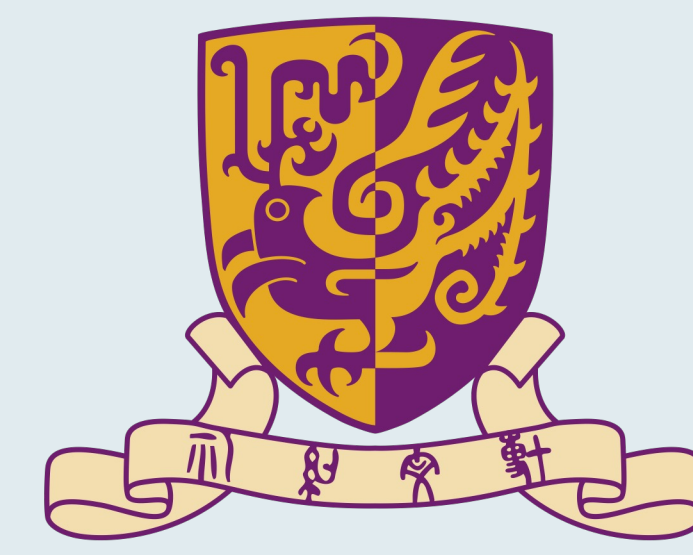


Understanding Constraint Inference in Safety-Critical Inverse Reinforcement Learning

Bo Yue¹, Shufan Wang², Ashish Gaurav^{3,4}, Jian Li², Pascal Poupart^{3,4}, Guiliang Liu^{1*}

¹School of Data Science, The Chinese University of Hong Kong, Shenzhen,

²Stony Brook University, ³University of Waterloo, ⁴Vector Institute



香港中文大學(深圳)
The Chinese University of Hong Kong, Shenzhen



Stony Brook
University

UNIVERSITY OF
WATERLOO



VECTOR
INSTITUTE | INSTITUT
VECTEUR

Abstract

Background: Constraint inference is crucial in safety-critical decision-making processes.

Literature: Existing methods, *Inverse Constrained Reinforcement Learning (ICRL)*, characterizes constraint learning as a inherently **complex** tri-level optimization problem.

Challenges: *Can we implicitly embed constraint signals into reward functions and effectively solve this problem using a classic reward inference algorithm?*

Methodology: *Inverse Reward Correction (IRC)* VS. *ICRL*

- IRC infers **a reward correction term**, which, when added to the reward function, ensures the optimality of the expert.
- ICRL infers **a cost function**, which, when serving as a constraint condition, ensures the optimality of the expert.

Takeaways:

- Training Efficiency:* **IRC** > **ICRL** (**IRC** learns constraint knowledge **faster**!)
- Cross-Environment Transferability:* **IRC** < **ICRL** (**IRC** fail to guarantee safety in target envs!)

Methods

Inverse Constraint Inference: Infer **constraint knowledge** followed by **expert policy**

IRC solver: $\pi^E = \max_{\pi} \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t [r + \Delta r](s_t, a_t)]$

(i) if $\pi^E(a|s) > 0$, then $Q_{\mathcal{M}U(r+\Delta r)}^{r+\Delta r, \pi^E}(s, a) = V_{\mathcal{M}U(r+\Delta r)}^{r+\Delta r, \pi^E}(s)$,

(ii) if $\pi^E(a|s) = 0$, then $Q_{\mathcal{M}U(r+\Delta r)}^{r+\Delta r, \pi^E}(s, a) \leq V_{\mathcal{M}U(r+\Delta r)}^{r+\Delta r, \pi^E}(s)$.

