

Generative Adversarial Imitation Learning

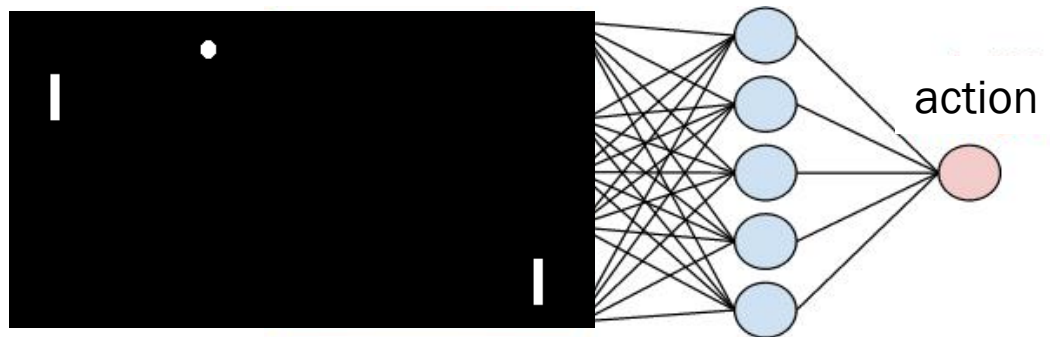
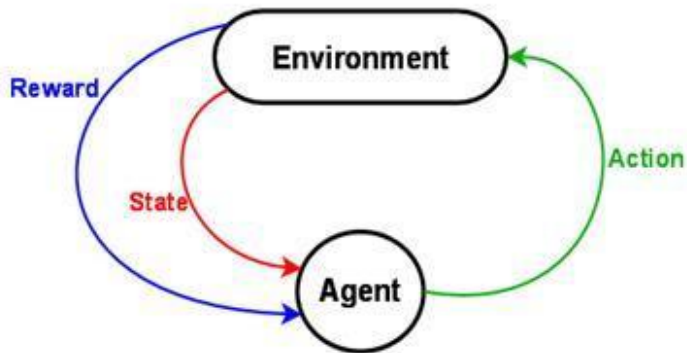
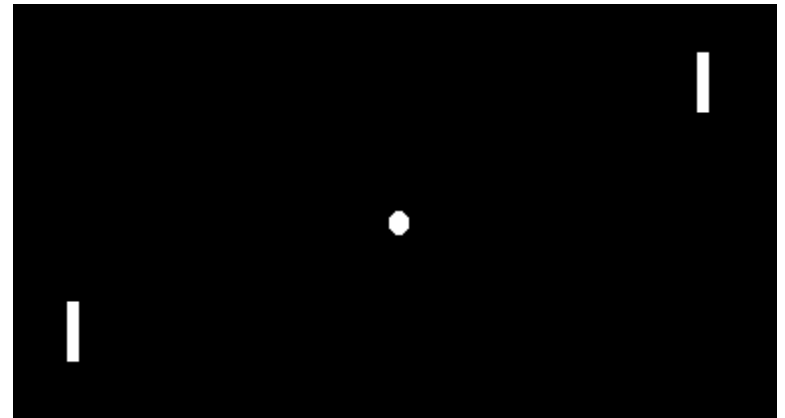
Stefano Ermon

Joint work with Jayesh Gupta, Jonathan Ho, Yunzhu Li,
Hongyu Ren, and Jiaming Song

Stanford University

Reinforcement Learning

- Goal: Learn policies
- High-dimensional, raw observations



Reinforcement Learning

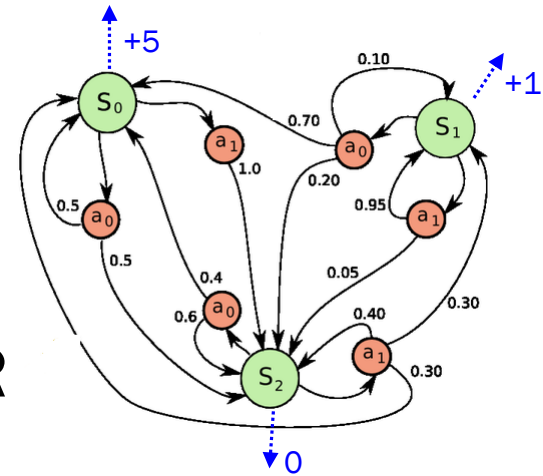
- **MDP**: Model for (stochastic) *sequential decision making* problems

- States **S**

- Actions **A**

- **Cost** function (immediate): **C**: $S \times A \rightarrow R$

- Transition Probabilities: $P(s' | s, a)$



- **Policy**: mapping from states to actions

– E.g., ($S_0 \rightarrow a_1$, $S_1 \rightarrow a_0$, $S_2 \rightarrow a_0$)

- Reinforcement learning: minimize total (expected, discounted) cost

$$\sum_{t=0}^{T-1} c(s_t)$$

Reinforcement Learning

$$RL(c) = \arg \min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_{\pi}[c(s, a)]$$



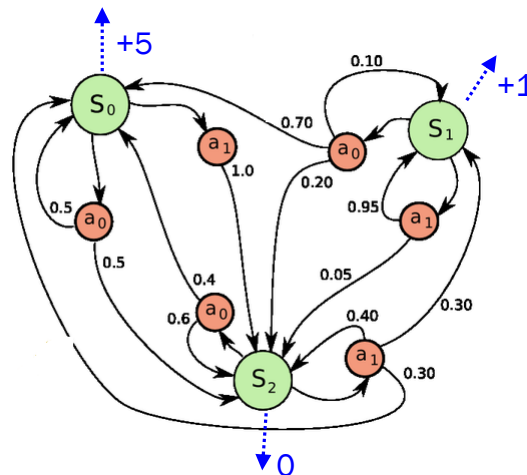
$C: S \times A \rightarrow R$

Cost

RL needs cost signal

Policy: mapping from states to actions

E.g., $(S_0 \rightarrow a_1,$
 $S_1 \rightarrow a_0,$
 $S_2 \rightarrow a_0)$

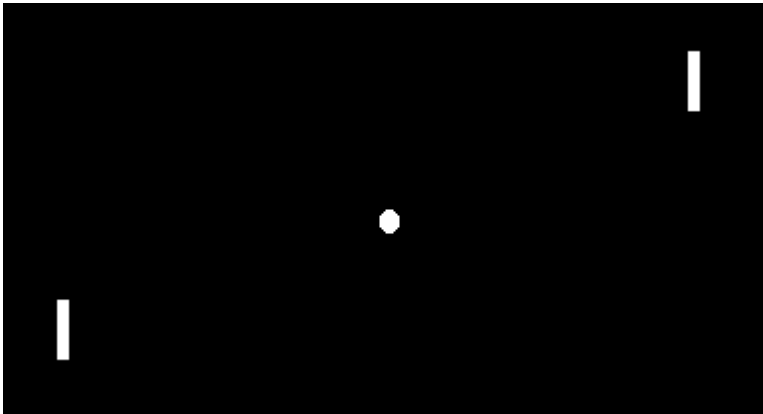


- States S
- Actions A
- Transitions: $P(s' | s, a)$

Imitation

Input: expert behavior generated by π_E

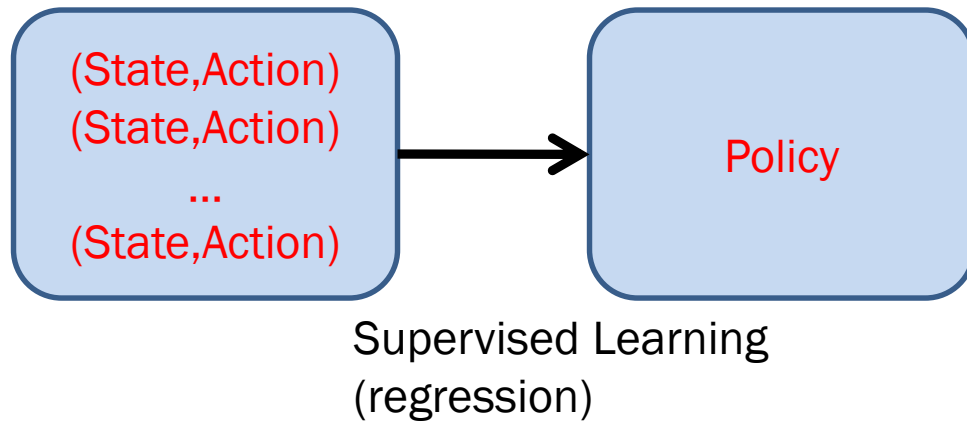
$$\{(s_0^i, a_0^i, s_1^i, a_1^i, \dots)\}_{i=1}^n \sim \pi_E$$



Goal: learn *cost function (reward) or policy*

(Ng and Russell, 2000), (Abbeel and Ng, 2004; Syed and Schapire, 2007), (Ratliff et al., 2006), (Ziebart et al., 2008), (Kolter et al., 2008), (Finn et al., 2016), etc.

