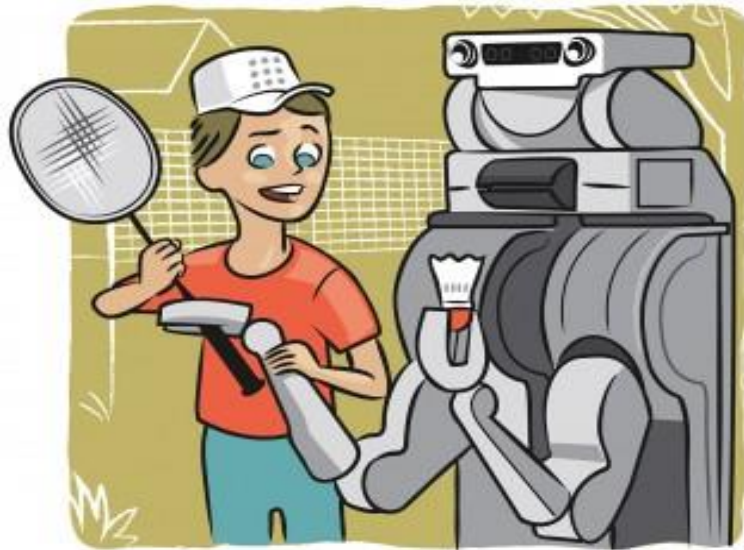


Inverse RL and Reward Learning from Preferences



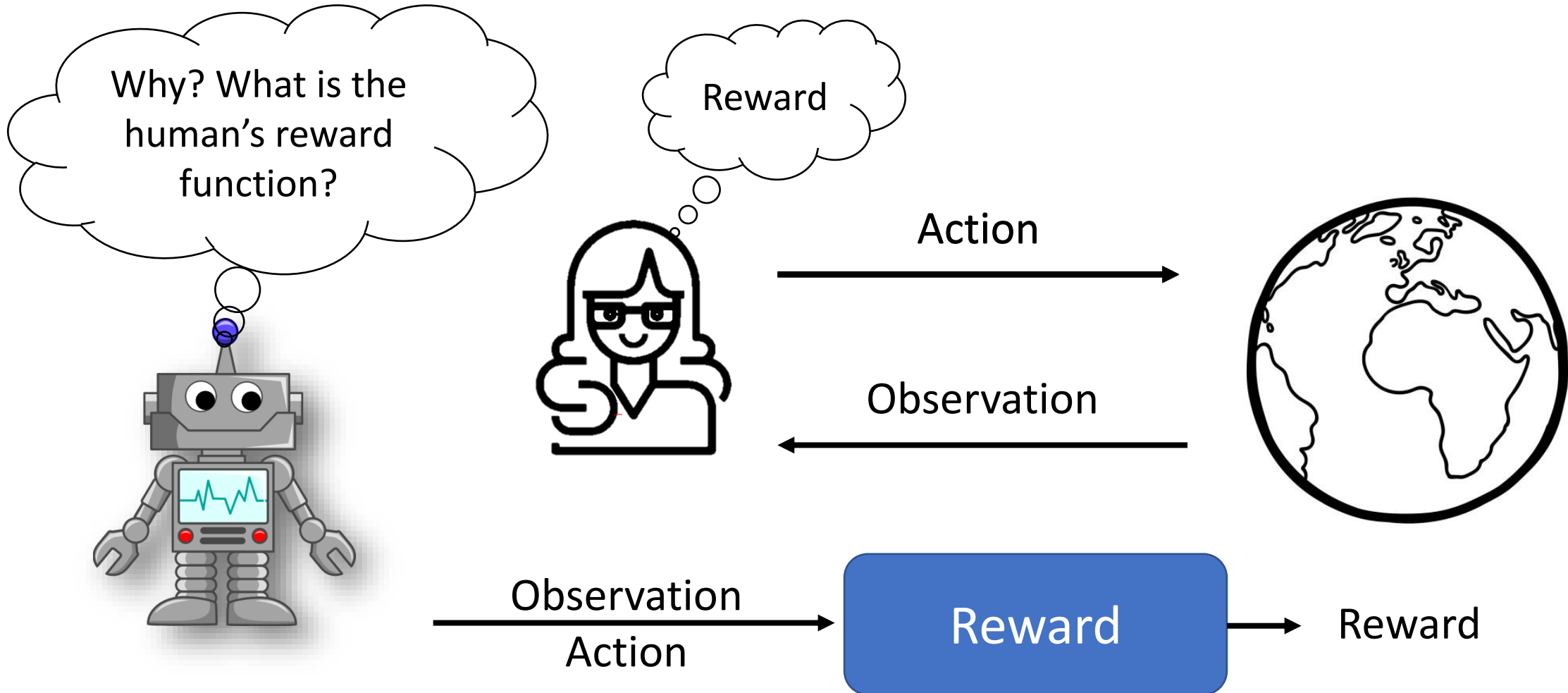
Instructor: Daniel Brown

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

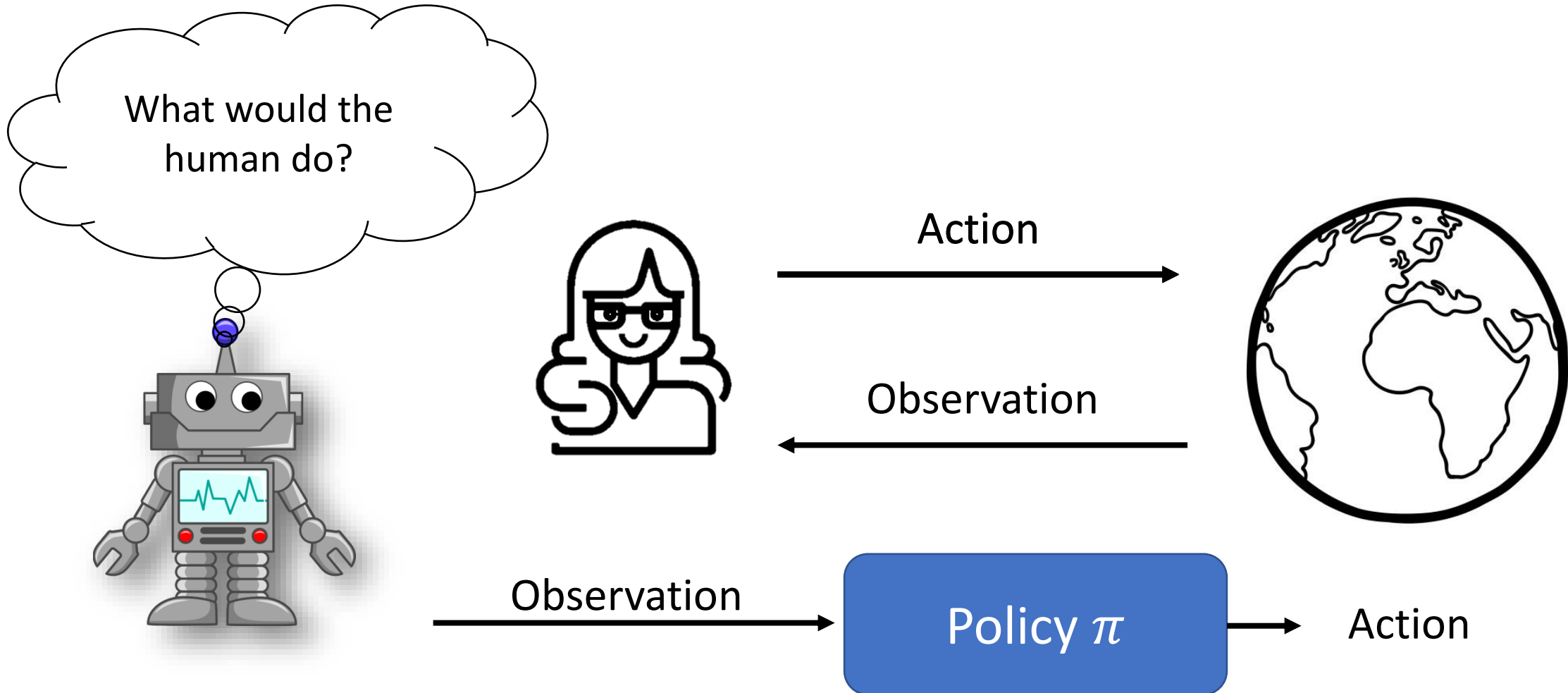
Course feedback is open

- Extra credit if class response rate is 70% or higher
 - Sliding scale if we reach 70%:
 - Extra credit points = $\text{response_rate_percentage} / 10$

Reward Learning (Inverse Reinforcement Learning)



Why not just imitate behavior? (Behavioral Cloning)





Human Intent Inference



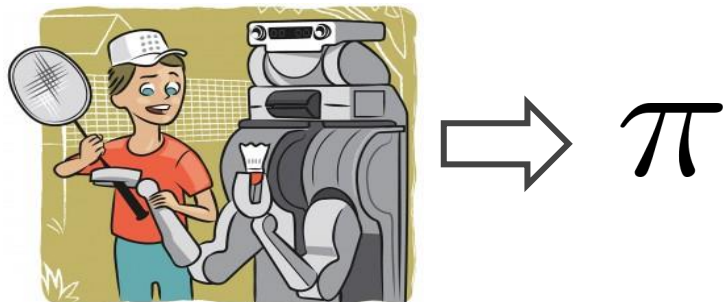
Inverse Reinforcement Learning

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

MDP/R

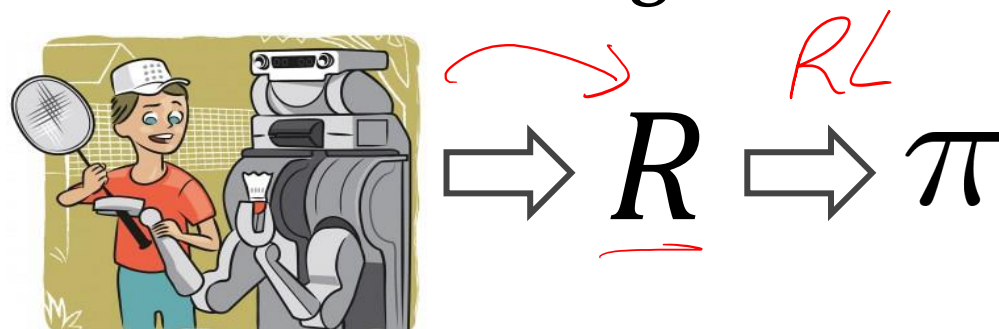
Imitation Learning

Behavioral Cloning



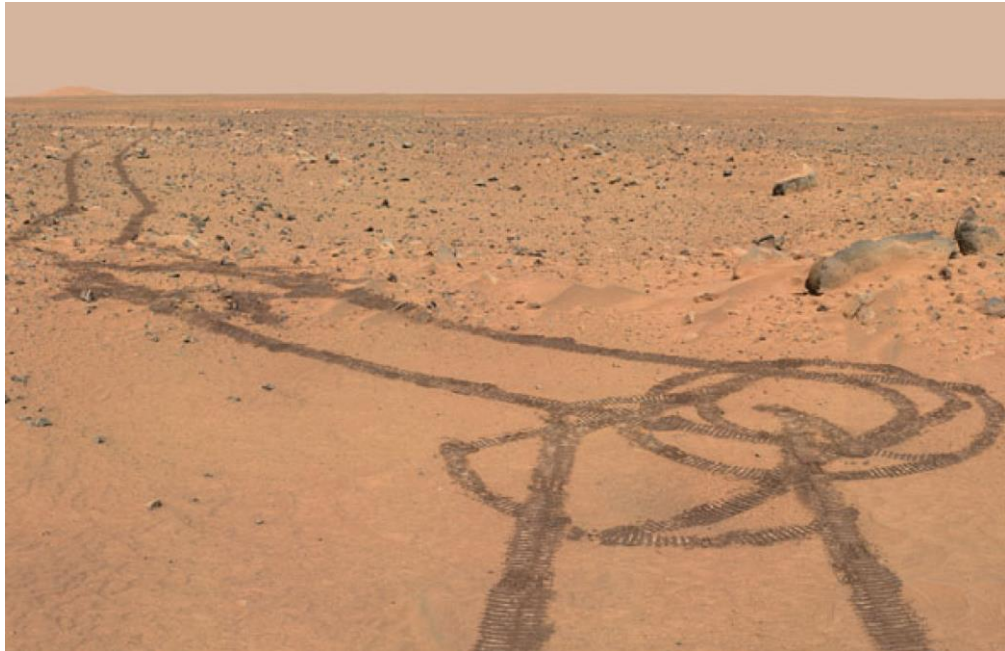
- Answers the “How?” question
- Mimic the demonstrator
- Learn mapping from states to actions
- Computationally efficient
- Compounding errors

Inverse Reinforcement Learning



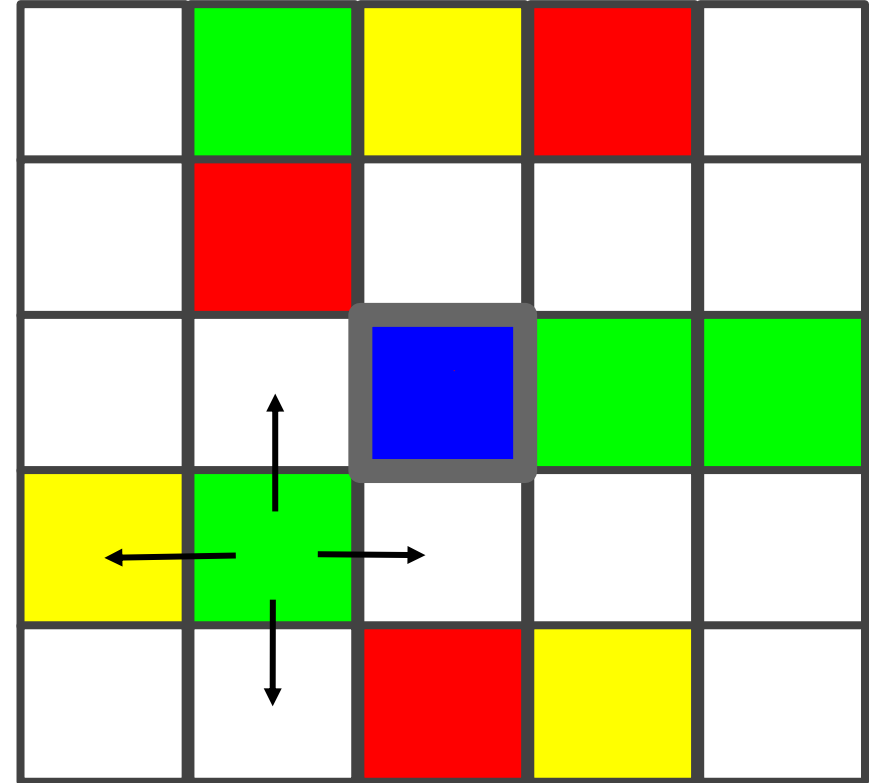
- Answers the “Why?” question
- Explain the demonstrator’s behavior
- Learn a reward function capturing the demonstrator’s intent
- Can require lots of data and compute
- Better generalization. Can recover from arbitrary states

IRL Example: Teaching a robot to navigate through demonstrations

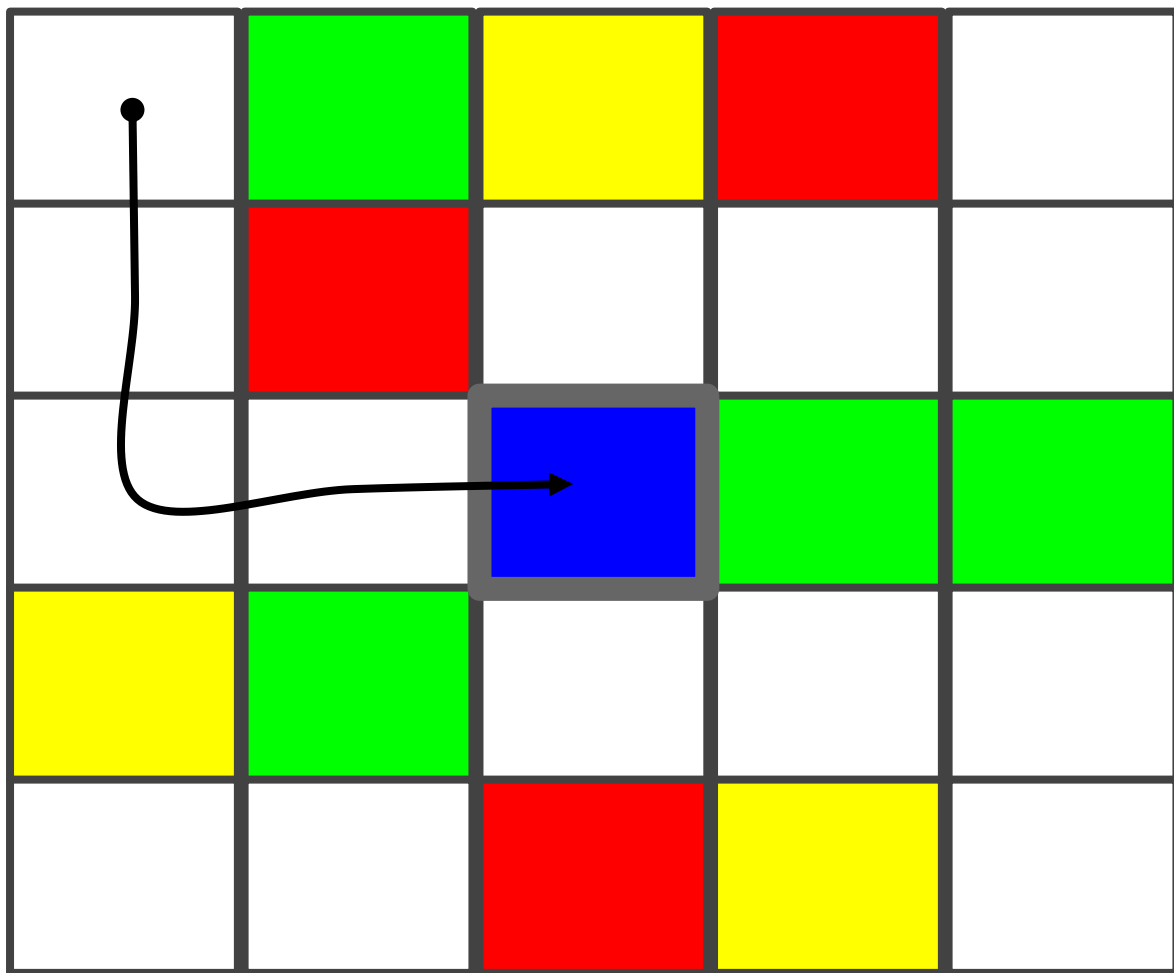


Toy version

T - deterministic



What is the reward?



