



Adversarial Imitation Learning via Random Search

Authors: **MyungJae Shin (Presenter)** and **Joongheon Kim**

School of Computer Science and Engineering, Chung-Ang University, Seoul, Republic of Korea

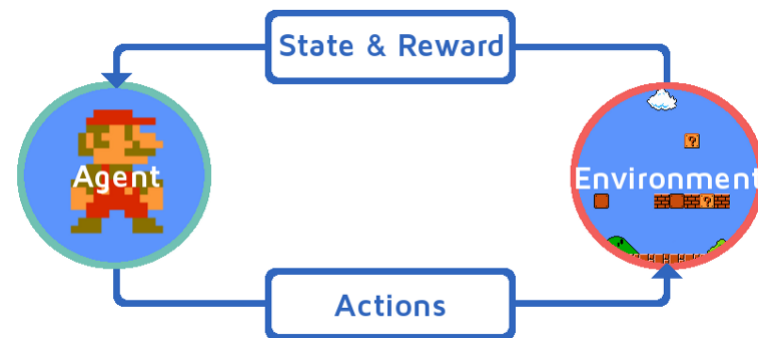
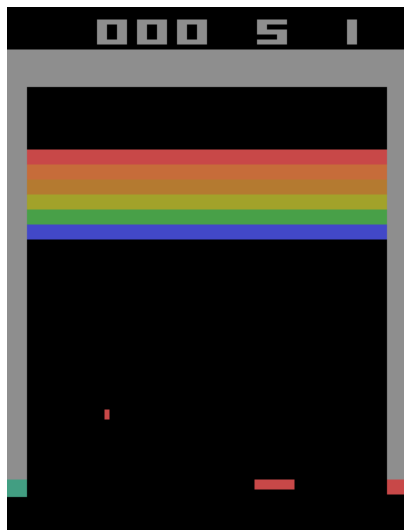
Emails: mjshin.cau@gmail.com , joongheon@gmail.com

Sites: github.com/170928 , <https://sites.google.com/site/joongheonkim/>

Reinforcement Learning



Goal : Learn policies
High-dimensional & raw observations



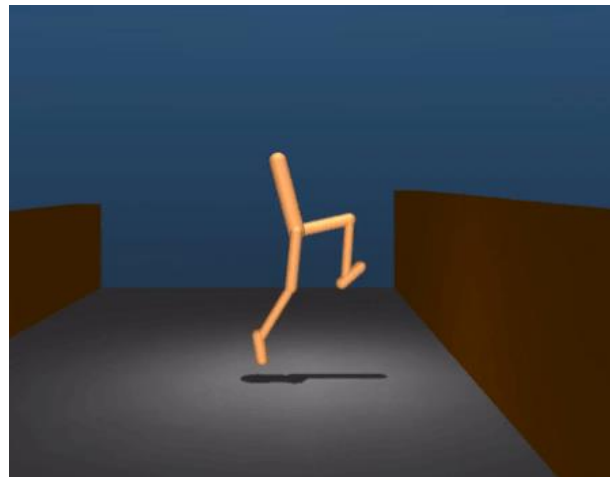
Imitation Learning



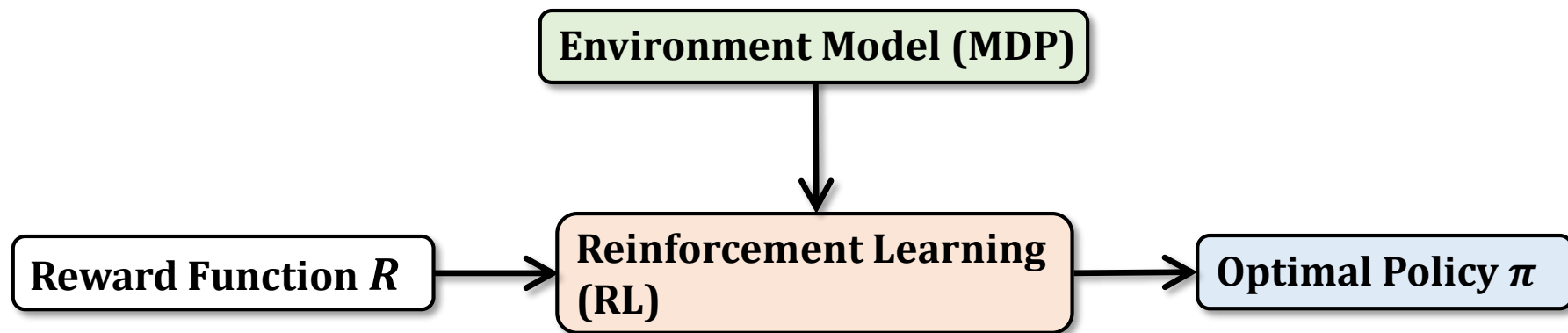
Input : expert behavior generated by expert π_E

