

Instructing **G**oal-**C**onditioned **R**einforcement **L**earning Agents with Temporal Logic Objectives

NeurIPS 2023

Rutgers University

Wenjie Qiu*

wq37@cs.rutgers.edu

Wensen Mao*

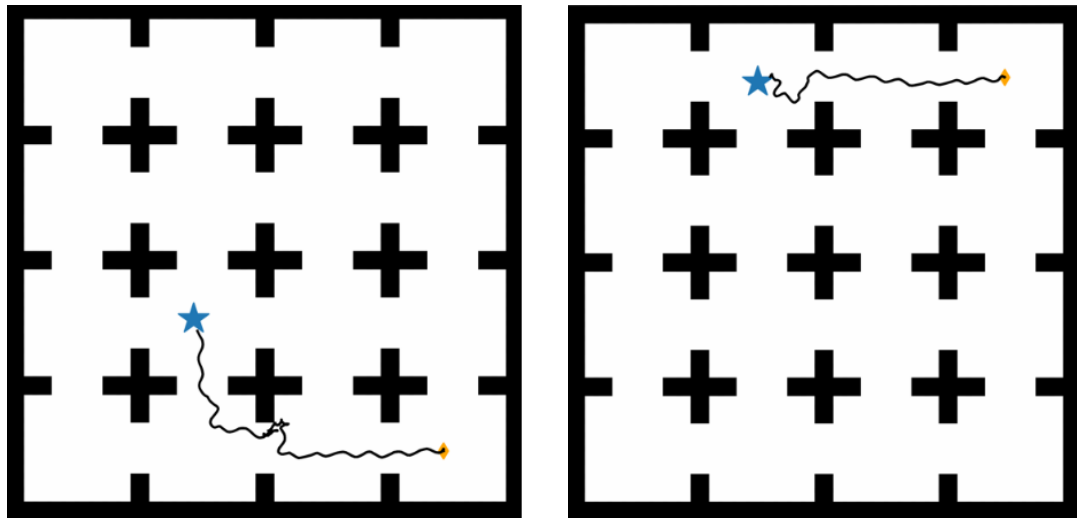
wm300@cs.rutgers.edu

He Zhu

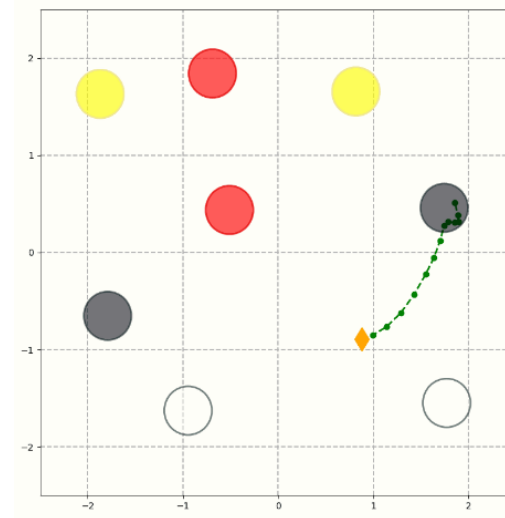
hz375@cs.rutgers.edu

Goal-Conditioned Reinforcement Learning (GCRL)

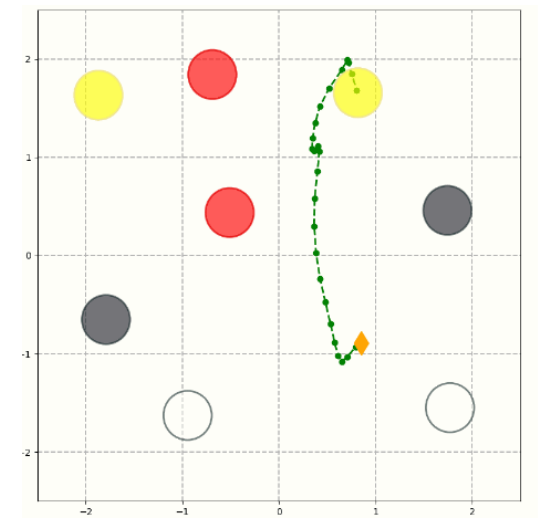
- Reach **arbitrary goals** in the goal space from initial environment states