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

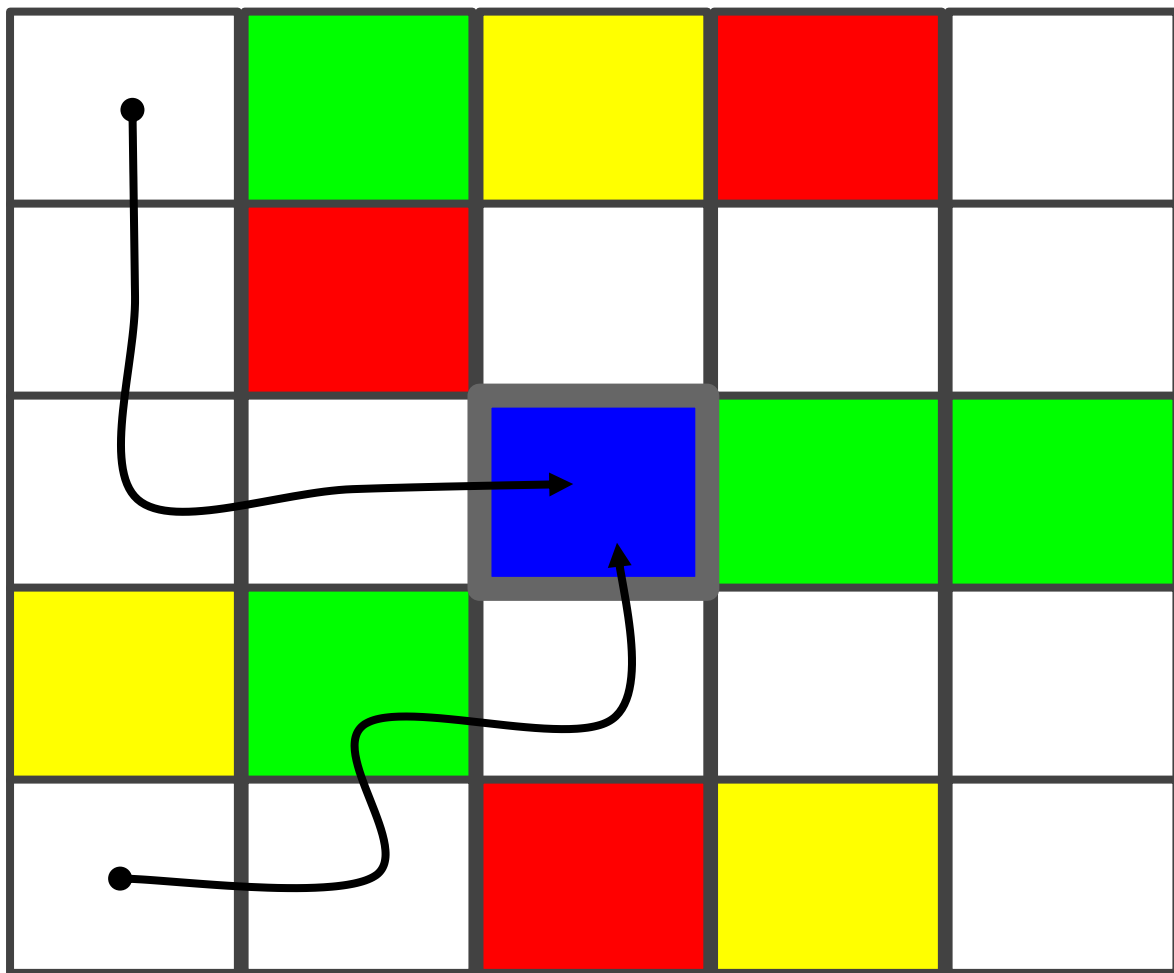
$$R(\text{Yellow}) = ?$$

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

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

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

What is the reward?



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

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

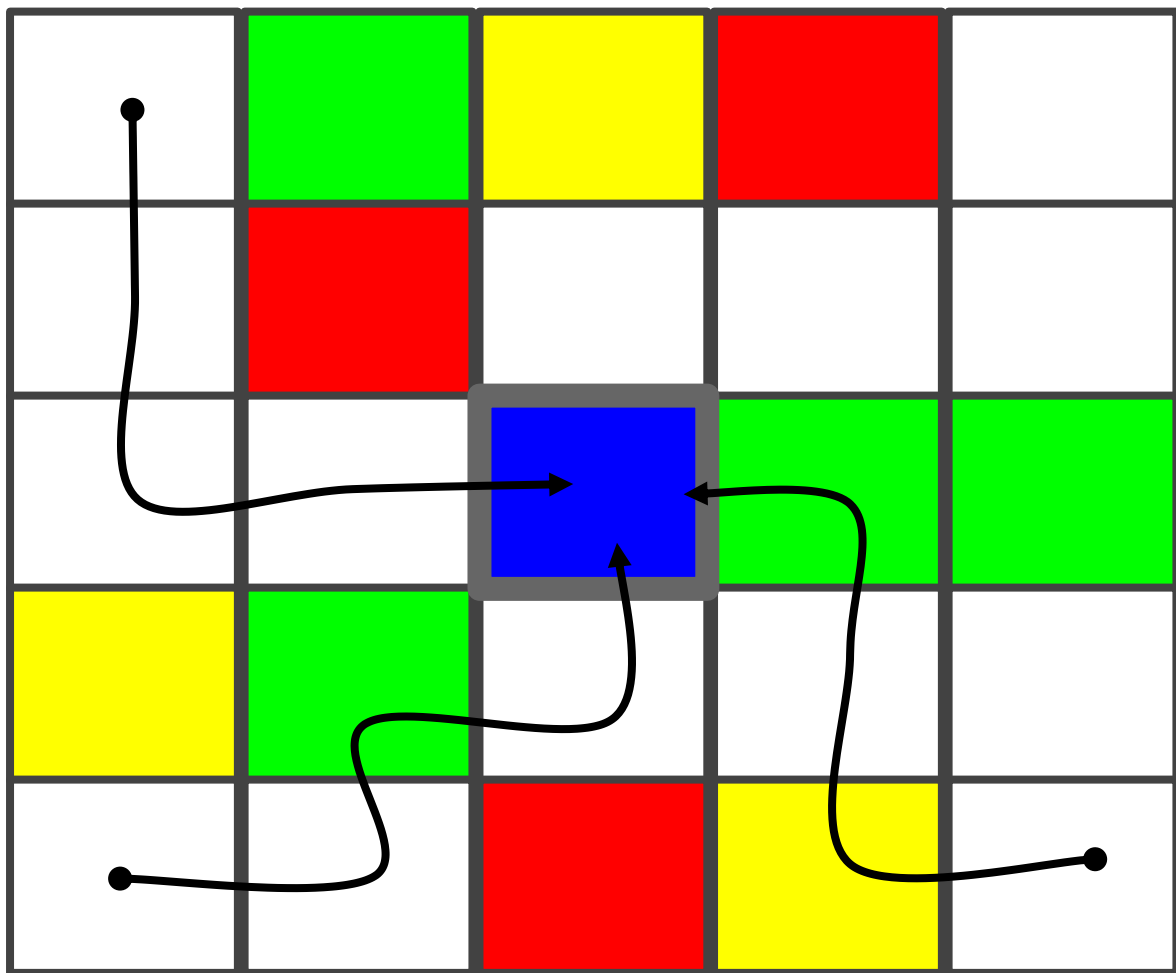
$$R(\text{Green}) = ?$$

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

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

What is the reward?

$\gamma = 1$



$$R(\text{blue square}) = ? \quad +10$$

$$R(\text{yellow square}) = ? \quad < 0 \quad +1$$

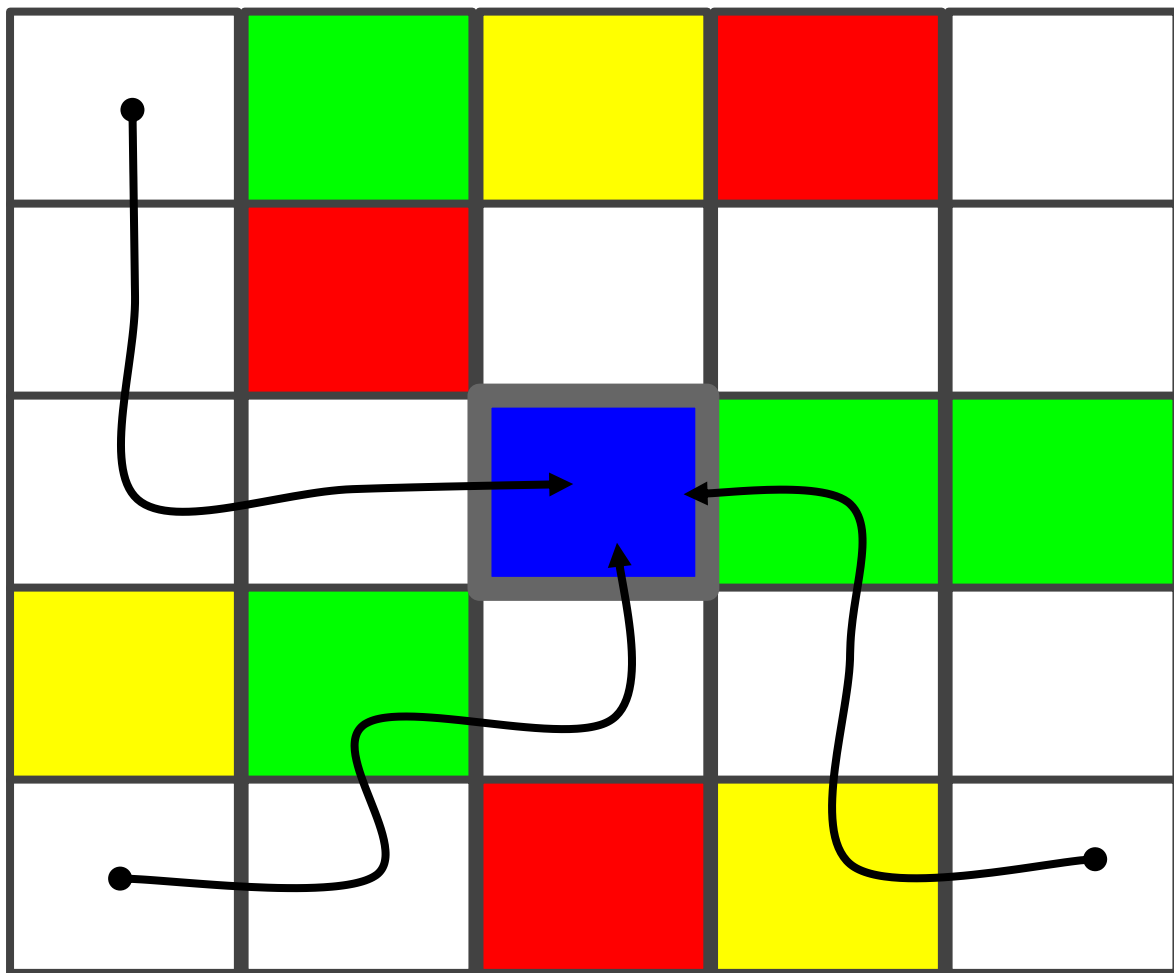
$$R(\text{green square}) = ? \quad < 0 \quad +1$$

$$R(\text{red square}) = ? \quad -10$$

$$R(\text{white square}) = ? \quad < 0 \quad +1$$

What is the reward?

$\gamma < 1$



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

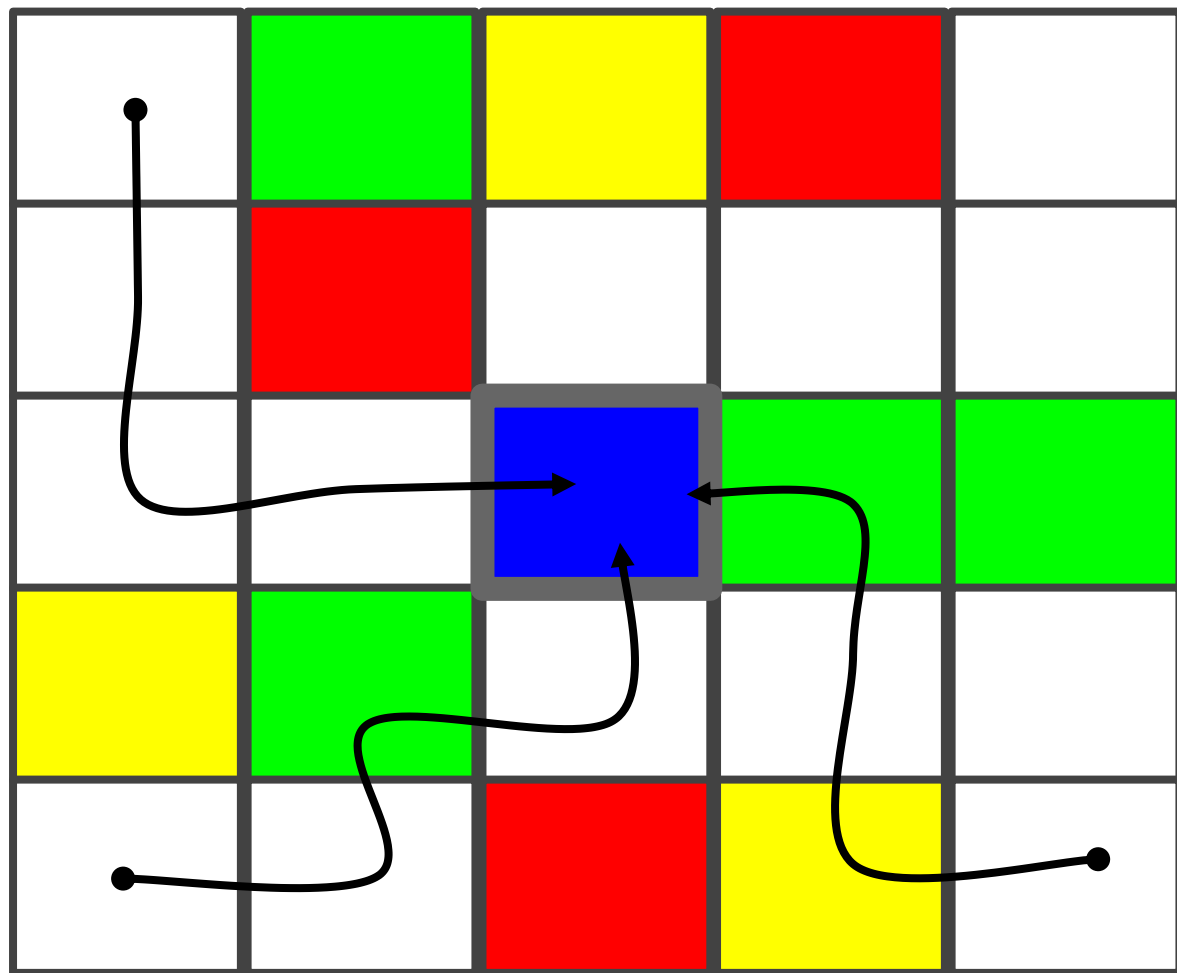
$$R(\text{Yellow}) = 0$$

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

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

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

What is the reward?