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

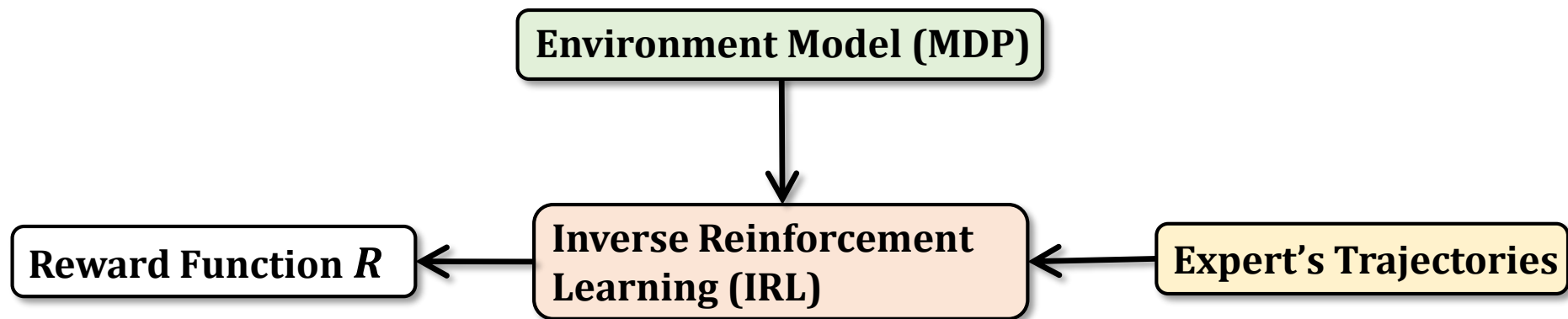
Goal : learn **cost function** or **policy**



Imitation Learning



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

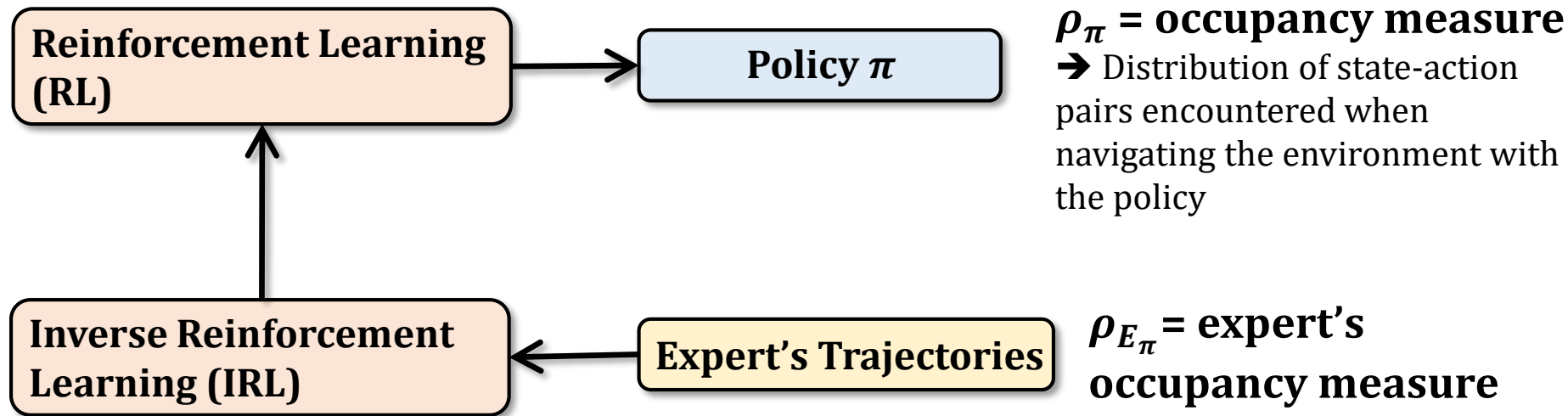


$$\max_R \left(\min_{\pi} \mathbb{E}_{\pi} [R(s, a)] - H(\pi) \right) - \mathbb{E}_{\pi_E} [R(s, a)]$$

Imitation Learning



$$\max_R -\psi(R) + \left(\min_{\pi} \mathbb{E}_{\pi} [R(s, a)] - H(\pi) \right) - \mathbb{E}_{\pi_E} [R(s, a)]$$



$$\min_{\pi} \psi^* (\rho_{\pi} - \rho_{E_{\pi}}) - H(\pi)$$

Imitation Learning



[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 finds a cost function such that the expert policy is uniquely optimal
- IRL as a procedure that tries to induce a policy that matches the expert's occupancy measure (generative model)

