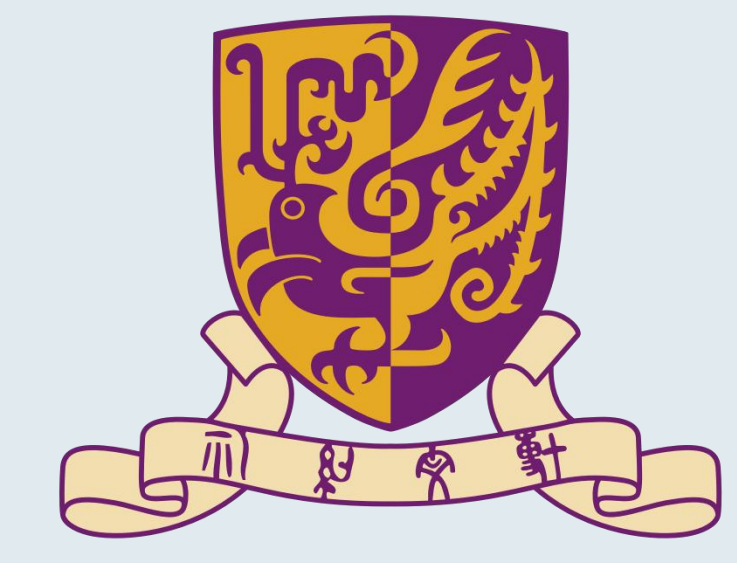# Provably Efficient Exploration in Inverse Constrained Reinforcement Learning

Bo Yue[1], Jian Li[2], Guiliang Liu[1]*
[1]School of Data Science, The Chinese University of Hong Kong, Shenzhen
[2]Stony Brook University

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

Stony Brook University

## Abstract

**Background:** Optimizing objective functions subject to *constraints* is fundamental in many real-world applications. However, ground-truth constraints are often *hard to specify, timevarying and context-dependent*.

**Literature:** *Inverse Constrained Reinforcement Learning (ICRL)*, recovers feasible constraints via training samples collected from *strategic exploration* in interactive environments.

**Challenges:** *The efficacy and efficiency of current sampling strategies remain unclear.*

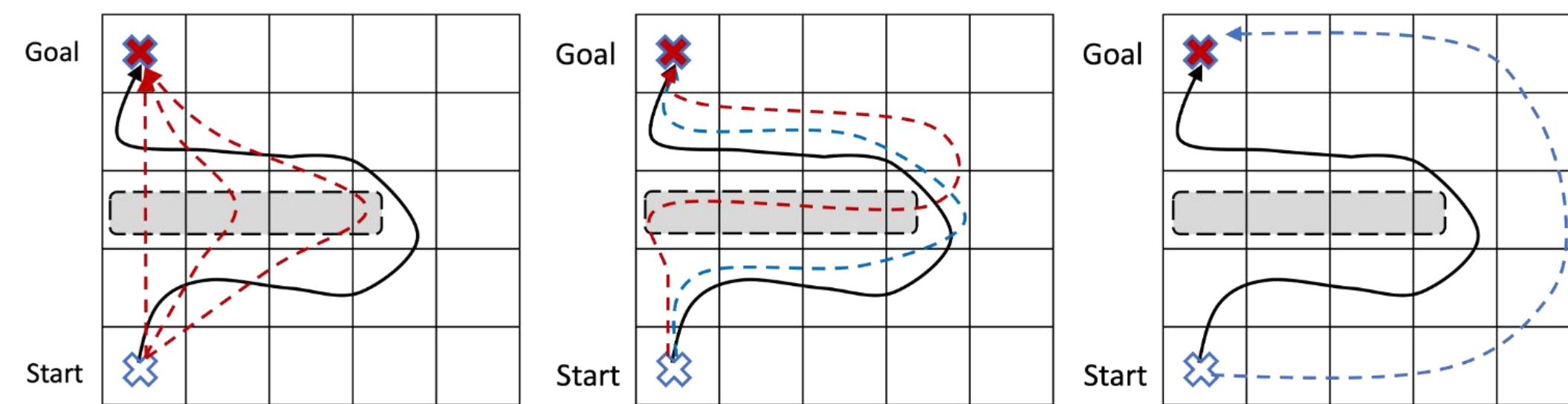**Methodology:** Two strategic exploration algorithms: BEAR and PCSE.