Behavioral Cloning



- Small errors compound over time (*cascading errors*)
- *Decisions are purposeful (require planning)*

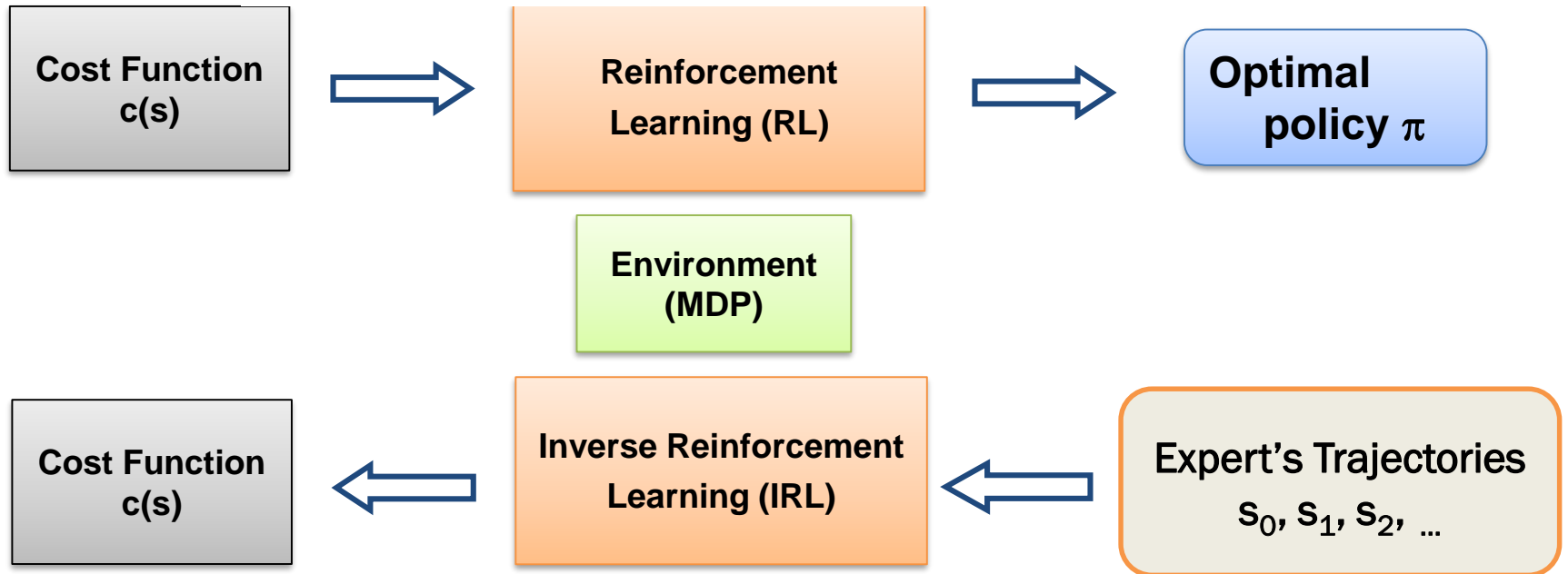
Inverse RL

- An approach to imitation
- Learns a cost c such that

$$\pi_E = \arg \min_{\pi \in \Pi} \mathbb{E}_{\pi} [c(s, a)]$$

Problem setup

$$RL(c) = \arg \min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_{\pi}[c(s, a)]$$



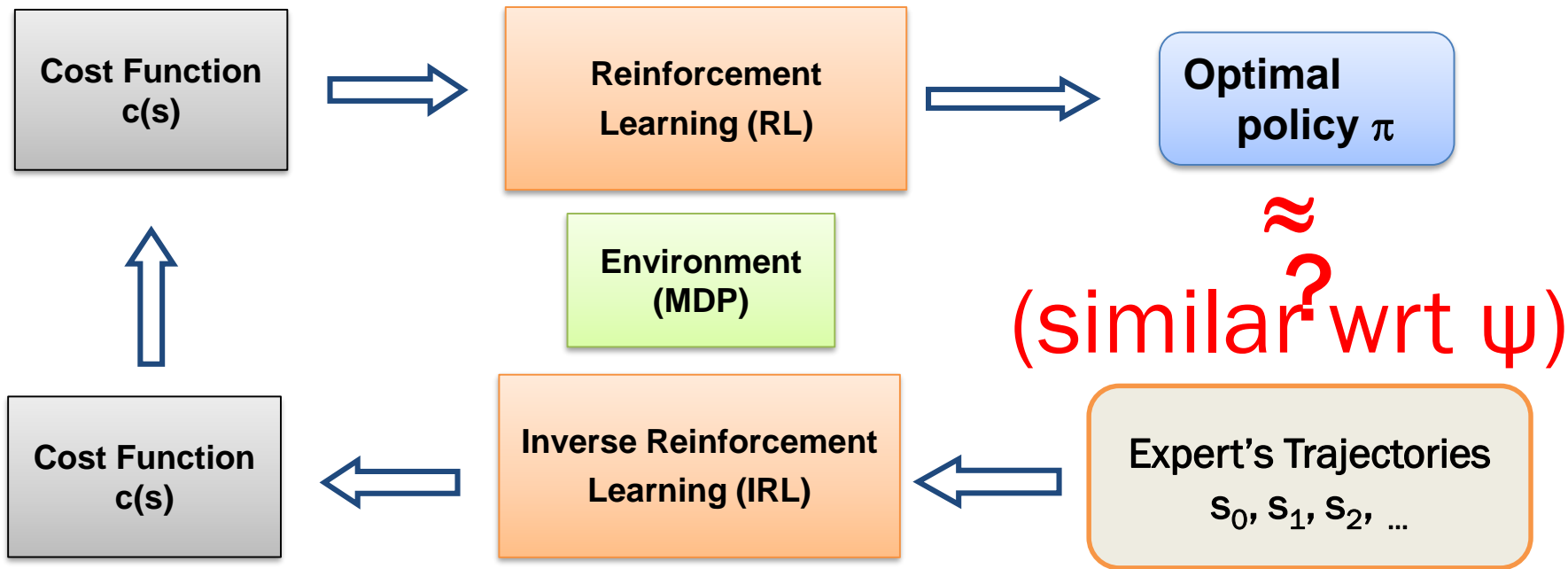
$$\underset{c \in \mathcal{C}}{\text{maximize}} \left(\min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_{\pi}[c(s, a)] \right) - \mathbb{E}_{\pi_E}[c(s, a)]$$

(Ziebart et al., 2010;
Rust 1987)

↑ Everything else
has high cost

↓ Expert has
small cost

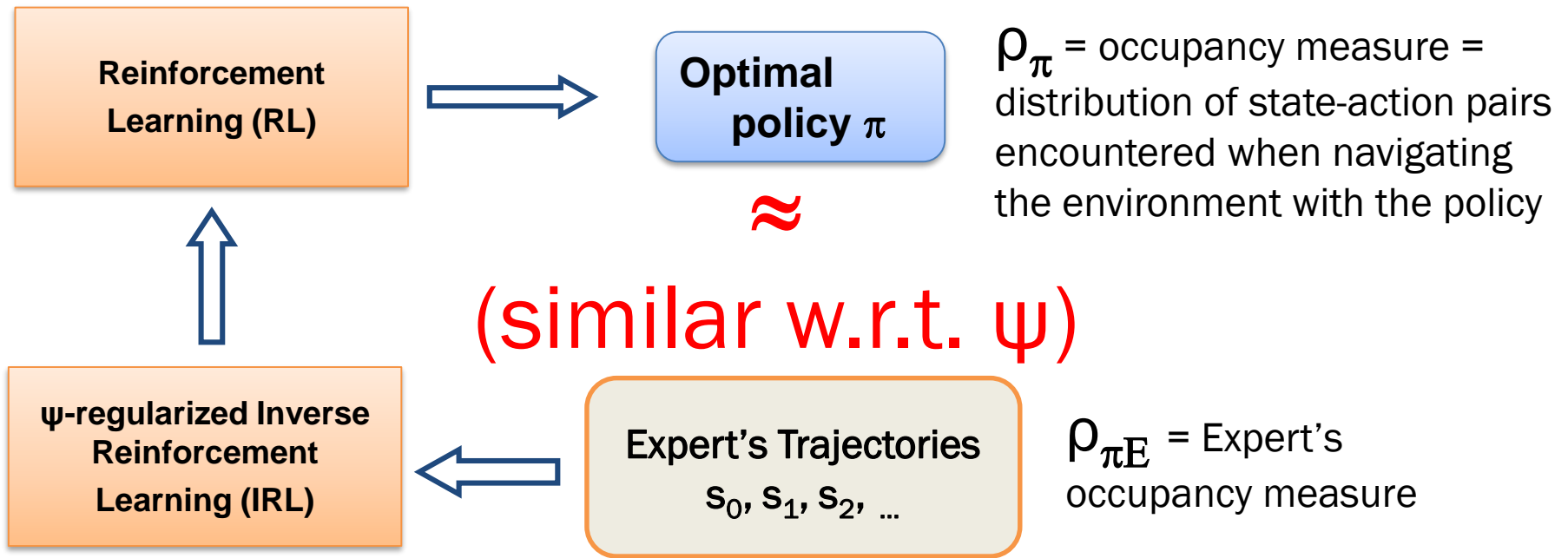
Problem setup



$$\text{IRL}_{\psi}(\pi_E) = \arg \max_{c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} \boxed{-\psi(c)} + \left(\min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_{\pi}[c(s, a)] \right) - \mathbb{E}_{\pi_E}[c(s, a)]$$

Convex cost regularizer

Combining RL ◦ IRL



Theorem: ψ -regularized inverse reinforcement learning, implicitly, **seeks a policy whose occupancy measure is close to the expert's**, as measured by ψ^* (convex conjugate of ψ)

$$\text{RL} \circ \text{IRL}_\psi(\pi_E) = \arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_\pi - \rho_{\pi_E})$$

Takeaway

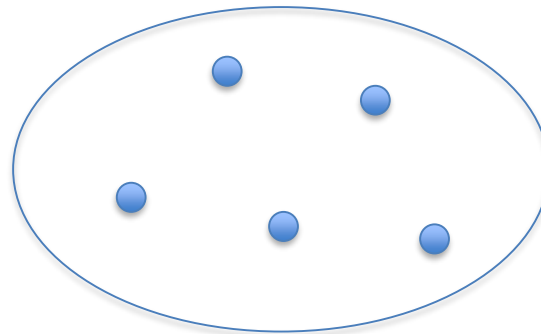
Theorem: ψ -regularized inverse reinforcement learning, implicitly, **seeks a policy whose occupancy measure is close to the expert's**, as measured by ψ^*

- Typical IRL definition: finding a cost function **c** such that the expert policy is uniquely optimal w.r.t. **c**
- Alternative view: IRL as a procedure that tries to induce a policy that matches the expert's occupancy measure (**generative model**)

Special cases

$$\text{RL} \circ \text{IRL}_\psi(\pi_E) = \arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_\pi - \rho_{\pi_E})$$