Generative Adversarial Imitation Learning (GAIL), *NIPS 2016*

Use this regularizer

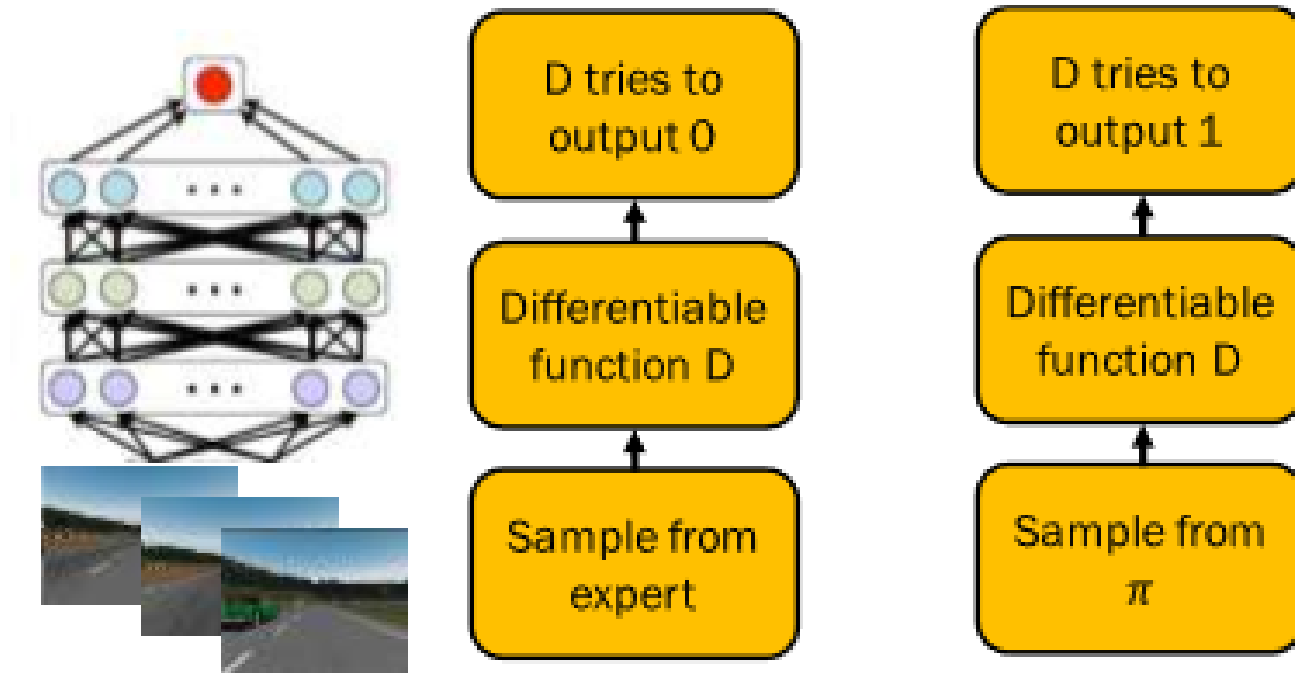
$$\psi_{GA}(R) = \begin{cases} \mathbb{E}_{\pi_E}[g(R(s, a))] & \text{if } R < 0 \\ +\infty & \text{otherwise} \end{cases}$$

Generative Adversarial Networks, Ian J. Goodfellow, *NIPS 2014*

Adversarial Imitation Learning via Random Search



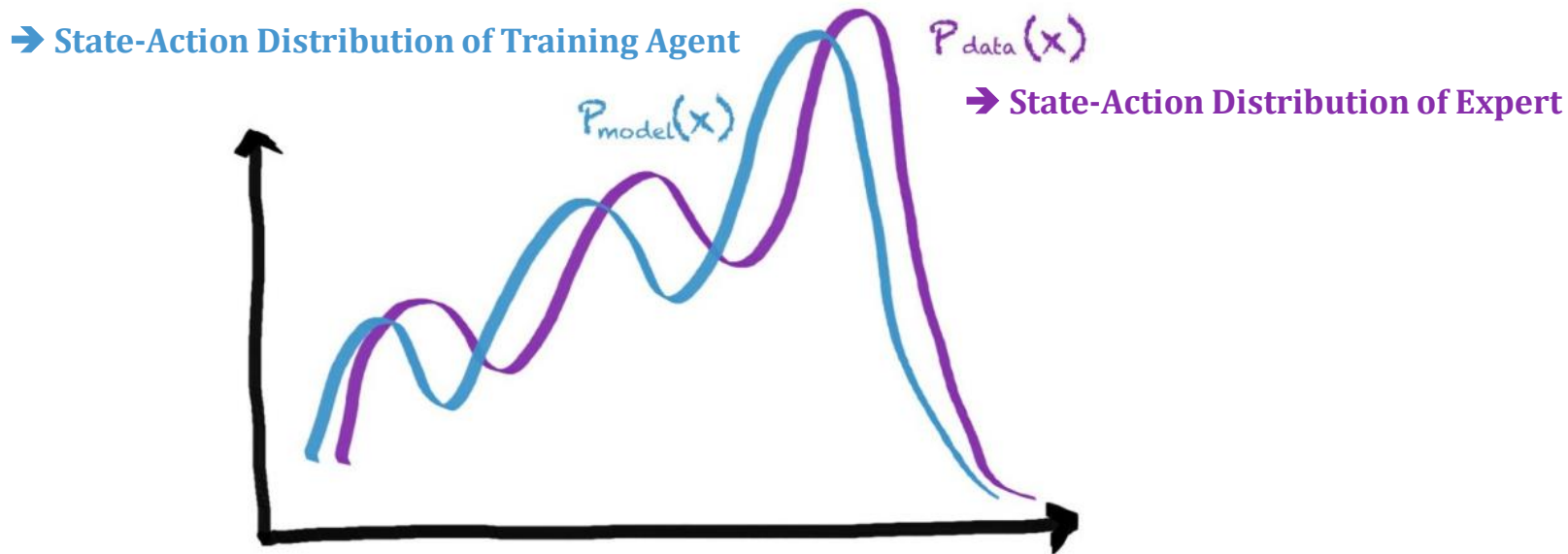
Generative Adversarial Imitation Learning (GAIL), *NIPS 2016*



