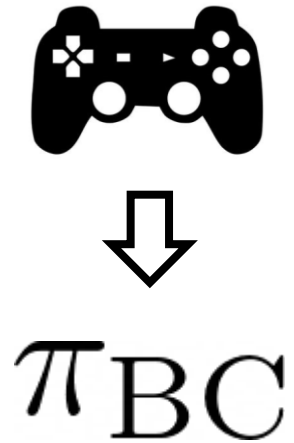
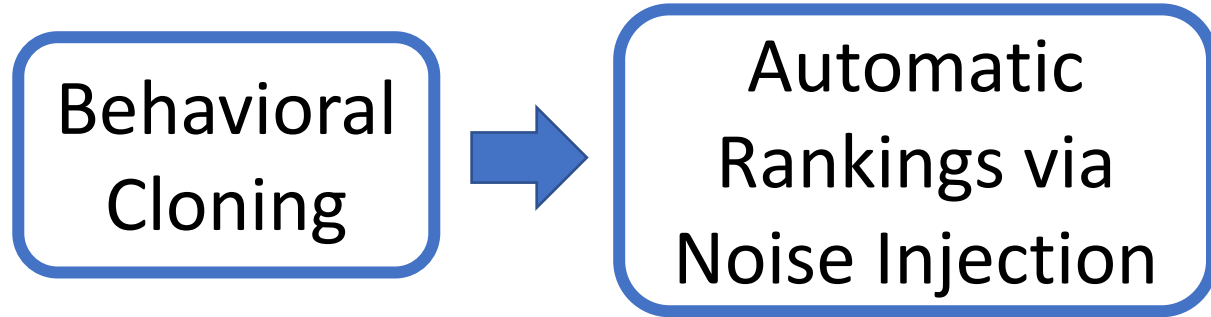
Behavioral
Cloning

$s \rightarrow a$



π_{BC}

Disturbance-based Reward Extrapolation (D-REX)



$\epsilon = 1.0$

\sim



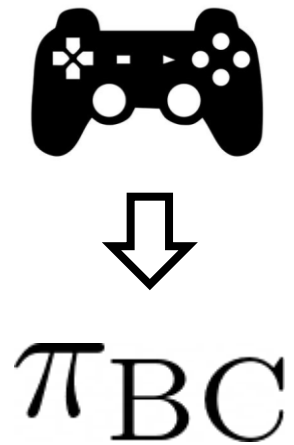
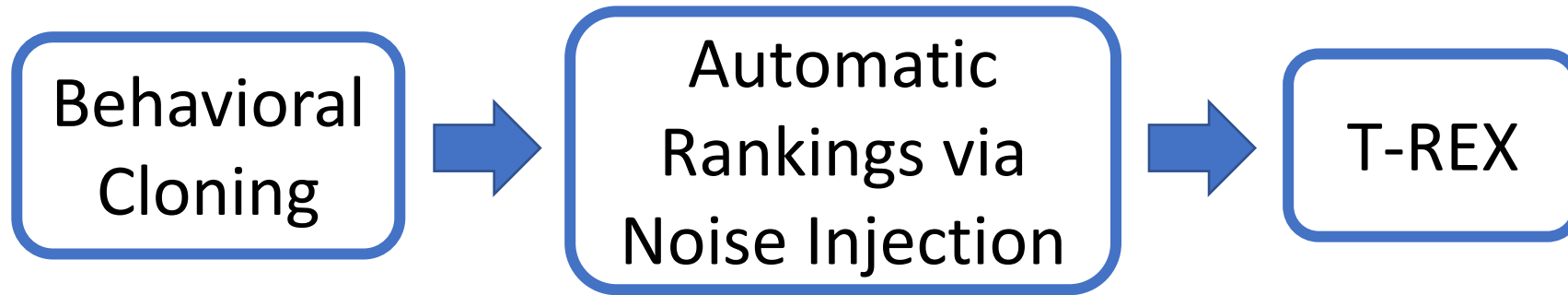
$\epsilon = 0.2$

\sim



$\epsilon = 0.01$

Disturbance-based Reward Extrapolation (D-REX)