- If $\psi(c)=\text{constant}$, then $\rho_{\tilde{\pi}} = \rho_{\pi_E}$
 - Not a useful algorithm. In practice, we only have sampled trajectories
- **Overfitting:** Too much flexibility in choosing the cost function (and the policy)



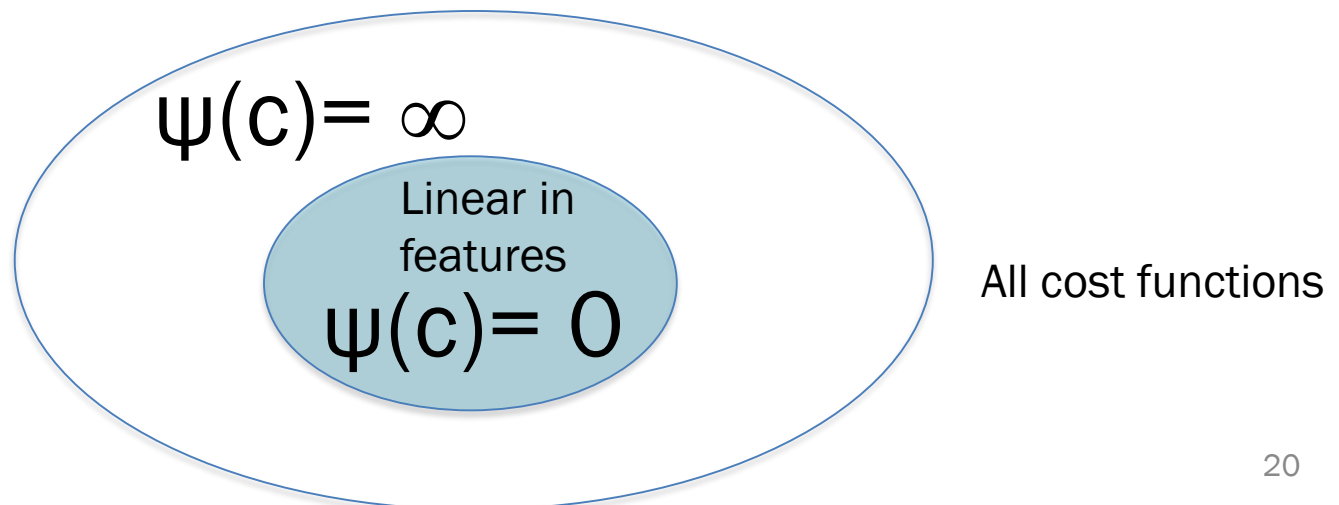
All cost functions
 $\psi(c)=\text{constant}$

Towards Apprenticeship learning

- Solution: use **features** $f_{s,a}$
- Cost $c(s,a) = \theta \cdot f_{s,a}$

$$\text{IRL}_\psi(\pi_E) = \arg \max_{c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} -\psi(c) + \left(\min_{\pi \in \Pi} -H(\pi) + \mathbb{E}_\pi[c(s, a)] \right) - \mathbb{E}_{\pi_E}[c(s, a)]$$

Only these “simple” cost functions are allowed



Apprenticeship learning

- For that choice of ψ , $\text{RL} \circ \text{IRL}_\psi$ framework gives apprenticeship learning

$$\text{RL} \circ \text{IRL}_\psi(\pi_E) = \arg \min_{\pi \in \Pi} -H(\pi) + \psi^*(\rho_\pi - \rho_{\pi_E})$$

- Apprenticeship learning: find π performing better than π_E over costs linear in the features
 - Abbeel and Ng (2004)
 - Syed and Schapire (2007)

Apprenticeship learning

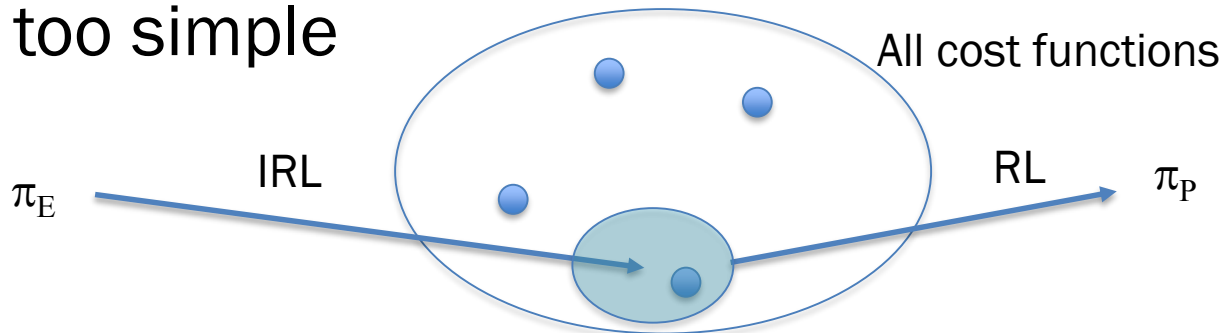
- Given $\{(s_0^i, a_0^i, s_1^i, a_1^i, \dots)\}_{i=1}^n \sim \pi_E$
- Goal: find π performing better than π_E over a class of costs

$$\underset{\pi}{\text{minimize}} \max_{c \in \mathcal{C}} \mathbb{E}_{\pi} [c(s, a)] - \mathbb{E}_{\pi_E} [c(s, a)]$$

Approximated using
demonstrations

Issues with Apprenticeship learning

- Need to craft features very carefully
 - unless the true expert cost function (assuming it exists) lies in C , there is no guarantee that AL will recover the expert policy
- $RL \circ IRL_{\psi}(\pi_E)$ is “encoding” the expert behavior as a cost function in C .
 - it might not be possible to decode it back if C is too simple

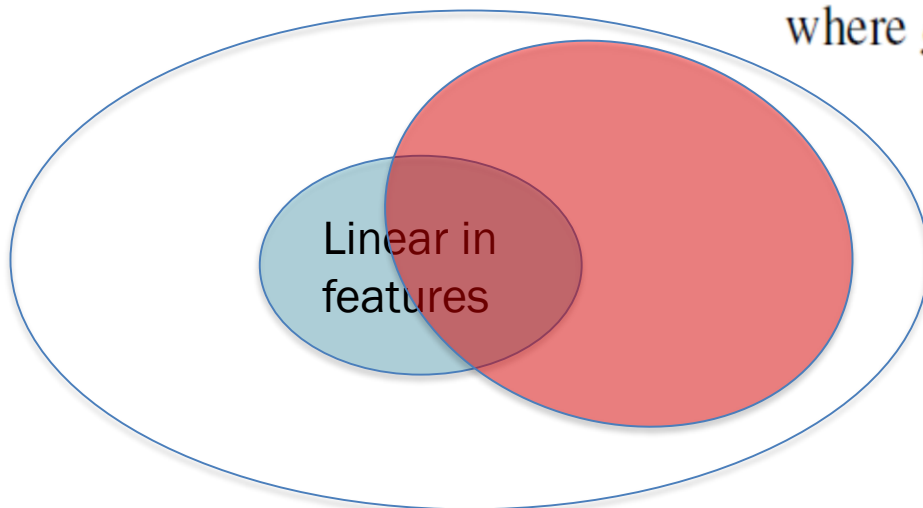


Generative Adversarial Imitation Learning

- **Solution:** use a more expressive class of cost functions

$$\psi_{\text{GA}}(c) \triangleq \begin{cases} \mathbb{E}_{\pi_E}[g(c(s, a))] & \text{if } c < 0 \\ +\infty & \text{otherwise} \end{cases}$$

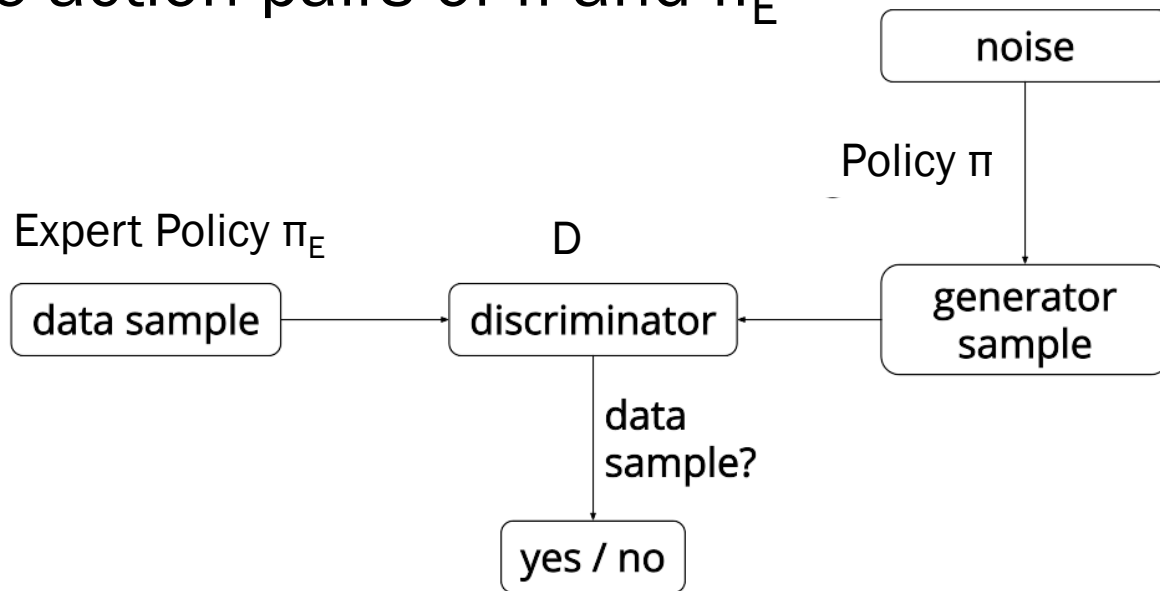
All cost functions



$$\text{where } g(x) = \begin{cases} -x - \log(1 - e^x) & \text{if } x < 0 \\ +\infty & \text{otherwise} \end{cases}$$