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

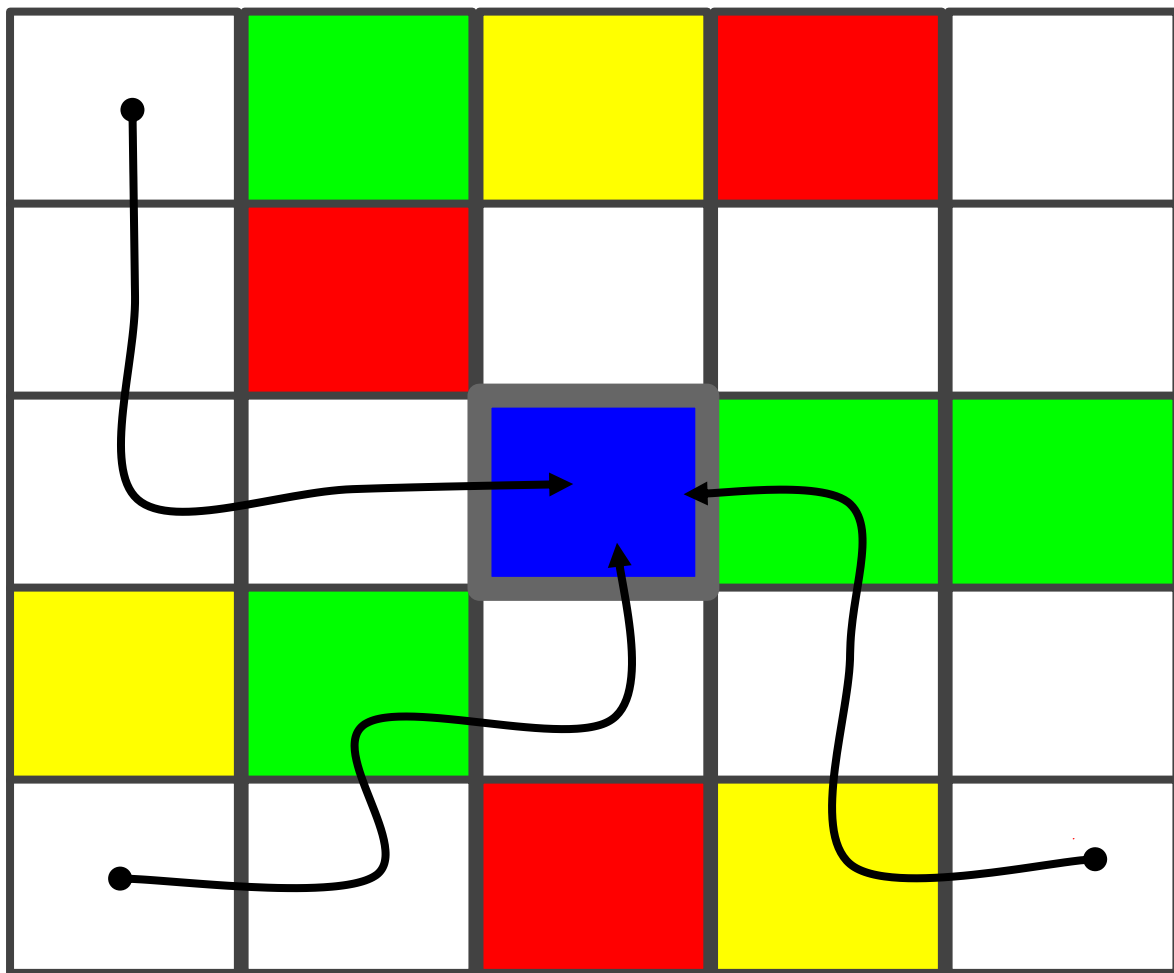
$$R(\text{Yellow}) = 0$$

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

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

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

What is the reward?



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

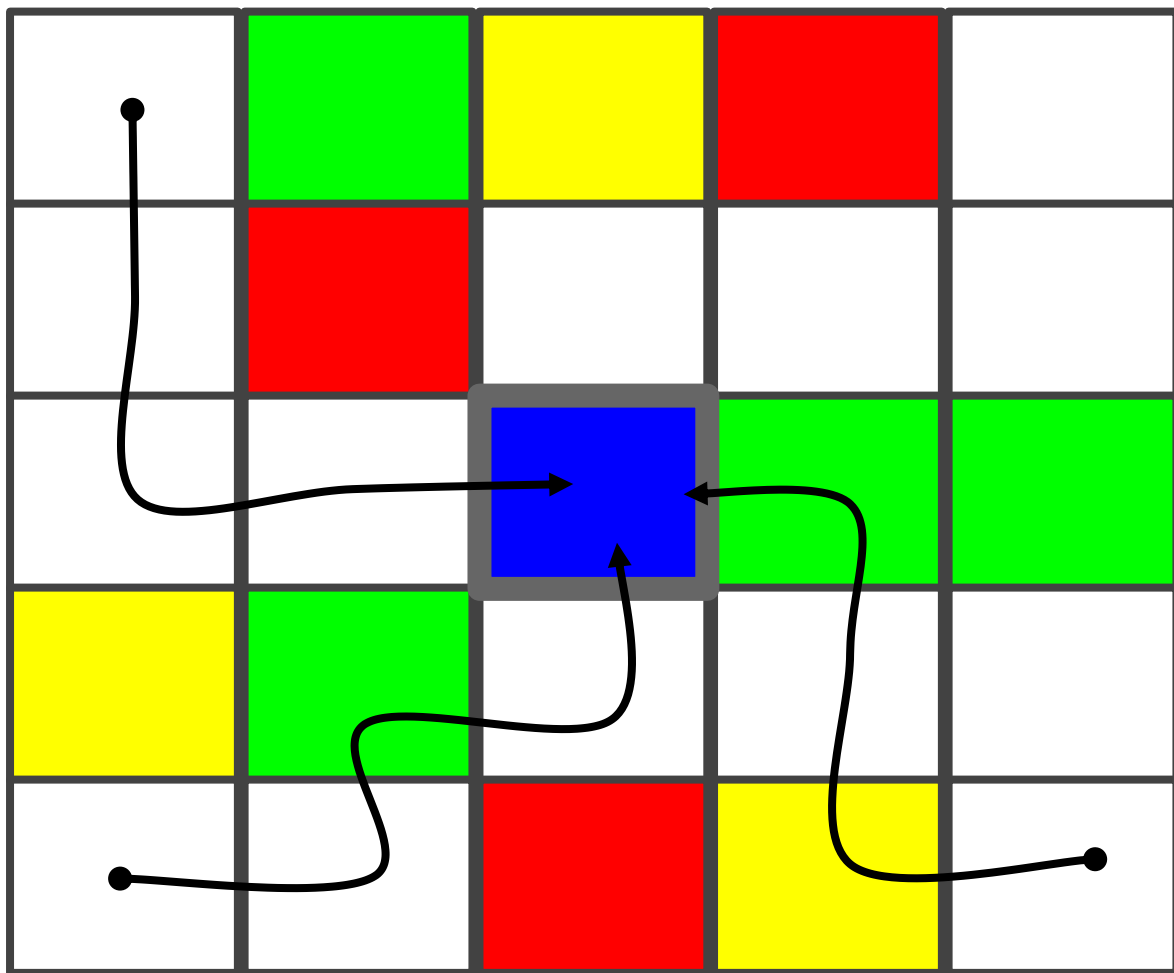
$$R(\text{Yellow square}) = -1$$

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

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

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

What is the reward?



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

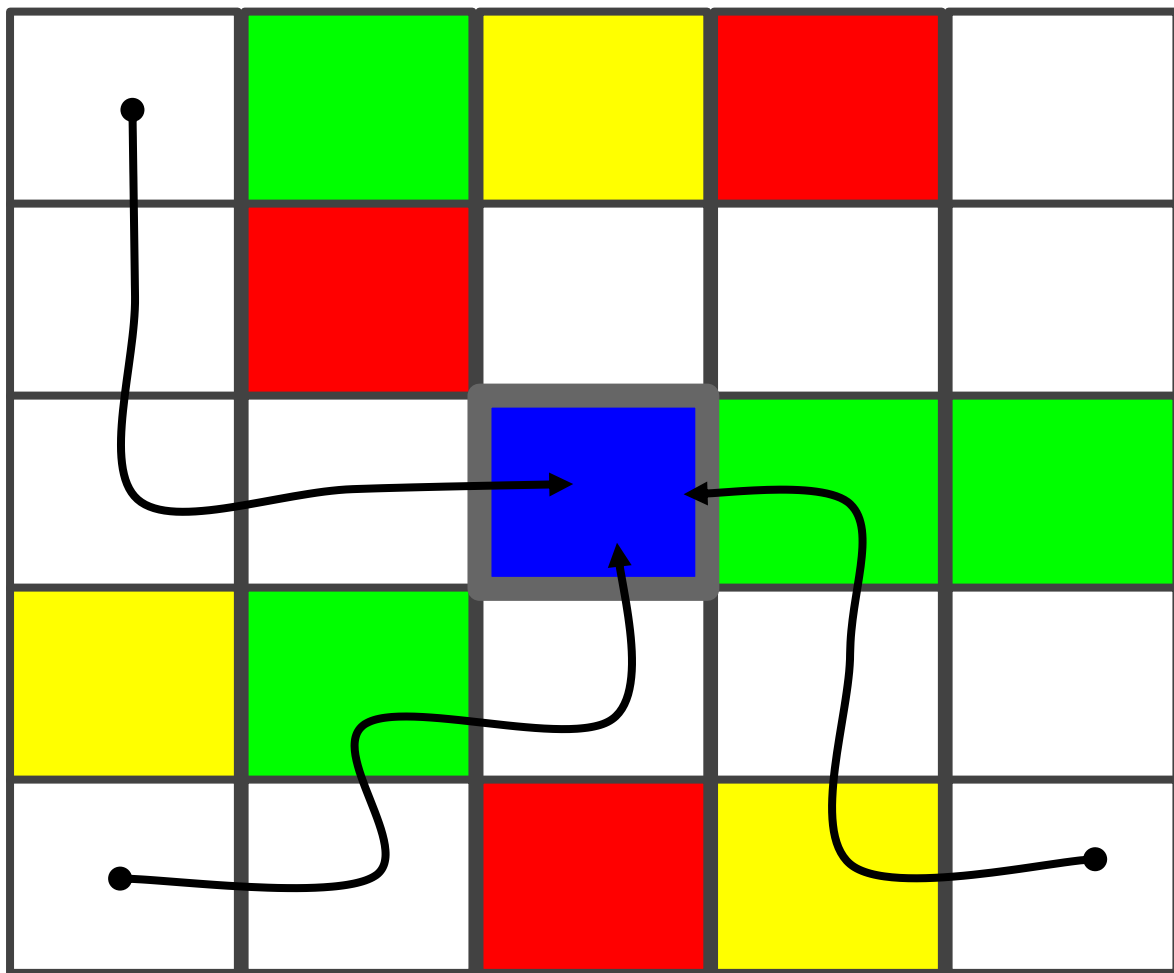
$$R(\text{Yellow}) = 0$$

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

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

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

What is the reward?



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

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

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

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

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

$C < 0$

Inverse Reinforcement Learning Formalism

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

How would you do this?

• maximize likelihood of D given \hat{R}

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

input

T, S, A, γ
MDP/R

computationally hard

2

Flashback: Approximate Q-Learning

$$Q(s, a) = w_1 f_1(s, a) + w_2 f_2(s, a) + \dots + w_n f_n(s, a)$$



- Q-learning with linear Q-functions:

$$\text{transition} = (s, a, r, s')$$

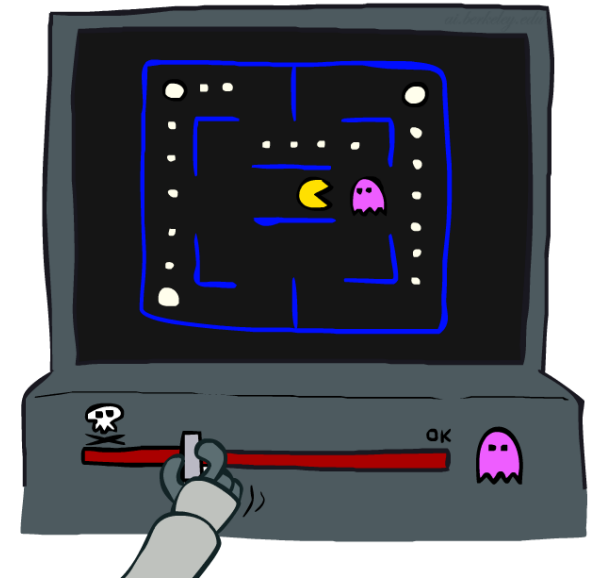
$$\text{difference} = \left[r + \gamma \max_{a'} Q(s', a') \right] - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [\text{difference}]$$

$$w_i \leftarrow w_i + \alpha [\text{difference}] f_i(s, a)$$

Exact Q's

Approximate Q's



- Intuitive interpretation:
 - Adjust weights of active features
 - E.g., if something unexpectedly bad happens, blame the features that were on: disprefer all states with that state's features
- Formal justification: online least squares

Feature count matching

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

$$R(s) = \mathbf{w}^T \phi(s) \subseteq \omega_1 \phi_1(s) + \omega_2 \phi_2(s) \dots$$

- 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:

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

Handwritten notes:

- Red arrows pointing to the word "features" in the first bullet point.
- Red arrow pointing to the term μ_π in the equation above.
- Red list of features: disc, # white, # red, # yellow, # blue.

