

# Algorithmic Human-Robot Interaction

---

(Shallow/Deep/Bottomless) Reinforcement Learning  
Inverse Reinforcement Learning

CSCI 7000

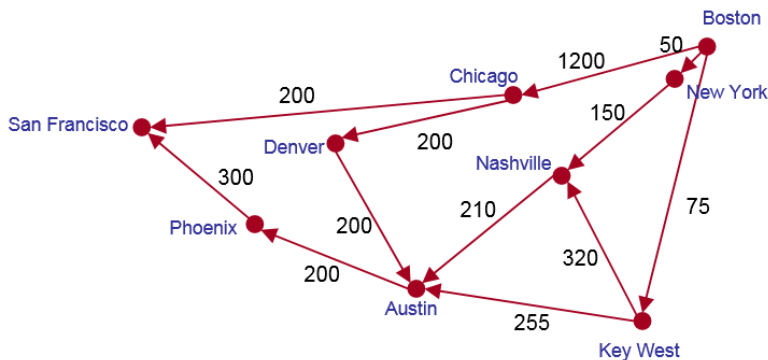
Prof. Brad Hayes

Computer Science Department

University of Colorado Boulder

# Planning Recap

# Search as a Problem-Solving Technique



Depth

0

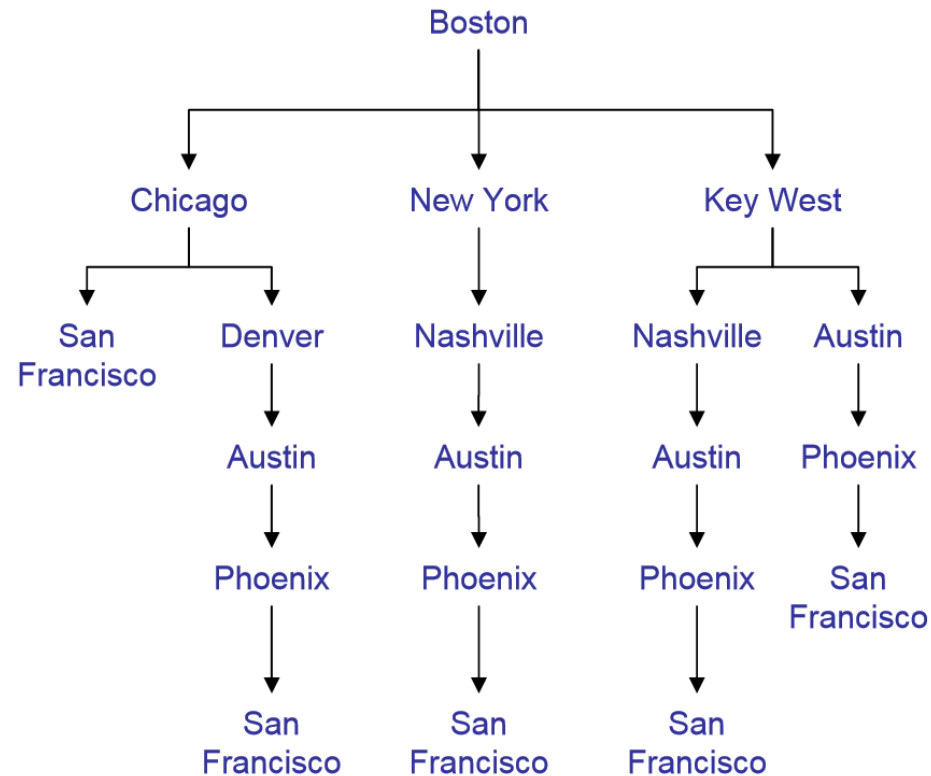
1

2

3

4

5



Branching Factor  $b=3$

# Lifted vs. Grounded Planning

Plan in this space

Execute in this space

Predicate Representation

On(a,b)  
Moving(d)  
Clear(Lhand)

...

...

...

...

...

...

1.0

1.0

0

0

0

-1.0

-1.0

1.0

Task  
Planning

Mapping Function

Compression  
Function

World State Vector

0.124

543.1

491.3

-50.11

194.2

...

...

...

...

Motion  
Planning



# Strategies

- Breadth-first search with sampled action parameterizations
- Pre-determined action parameterizations
  - e.g., `Place(object,x,y,z)`
- Plan for abstract 'areas'
  - **Lazy Grounding:**  
Sample parameterizations and add that action to plan if a valid parameterization is found
  - **Greedy Grounding:**  
Sample parameterizations and add that specific parameterization to the plan once a valid one is found

# Planning Heuristics

- FastForward Heuristic

- Remove deletions from all operators and run planner
- Resulting plan gives a heuristic for # steps from current state to goal
- FF Plan can provide ordering constraints on actions, even within partially specified plans

- Learned Cost Function

- Learning Real-Time A\* (LRTA\*)
- Update heuristic dynamically: Value function  $V(s)$
- Connection back to reinforcement learning:  $V(s)$  and  $Q(s,a)$  for search heuristic