Generative Adversarial Imitation Learning

- ψ^* = optimal negative log-loss of the binary classification problem of distinguishing between state-action pairs of π and π_E



$$\psi_{\text{GA}}^*(\rho_\pi - \rho_{\pi_E}) = \sup_{D \in (0,1)^{\mathcal{S} \times \mathcal{A}}} \mathbb{E}_\pi[\log(D(s, a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s, a))]$$

Generative Adversarial Networks

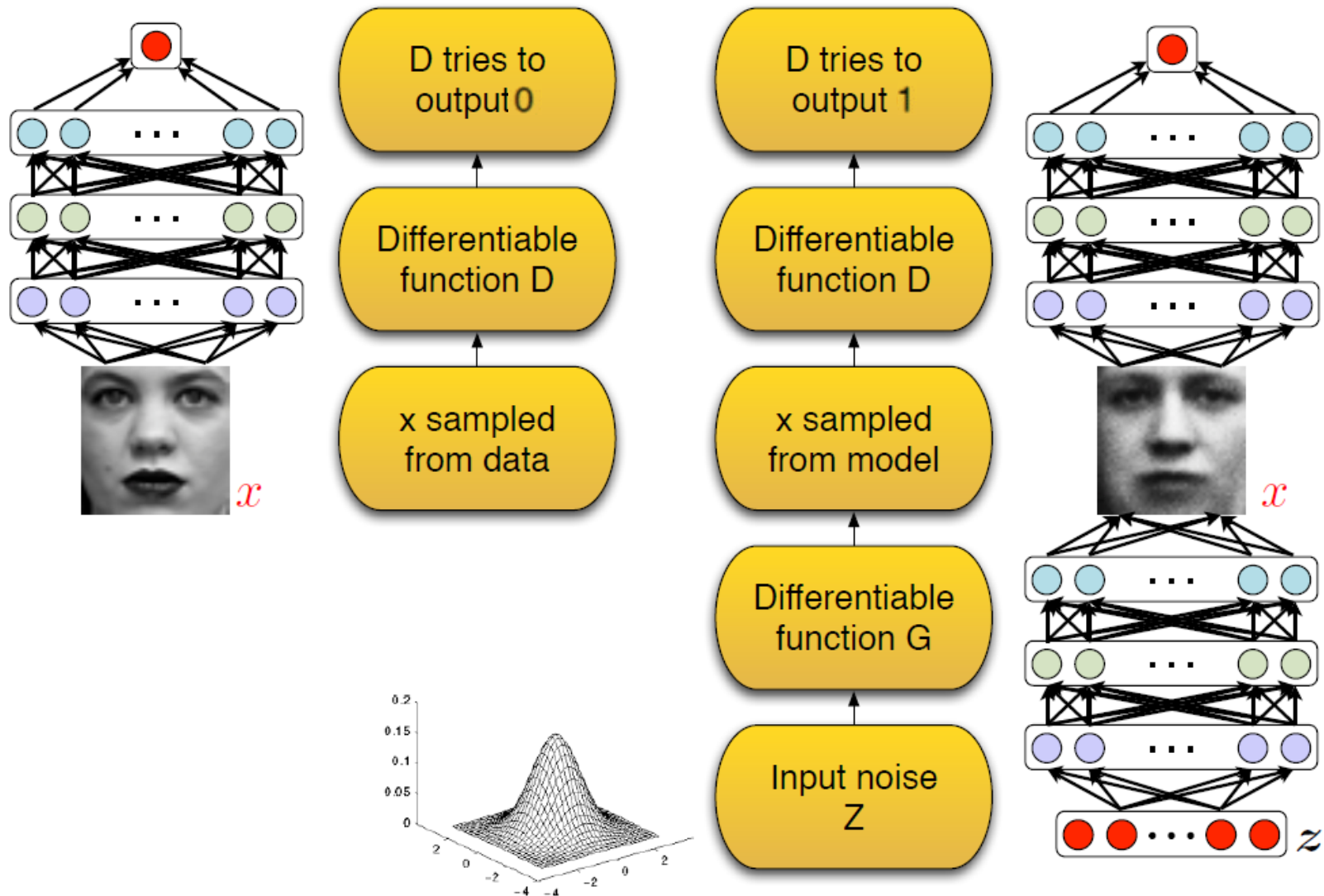
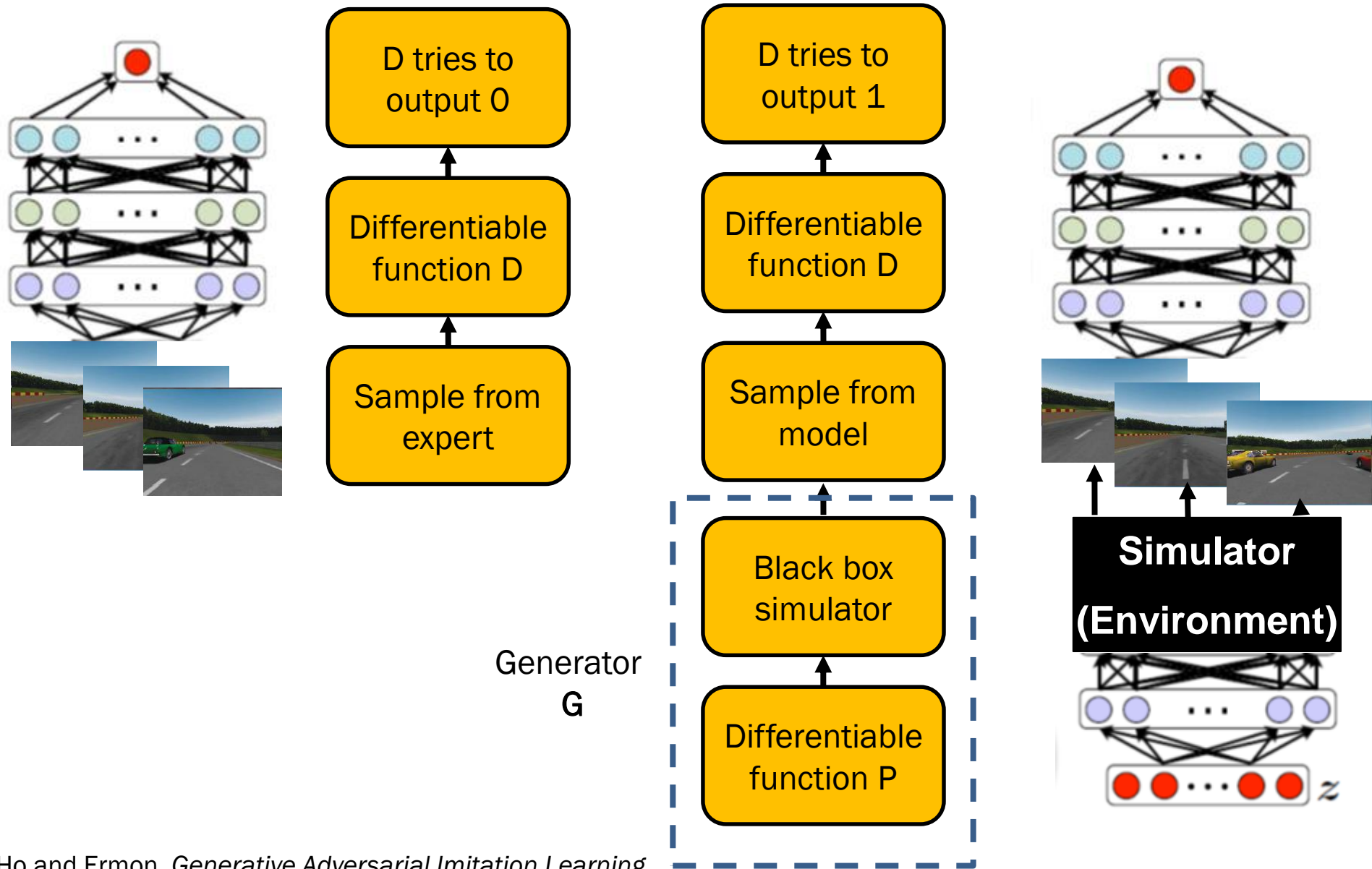


Figure from Goodfellow et al, 2014

GAIL



How to optimize the objective

- Previous Apprenticeship learning work:
 - Full dynamics model
 - Small environment
 - Repeated RL
- We propose: gradient descent over policy parameters (and discriminator)

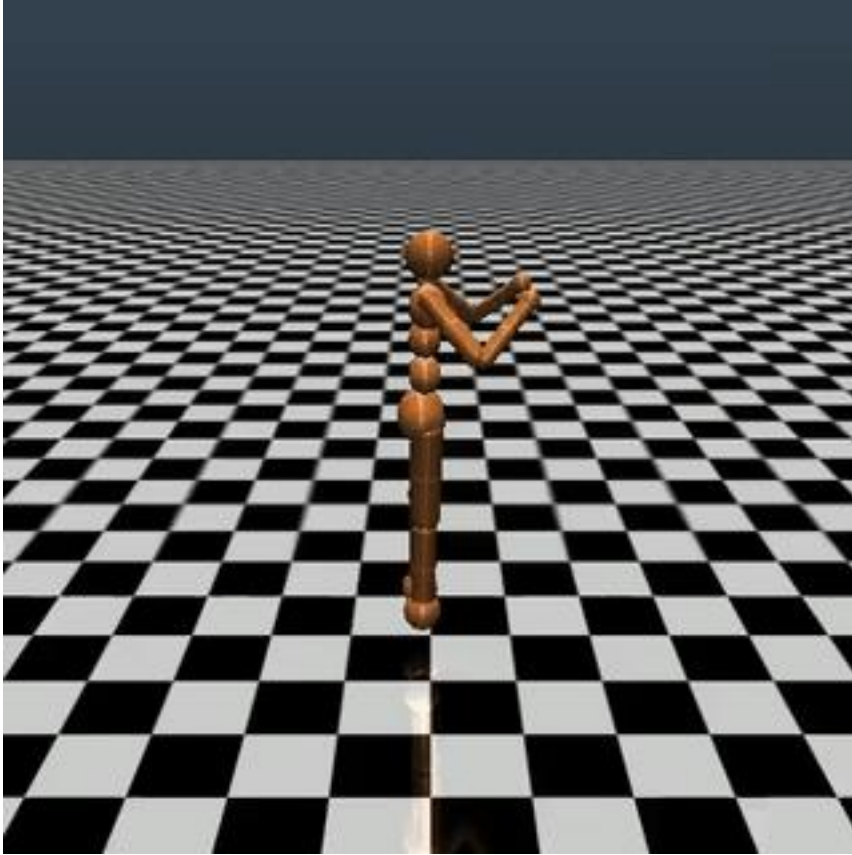
Properties

- Inherits pros of policy gradient
 - Convergence to local minima
 - Can be model free
- Inherits cons of policy gradient
 - High variance
 - Small steps required