$\epsilon = 1.0$

\simeq

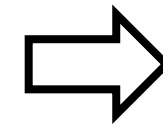


$\epsilon = 0.2$

\simeq

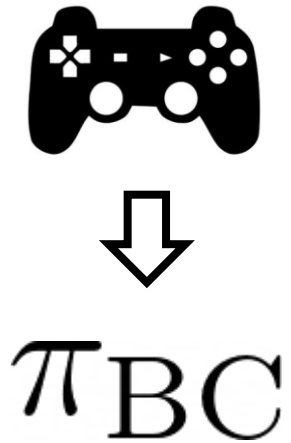
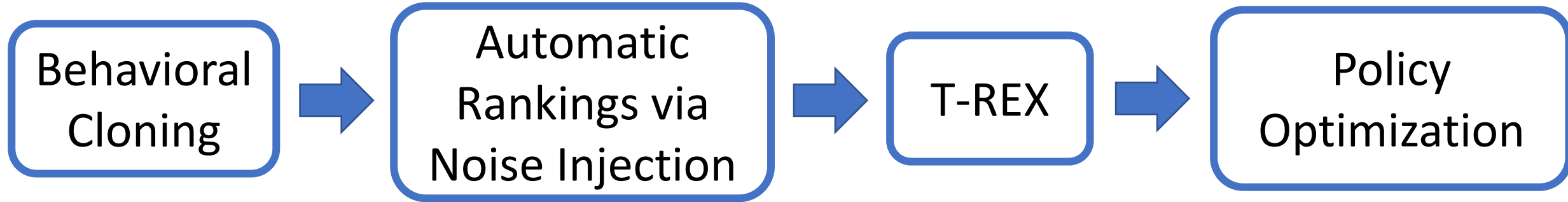


$\epsilon = 0.01$



Reward
Function
 $R(s)$

Disturbance-based Reward Extrapolation (D-REX)



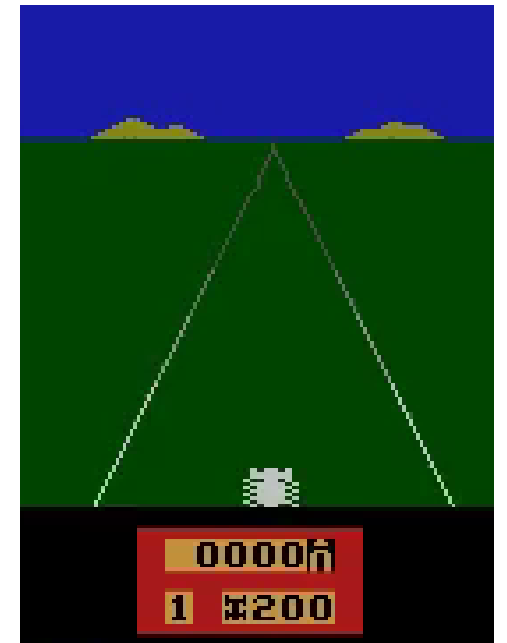
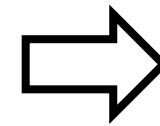
$\epsilon = 1.0$