# Admissible Heuristics and Pruning Strategies

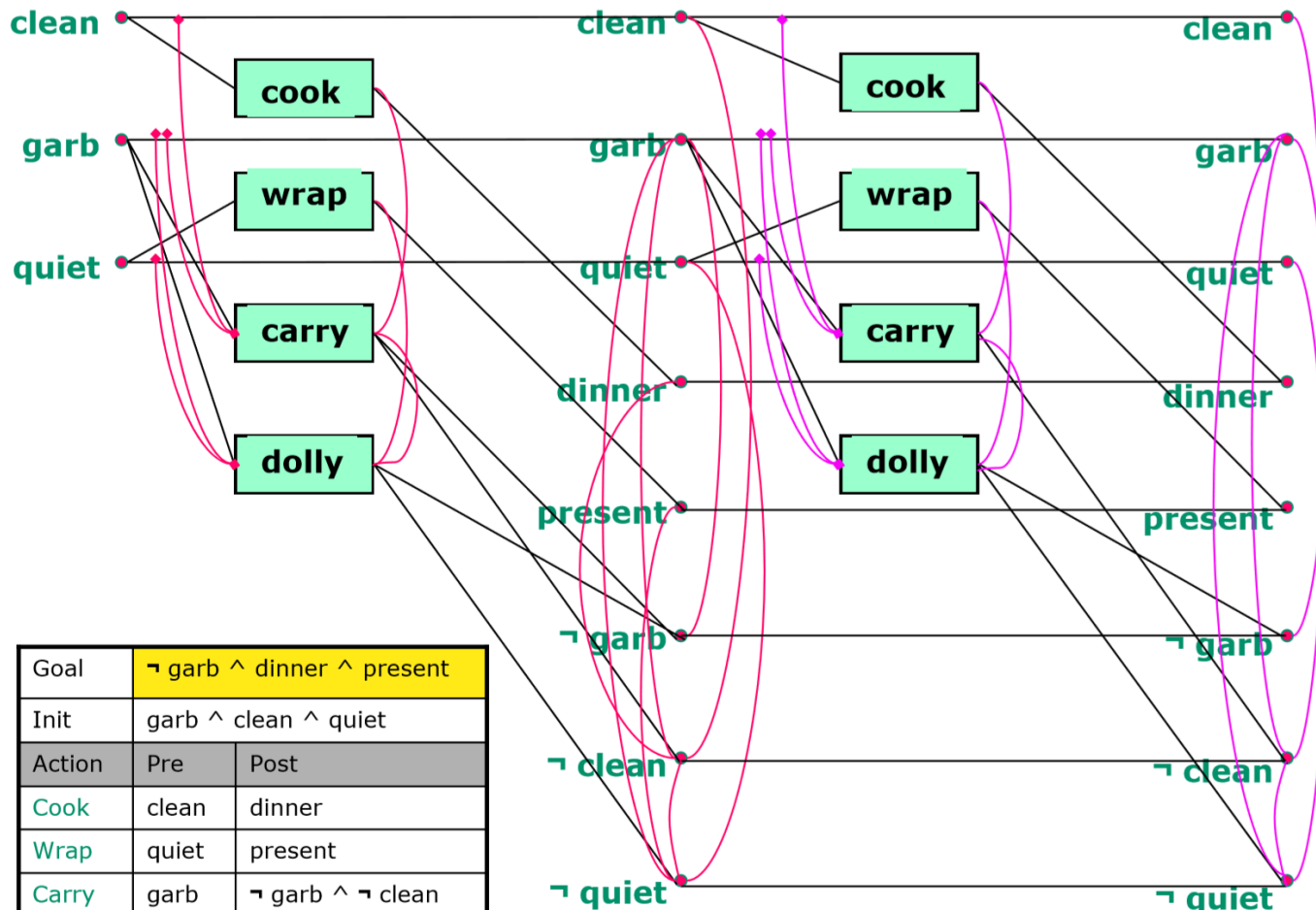
## Admissible heuristics from delete-relaxation:

- $h_{max} = i$  iff goal  $g$  is reachable in  $i$  steps from  $s$ 
  - Admissible because of GRAPHPLAN Algorithm
- $h_{FF} =$  number of **different** actions in  $\pi(g)$