Properties

- Inherits pros of policy gradient
 - Convergence to local minima
 - Can be model free
- Inherits cons of policy gradient
 - High variance
 - Small steps required
- **Solution: trust region policy optimization**

Results



Results

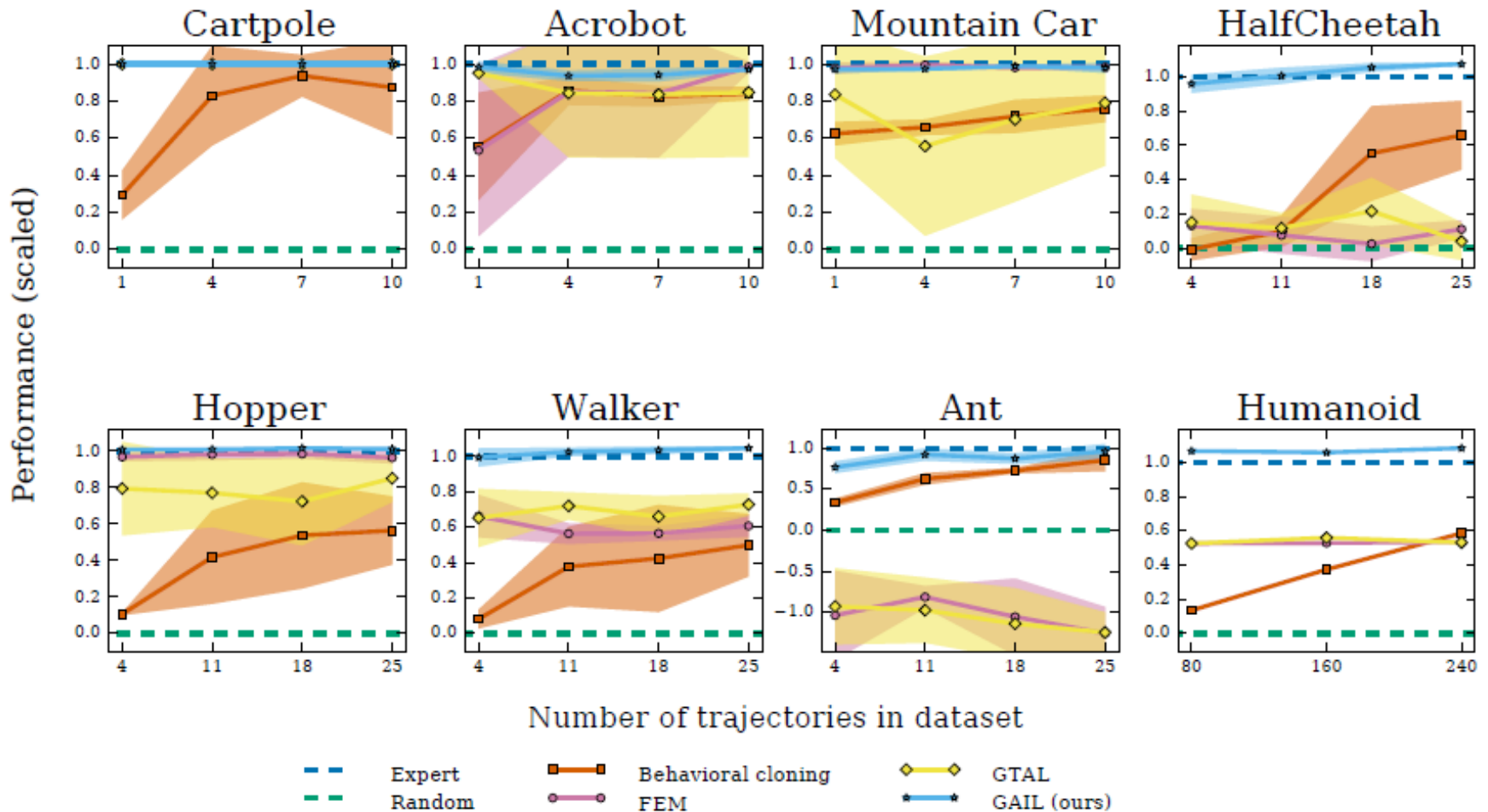
Input: driving demonstrations (Torcs)

Output policy:



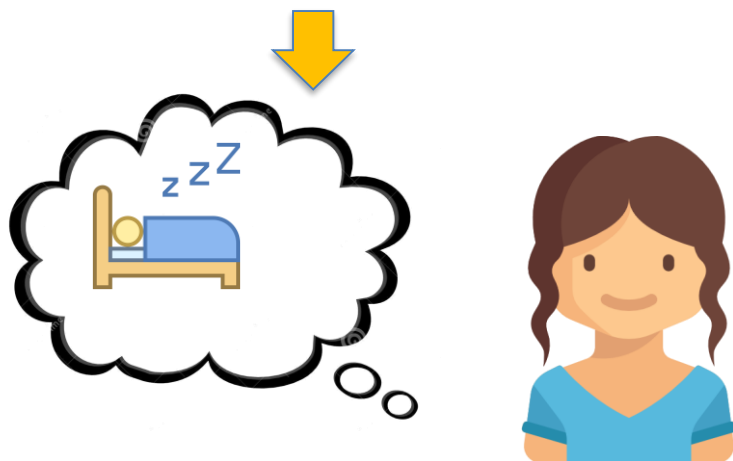
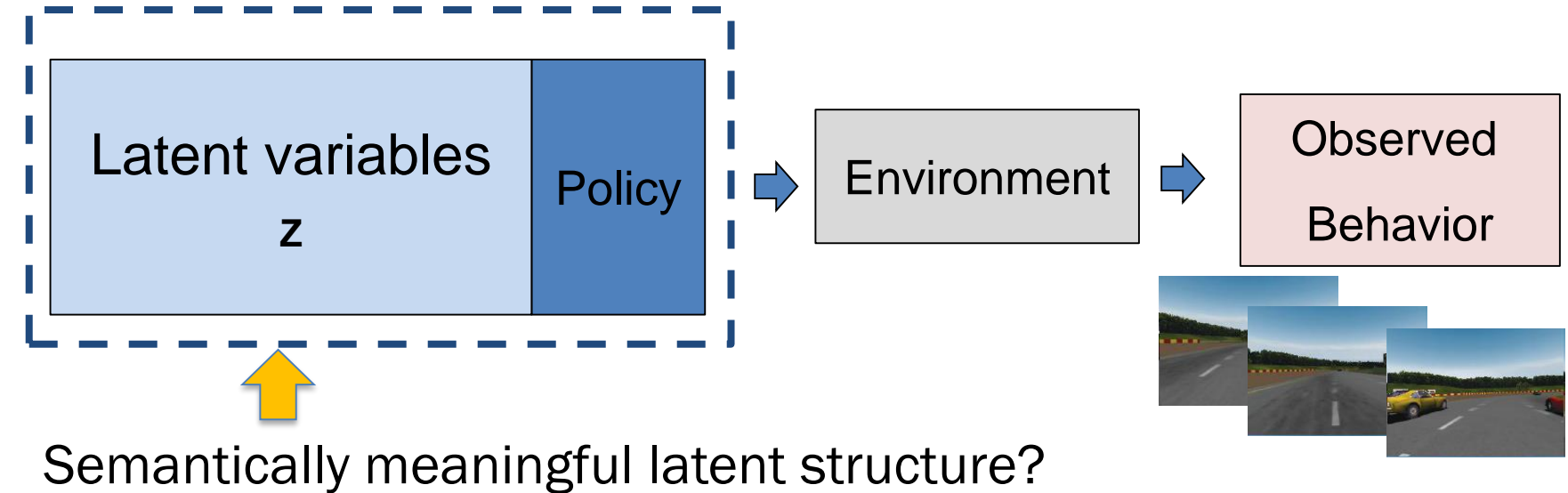
From raw visual inputs

Experimental results



Latent structure in demonstrations

Human model



InfoGAIL

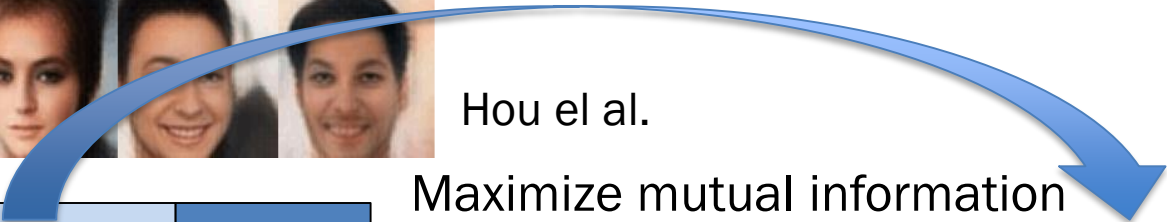
Latent structure



Observed data

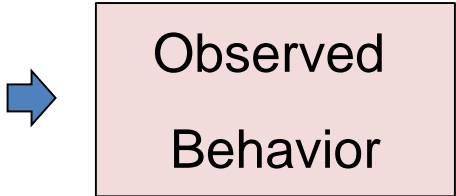
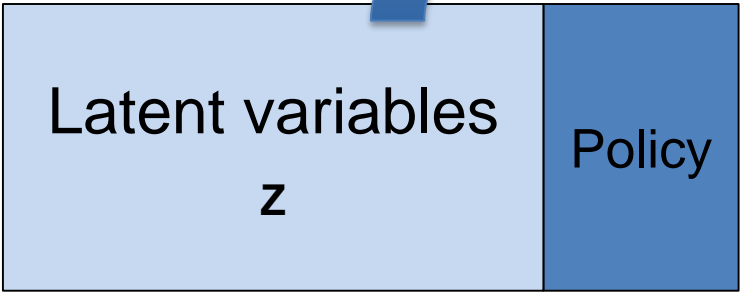