$\epsilon = 0.2$



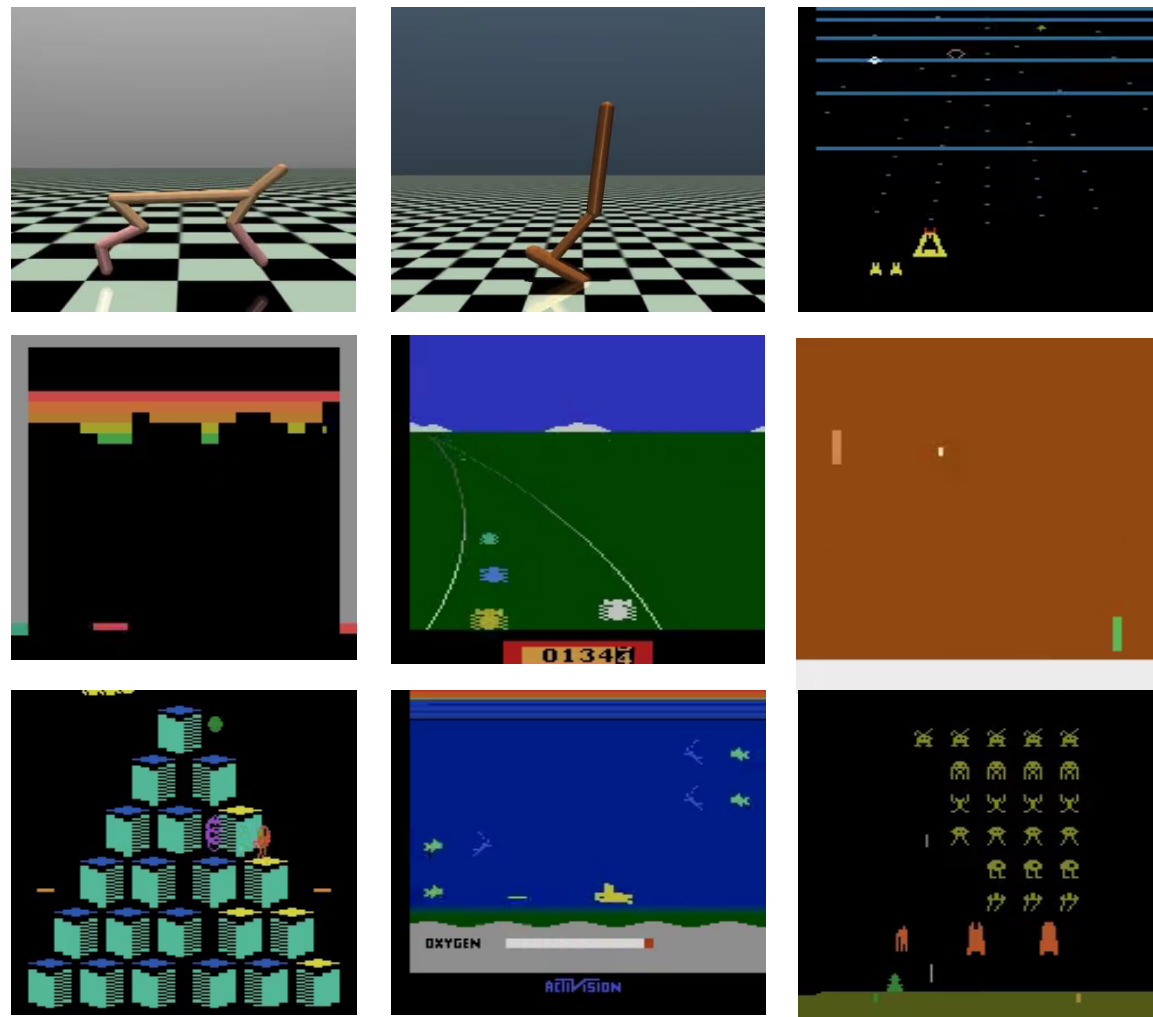
$\epsilon = 0.01$

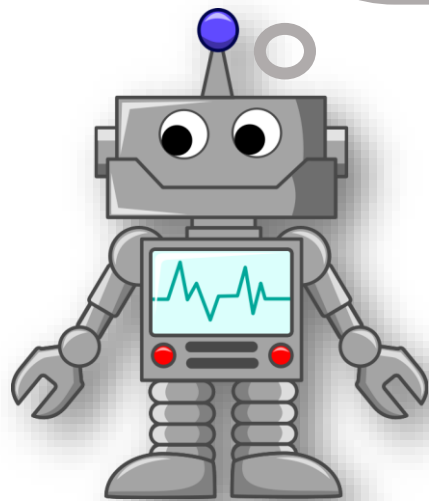
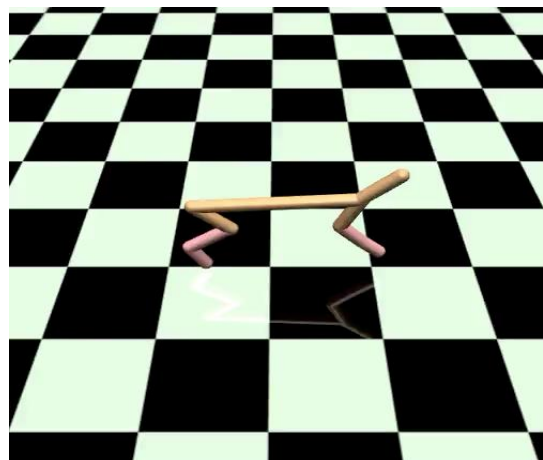