## Search Pruning Strategies:

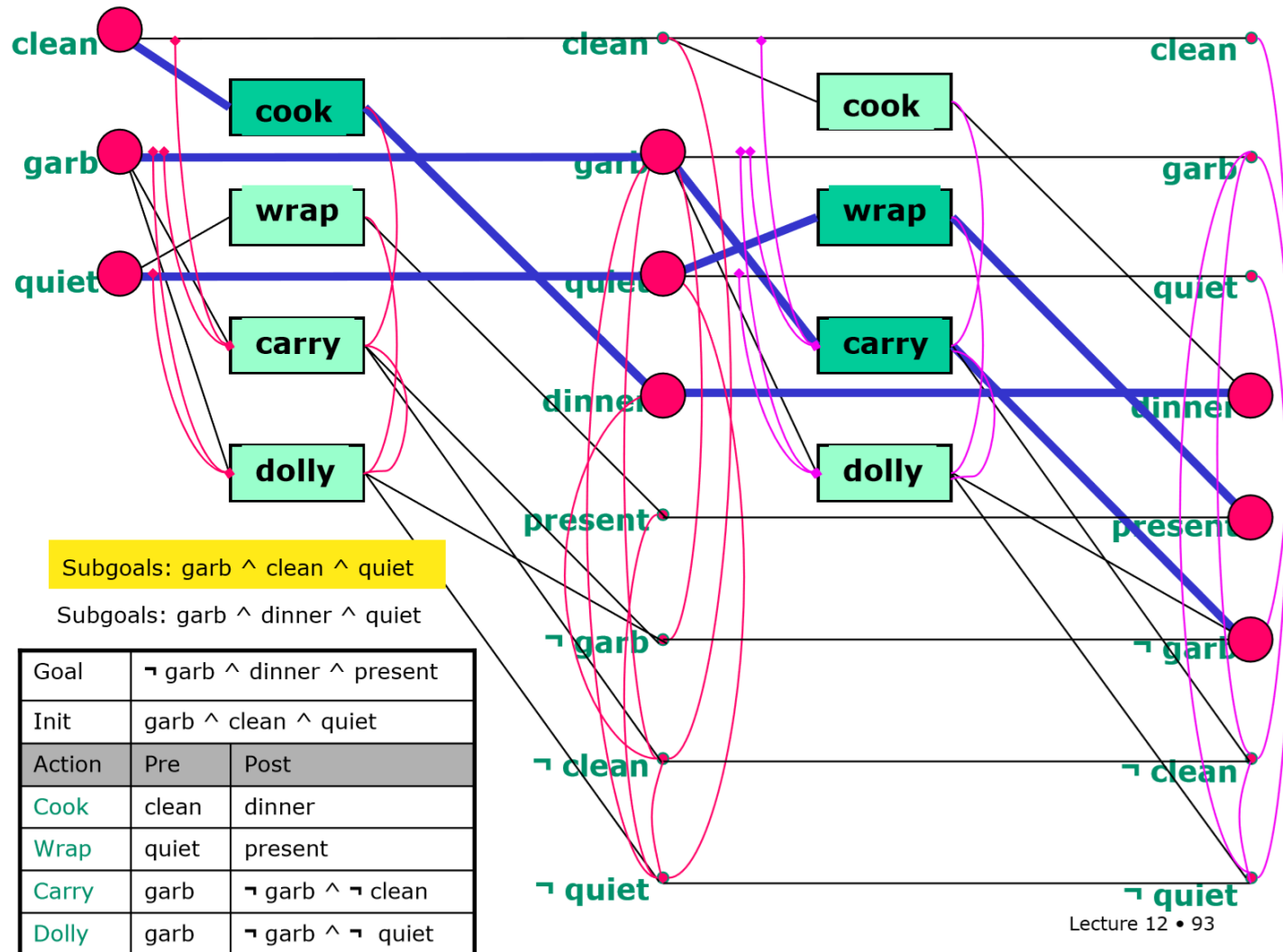
- Novelty: Assign each state  $novelty(s) = i$  if  $i$  predicates take on a value never seen in path from  $s_0$  to  $s$
- IW-search (Iterative Width): Ignore any nodes at the search frontier with  $novelty(s) > i$ .
  - Call planner for increasing  $i$  until problem solved or  $i$  exceeds number of variables in the problem.