IRC should be **optimal** regarding $r + \Delta r$

ICRL solver: $\pi^E = \max_{\pi} \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t [r - \lambda^* c](s_t, a_t)]$

(i) if $\pi^E(a|s) > 0$, $Q_{\mathcal{M}Uc}^{c, \pi^E}(s, a) - V_{\mathcal{M}Uc}^{c, \pi^E}(s) = 0$;

(ii) if $\pi^E(a|s) = 0$ and $A_{\mathcal{M}Uc}^{r, \pi^E}(s, a) > 0$, $Q_{\mathcal{M}Uc}^{c, \pi^E}(s, a) - V_{\mathcal{M}Uc}^{c, \pi^E}(s) > 0$;

(iii) if $\pi^E(a|s) = 0$ and $A_{\mathcal{M}Uc}^{r, \pi^E}(s, a) \leq 0$, $Q_{\mathcal{M}Uc}^{c, \pi^E}(s, a) - V_{\mathcal{M}Uc}^{c, \pi^E}(s) \leq 0$.

ICRL should be **optimal** regarding r under constraint condition $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)] \leq \epsilon$

Implicit
Model

Explicit
Model

Theoretical Findings

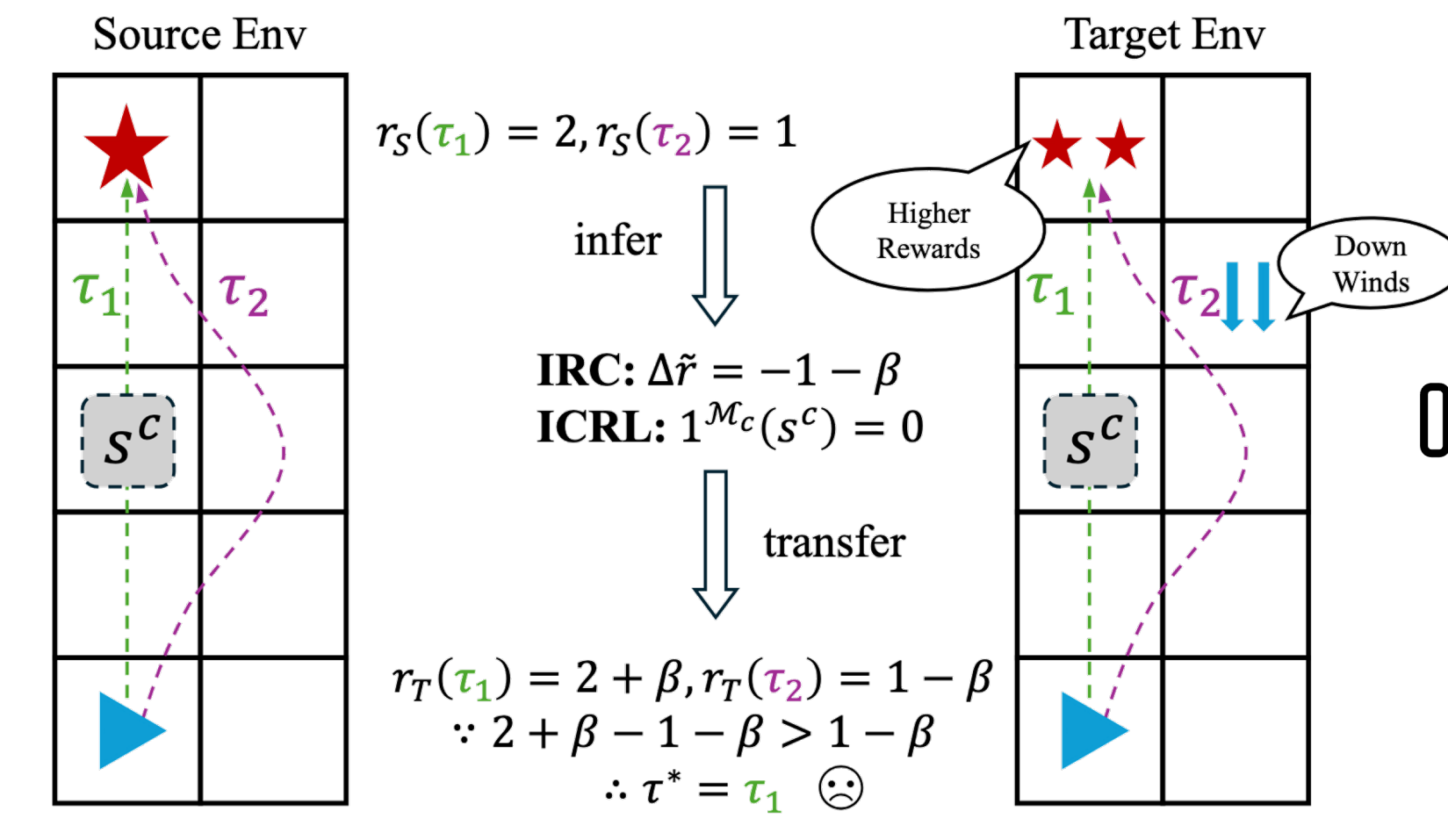
Sample Complexity

IRC solves $\max_{\Delta r} \min_{\pi} \mathcal{J}(\pi^E, r + \Delta r) - \mathcal{J}(\pi, r + \Delta r)$.