D-REX Policy

Experiments

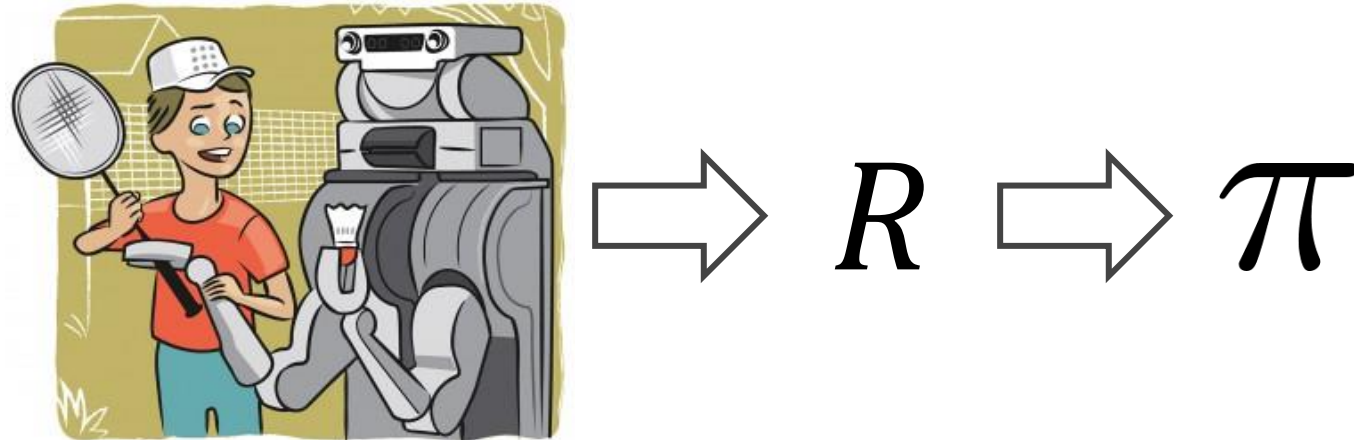
D-REX consistently outperforms the best demonstration as well as outperforming BC and GAIL.



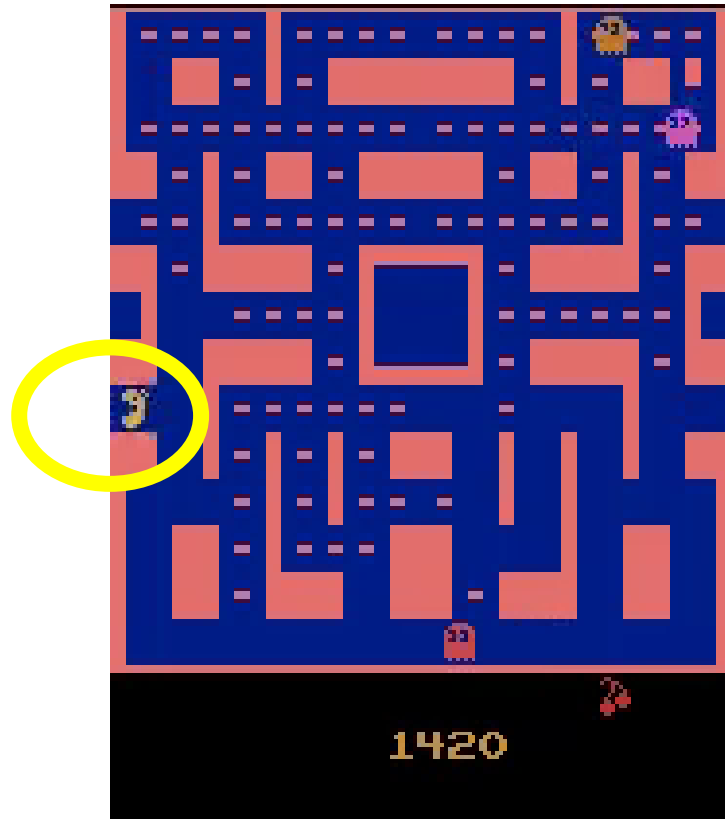


AI systems can **efficiently** infer human intent from **suboptimal demonstrations**.