# GRAPHPLAN Graph



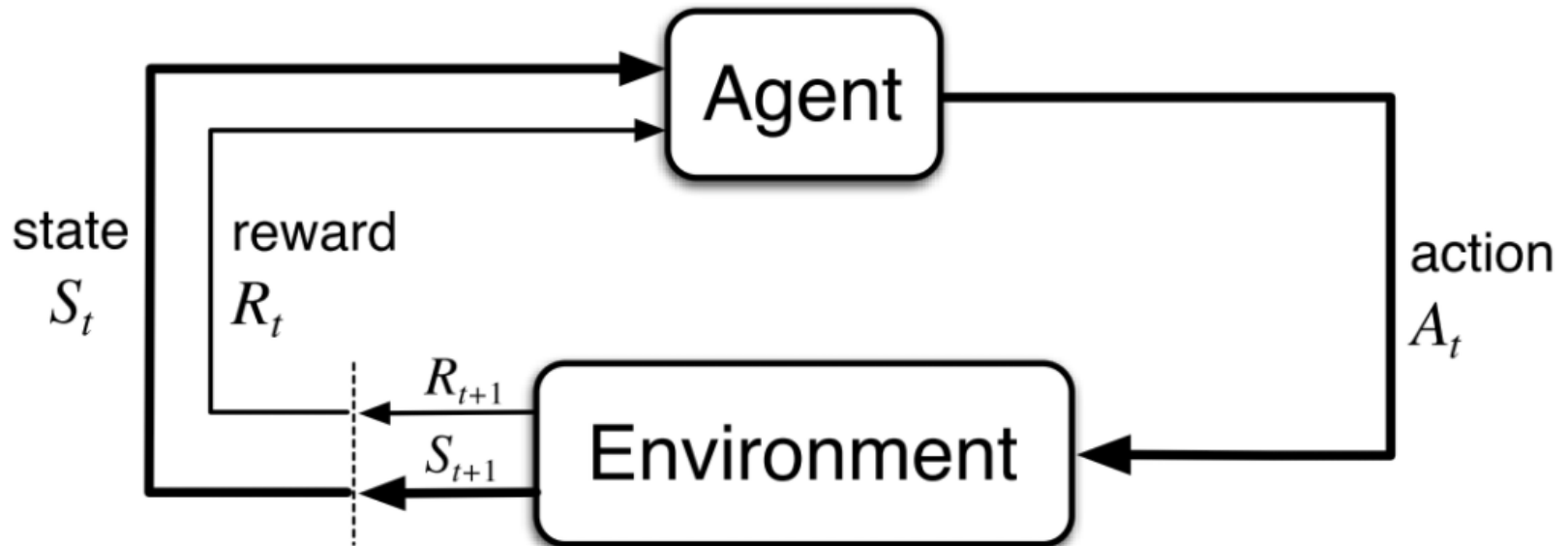


# GRAPHPLAN Graph Solution



# Inverse Reinforcement Learning

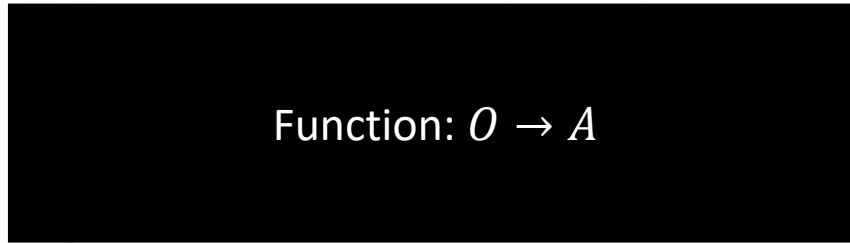
# Reinforcement Learning



# Reinforcement Learning



$\mathbf{o}$



Function:  $O \rightarrow A$

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



$\mathbf{a}$

$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action

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

$\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$  – policy (fully observed)



$\mathbf{o}_t$  – observation



$\mathbf{s}_t$  – state

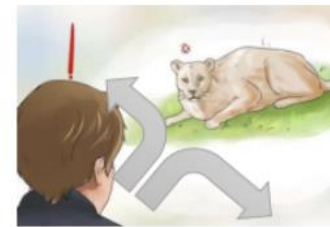
# Reinforcement Learning



$\mathbf{o}$

Function:  $O \rightarrow A$

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



$\mathbf{a}$

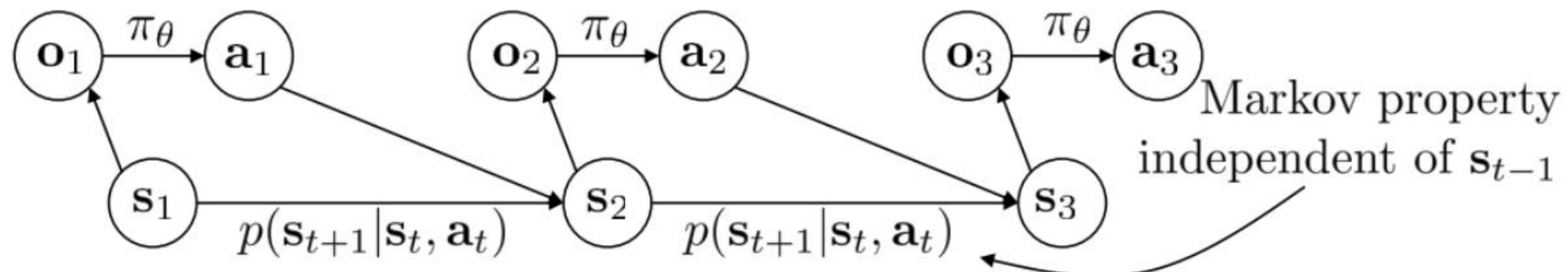
$\mathbf{s}_t$  – state

$\mathbf{o}_t$  – observation

$\mathbf{a}_t$  – action

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

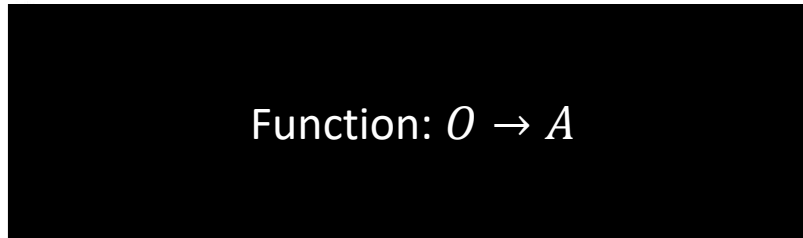
$\pi_{\theta}(\mathbf{a}_t|\mathbf{s}_t)$  – policy (fully observed)



# Reward Functions



$\mathbf{o}_t$



Function:  $O \rightarrow A$

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



$\mathbf{a}_t$

which action is better or worse?