ICRL solves $\max_c \max_{\lambda} \min_{\pi} \mathcal{J}(\pi^E, r - \lambda c) - \mathcal{J}(\pi, r - \lambda c)$.

Bi-level to tri-level
A higher complexity of
 $1/(1-\gamma)^2$

Safety



Safety	IRC	ICRL
Hard Constraint	✗	✓
Soft Constraint	influenced by different rewards and transitions	influenced by different transitions

Source: $\mathcal{M}_c = \mathcal{M} \cup c^E = (\mathcal{S}, \mathcal{A}, P_{\mathcal{T}}, r, c^E)$ Target: $\mathcal{M}'_c = \mathcal{M}' \cup (c')^E = (\mathcal{S}, \mathcal{A}, P'_{\mathcal{T}}, r', (c')^E)$

$$Q^{r'}(s, a) - Q^r(s, a) = \left[\underbrace{(Y')^{-1}(r' - r)}_{\text{reward transfer shift}} + \underbrace{(Y')^{-1}AQ^r}_{\text{transition transfer shift}} + \underbrace{(Y')^{-1}BQ^r}_{\text{expert policy transfer shift}} \right](s, a).$$

$$Q^{c'}(s, a) - Q^c(s, a) = \left[\underbrace{(Y')^{-1}AQ^c}_{\text{transition transfer shift}} + \underbrace{(Y')^{-1}BQ^c}_{\text{expert policy transfer shift}} \right](s, a)$$