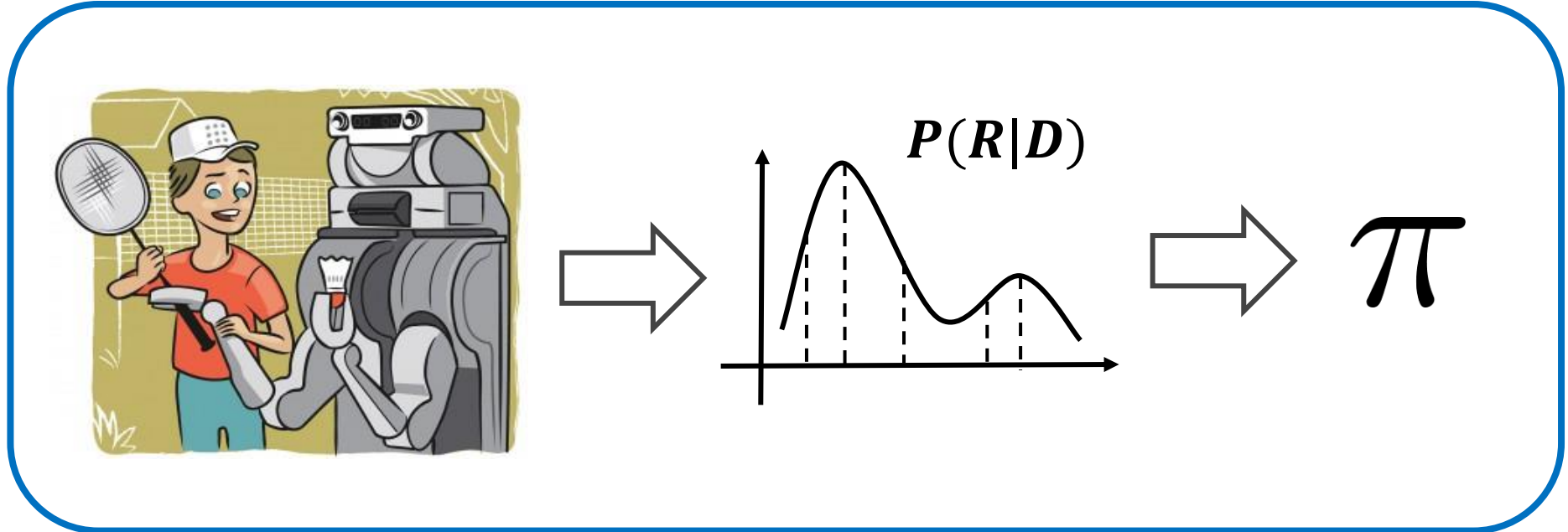
T-REX only learns a maximum likelihood estimate of the reward function.



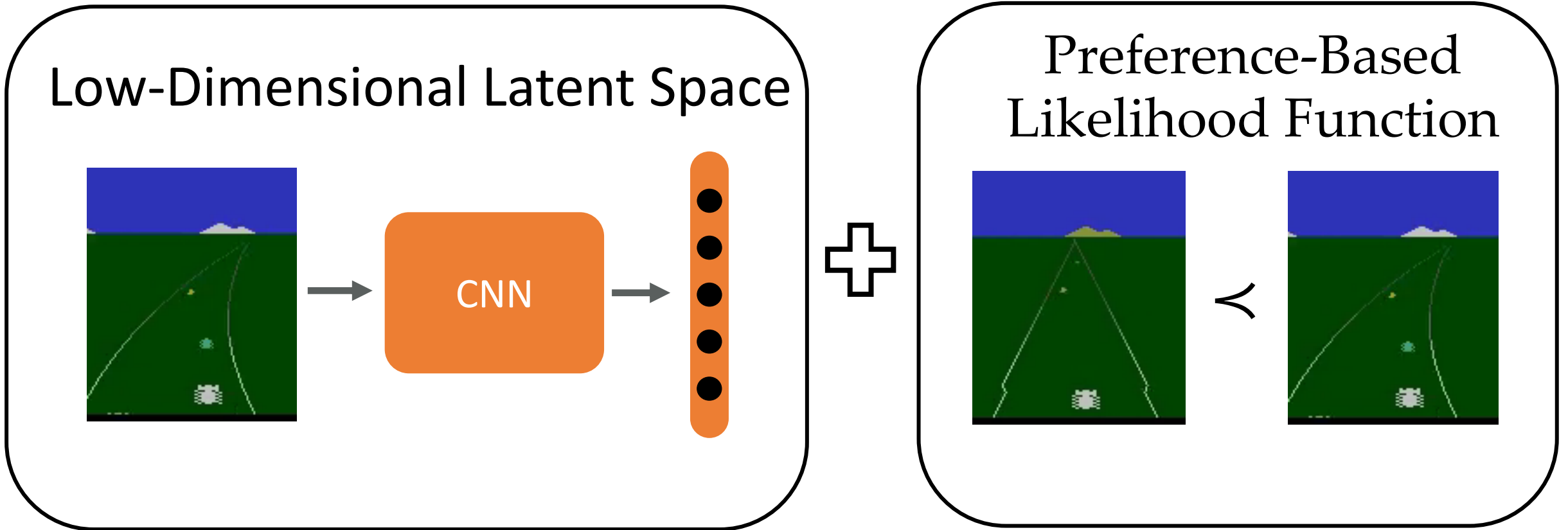
Reward Hacking



- Overfit to spurious correlations
- No consideration of alternative hypotheses



Idea: Fast Bayesian Inference



Next time: LLMs and ChatGPT