$r(\mathbf{s}, \mathbf{a})$ : reward function

tells us which states and actions are better

$\mathbf{s}$ ,  $\mathbf{a}$ ,  $r(\mathbf{s}, \mathbf{a})$ , and  $p(\mathbf{s}' | \mathbf{s}, \mathbf{a})$  define

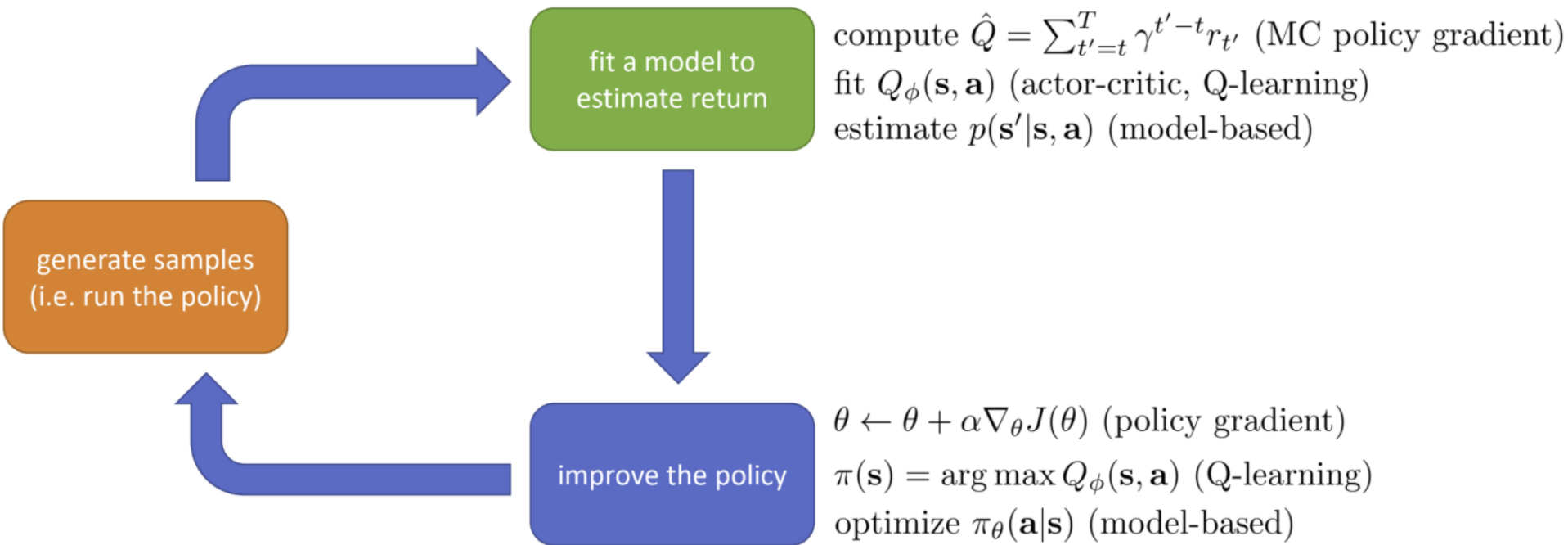
Markov decision process



high reward



low reward



# Anatomy of an RL Algorithm

# Policy Differentiation

Trajectory:  $\tau$   
Trajectory length:  $T$   
Policy:  $\pi$   
Reward:  $r$   
Parameters:  $\theta$   
Gradient:  $\nabla$   
 $J$ : Expected reward

$$\theta^* = \underbrace{\text{Best } \theta \text{ based on expected reward over the } T\text{-length trajectory}}_{J(\theta)}$$

$$\underbrace{\pi_\theta(s_1, a_1, \dots, s_T, a_T)}_{\pi_\theta(\tau)} = \text{Probability of given trajectory using policy } \pi_\theta$$

$$J(\theta) = \text{Expected reward value given trajectory } (\tau) \text{ sampled from } \pi_\theta$$

Gradient of policy w.r.t.  $\theta$  is equal to policy \* gradient of log policy

$$\nabla_\theta J(\theta) = \text{Take gradient of } J(\theta) = \text{Use convenient identity to get rid of non-log policy} = \text{Sample from policy instead of enumerating all possible trajectories}$$

Rewrite substituting in worked-out gradient

Expand policy(trajectory) and take log gradient



# Policy Differentiation

Trajectory:  $\tau$   
 Trajectory length:  $T$   
 Policy:  $\pi$   
 Reward:  $r$   
 Parameters:  $\theta$   
 Gradient:  $\nabla$   
 $J$ : Expected reward

$$\theta^* = \arg \max_{\theta} \underbrace{\sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim p_{\theta}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)]}_{J(\theta)}$$

$$\underbrace{\pi_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)}_{\pi_{\theta}(\tau)} = p(\mathbf{s}_1) \underbrace{\prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}_{\text{green underline}}$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [r(\tau)] = \int \pi_{\theta}(\tau) r(\tau) d\tau$$

$$\underbrace{\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)}_{\text{orange underline}} = \pi_{\theta}(\tau) \frac{\nabla_{\theta} \pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \underbrace{\nabla_{\theta} \pi_{\theta}(\tau)}_{\text{blue underline}}$$

$$\nabla_{\theta} J(\theta) = \int \underbrace{\nabla_{\theta} \pi_{\theta}(\tau)}_{\text{blue underline}} r(\tau) d\tau = \int \underbrace{\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau)}_{\text{orange underline}} r(\tau) d\tau = E_{\tau \sim \pi_{\theta}(\tau)} [\underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau)}_{\text{green underline}} r(\tau)]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$\nabla_{\theta} \left[ \cancel{\log p(\mathbf{s}_1)} + \sum_{t=1}^T \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) + \cancel{\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} \right]$$

# Using the Policy Gradient

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left( \sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