Adversarial Imitation Learning via Random Search



Generative Adversarial Imitation Learning (GAIL), *NIPS 2016*

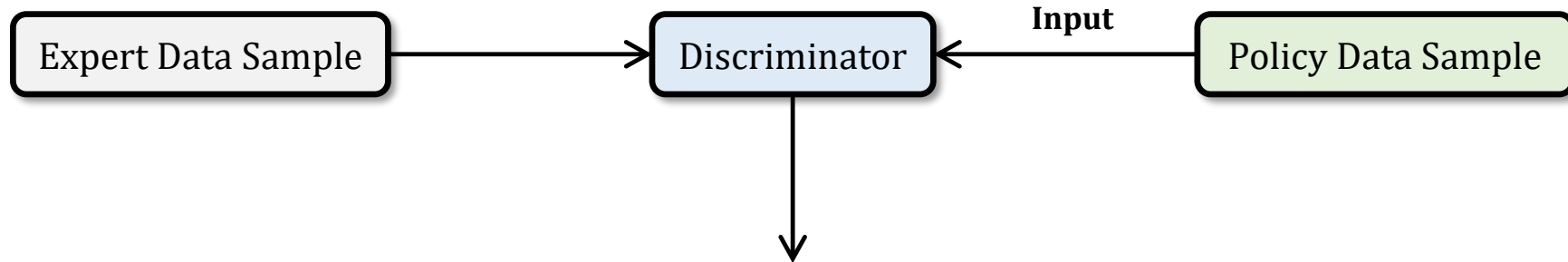


Based on the output of the discriminator (Generative Adversarial Networks, [Ian J. Goodfellow](#), 2014), we could know the difference between the distribution of expert data and that of agent.

Adversarial Imitation Learning via Random Search



$$\text{minimize } \mathbb{E}_{\pi}[\log(D(s, a))] + \mathbb{E}_{\pi_E}[\log(1 - D(s, a))]$$



$D(s, a)$: Probability between 0 and 1

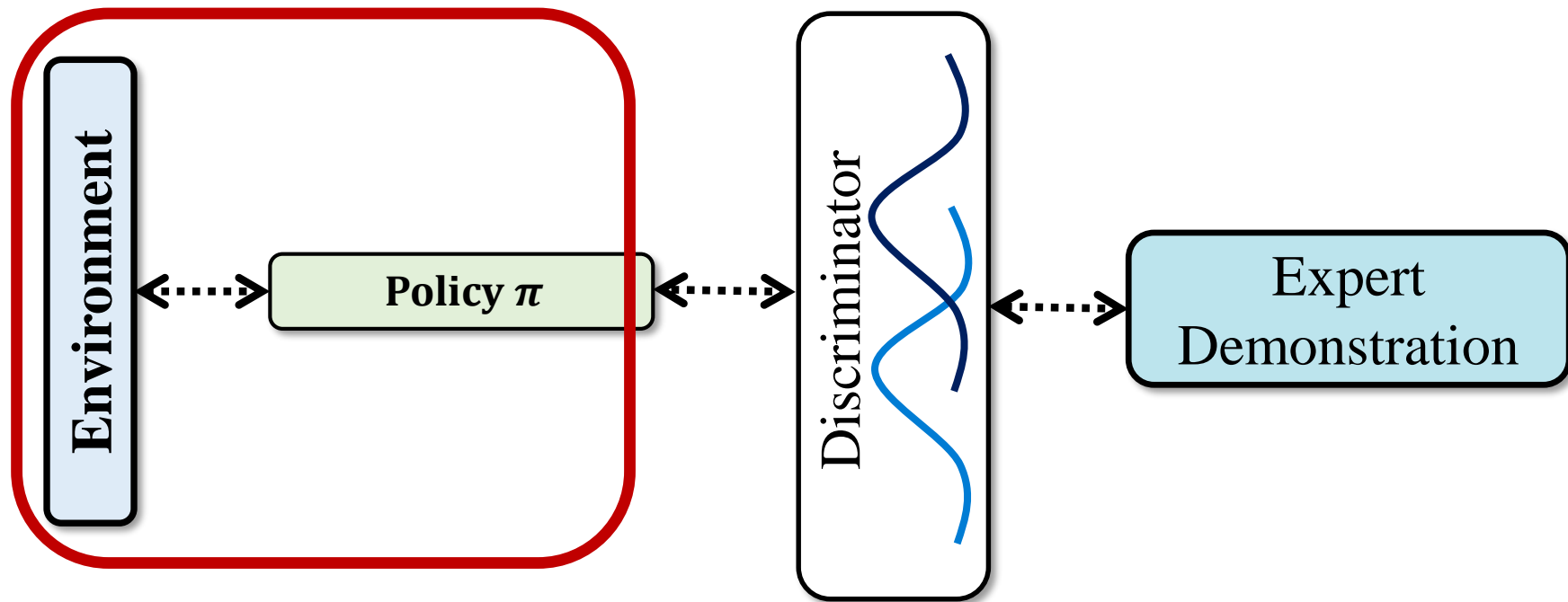
The probability that the input data sample is the expert data sample