Inverse reinforcement learning: feature matching

(Abbeel and Ng 2004, Syed and Schapire 2007)

$$\|x\|_\infty = \max_i |x_i|$$

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

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

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

difference in features

- 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)

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

- 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 $\mathbf{w} \leftarrow \mathbf{w} + \alpha(\mu_D - \mu_\pi)$

↑
step size

$$\mu_D[i] > \mu_\pi[i]$$

Soft Value Iteration

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

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

Soft Maximum



Policy is a softmax policy.

$$\begin{aligned} & e^{Q(A, S) - \log \sum_b e^{Q(b, S)}} \\ &= \frac{e^{Q(A, S)}}{\sum_b e^{Q(b, S)}} \end{aligned}$$

Soft Maximum

$$\log(e^{-b} \cdot e^b) = \log(e^0) = 0$$

- Let $a > b$

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

$$\log(e^a + e^b) + \log e^{-b} + \log e^b$$

$$\log((e^a + e^b)e^{-b}) + b$$

$$= \log(e^{a-b} + 1) + b$$

- If $a = b$

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

$$\log(2e^a)$$

$$= \log(2) + \log e^a = \log(2) + a$$

- In general $\max\{x_1, x_2, \dots, x_n\} \leq \log \sum_i e^{x_i} \leq \max\{x_1, \dots, x_n\} + \log n$

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 value of terminal states to 0 and other values to $-\infty$
- Repeat:
 - 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

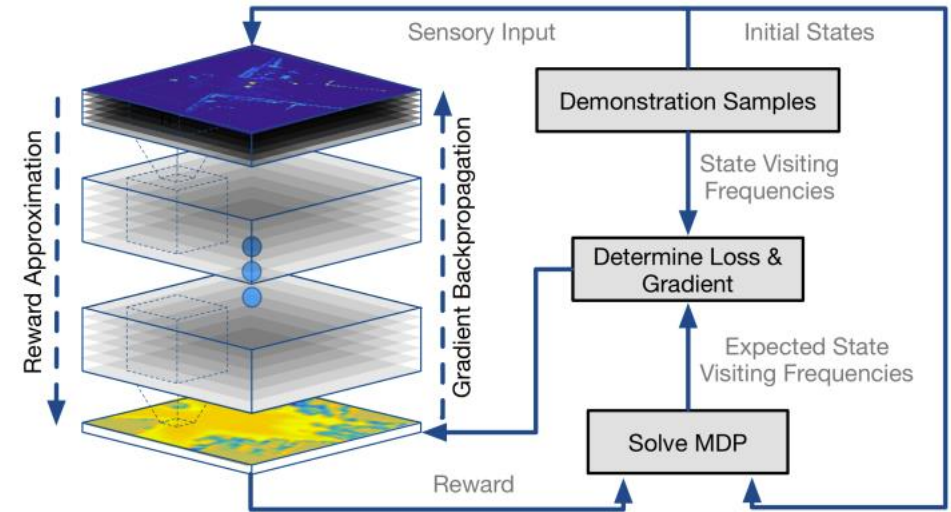
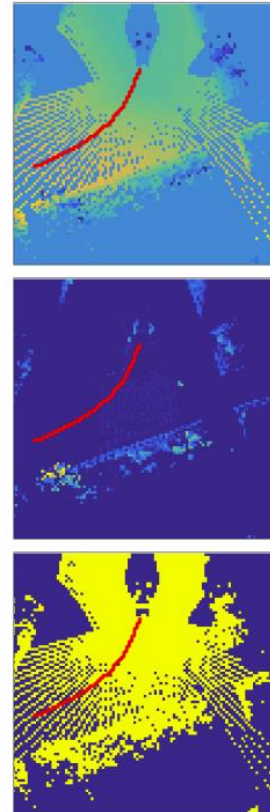


Fig. 1: Schema for training neural networks in the Maximum Entropy paradigm for IRL.

Another way to look at MaxEnt IRL

$$\int_{\tau} e^{R_{\theta}(\tau)} d\tau$$

$$P(\tau) \propto \frac{e^{R_{\theta}(\tau)}}{Z}$$

$$\begin{aligned} & \max_{\theta} \log P(\tau) \\ Z &= \int e^{R_{\theta}(\tau)} d\tau \\ & \max_{\theta} \log \frac{e^{R_{\theta}(\tau)}}{Z} \end{aligned}$$

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

$$|D| = N$$

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

$$\max_{\theta} \log e^{R_{\theta}(\tau)} - \log Z$$

How to avoid fully solving MDP

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

- Estimate Z with a finite set of trajectories Z_{τ} .
- Loop:
 - Update parameters θ 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

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

Uniform sampling to approximate Z.

Noisy perturbations of demonstrations to approximate Z

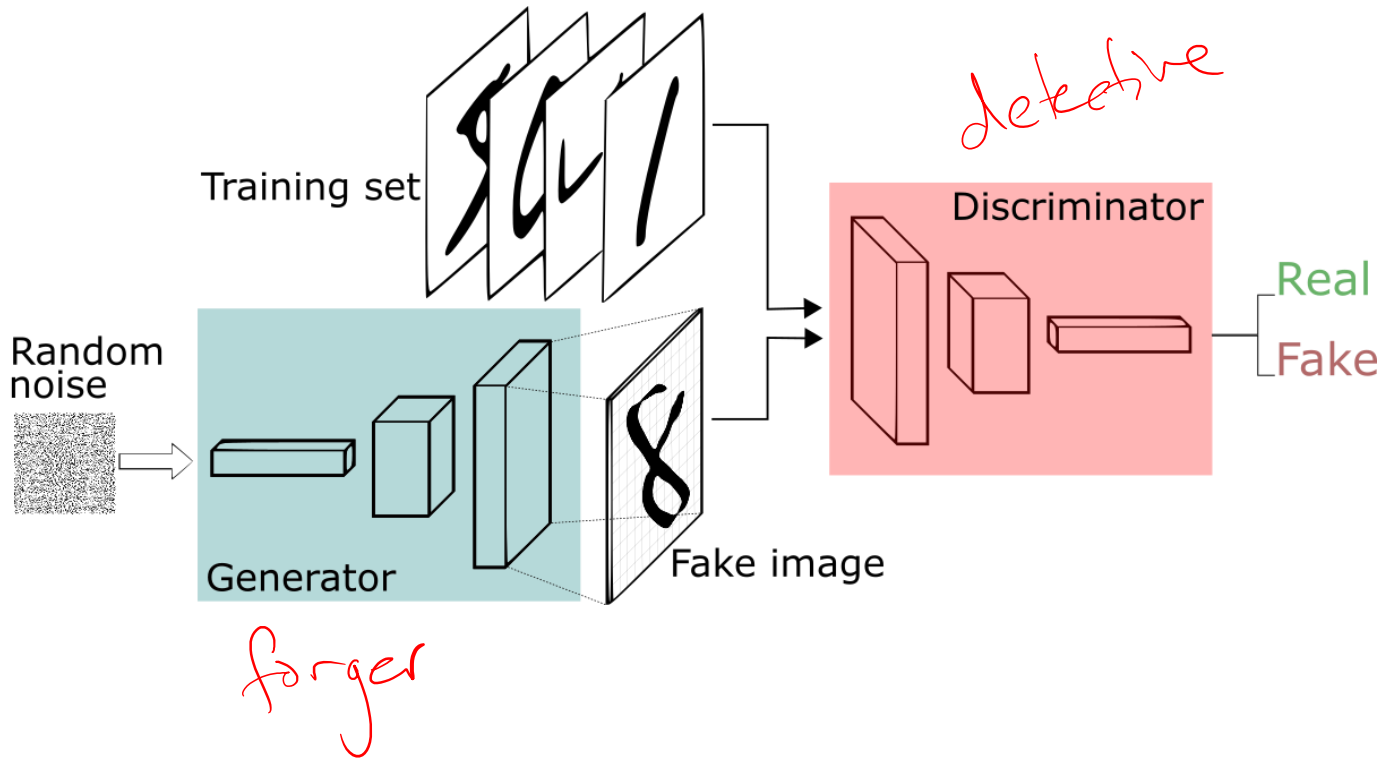
Use current policy to approximate Z.

~~✗~~ Alternate between a few steps of reward updates and a few steps of policy updates.

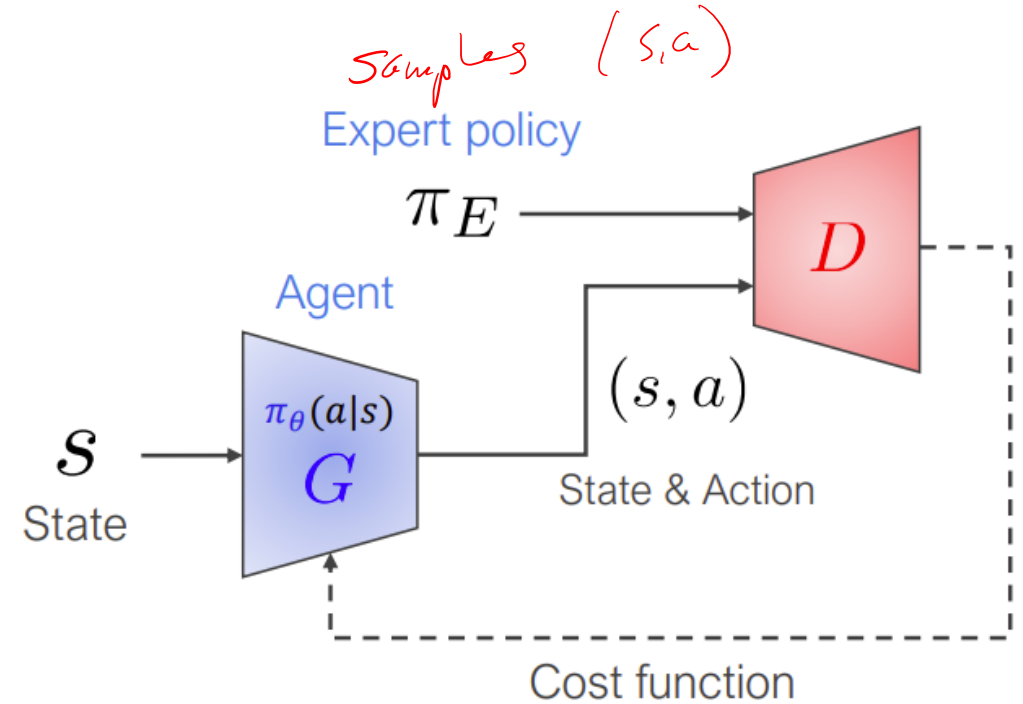
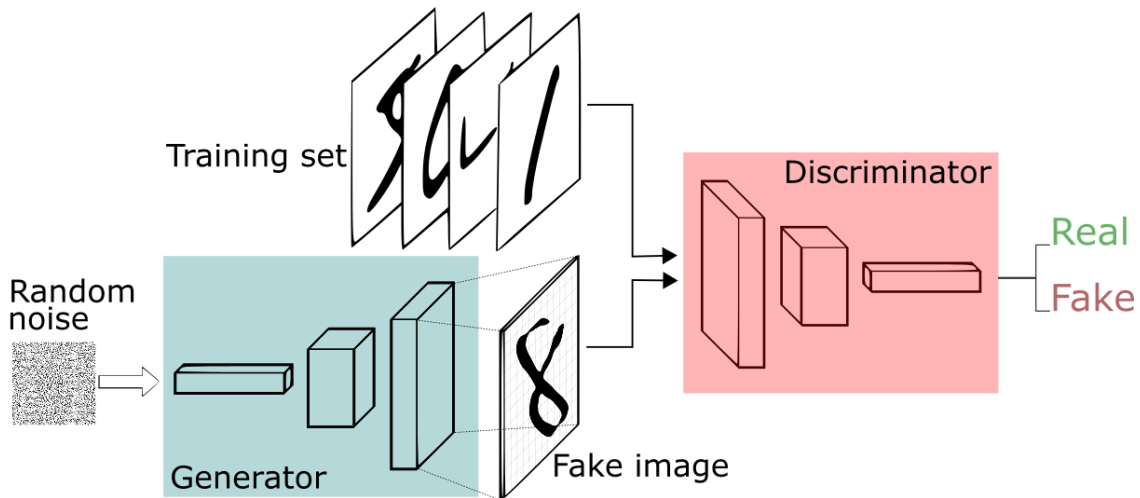
Finn et al. "Guided Cost Learning." 2016



GANs (Generative Adversarial Networks)



GAIL (Generative Adversarial Imitation Learning)



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

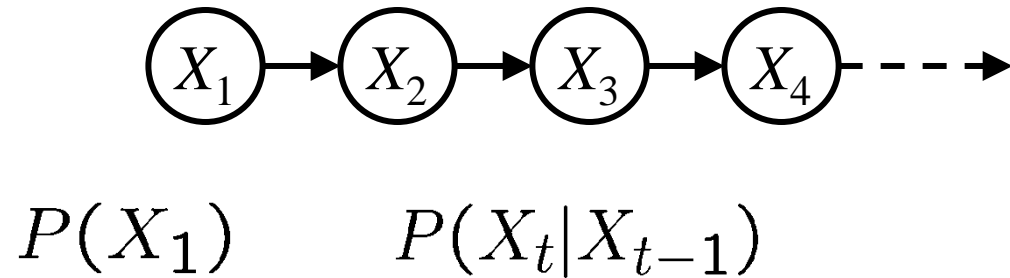
- Can we get a samples from the full Bayesian posterior?

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

likelihood *prior*

Markov Chain Monte Carlo (MCMC)

Markov chain:



Stationary Distribution: $P_\infty(X) = P_{\infty+1}(X) = \sum_x P(X|x)P_\infty(x)$

MCMC is a sampling approach for Bayesian inference where we construct a Markov chain such that the stationary distribution is the posterior distribution we care about.

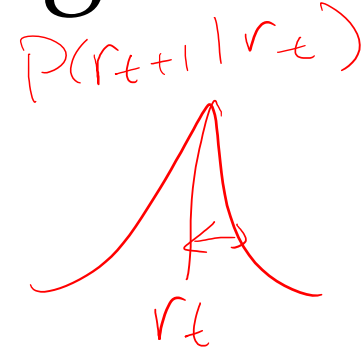
$$r_0 \rightarrow r_1 \rightarrow r_2 \rightarrow \dots$$

MCMC (Metropolis Hastings Algorithm)

- We want to sample from $P(R|D)$
- Start with random sample r_0
- Loop
 - Sample $r' \sim q(R_{t+1}|r_t)$
 - With probability $\min\left\{1, \frac{P(r'|D)}{P(r_t|D)}\right\}$ set $r_{t+1} = r'$
 - Else set $r_{t+1} = r_t$

proposal

Assume q is symmetric. For example, a Gaussian distribution with mean x_t and standard deviation σ



$r_{t+1} = r'$
 ~~$r_{t+1} = r_t$~~
 posterior ratio

Accept!
 Reject!



Normalizing constant cancels in the ratio!

$$\frac{P(r'|D)}{P(r|D)} = \frac{\frac{P(D|r')P(r')}{P(D)}}{\frac{P(D|r)P(r)}{P(D)}} = \frac{P(D|r')P(r')}{P(D|r)P(r)}$$

Bayesian Inverse Reinforcement Learning

(Ramachandran and Amir 2007)

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

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

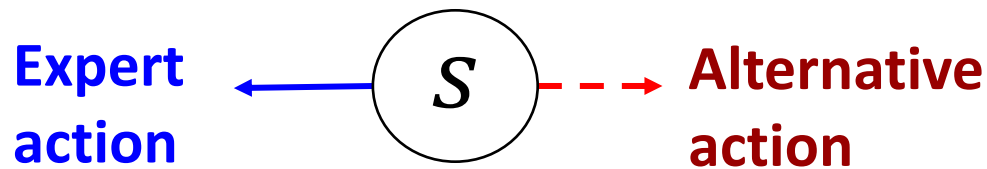
$Q^*(s, a, R) =$ How much reward will I expect to see if I take action a in state s and act optimally thereafter.