Reach ★ with random start ♦



Go to JetBlack ● zone

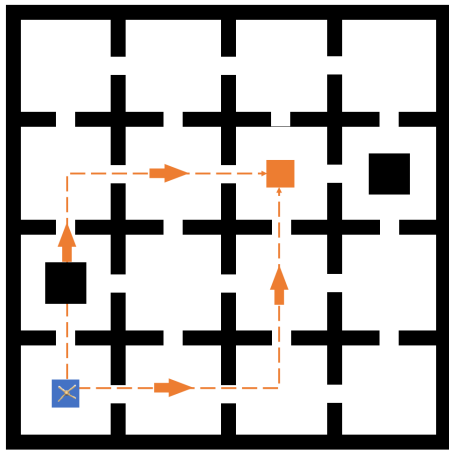


Go to Yellow ● zone

Linear Temporal Logic (LTL)

- Define a sequence of tasks in the chronological order
- May include **infinite loop** task (ω -regular property)

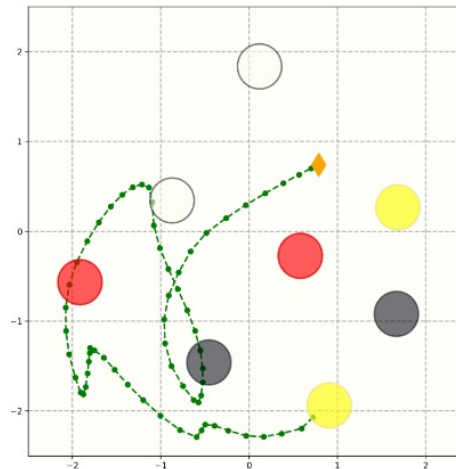
Ant Maze



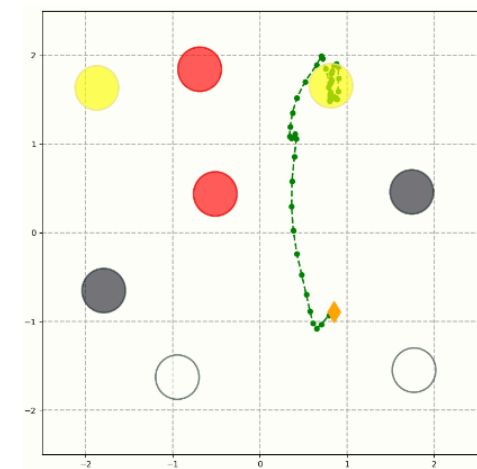
$$F \left((room(0,2) \vee room(2,0)) \wedge F room(2,2) \right)$$

Reach **orange room** via one possible **orange path**

ZoneEnv

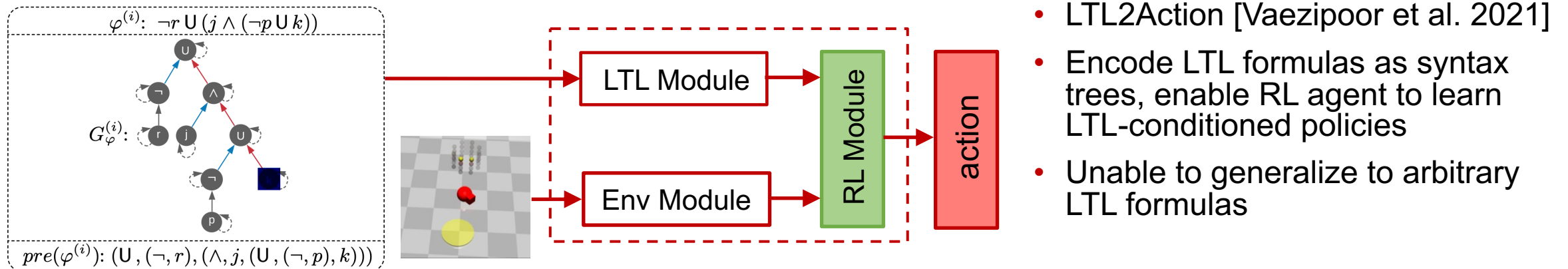
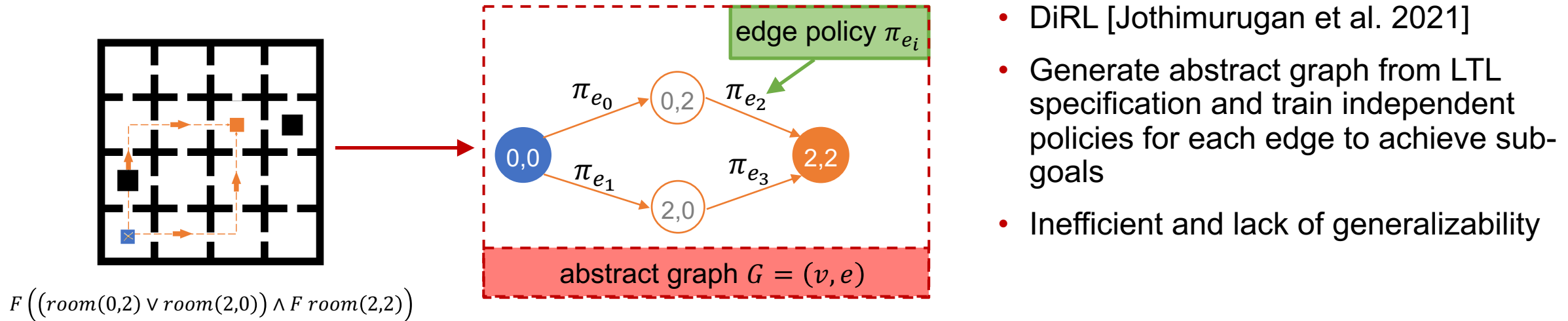