Adversarial Imitation Learning via Random Search



Challenge

A lot of interaction with the environment is required to optimize the policy through GAIL framework





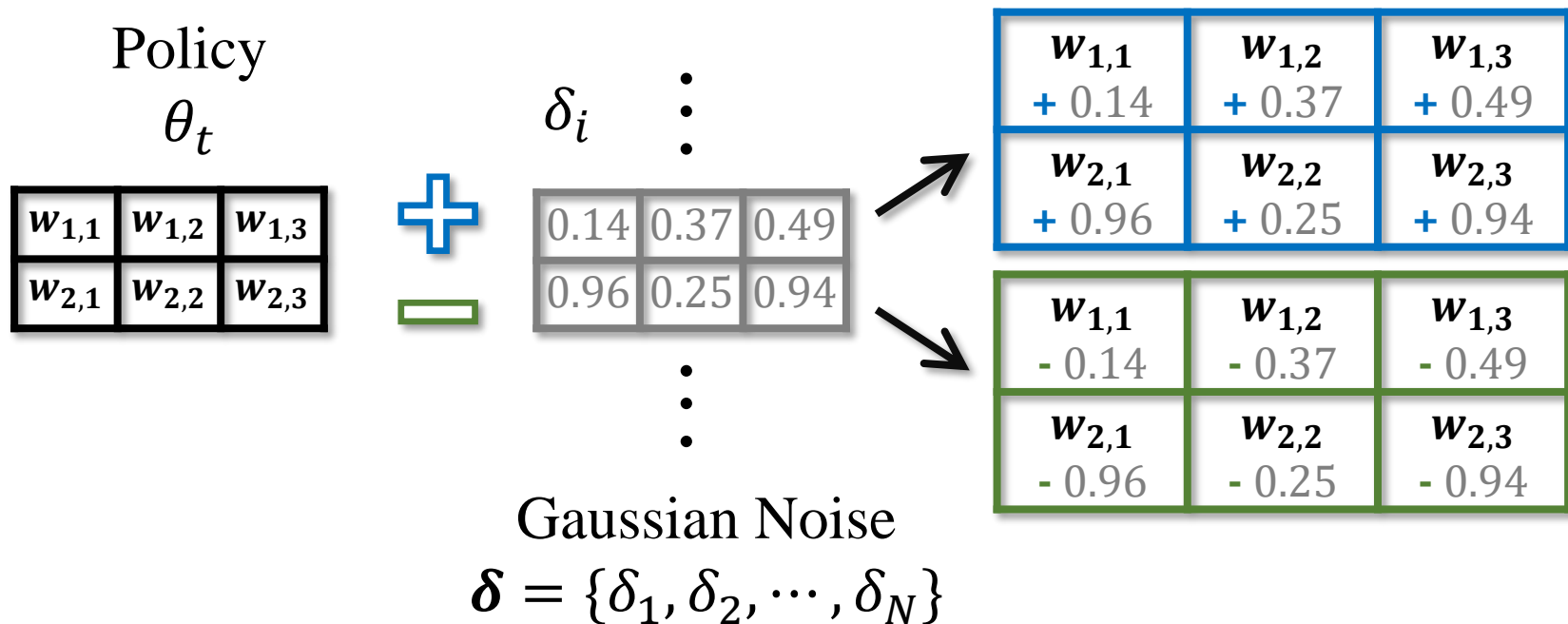
Usually in AI:

$$f'(x) = \frac{df}{dx}$$

Proposed method

$$f'(a) = \frac{f(a+h) - f(a)}{h}$$

Adversarial Imitation Learning via Random Search




Adversarial Imitation Learning via Random Search




R_{d-pos}

$w_{1,1}$ $+d_{1,1}$	$w_{1,2}$ $+d_{1,2}$	$w_{1,3}$ $+d_{1,3}$
$w_{2,1}$ $+d_{2,1}$	$w_{2,2}$ $+d_{2,2}$	$w_{2,3}$ $+d_{2,3}$




R_{e-pos}

$w_{1,1}$ $+e_{1,1}$	$w_{1,2}$ $+e_{1,2}$	$w_{1,3}$ $+e_{1,3}$
$w_{2,1}$ $+e_{2,1}$	$w_{2,2}$ $+e_{2,2}$	$w_{2,3}$ $+e_{2,3}$




R_{f-pos}

$w_{1,1}$ $+f_{1,1}$	$w_{1,2}$ $+f_{1,2}$	$w_{1,3}$ $+f_{1,3}$
$w_{2,1}$ $+f_{2,1}$	$w_{2,2}$ $+f_{2,2}$	$w_{2,3}$ $+f_{2,3}$




R_{g-pos}

$w_{1,1}$ $+g_{1,1}$	$w_{1,2}$ $+g_{1,2}$	$w_{1,3}$ $+g_{1,3}$
$w_{2,1}$ $+g_{2,1}$	$w_{2,2}$ $+g_{2,2}$	$w_{2,3}$ $+g_{2,3}$