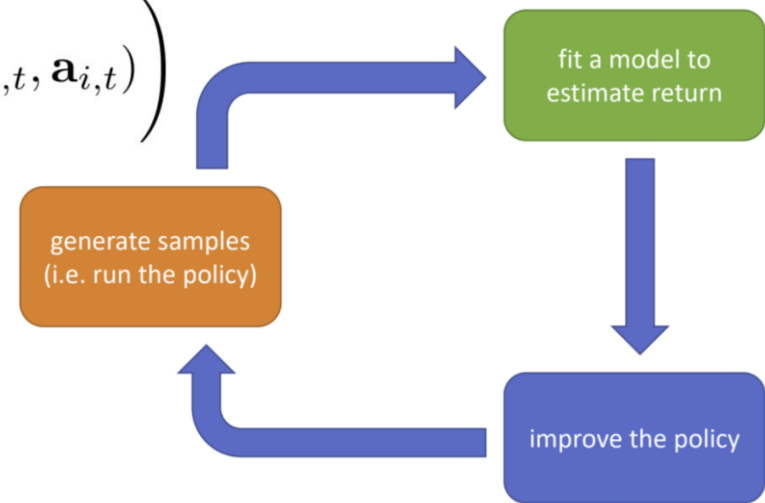
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)$$

REINFORCE algorithm:

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$  (run it on the robot)
2.  $\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



# Where does reward come from?

Computer Games

reward



Mnih et al. '15

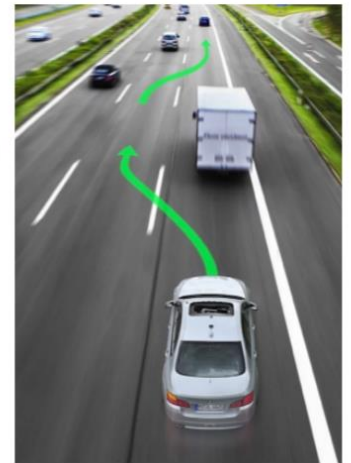
robotics



dialogue



autonomous driving



What's the reward here?  
Usually we use a proxy.

It's generally easier to provide expert data  
than to specify a useful reward function.  
Consider two cases:

- Autonomous helicopter tricks?
- Growing a plant?

# Inverse Optimal Control / Inverse RL

Given:

- State and action space
- Roll-outs from  $\pi^*$
- Dynamics Model (sometimes)

Goal:

- Recover reward function
- Use reward function to derive  $\pi \approx \pi^*$

## Challenges:

- Problem is underspecified
- Hard to evaluate your learned reward
- Demonstrations might be good, but not optimal

# Maximum Entropy IRL

(Ziebart et al. 2008)

## Notation:

$$\tau = \{s_1, a_1, \dots, s_t, a_t, \dots, s_T\}$$

trajectory

$$R_\psi(\tau) = \sum_t r_\psi(s_t, a_t)$$

learned reward

$$\mathcal{D} : \{\tau_i\} \sim \pi^*$$

expert demonstrations

$$p(\tau) = \frac{1}{Z} \exp(R_\psi(\tau))$$

Probability is exponential in the reward  
(good trajectories are more likely)

$$\max_\psi \sum_{\tau \in \mathcal{D}} \log p_{r_\psi}(\tau)$$

Maximize likelihood of a given reward  
function parameterization

$$Z = \int \exp(R_\psi(\tau)) d\tau$$

Integral over all possible trajectories ☹️

# MaxEnt IRL

1. Initialize  $\psi$  (reward function params),  
gather demonstrations  $D$
2. Solve for optimal policy  $\pi(a|s)$  with reward  $r_\psi$
3. Solve for state visitation frequencies  $p(s|\psi)$
4. Compute gradient:  $\nabla_\psi L = -\frac{1}{|D|} \sum_{\tau_d \in D} \frac{dr_\psi}{d\psi}(\tau_d) - \sum_s p(s|\psi) \frac{dr_\psi}{d\psi}(s)$
5. Update  $\psi$  with one gradient step using  $\nabla_\psi L$
6. GOTO 2

Must solve the whole MDP in the inner loop of  
finding the reward function!

# Making IRL work for complex problems

Must handle:

- (1) Unknown dynamics
- (2) Solving MDP in an inner loop

$$p(\tau) = \frac{1}{Z} \exp(R_\psi(\tau))$$

$$\max_{\psi} \sum_{\tau \in D} \log p_{r_\psi}(\tau)$$

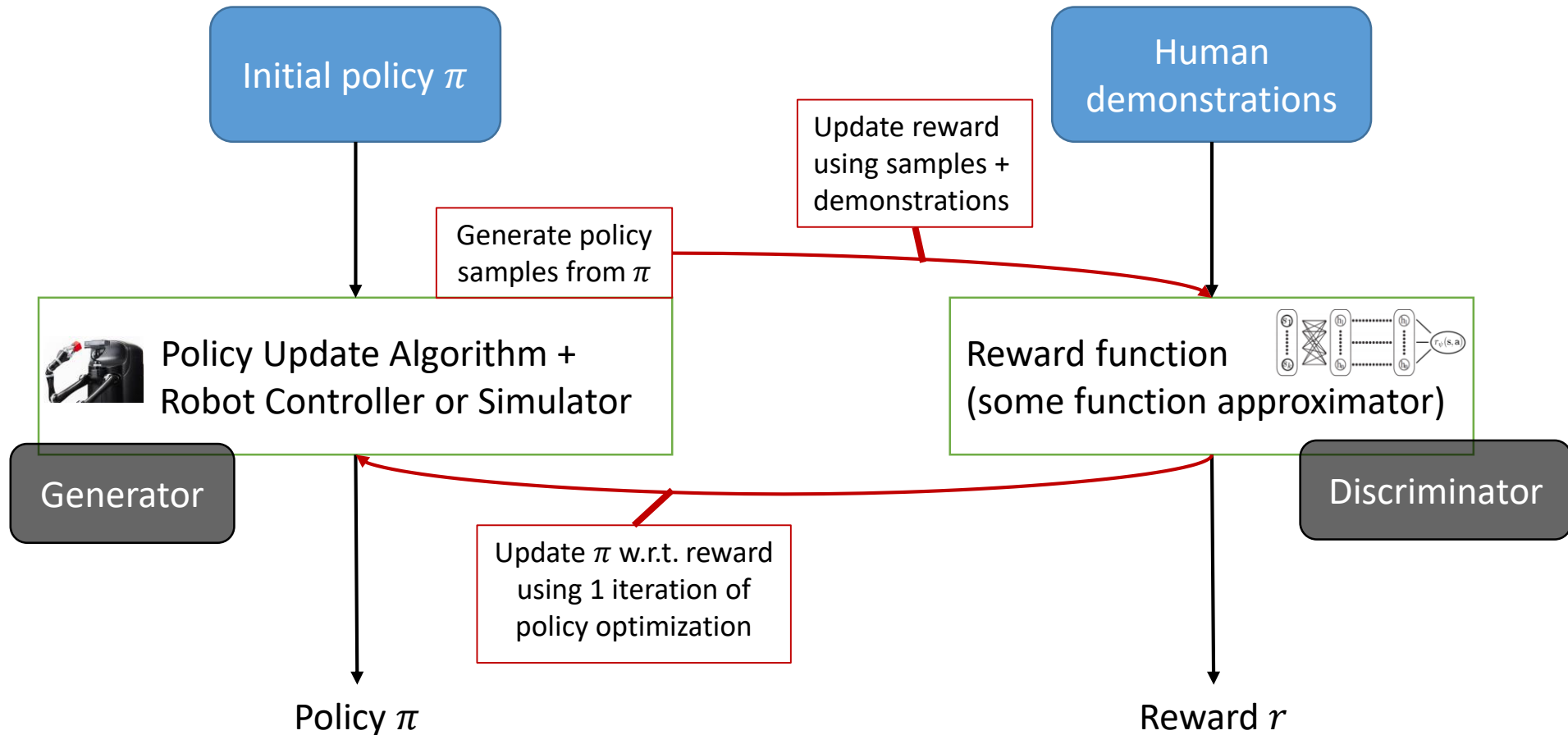
$$Z = \int \exp(R_\psi(\tau)) d\tau$$

Adaptively sample increasingly close to  $\psi$  by constructing a policy to do it for us

Sample from what?  
Can't sample from policies near  $\psi$  because we're solving for  $\psi$ !

Can sample to approximate Z!

# Guided Cost Learning (Generative Adversarial Imitation)





# Brief Aside: Generative Adversarial Networks



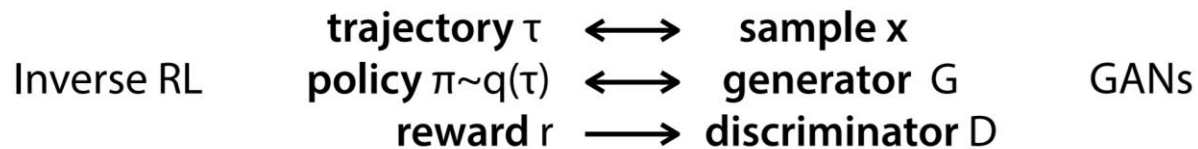
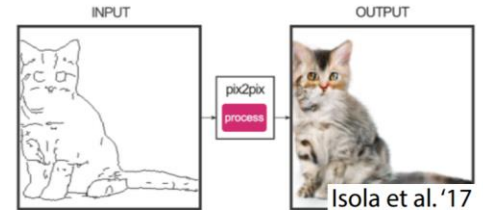
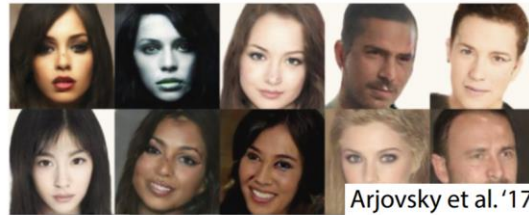
Unpaired Image-to-image Translation  
using Cycle-consistent  
Adversarial Networks

