

# **Human-Interactive Subgoal Supervision for Efficient Inverse Reinforcement Learning**

**Xinlei Pan**

University of California, Berkeley  
Berkeley, California, USA  
[xinleipan@berkeley.edu](mailto:xinleipan@berkeley.edu)

**Yan Xu**

Carnegie Mellon University  
Pittsburgh, Pennsylvania, USA  
[yxu2@andrew.cmu.edu](mailto:yxu2@andrew.cmu.edu)

**Eshed Ohn-Bar**

Carnegie Mellon University  
Pittsburgh, Pennsylvania, USA  
[eshedohnbar@gmail.com](mailto:eshedohnbar@gmail.com)

**Yilin Shen**

Samsung Research America  
Mountain View, California, USA  
[yilin.shen@samsung.com](mailto:yilin.shen@samsung.com)

**Nicholas Rhinehart**

Carnegie Mellon University  
Pittsburgh, Pennsylvania, USA  
[nrhineha@cs.cmu.edu](mailto:nrhineha@cs.cmu.edu)

**Kris M. Kitani**

Carnegie Mellon University  
Pittsburgh, Pennsylvania, USA  
[kkitani@cs.cmu.edu](mailto:kkitani@cs.cmu.edu)



# Maximum Entropy IRL

$$\begin{aligned}\theta^* &= \arg \max_{\theta} - \sum_{d_i} p(d_i|\theta) \log(p(d_i|\theta)) \\ s.t. \quad f^\pi &= \tilde{f}^D \text{ (feature matching)}\end{aligned}$$

其中  $\tilde{f}^D = \frac{1}{N} \sum_{d_i \in D} \sum_{t=0}^k f_{it}$   
 $f^\pi = \sum_{d_i} p(d_i|\theta) f_{d_i}$

## Notation

$\theta$  : parameter of reward function

$\pi$  : optimal policy

$d_i$  : trajectory

$f_s$  : feature of state  $s$

$f_{d_i}$  : feature of trajectory

$f^\pi$  : feature expectation of trajectory generated by policy

$\tilde{f}^D$  : feature expectation of trajectory given by human

# IRL from Failure

$$\max_{\pi, w, z} H(D) + wz - \frac{\lambda}{2} \|w\|^2$$

$$s.t. \quad f^\pi = \tilde{f}^D$$

$$f^\pi - \tilde{f}^F = z$$

Optimization step

$$\theta_d = \theta_d - \alpha(f^\pi - \tilde{f}^D)$$

$$\theta_f = \frac{f^\pi - \tilde{f}^F}{\lambda}$$

其中  $H(D)$  causal entropy is defined as:

$$H(D) = - \sum_t \sum_{s_{1:t}, a_{1:t}} p(a_{1:t}, s_{1:t}) \log(p(a_t | s_t))$$

## Notation

$\tilde{f}^F$  : feature expectation of failure experience

$z$  : difference between feature expectation of failure experience and feature expectation following  $\pi$

$w$  : lagrange multiplier of  $z$

$\theta_d, \theta_f$ : parameter of reward function learning from demonstration and failure experience respectively.



# HI-IRL

- Interative setting
- Provide subgoal supervision-breaking task into subtasks.

## Step1

Human expert provides several full demonstrations and define subgoals

## Step2

Agents tries the defined subtasks.  $s_r \rightarrow s_{sub_1} \rightarrow \dots \rightarrow s_{sub_i} \rightarrow s_{goal}$ .

## Step3

Human provides further demonstrations if needed.

## Step4

Learning reward functions from both failure experiences and expert demonstrations

# Algorithm

---

**Algorithm 2** Human-Interactive Inverse Reinforcement Learning (HI-IRL)

---

**Require:** Set of initial demonstrations  $d_0$ ,  $T$ , State Transition Matrix  $\mathcal{T}$ ,  $\theta^0$ , all state raw feature  $f$ , and human  $\mathcal{H}$ .