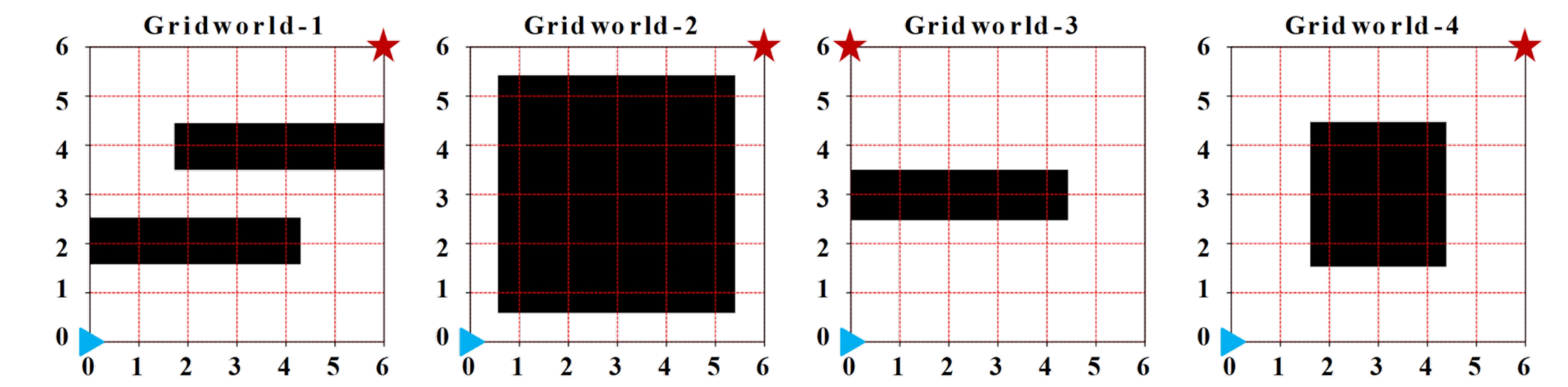
Optimality

$$\epsilon = 2 \max \left\{ d_1^2 \sin(\theta_{\max}(P'_{\mathcal{T}}, P_{\mathcal{T}}))^2 / 2, 2\epsilon_1 / \sigma_{\mathcal{R}} \right\} / \eta, \quad d_1 = \|[c^E - \hat{c}]_{\mathcal{U}_{P'_{\mathcal{T}}}}\|_2$$

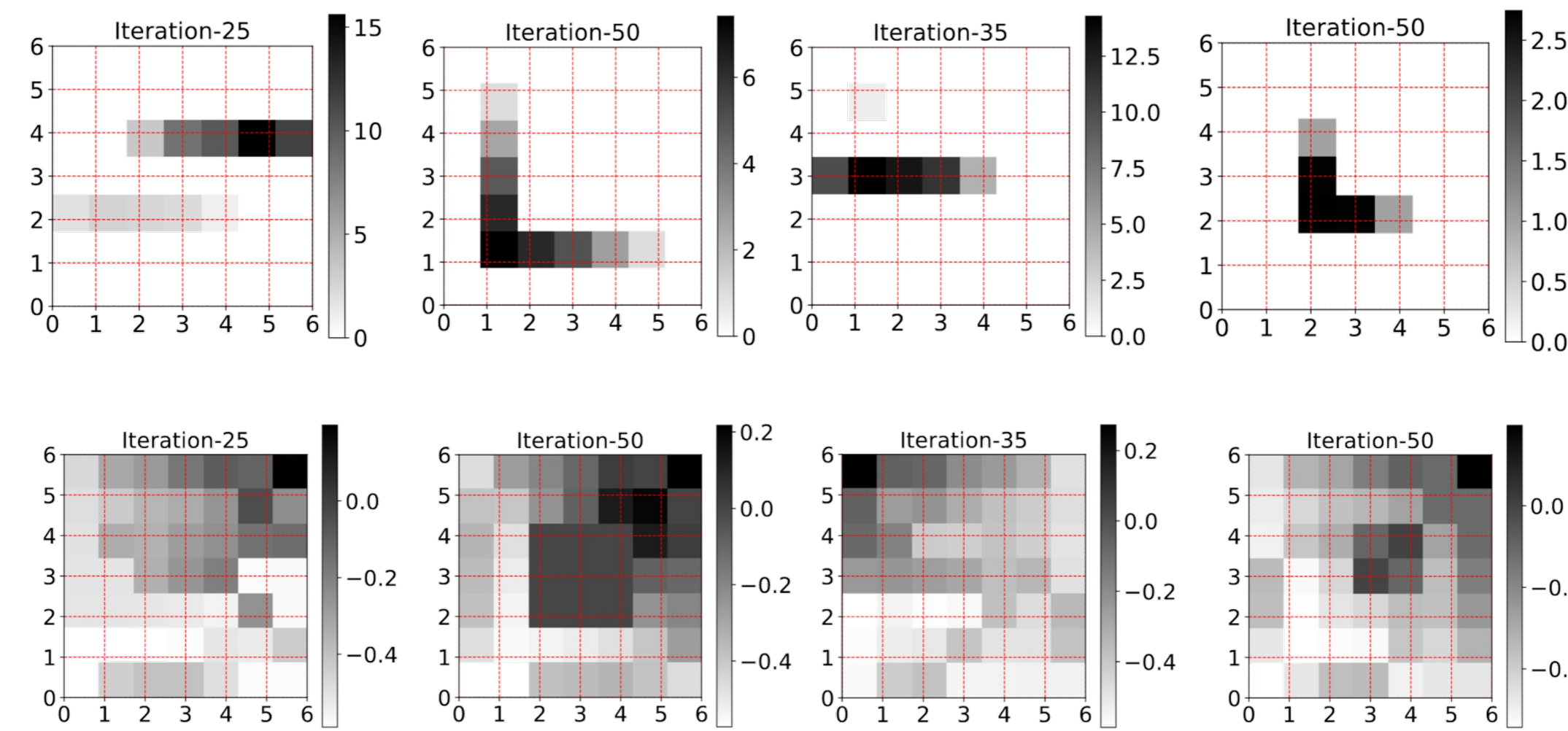
If the two transition laws are close and the recovered cost has a small suboptimality gap in the target environment, then ϵ -optimality of the recovered cost is guaranteed