R_{d-neg}

$w_{1,1}$ $-d_{1,1}$	$w_{1,2}$ $-d_{1,2}$	$w_{1,3}$ $-d_{1,3}$
$w_{2,1}$ $-d_{2,1}$	$w_{2,2}$ $-d_{2,2}$	$w_{2,3}$ $-d_{2,3}$




R_{e-neg}

$w_{1,1}$ $-e_{1,1}$	$w_{1,2}$ $-e_{1,2}$	$w_{1,3}$ $-e_{1,3}$
$w_{2,1}$ $-e_{2,1}$	$w_{2,2}$ $-e_{2,2}$	$w_{2,3}$ $-e_{2,3}$




R_{f-neg}

$w_{1,1}$ $-f_{1,1}$	$w_{1,2}$ $-f_{1,2}$	$w_{1,3}$ $-f_{1,3}$
$w_{2,1}$ $-f_{2,1}$	$w_{2,2}$ $-f_{2,2}$	$w_{2,3}$ $-f_{2,3}$



R_{g-neg}

$w_{1,1}$ $-g_{1,1}$	$w_{1,2}$ $-g_{1,2}$	$w_{1,3}$ $-g_{1,3}$
$w_{2,1}$ $-g_{2,1}$	$w_{2,2}$ $-g_{2,2}$	$w_{2,3}$ $-g_{2,3}$