**Return:** Reward function  $\theta^{T+1}$

**Define:**  $\mathcal{D}$ : positive demonstrations;  $\mathcal{F}$ : failure experience;  $\mathcal{E}$ : agent experience;  $\mathcal{S}_{sub}$ : set of subgoal states

**Start:**

$\mathcal{S}_{sub} = \text{specify\_subgoals}(\mathcal{H})$

$\mathcal{D} = d_0;$

$\theta^1 = \text{MaxEntIRL}(\mathcal{D}, \theta^0)$

**for**  $t \in 1, 2, \dots, T$

$\mathcal{E} = \text{ROLLOUT}(\theta^t, \mathcal{S}_{sub})$

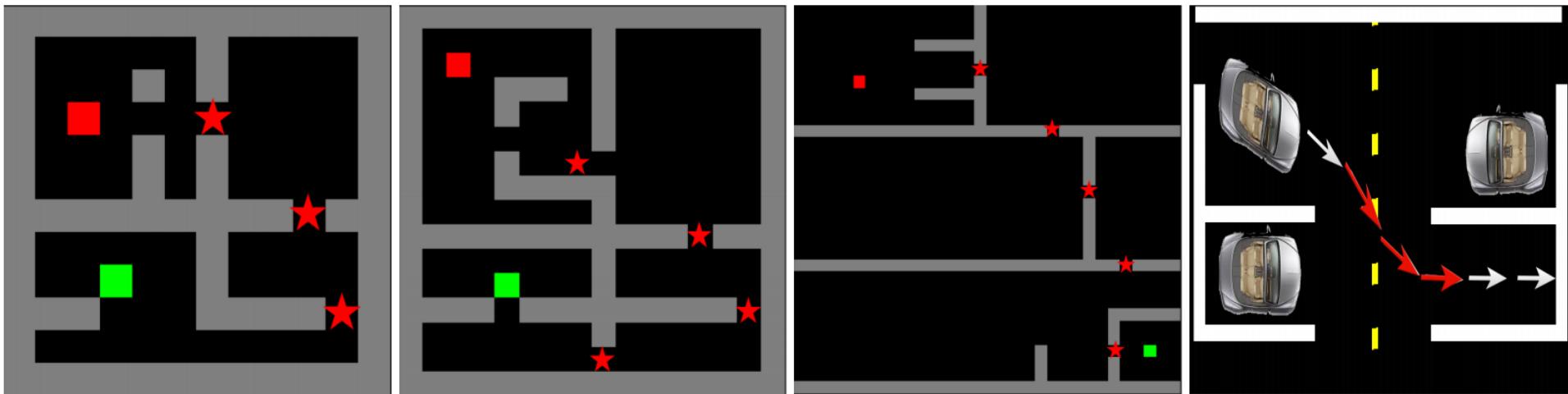
**for**  $e$  in  $\mathcal{E}$

$\mathcal{F}, \mathcal{D} = \text{UPDATEDEMO}(e, \theta^t, \mathcal{D})$

$\theta_d^{t+1}, \theta_f^{t+1} = \text{IRLFF}(\mathcal{F}, \mathcal{D}, \mathcal{T}, \theta_d^t, f)$  (Alg. 1)