Jun-Yan Zhu\* Taesung Park\* Phillip Isola Alexei A. Efros  
UC Berkeley

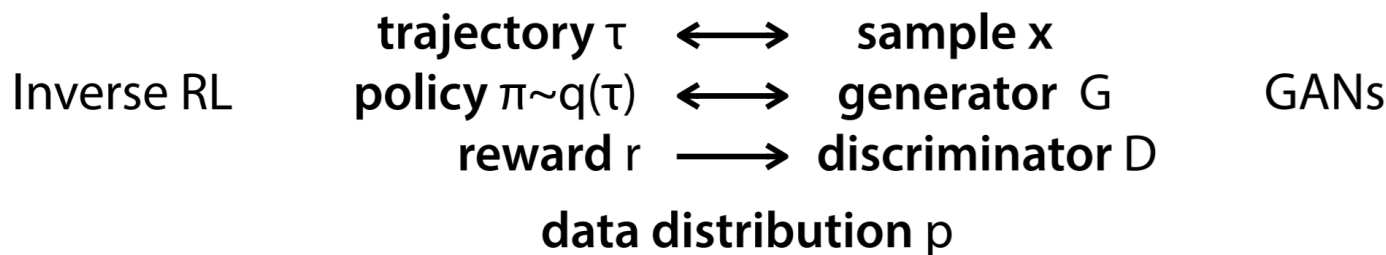


# Inverse RL $\Leftrightarrow$ GANs

Similar to inverse RL, **GANs** learn an objective for generative modeling.



# Discriminators and Loss Functions



Guided Cost Learning:

$$D^*(\tau) = \frac{p(\tau)}{p(\tau) + q(\tau)}$$

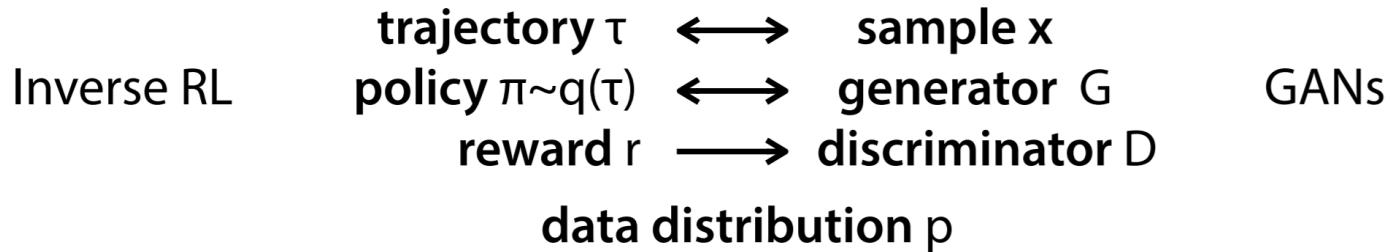
Guided Adversarial Imitation Learning:

$$D_\psi(\tau) = \text{some classifier}$$

$$D_\psi(\tau) = \frac{\frac{1}{Z} \exp(R_\psi)}{\frac{1}{Z} \exp(R_\psi) + q(\tau)}$$

$$\mathcal{L}_{\text{discriminator}}(\psi) = \mathbb{E}_{\tau \sim p}[-\log D_\psi(\tau)] + \mathbb{E}_{\tau \sim q}[-\log(1 - D_\psi(\tau))]$$

# Generators and Loss Functions



$$\begin{aligned}\mathcal{L}_{\text{generator}}(\theta) &= \mathbb{E}_{\tau \sim q} [\log(1 - D_{\psi}(\tau)) - \log D_{\psi}(\tau)] \\ &= \mathbb{E}_{\tau \sim q} [\log q(\tau) + \log Z - R_{\psi}(\tau)]\end{aligned}$$



Entropy-regularized RL

Must train generator/policy with RL because we don't know the dynamics of the environment!

(otherwise we could just train using the discriminator signal straight through to the policy)

# IRL Recap

**Goal: Infer reward function underlying expert demonstrations**

**Evaluating the partition function ( $Z$ ):**

- Initial approaches solve the MDP in the inner loop of IRL (or assume known dynamics).
- Can estimate  $Z$  using sampling!

**Connection to Generative Adversarial Networks:**

- Sampling-based MaxEnt IRL is a GAN with a special form of discriminator, using RL to optimize the generator.

## Further Reading

### Classic Papers

- **Abbeel & Ng ICML '04.** *Apprenticeship Learning via Inverse Reinforcement Learning.* Good introduction to inverse reinforcement learning
- **Ziebart et al. AAAI '08.** *Maximum Entropy Inverse Reinforcement Learning.* Introduction of probabilistic method for inverse reinforcement learning

### Modern Papers

- **Wulfmeier et al. arXiv '16.** *Deep Maximum Entropy Inverse Reinforcement Learning.* MaxEnt IRL using deep reward functions
- **Finn et al. ICML '16.** *Guided Cost Learning.* Sampling-based method for MaxEnt IRL that handles unknown dynamics and deep reward functions
- **Ho & Ermon NIPS '16.** *Generative Adversarial Imitation Learning.* IRL method building on Abbeel & Ng '04 using generative adversarial networks

# Designing Your Evaluation

## **Evaluation Design:**

What are your hypotheses about your system?

How will you test them?

What are you trying to prove with this work?

## **Experiment Design:**

Do you need human subjects?

Are your conditions likely to test your hypotheses?

Within-subjects or between-subjects?

## **Protocol design:**

Someone not on your project should be able to run your experiment with this script!