Bayesian Inverse Reinforcement Learning

(Ramachandran and Amir 2007)

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

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



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

Bayesian Inverse Reinforcement Learning

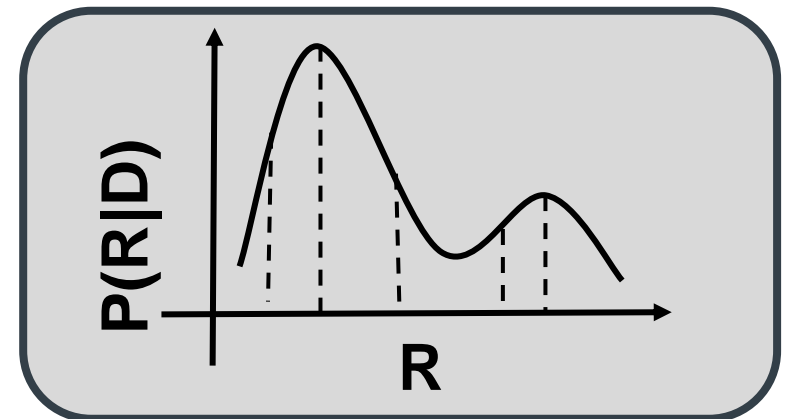
(Ramachandran and Amir 2007)

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

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

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

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



Applications of Bayesian IRL

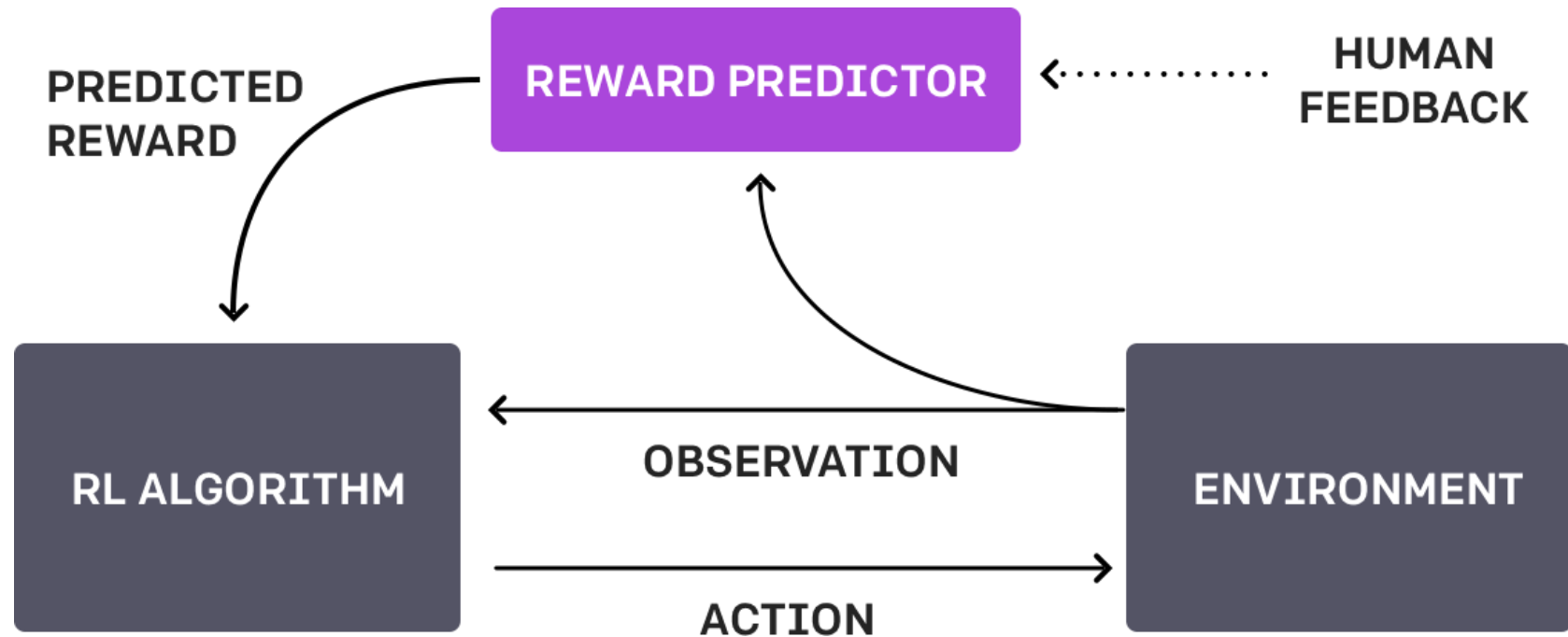
$P(R|D)$

- Active Learning — the algorithm picks its training data
- Uncertainty Estimation
- Demonstration Sufficiency

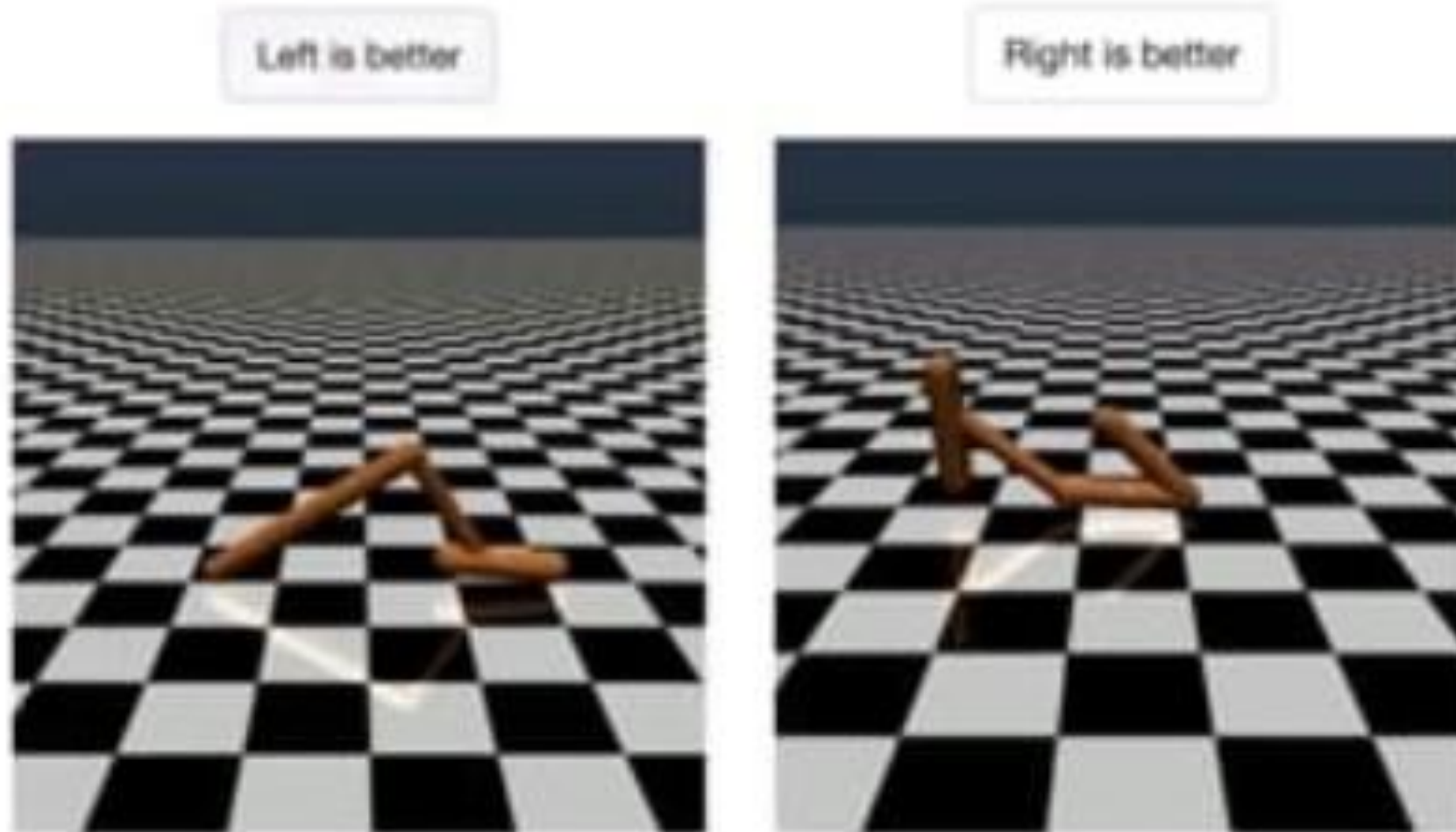
Autonomous Assessment of Demonstration Sufficiency via Bayesian Inverse Reinforcement Learning



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

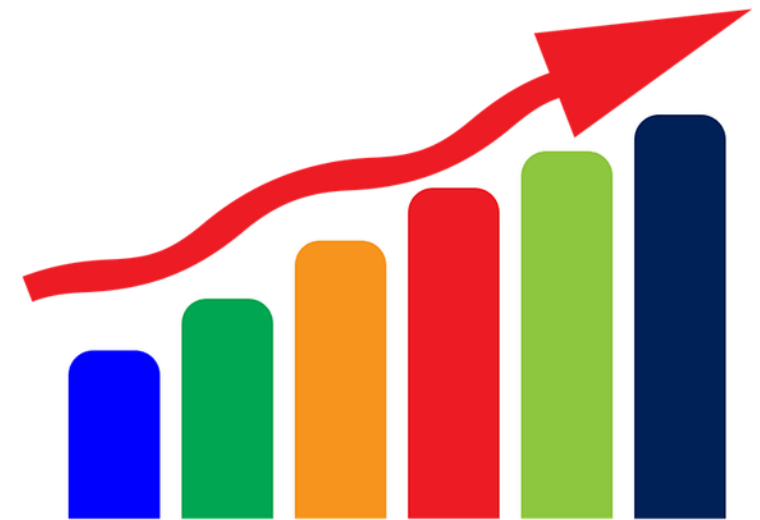
Prior approaches ...

~~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 ...

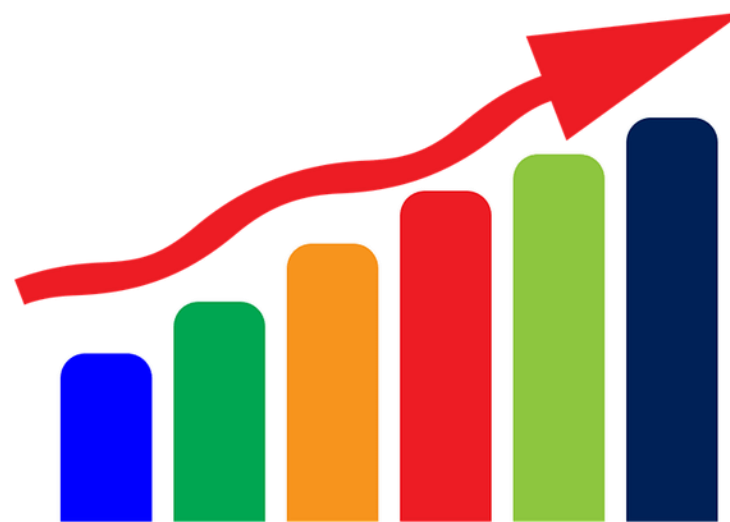
~~1. Typically couldn't do much better than the demonstrator.~~

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.

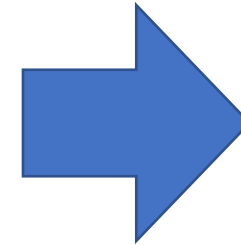
Pre-Ranked
Demonstrations



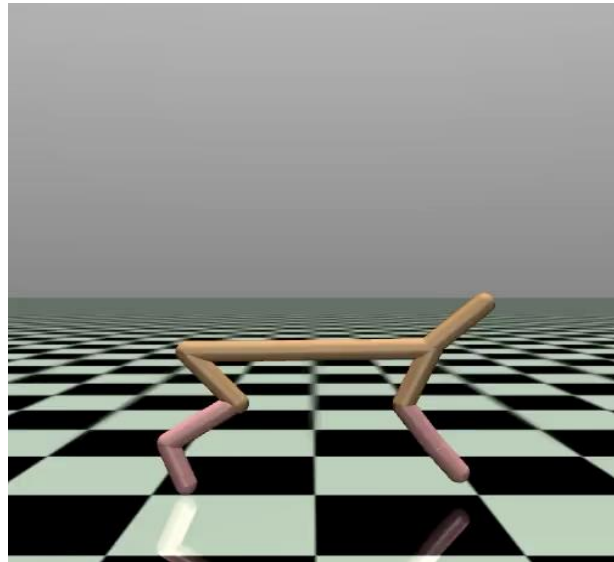
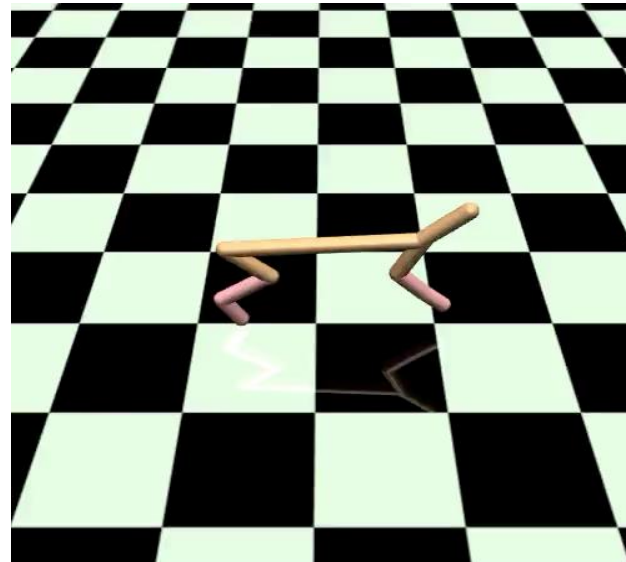
Trajectory-ranked Reward Extrapolation (T-REX)



Reward
Function

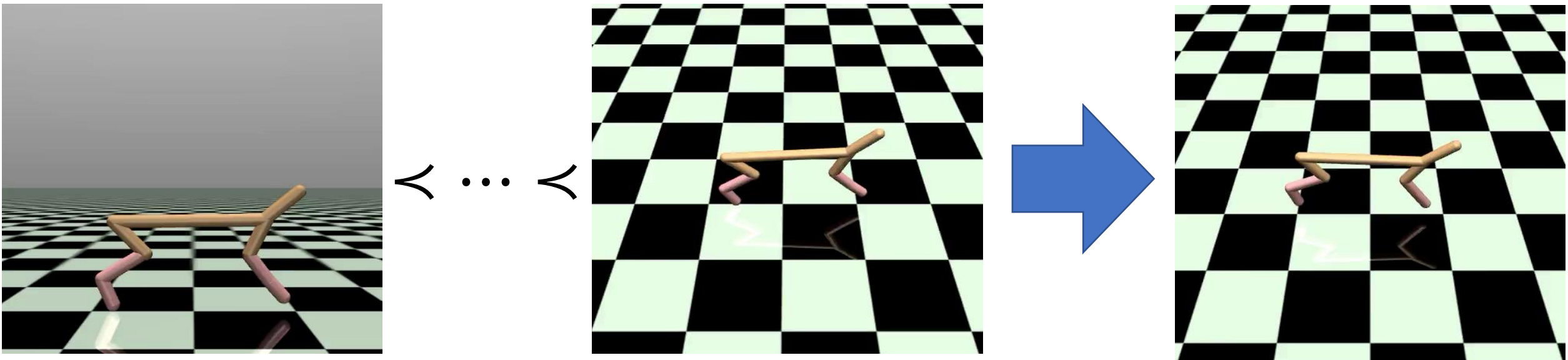


$\wedge \dots \wedge$



Pre-ranked demonstrations

Trajectory-ranked Reward Extrapolation (T-REX)



Pre-ranked demonstrations

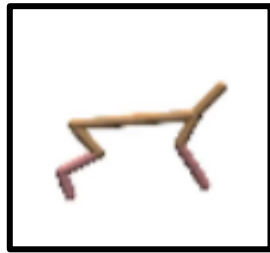
T-REX Policy

Reward Function

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

Examples of \mathcal{S} :

Current Robot Joint
Angles and Velocities



→ 0.5



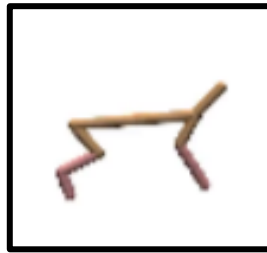
→ -0.7

Reward Function

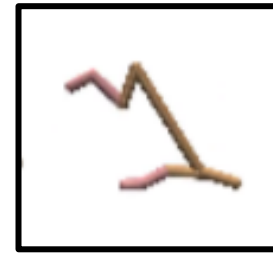
$$R_{\theta}: \mathcal{S} \rightarrow \mathbb{R}$$

Examples of \mathcal{S} :