$\theta^{t+1} = (\theta_d^{t+1}, \theta_f^{t+1})$

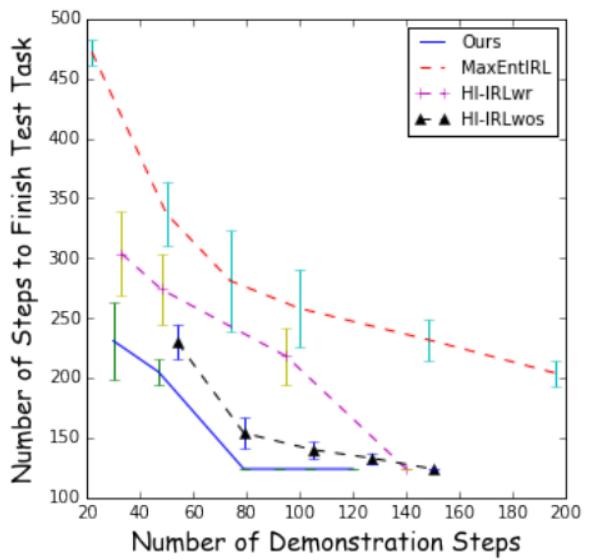
# Experiment



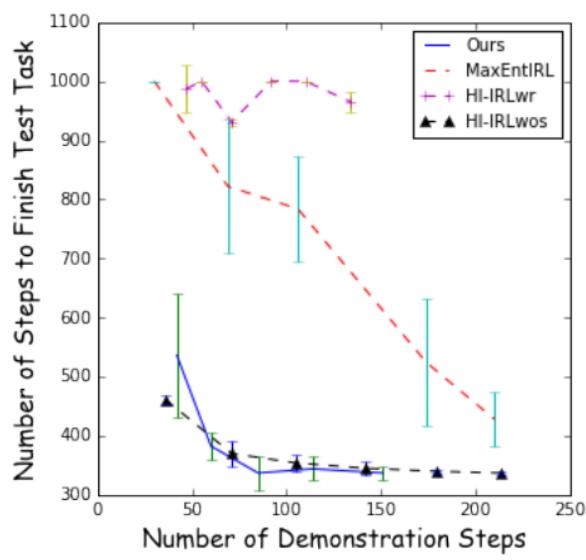
Baseline

1. MaxEntIRL
2. HI-IRLwos. The agent will be required to complete entire task and expert provide full demonstration.
3. HI-IRLwr. Randomly select states as subgoals

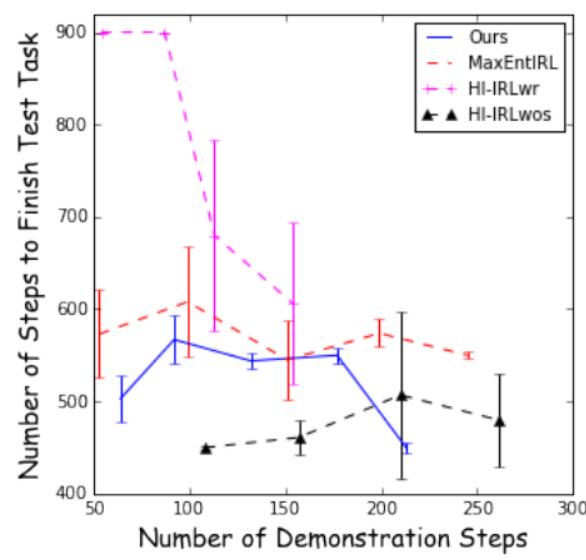
# Result



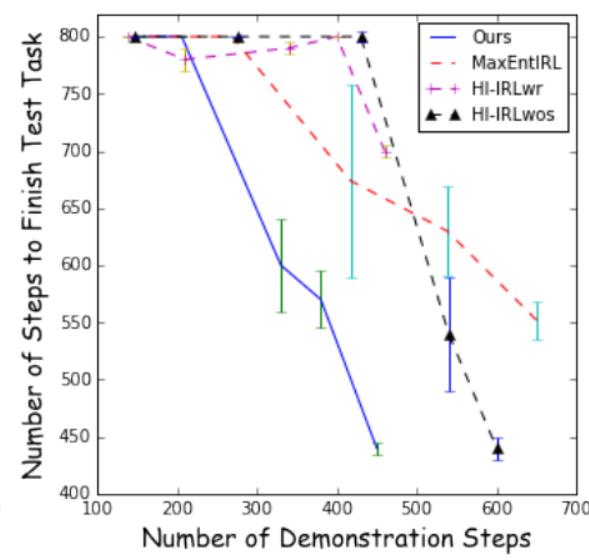
(a)



(b)



(c)



(d)