Infer structure



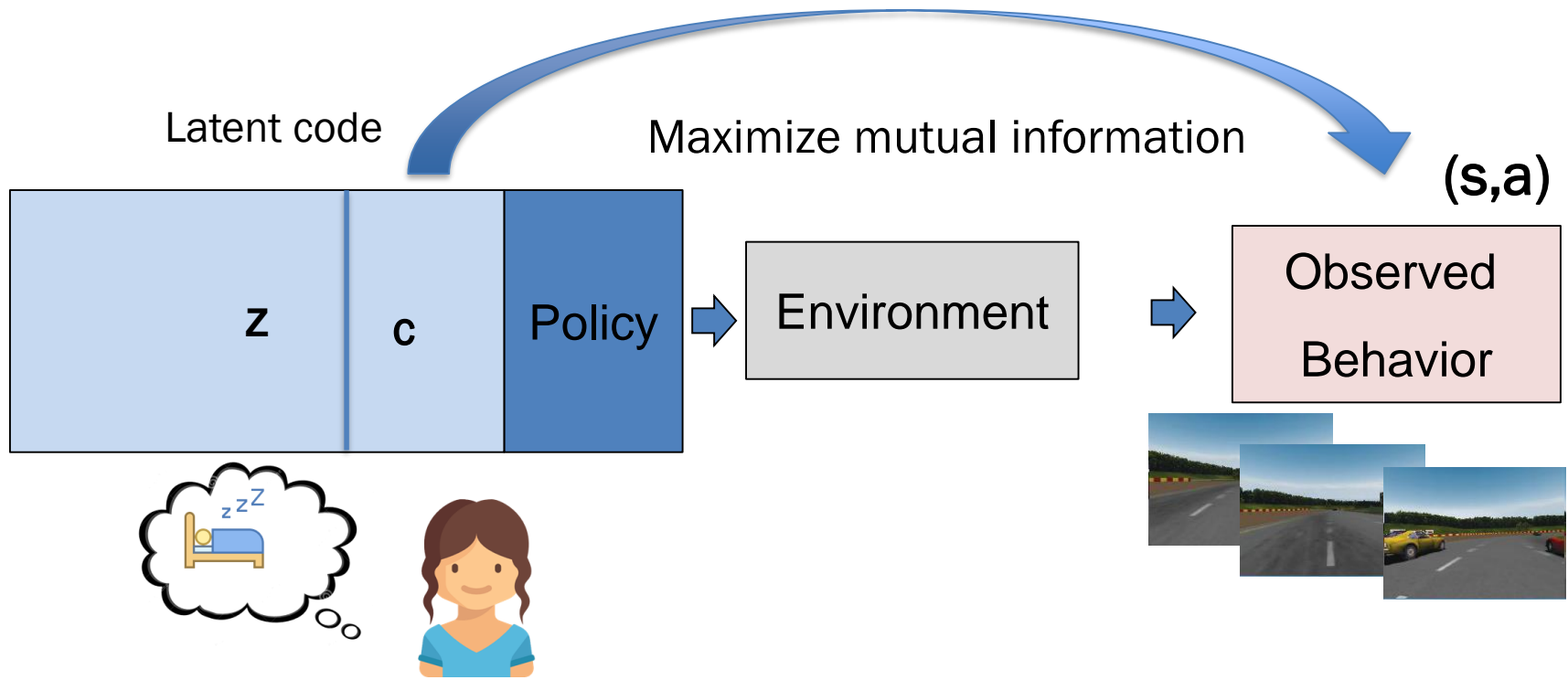
Hou et al.

Maximize mutual information

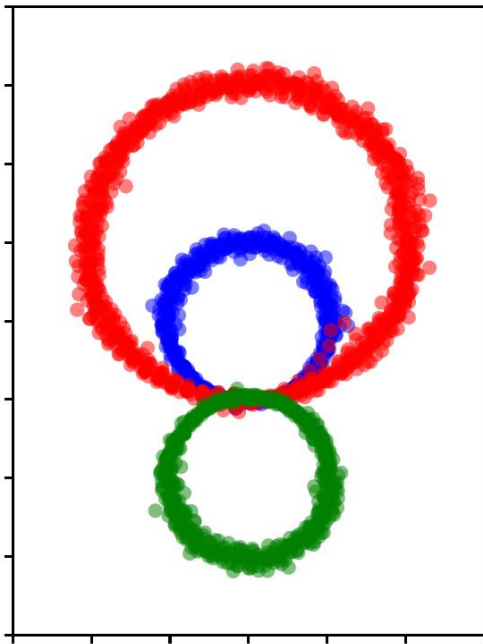


InfoGAIL

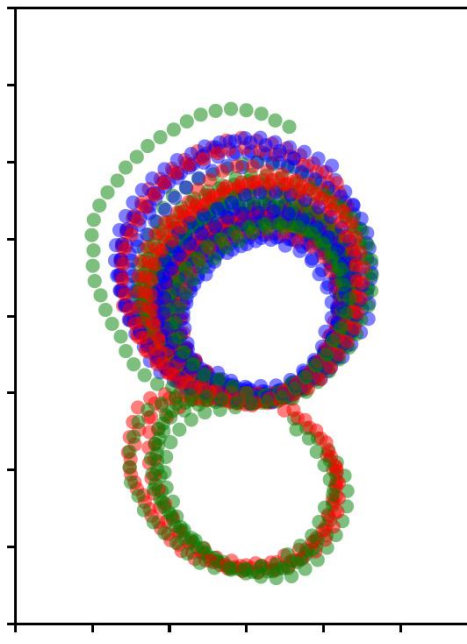
$$L_I(\pi_\theta, Q_\psi) = \mathbb{E}_{c \sim p(c), a \sim \pi_\theta(\cdot | s, c)} [\log Q_\psi(c | s, a)] + H(c) \\ \leq I(c; s, a)$$