Current Robot Joint
Angles and Velocities



→ 0.5



→ -0.7

Short
Sequence of
Images

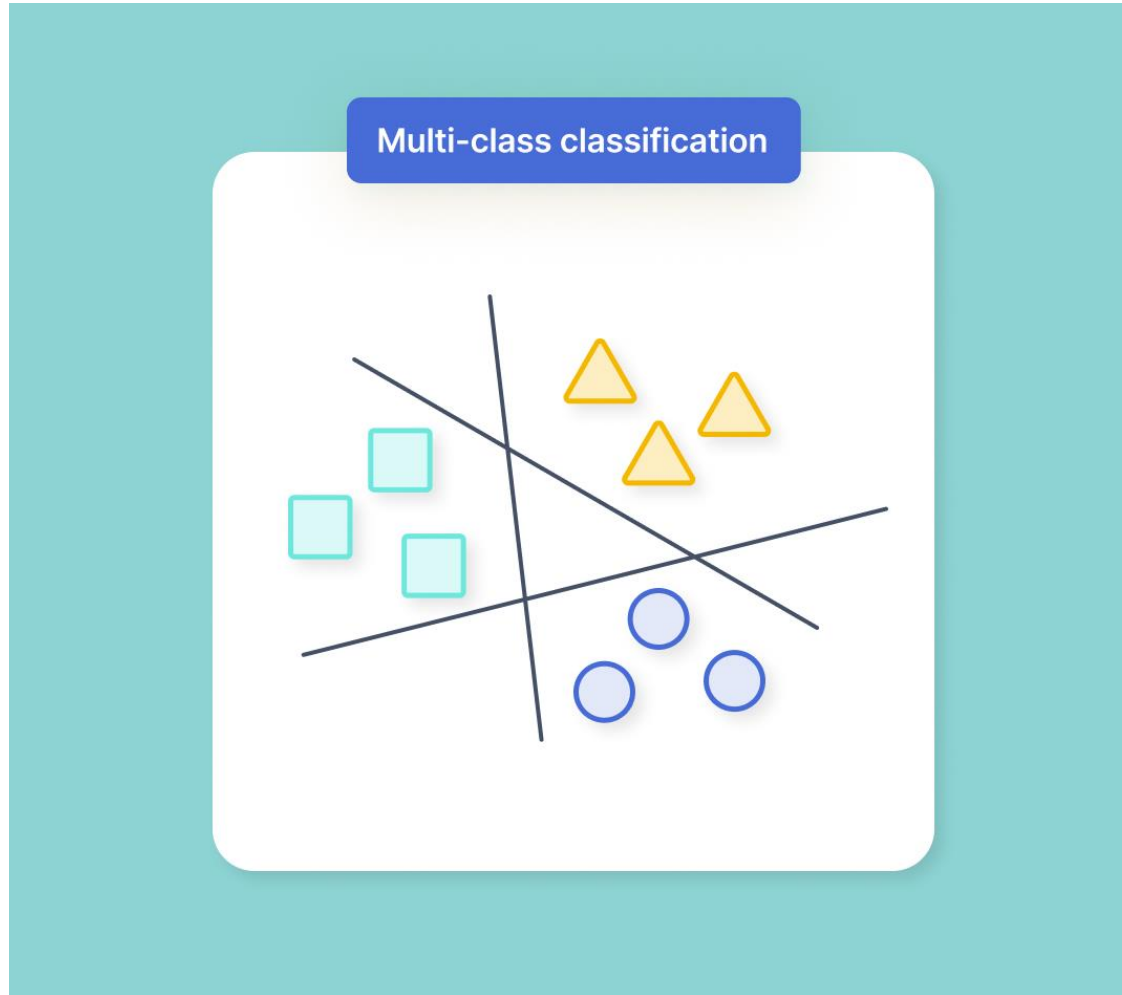


→ 0.9



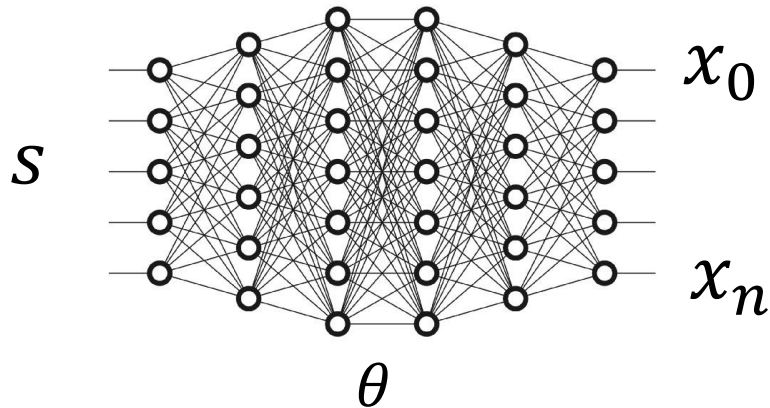
→ -1.2

Binary Classification and the Cross Entropy Loss



Flashback: How should we parameterize our policy?

- We need to be able to do two things:
 - Sample actions $a_t \sim \pi_\theta(\cdot | s_t)$
 - Compute log probabilities $\log \pi_\theta(a_t | s_t)$
- Categorical (classifier over discrete actions)
 - Typically, you output a value x_i for each action (class) and then the probability is given by a softmax equation



$$\pi_\theta(a_i | s) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

Cross Entropy

$$H(p, q) = - \sum_{x \text{ classes}} p(x) \log q(x)$$

True probability distribution (one-shot)

Your model's predicted probability distribution

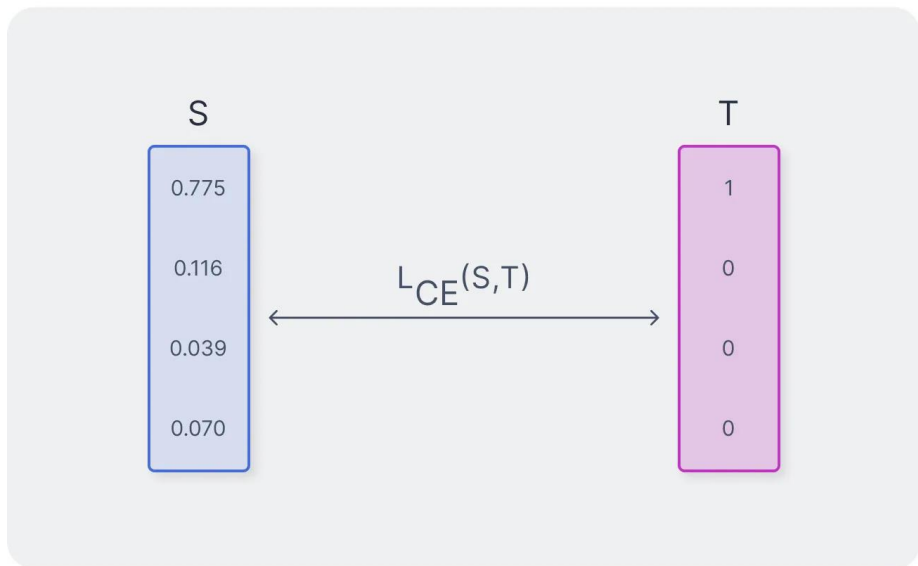
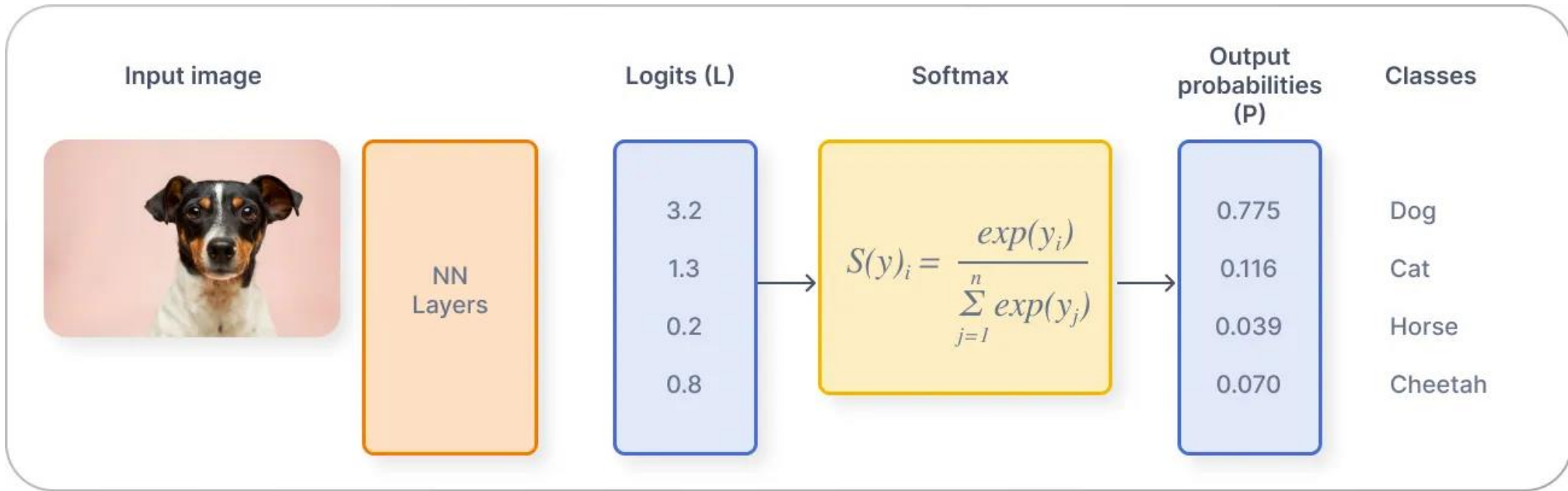
S

Softmax

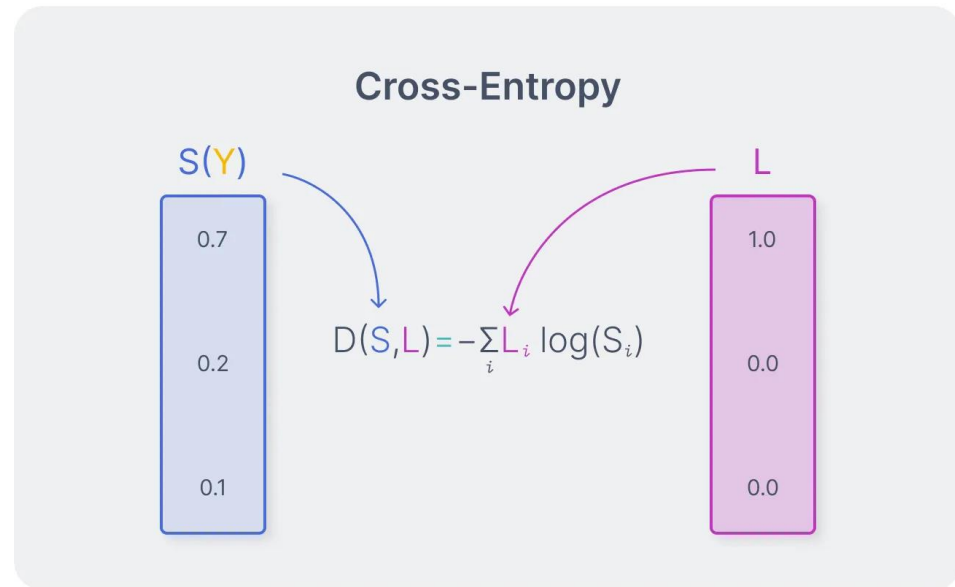
Cross-Entropy
Loss

$$f(s)_i = \frac{e^{s_i}}{\sum_j^C e^{s_j}}$$

$$CE = - \sum_i^C t_i \log(f(s)_i)$$



v7



v7

Trajectory-ranked Reward Extrapolation (T-REX)

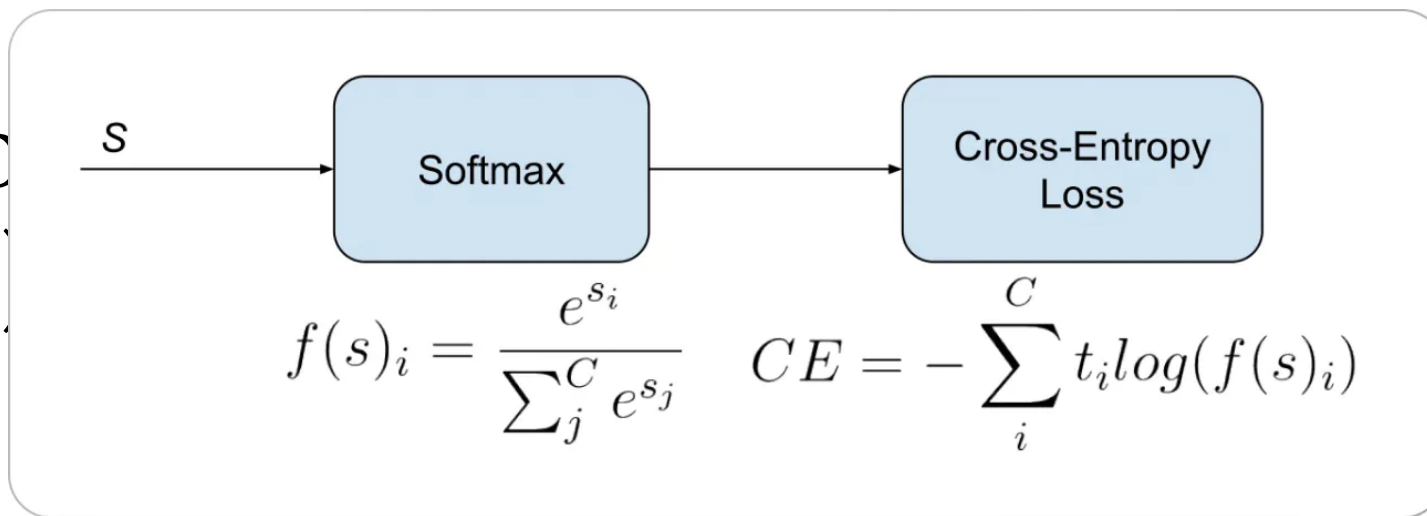
$$\tau_1 \prec \tau_2 \prec \dots \prec \tau_T$$

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

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

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

Trajectory-ranked Reward Extrapolation (T-REX)

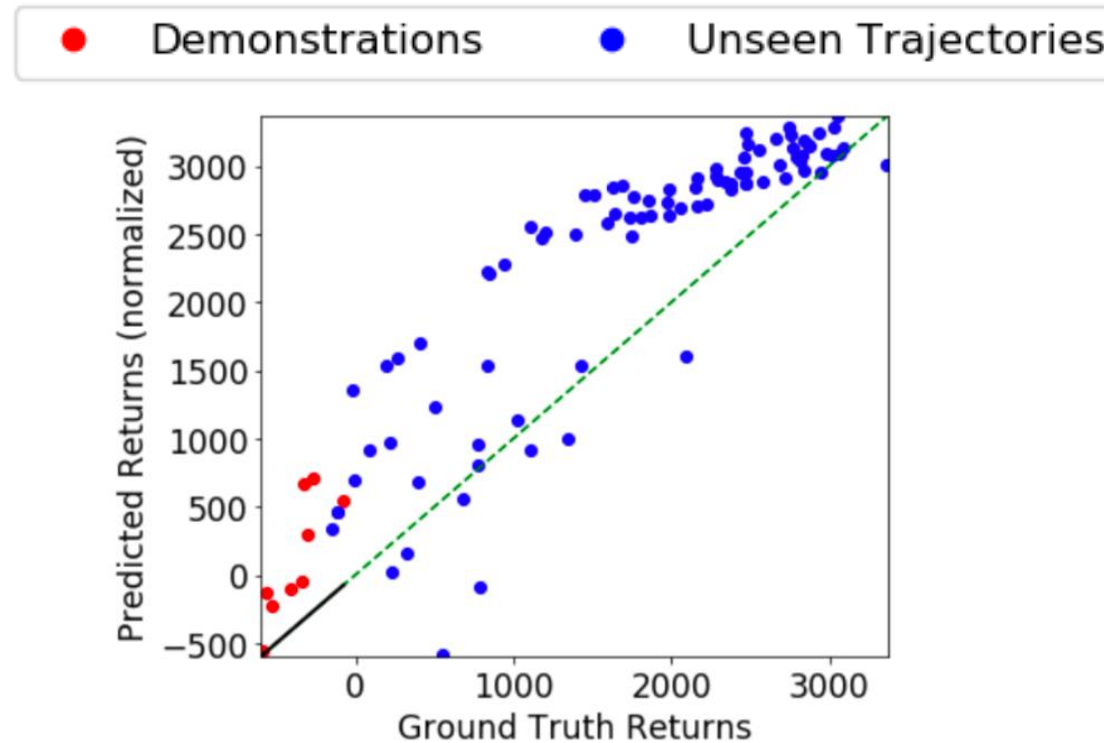
$$\tau_1 \prec \tau_2 \prec \dots \prec \tau_T$$