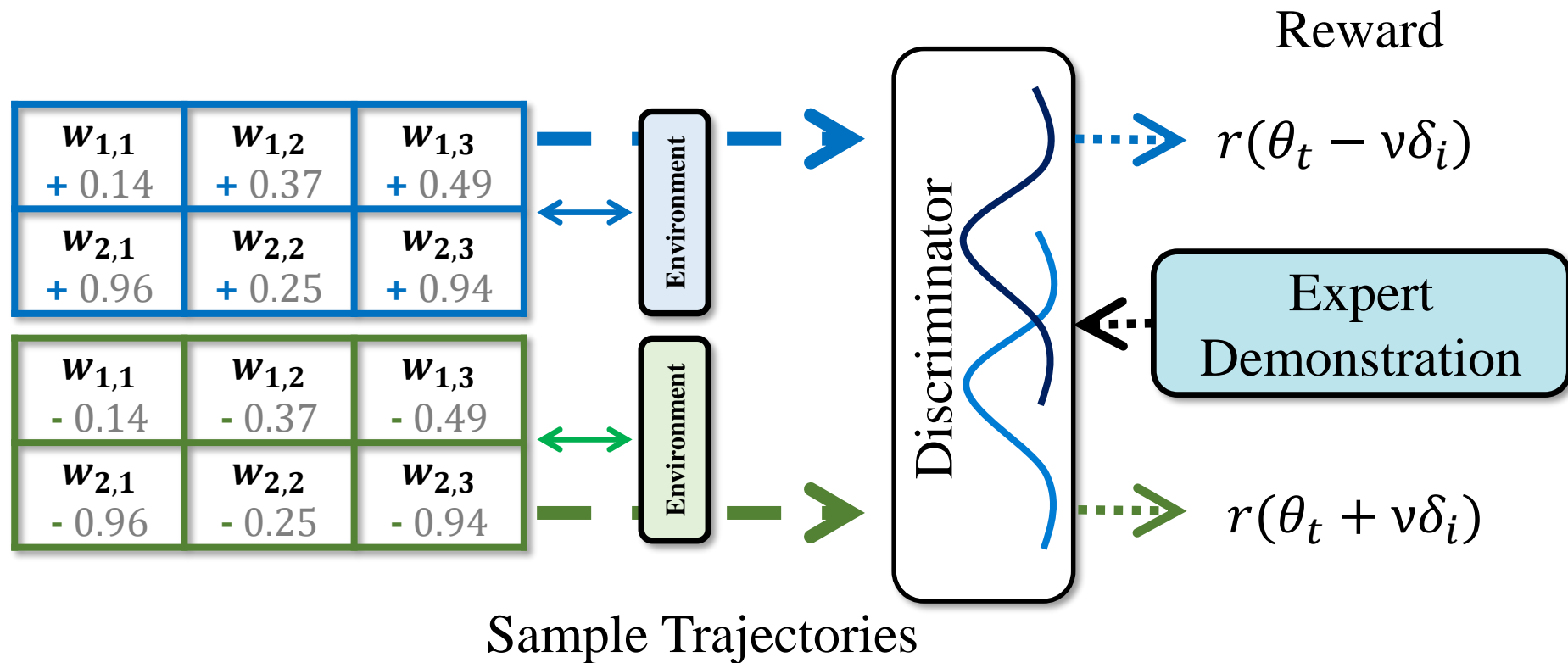


Adversarial Imitation Learning via Random Search

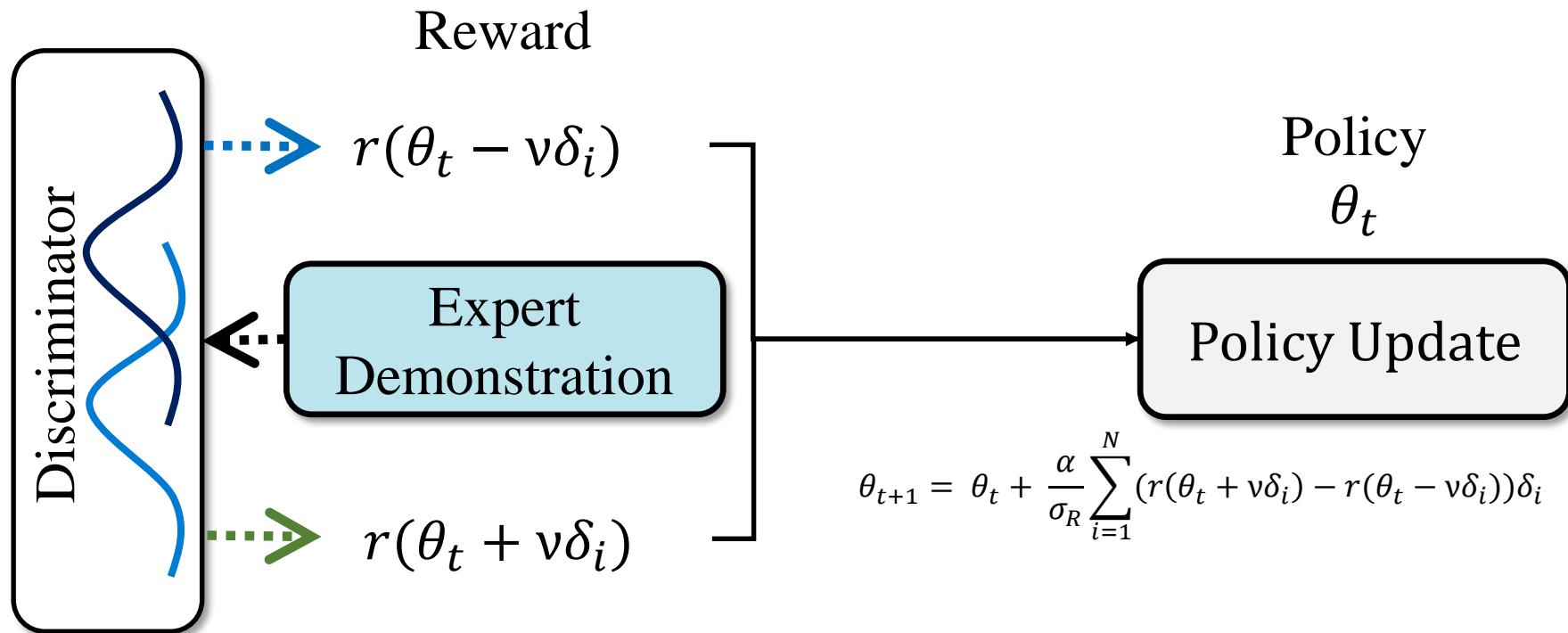


$$\begin{array}{c} \text{Policy} \\ \theta_{t+1} \\ \begin{array}{|c|c|c|} \hline w_{1,1} & w_{1,2} & w_{1,3} \\ \hline w_{2,1} & w_{2,2} & w_{2,3} \\ \hline \end{array} \end{array} = \begin{array}{c} \text{Policy} \\ \theta_t \\ \begin{array}{|c|c|c|} \hline w_{1,1} & w_{1,2} & w_{1,3} \\ \hline w_{2,1} & w_{2,2} & w_{2,3} \\ \hline \end{array} \end{array} + \left[\begin{array}{c} (R_{d-pos} - R_{d-neg}) * \begin{array}{|c|c|c|} \hline d_{1,1} & d_{1,2} & d_{1,3} \\ \hline d_{2,1} & d_{2,2} & d_{2,3} \\ \hline \end{array} + \\ (R_{e-pos} - R_{e-neg}) * \begin{array}{|c|c|c|} \hline e_{1,1} & e_{1,2} & e_{1,3} \\ \hline e_{2,1} & e_{2,2} & e_{2,3} \\ \hline \end{array} + \\ \vdots \end{array} \right]$$