**Key Takeaways:**
- Both BEAR and PCSE are theoretically grounded with tractable sample complexity. (efficiency)
- PCSE outperforms six other baselines in recovering ground-truth constraints. (efficacy)

Sampling Trajectory: $\mathcal{D}_{\pi_{i+1}}$

Forward Step: Poliy Upate $(\pi_i \to \pi_{i+1})$

Constraint Function: $f_c(\phi_{i+1})$

*I* Iterations

Inverse Step: Constraint Upate $(\phi_i \to \phi_{i+1})$

Expert Trajectory: $\mathcal{D}_e$

## Recovery of Feasible Constraints

Since the expert policy maximizes rewards within a limited budget, two key insights emerge
1) if a policy achieves underline{higher rewards} than the expert policy, **the underlying constraints must be violated**, unsafe state-action pairs by examining these infeasible trajectories (left figure ⬇)
2) if a policy achieves underline{the same or lower rewards} than the expert policy, this suggests an absence of notable constraint-violating actions, implying that **the underlying constraints may or may not be violated** (middle & right figure ⬇)



Trajectories of the expert policy (black) and exploratory policies (red and blue) in a Gridworld. The constraint (gray) is not observable.

Feasible cost function

$c = A_{\mathcal{M}}^{r,\pi^E}\zeta + (E - \gamma P_{\mathcal{T}})V^c$

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

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

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

## Sample Complexity Analysis

### Estimation of Expert Policy and Transitions

Use visitation counts: $\widehat{P}_{\mathcal{T}k}(s'|s,a) = \dfrac{N_k(s,a,s')}{N_k^+(s,a)}, \quad \widehat{\pi}_k^E(a|s) = \dfrac{N_k^E(s,a)}{N_k^{E+}(s)}$

### Error Propagation to Constraint Estimation

$|c - \widehat{c}|(s,a) \le \dfrac{2(\chi(s,a) + \chi)}{1 + (\chi(s,a) + \chi)/C_{\max}}, \quad \chi = \max \chi(s,a)$

$\chi(s,a) = \gamma \left\| (P_{\mathcal{T}} \underset{①}{-} \widehat{P}_{\mathcal{T}}) V^c \right\|(s,a) + \left| A_{\mathcal{M}}^{r,\pi^E} - A_{\widehat{\mathcal{M}}}^{r,\widehat{\pi}^E} \right| \zeta(s,a)$

$\left| A_{\mathcal{M}}^{r,\pi} - A_{\widehat{\mathcal{M}}}^{r,\widehat{\pi}} \right| \le \dfrac{2\gamma}{1-\gamma} \left\| (\widehat{P}_{\mathcal{T}} \underset{②}{-} P_{\mathcal{T}}) V_{\widehat{\mathcal{M}}}^{r,\widehat{\pi}^E} \right\| + \dfrac{\gamma(1+\gamma)}{1-\gamma} \left| (\pi \underset{③}{-} \widehat{\pi}) P_{\mathcal{T}} V_{\mathcal{M}}^{r,\pi^E} \right|.$

### Strategic Exploration (BEAR and PCSE)

**BEAR**: guides the exploration policy to visit (s,a) to **minimize the upper bound of constraint estimation errors;**
**PCSE**: further **restricts exploration over plausibly optimal policies**

Sample complexity of BEAR:

$n \le \widetilde{\mathcal{O}}\left( \dfrac{\check{\sigma}^2(2C_{\max} - \varepsilon_K(1-\gamma))^2}{(1-\gamma)^2 \varepsilon_K^2 C_{\max}^2} \right)$

Sample complexity of PCSE:

$n \le \widetilde{\mathcal{O}}\left( \min\left\{ \widetilde{\mathcal{O}}\left( \dfrac{\check{\sigma}^2(2C_{\max} - \varepsilon_K(1-\gamma))^2}{(1-\gamma)^2\varepsilon_K^2 C_{\max}^2} \right), \dfrac{\sigma^2(6\varepsilon_{K-1}+\epsilon)^2 SA}{\min_{(s,a)}\left( A_{\widehat{\mathcal{M}}\cup\widehat{c}}^{c,*}(s,a) \right)^2 \varepsilon_K^2} \right\} \right)$