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

Minimize cross-entropy loss

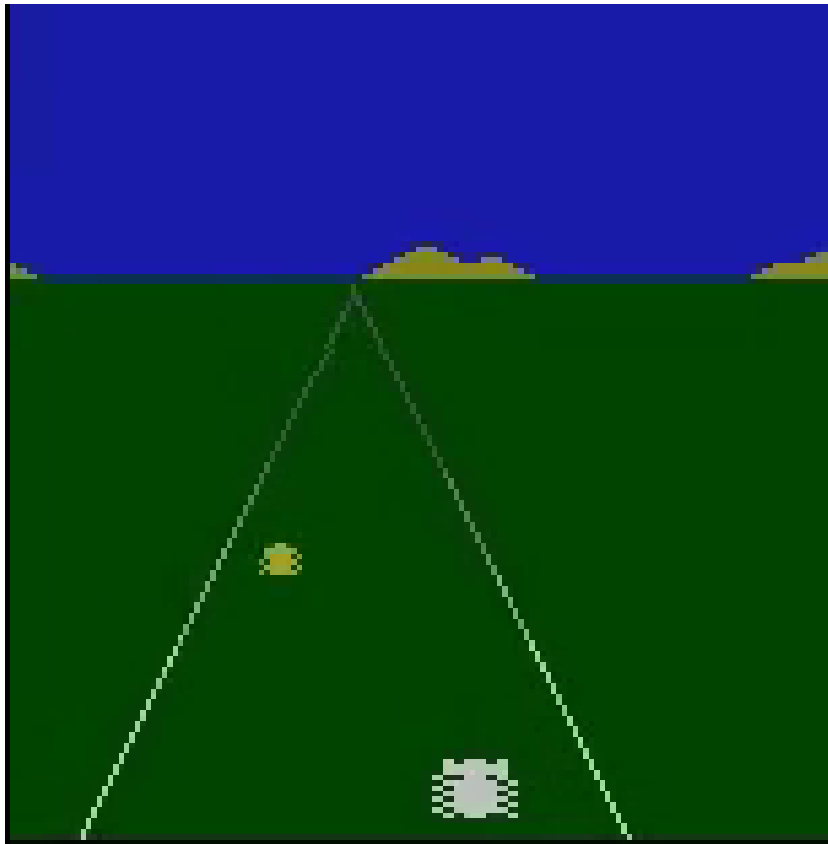
$$\mathcal{L}(\theta) = - \sum_{\tau_i \prec \tau_j} \frac{\exp \sum_{s \in \tau_j} R_\theta(s)}{\exp \sum_{s \in \tau_i} R_\theta(s) + \exp \sum_{s \in \tau_j} R_\theta(s)}$$

Reward Extrapolation

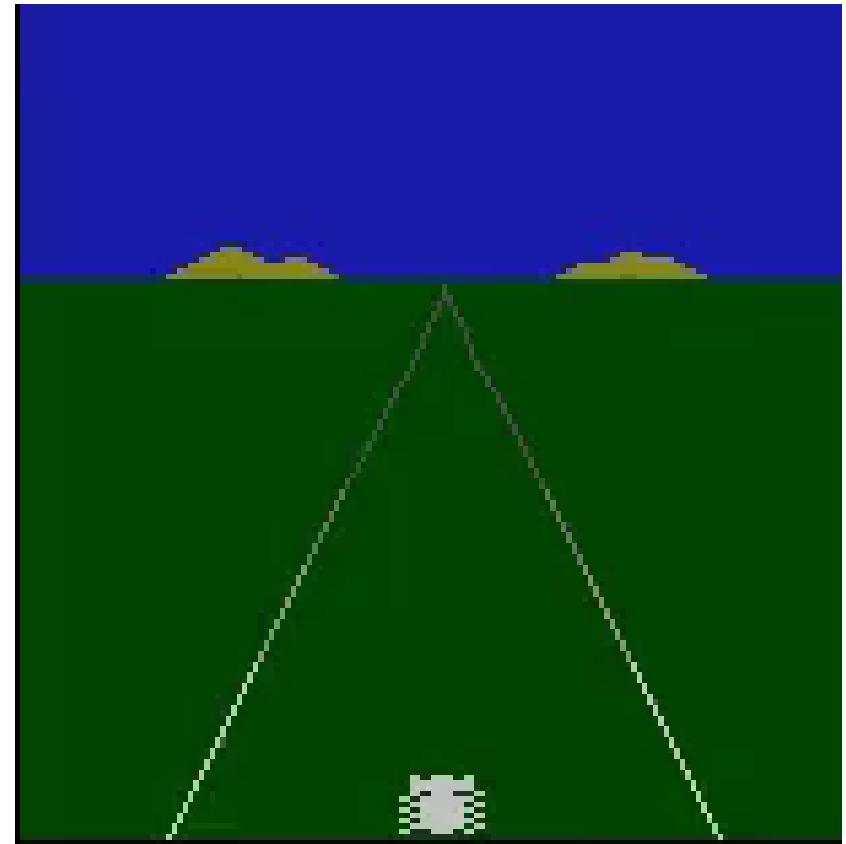


T-REX can extrapolate beyond the performance of the best demo

“Autonomous Driving” in Atari



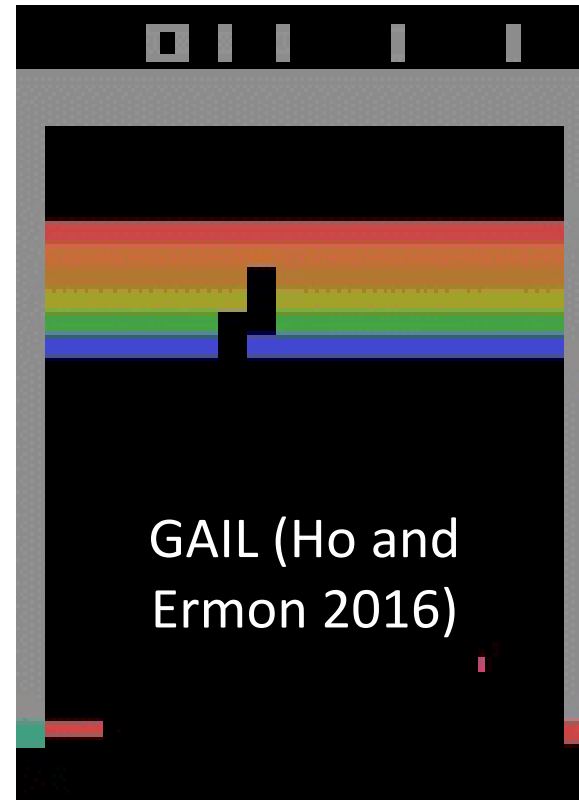
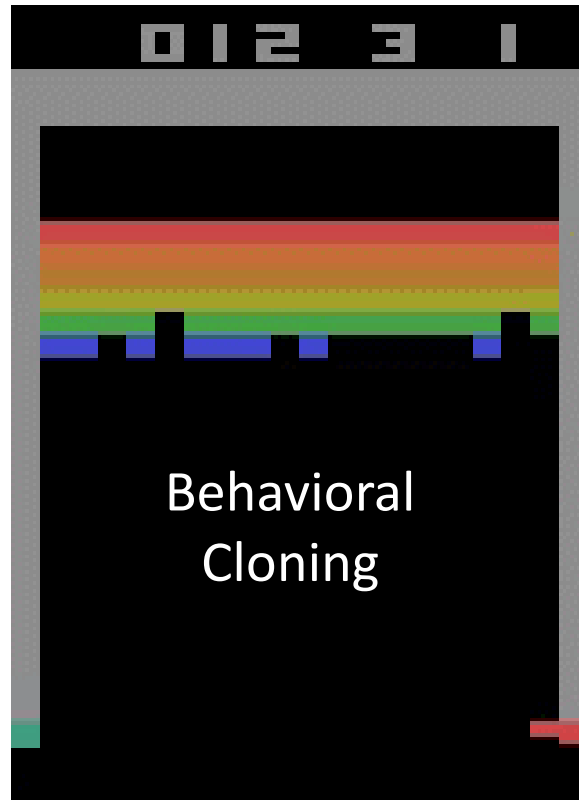
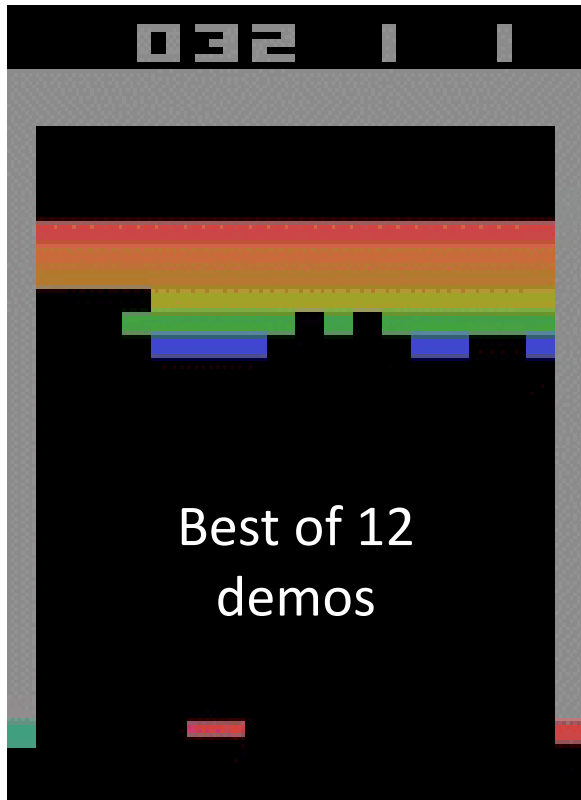
Best demo (Score = 84)



T-REX (Score = 520)

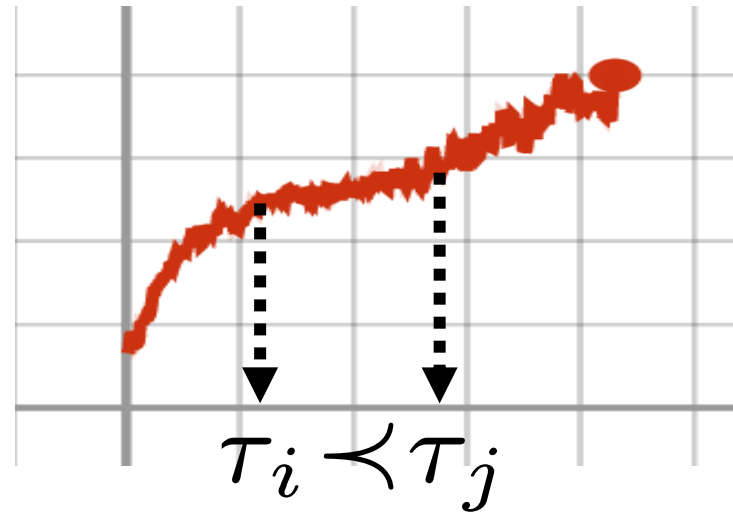
Uses only 12 ranked demonstrations

Atari Breakout

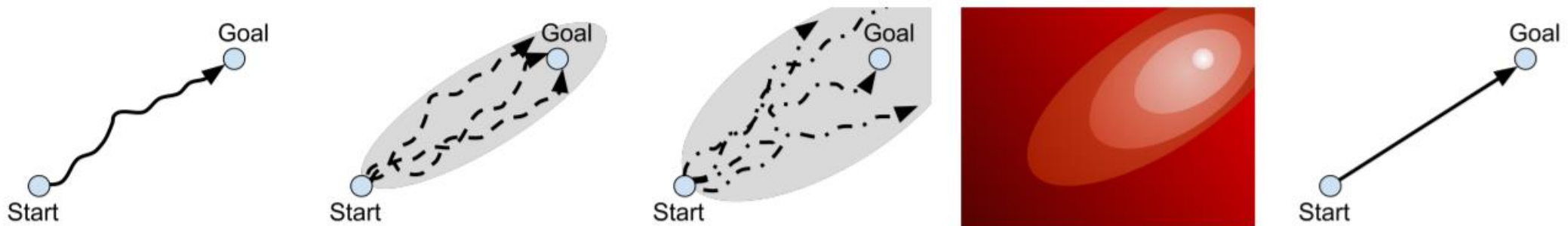


What if you don't have explicit preference labels?

Learning from a learner [ICML'19]

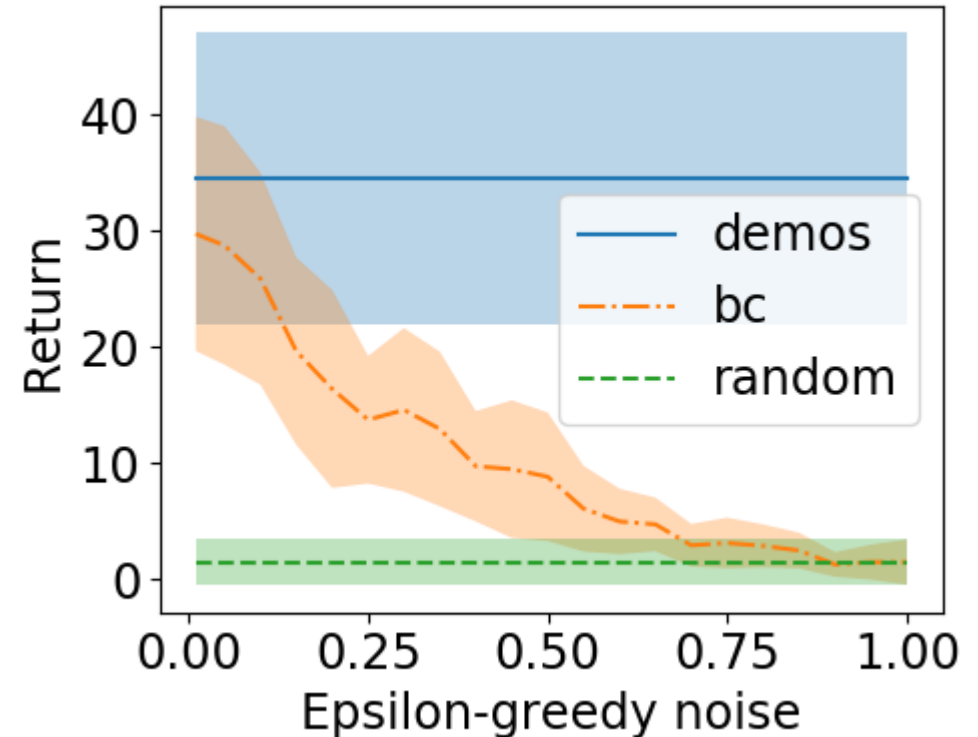


Automatic preference label generation [CoRL'20]



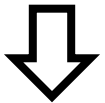
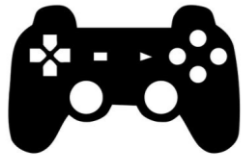
Automatic Rankings via Noise Injection

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



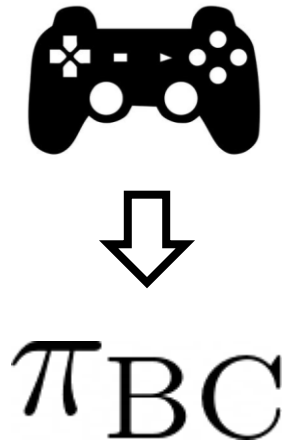
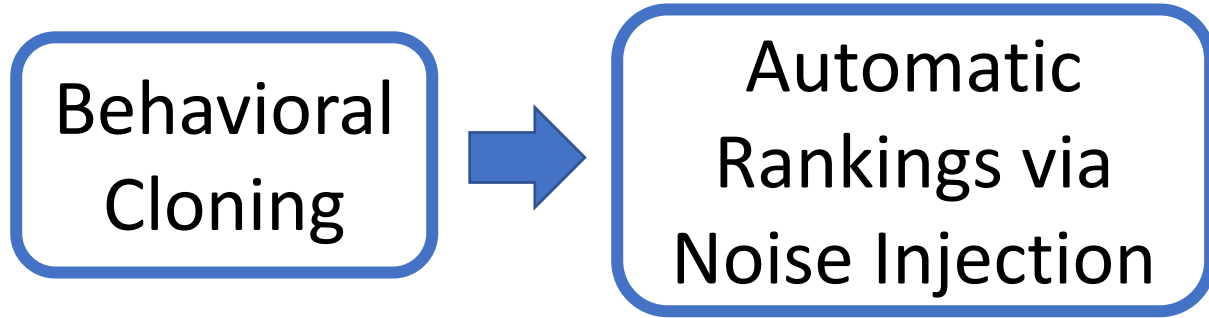
Disturbance-based Reward Extrapolation (D-REX)

Behavioral
Cloning



π_{BC}

Disturbance-based Reward Extrapolation (D-REX)