Results

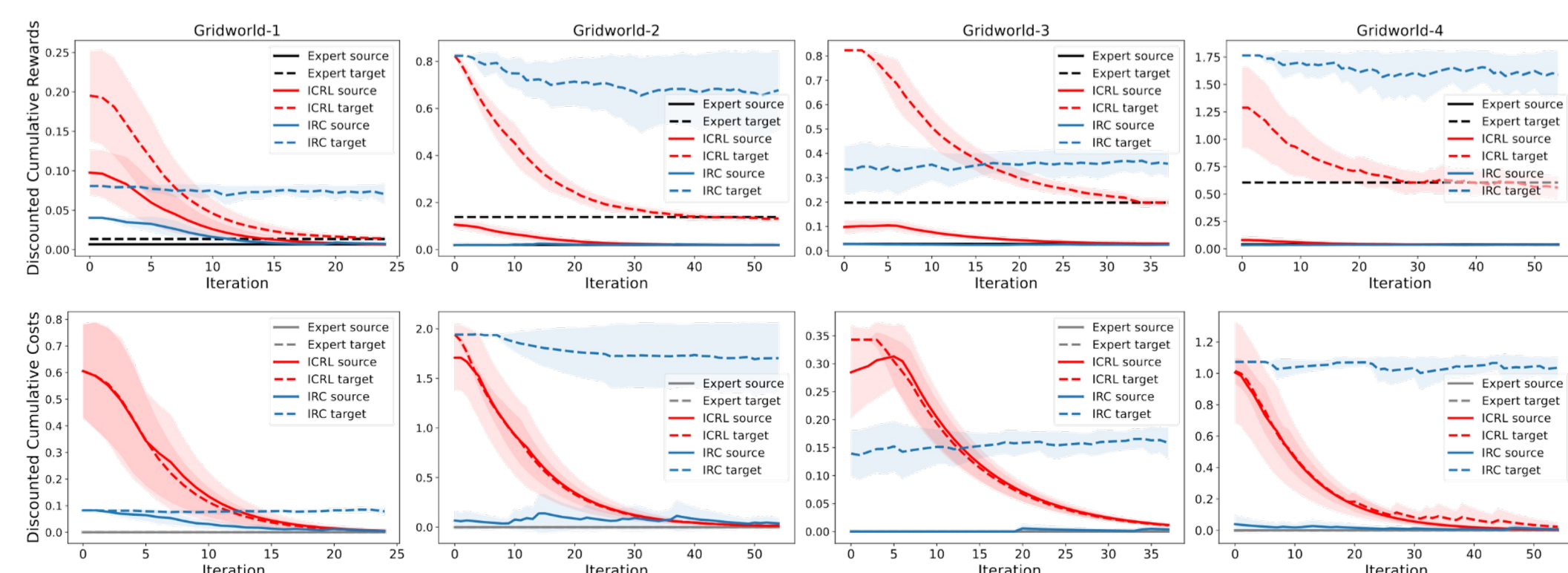
Discrete Envs



Constraint Knowledge (Up: ICRL; Bottom: IRC)



Learning Curves



More details

Group Info