**Algorithm 1** BEAR and PCSE for ICRL in an unknown environment
**Input:** significance $\delta \in (0,1)$, target accuracy $\varepsilon$, maximum number of samples per iteration $n_{\max}$;
Initialize $k \leftarrow 0$, $\varepsilon_0 = \frac{1}{1-\gamma}$;
**while** $\varepsilon_k > \varepsilon$ **do**
  Solve RL problem defined by $\mathcal{M}^{c_k}$ to obtain the exploration policy $\pi_k$;
  Solve optimization problem in (15) to obtain the exploration policy $\pi_k$;
  Explore with $\pi_k$ for $n_e$ episodes;
  For each episode, collect $n_{\max}$ samples from $\mathcal{S} \times \mathcal{A}$;
  Update accuracy $\varepsilon_{k+1} = $
    $\max_{(s,a)\in\mathcal{S}\times\mathcal{A}} \mathcal{C}_{k+1}(s,a)/(1-\gamma)$;
  Update accuracy $\varepsilon_{k+1} = $
    $\|\mu_0^\top(I_{\mathcal{S}\times\mathcal{A}} - \gamma P_{\mathcal{T}}\pi)^{-1}\mathcal{C}_k\|_\infty$;
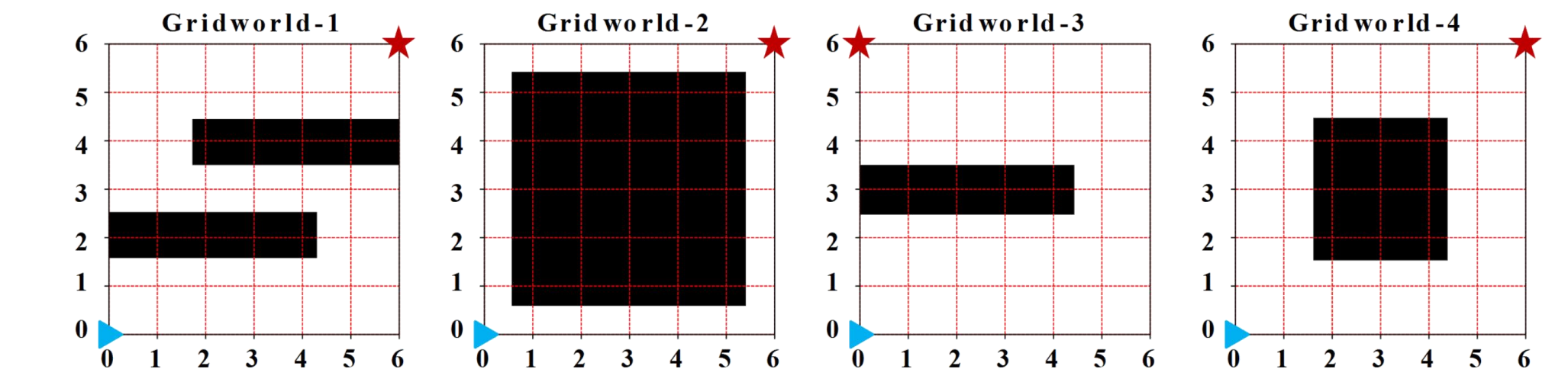  Update $\widehat{\pi}_{k+1}^E$ and $\widehat{P}_{\mathcal{T}k+1}$ in (7);
  $k \leftarrow k+1$.
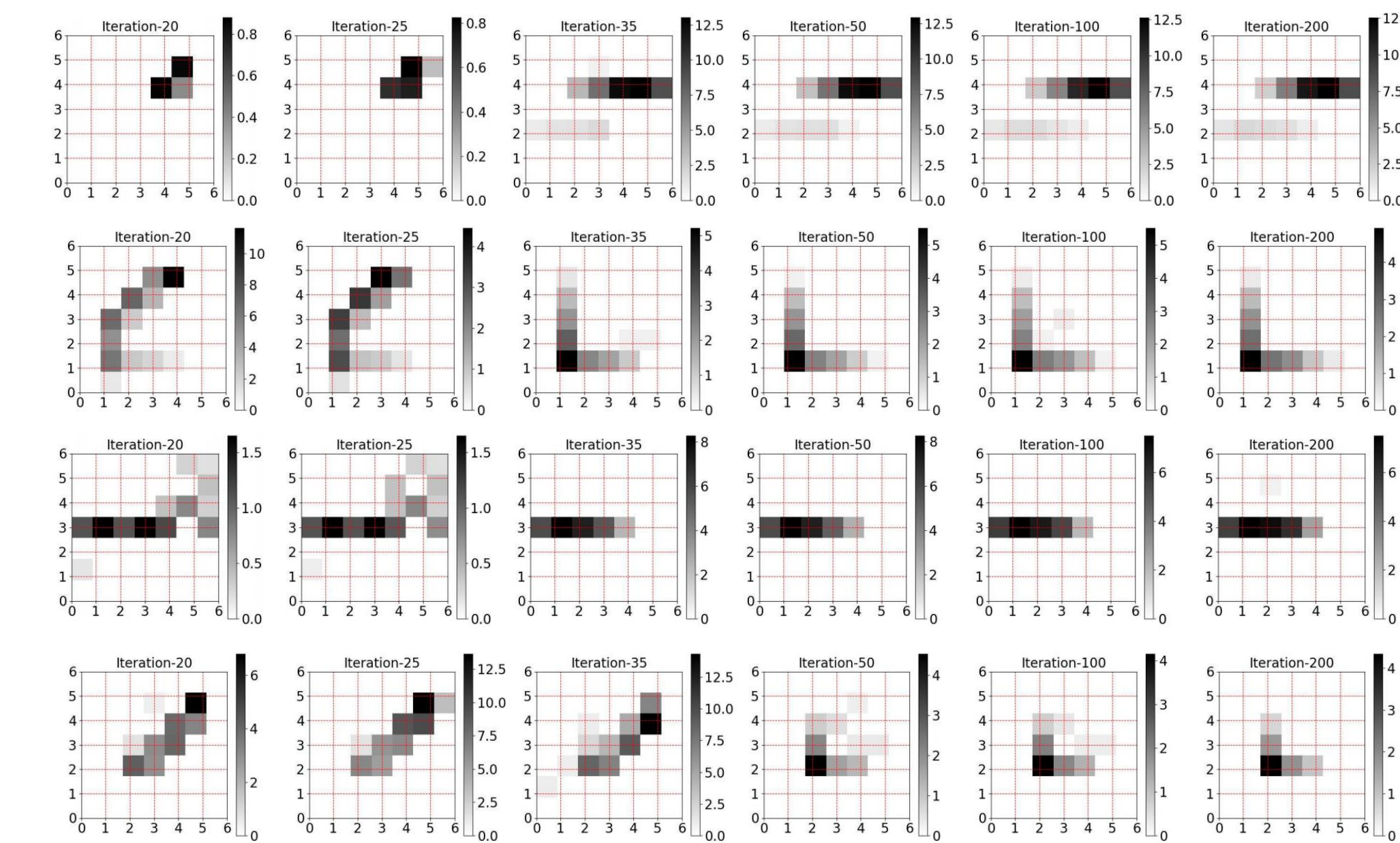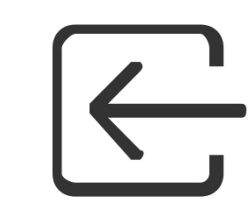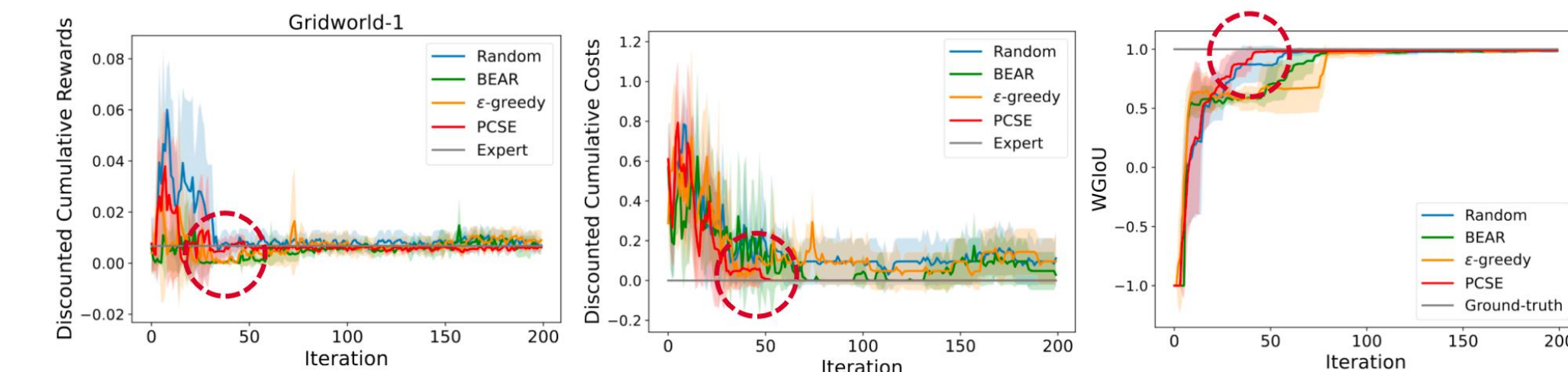**end while**

## Empirical Results

### Gridworld Envs



### Constraint Recovery Visualization (PCSE)



### Learning Curves



More details

Group Info