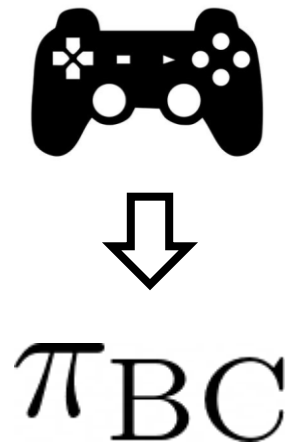
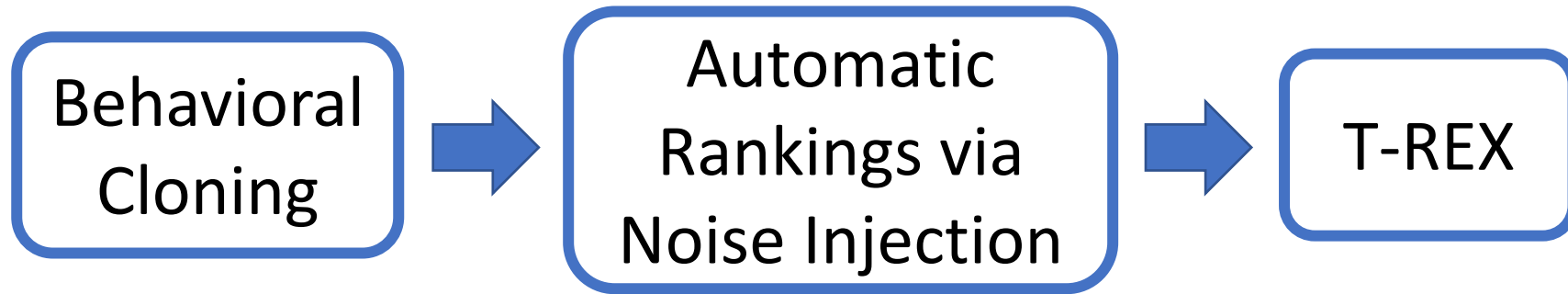
\sim



\sim



Disturbance-based Reward Extrapolation (D-REX)



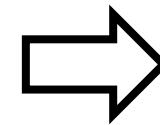
$\epsilon = 1.0$



$\epsilon = 0.2$

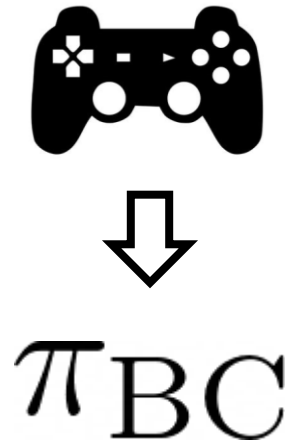
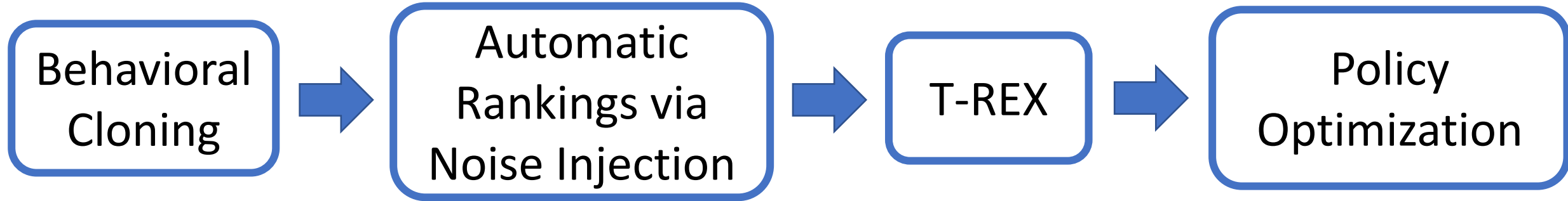


$\epsilon = 0.01$



Reward
Function
R(s)

Disturbance-based Reward Extrapolation (D-REX)



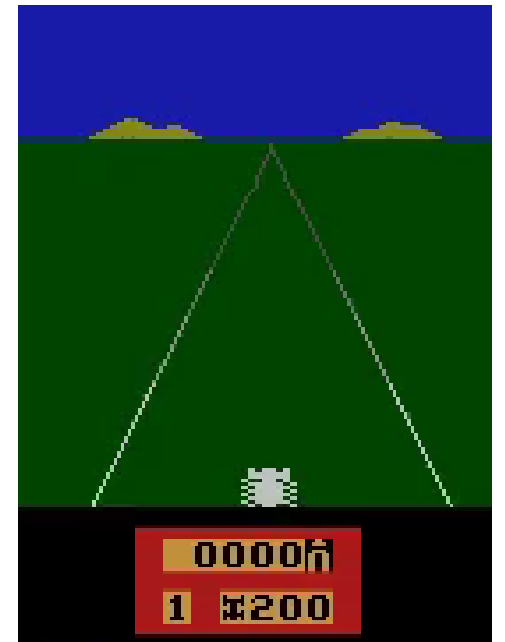
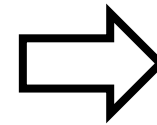
$\epsilon = 1.0$



$\epsilon = 0.2$



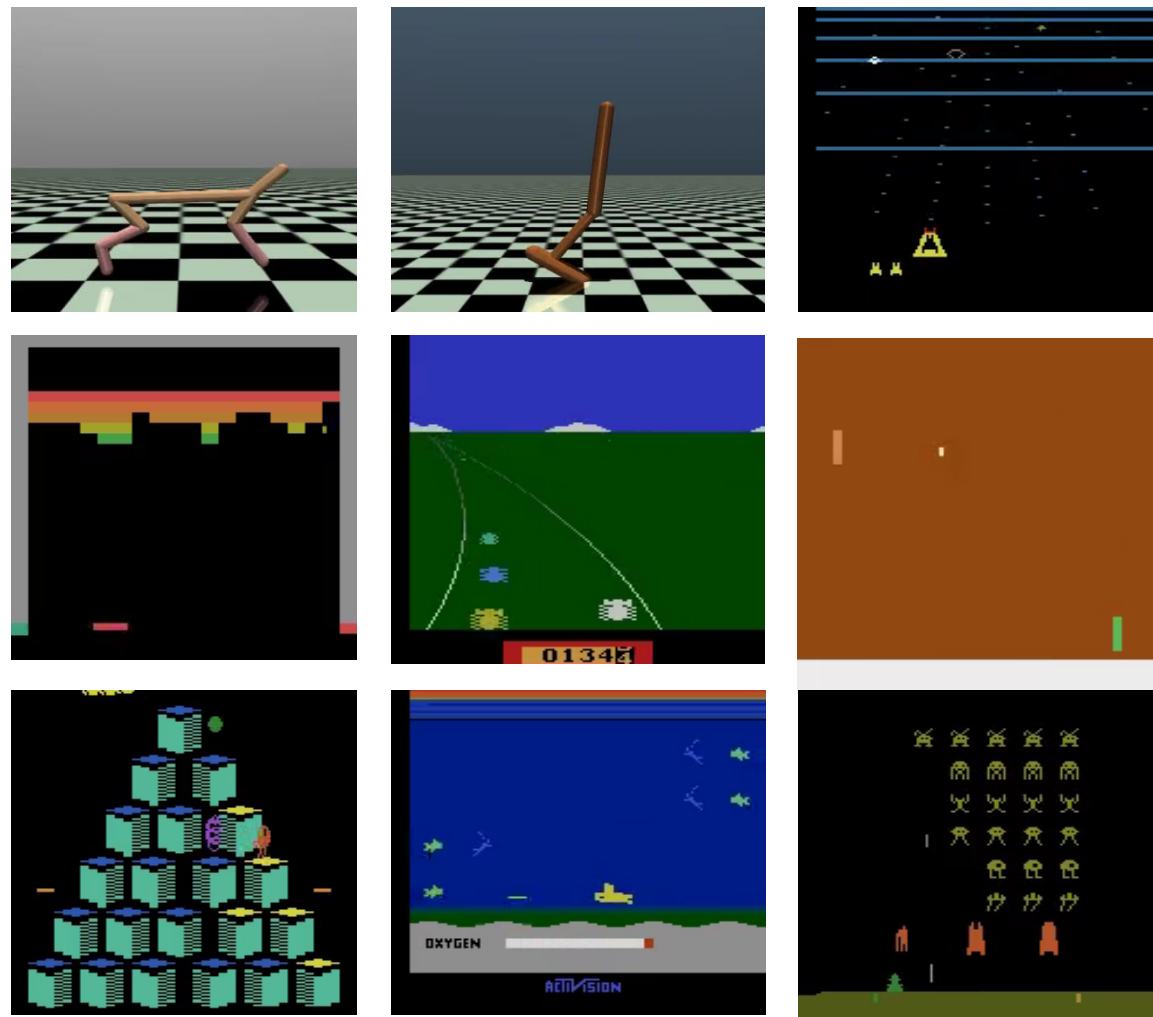
$\epsilon = 0.01$

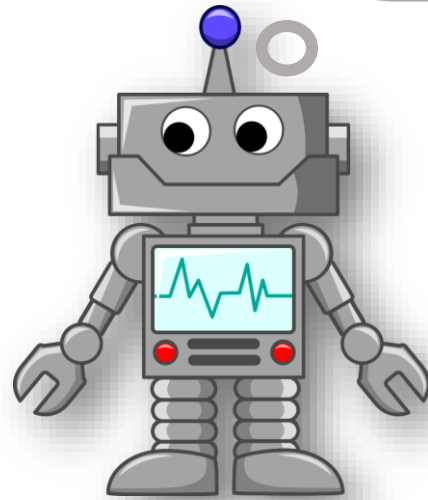
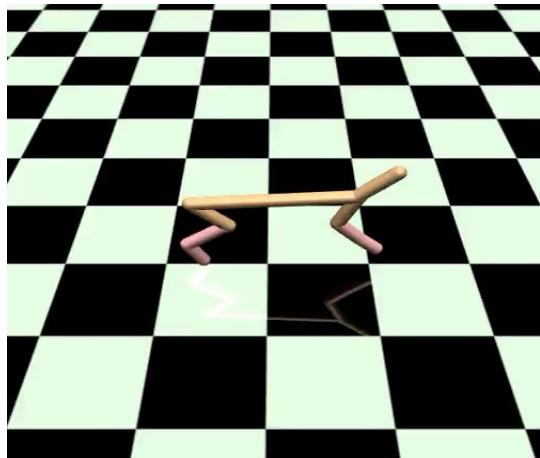


D-REX Policy

Experiments

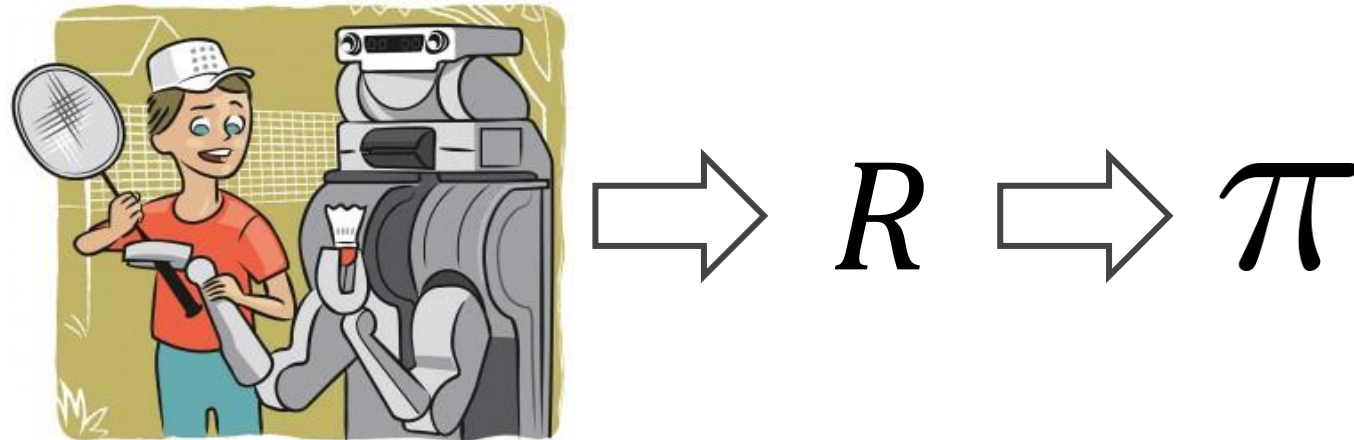
D-REX consistently outperforms the best demonstration as well as outperforming BC and GAIL.



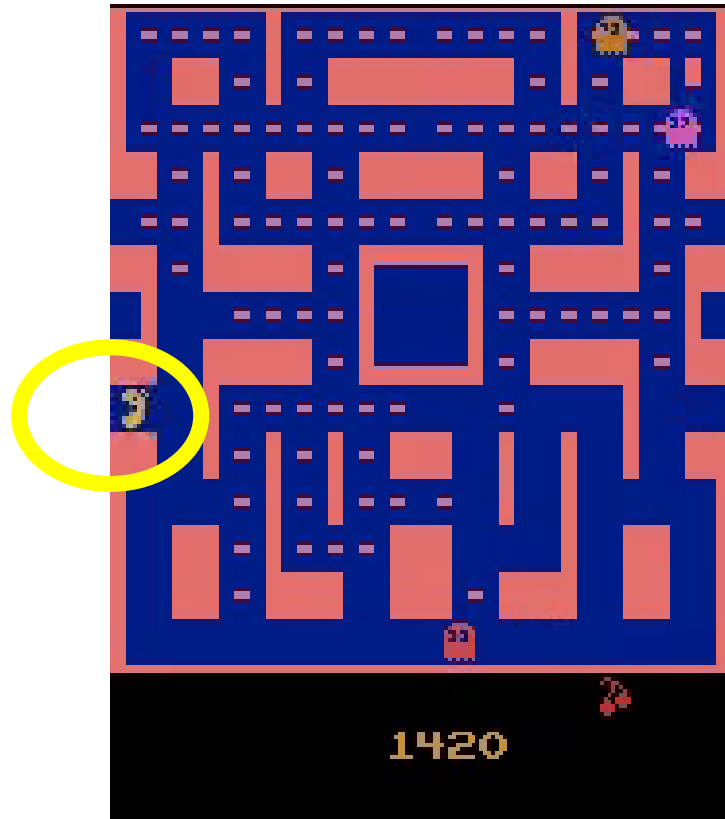


AI systems can **efficiently** infer human intent from **suboptimal demonstrations**.

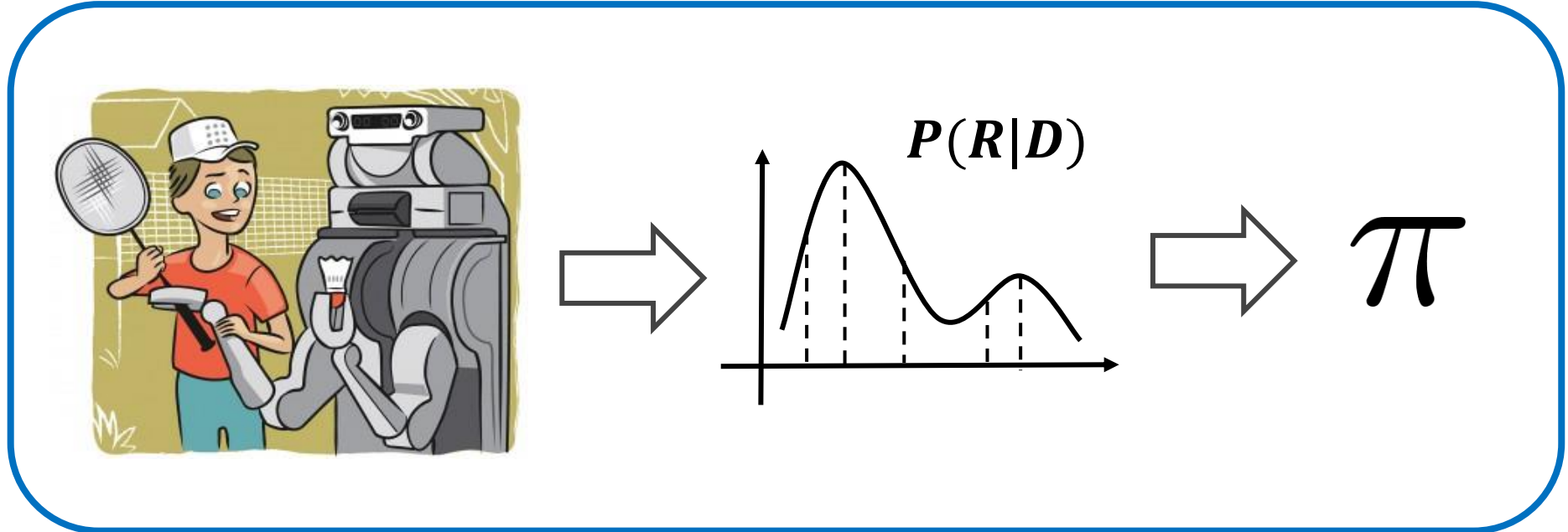
T-REX only learns a maximum likelihood estimate of the reward function.



Reward Hacking



- Overfit to spurious correlations
- No consideration of alternative hypotheses



Next time: LLMs and ChatGPT