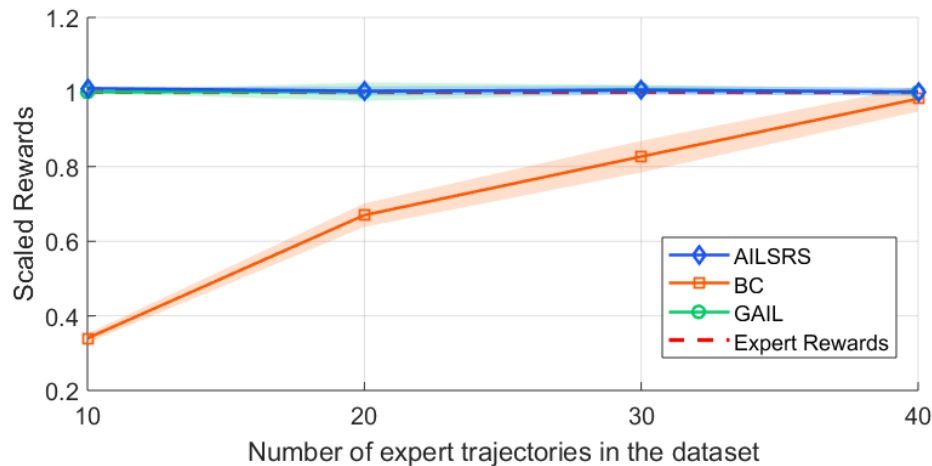
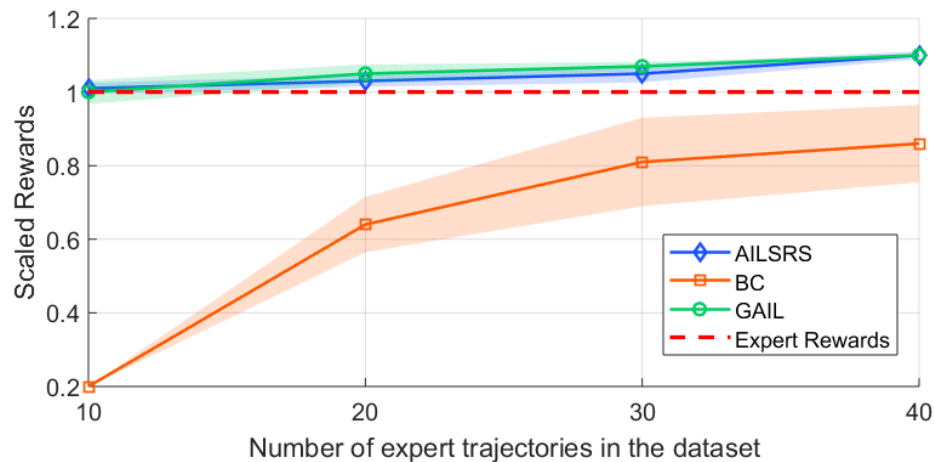
Adversarial Imitation Learning via Random Search



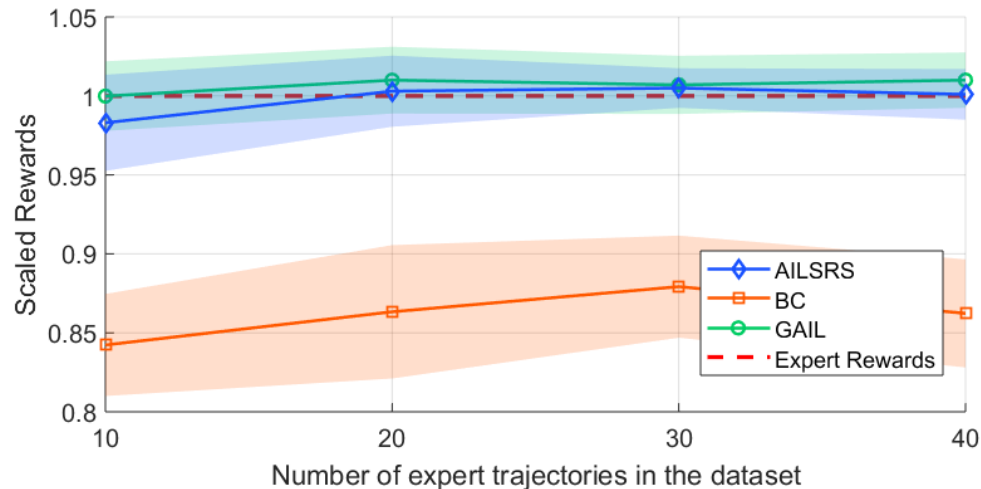
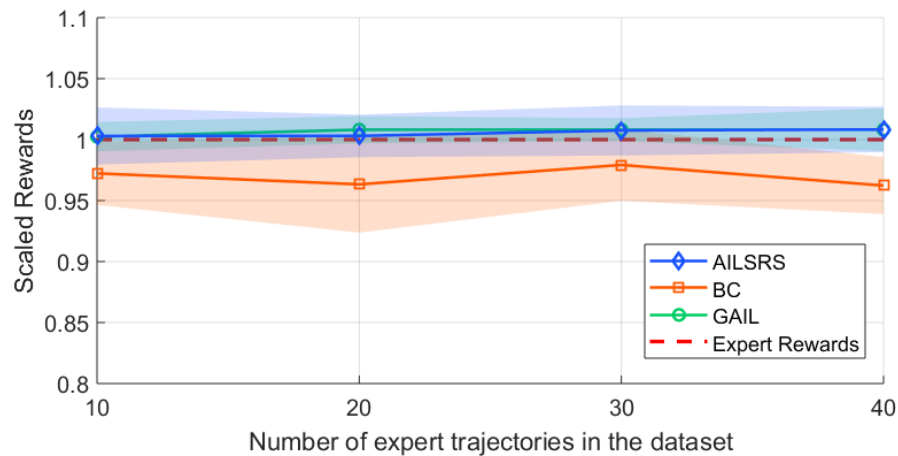
Adversarial Imitation Learning via Random Search



Adversarial Imitation Learning via Random Search



Adversarial Imitation Learning via Random Search



Thank You