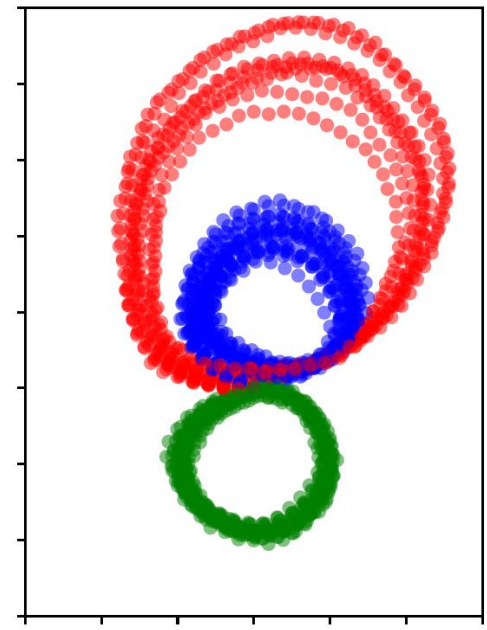
Synthetic Experiment



Demonstrations

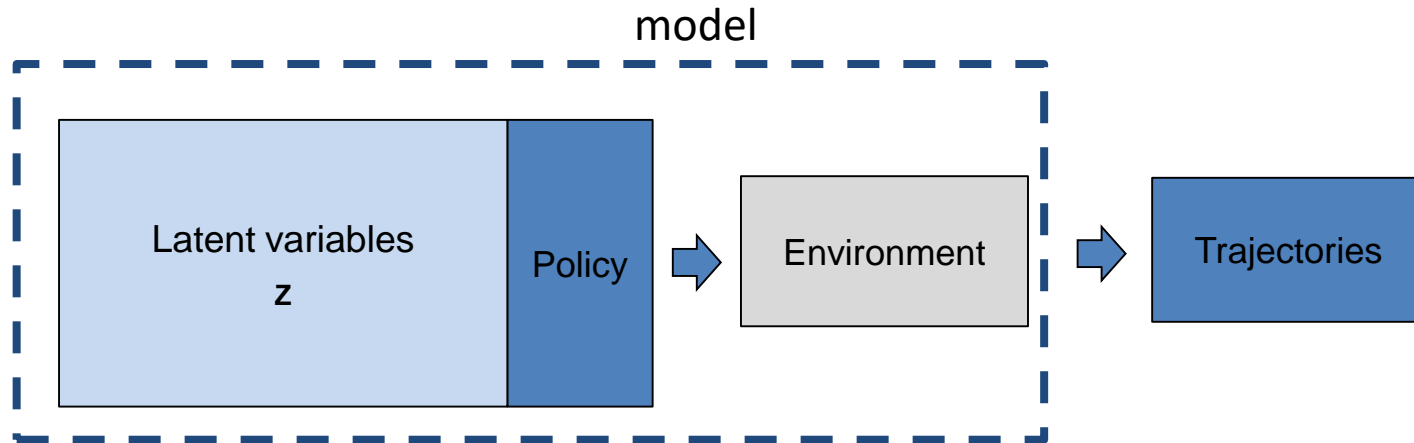


GAIL



Info-GAIL

InfoGAIL



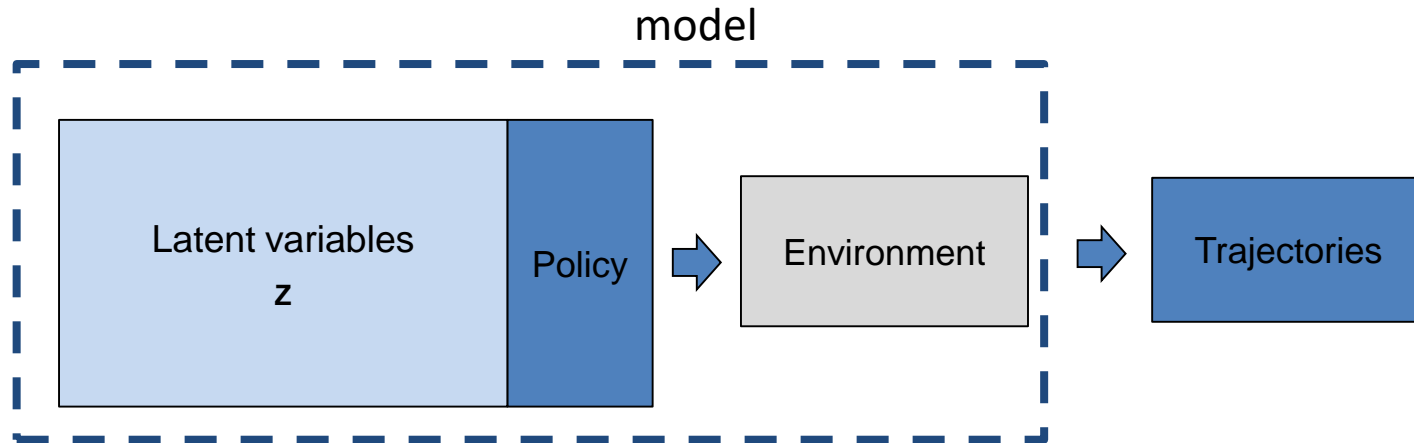
Pass left ($z=0$)



Pass right ($z=1$)



InfoGAIL



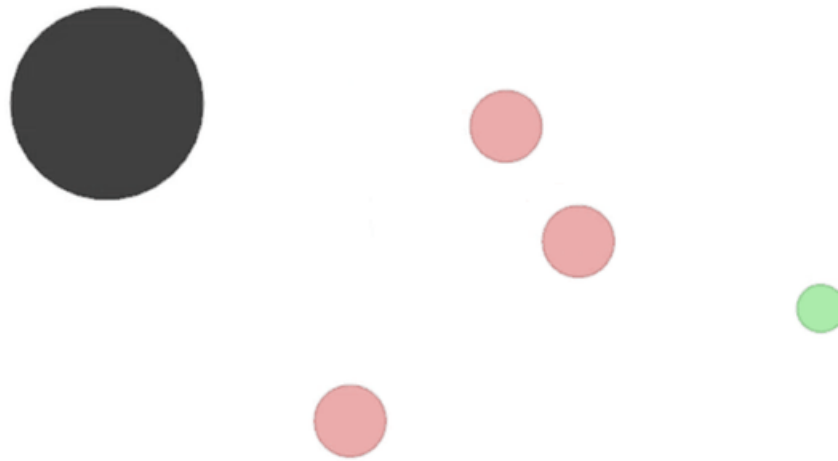
Turn inside ($z=0$)



Turn outside ($z=1$)

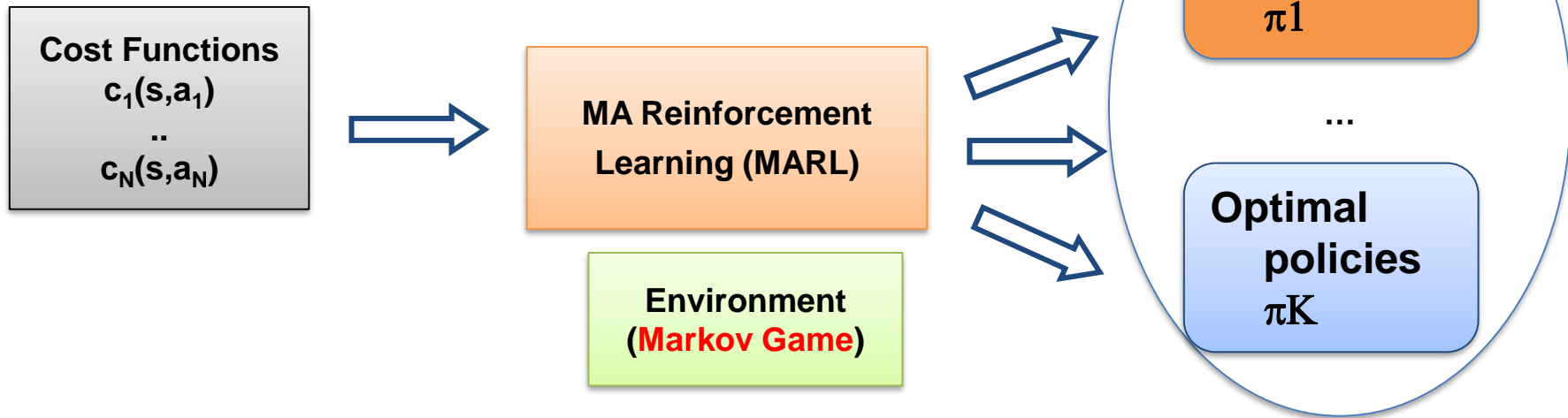


Multi-agent environments



What are the goals of these 4 agents?

Problem setup



	R	L
R	0,0	10,10
L	10,10	0,0

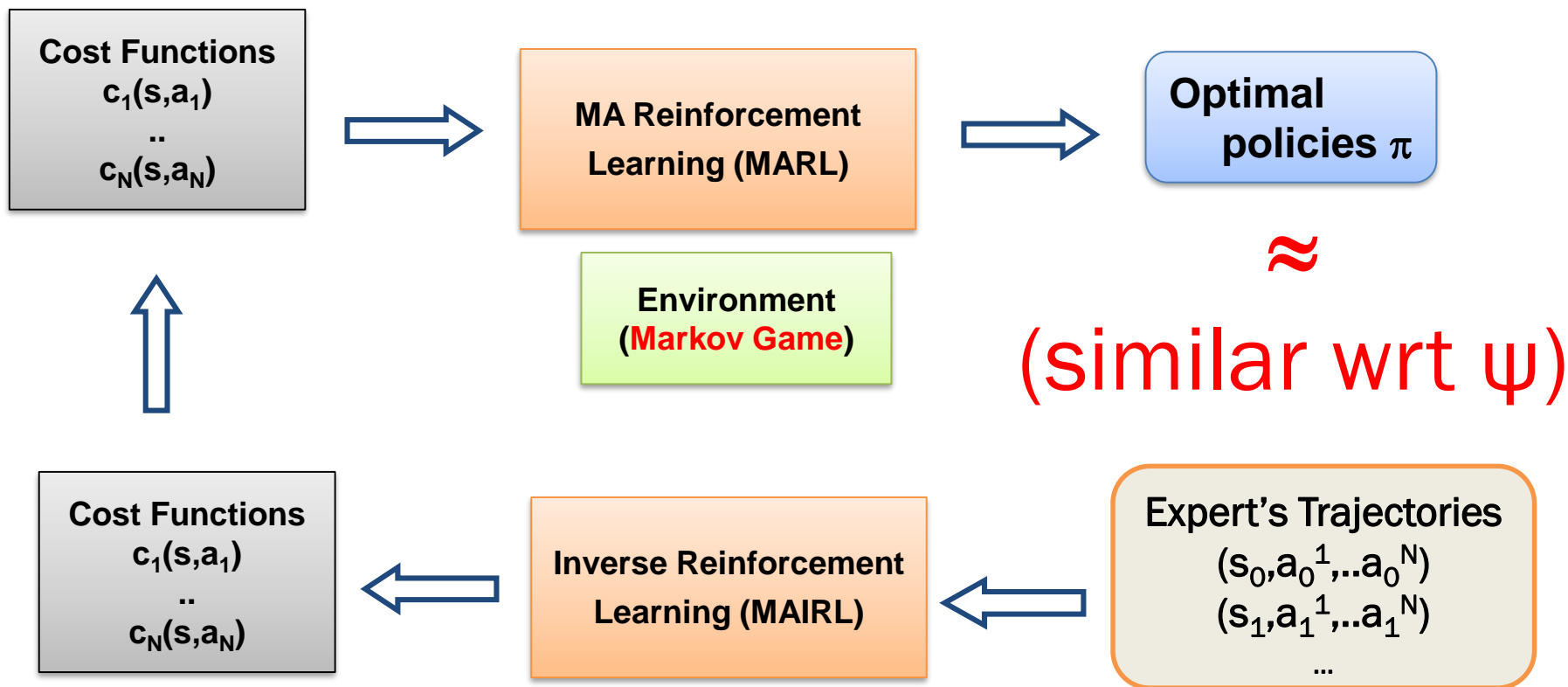
DRIVE ON LEFT



DRIVE ON RIGHT



Problem setup

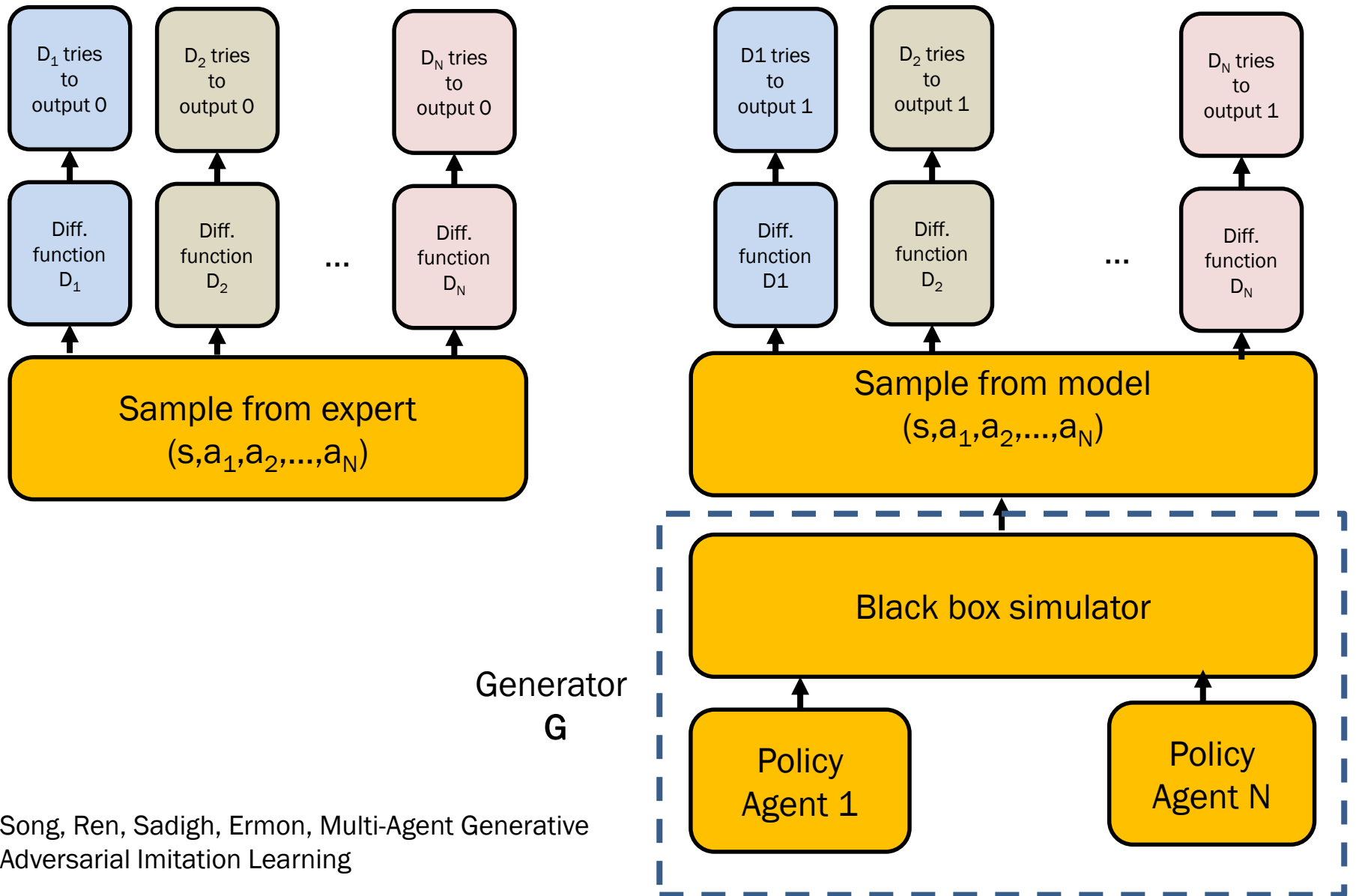


$$\text{MIM}_\psi(\pi_E) = \arg \max_{\pi \in \Pi} \max_v \min_{r \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}} \mathcal{L}_\psi(\pi, v)$$

$$\mathcal{L}_\psi(\pi, v) = -f_r(\pi, v) + f_r(\pi_E, v) + \psi(r)$$

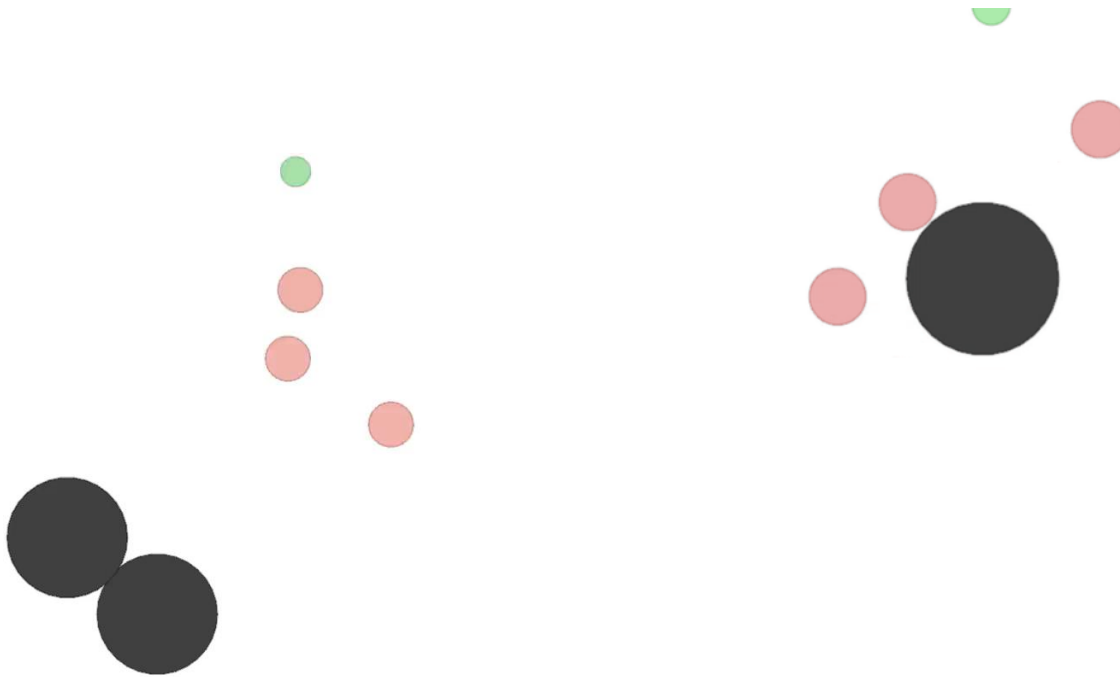
$$r \in \text{MAIRL}(\pi_E)$$

MAGAIL



Generator
G

Environments



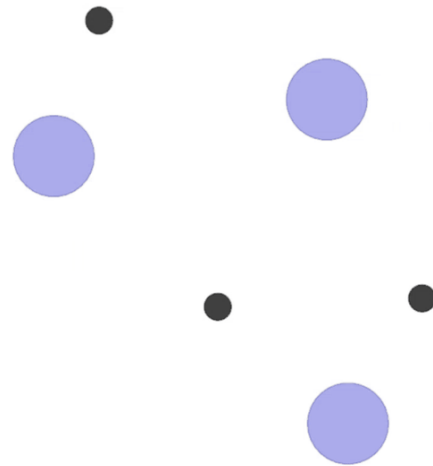
Demonstrations

MAGAIL

Environments

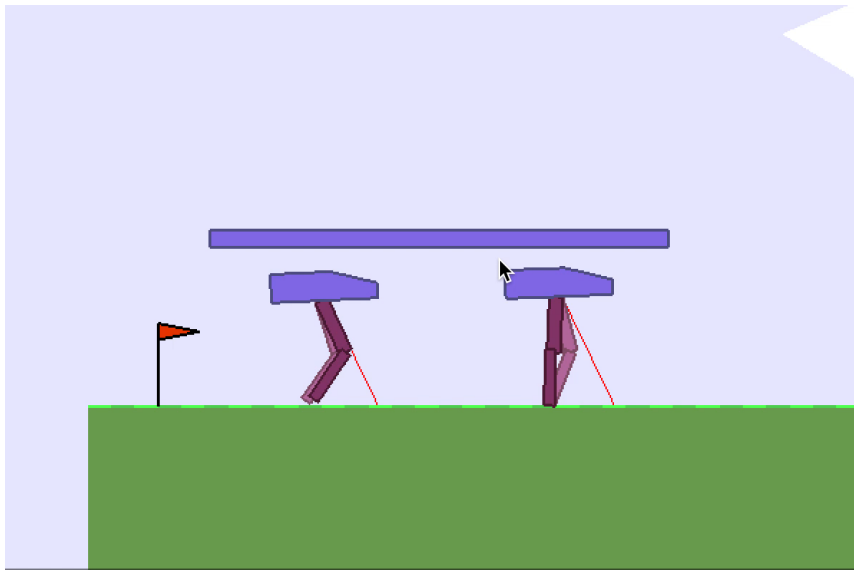


Demonstrations

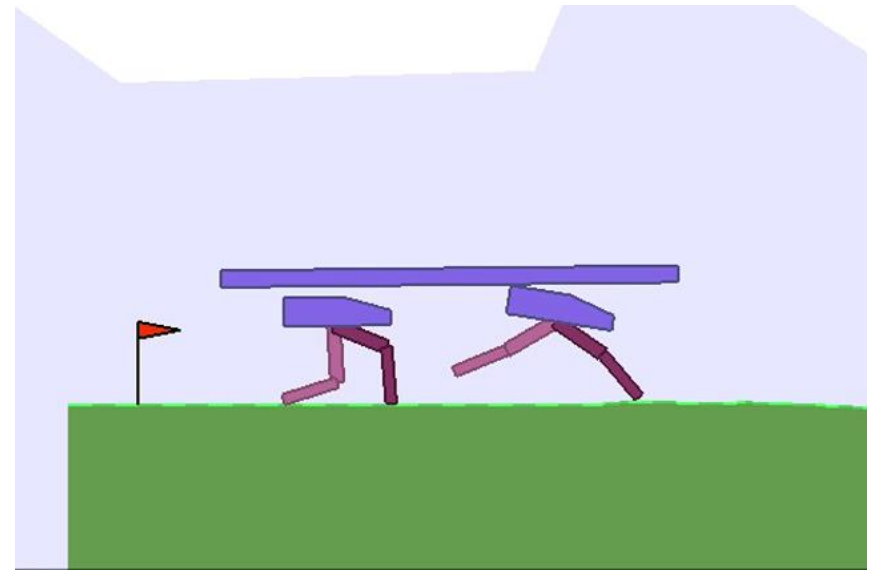


MAGAIL

Suboptimal demos



Expert



MAGAIL

lighter plank + bumps on ground

Conclusions

- IRL is a dual of an occupancy measure matching problem (generative modeling)
- Might need flexible cost functions
 - GAN style approach
- Policy gradient approach
 - Scales to high dimensional settings
- Towards unsupervised learning of latent structure from demonstrations