$$F(\text{grey} \wedge F(\text{white} \wedge F(\text{red} \wedge F \text{yellow})))$$

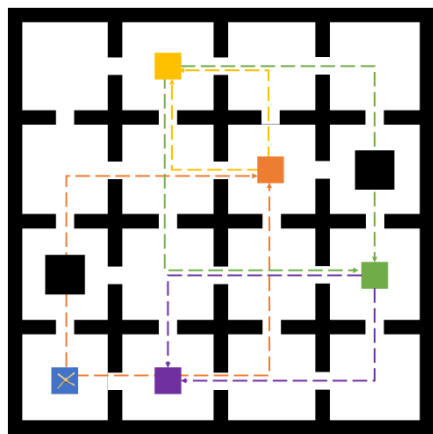


$$GF \text{ yellow}$$

Prior approaches

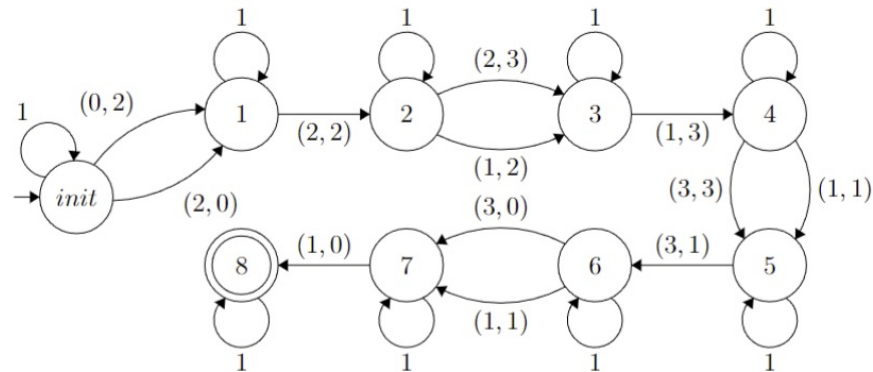


GCRL-LTL algorithm framework

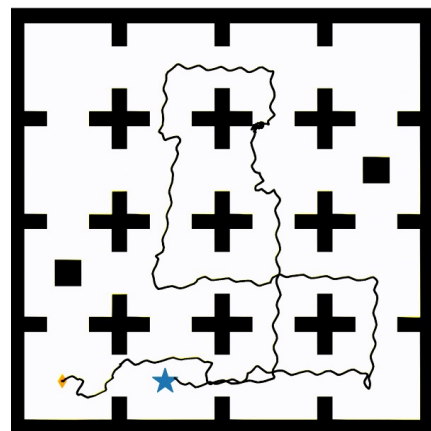


LTL specification ϕ

Buchi automata
abstract graph

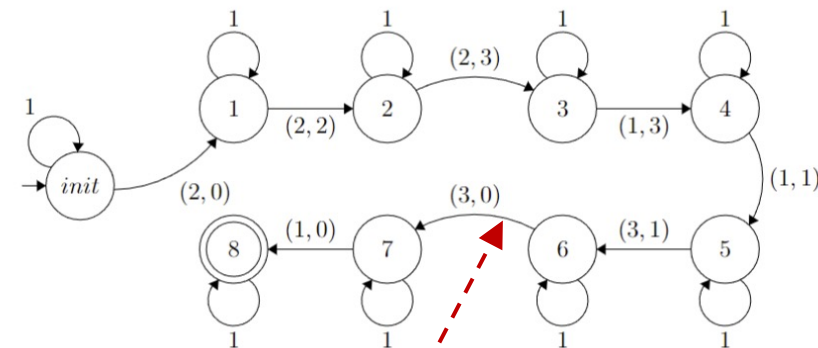


optimal path finding



execute with goal-
conditioned policy

$$\pi(a|s, \varphi)$$



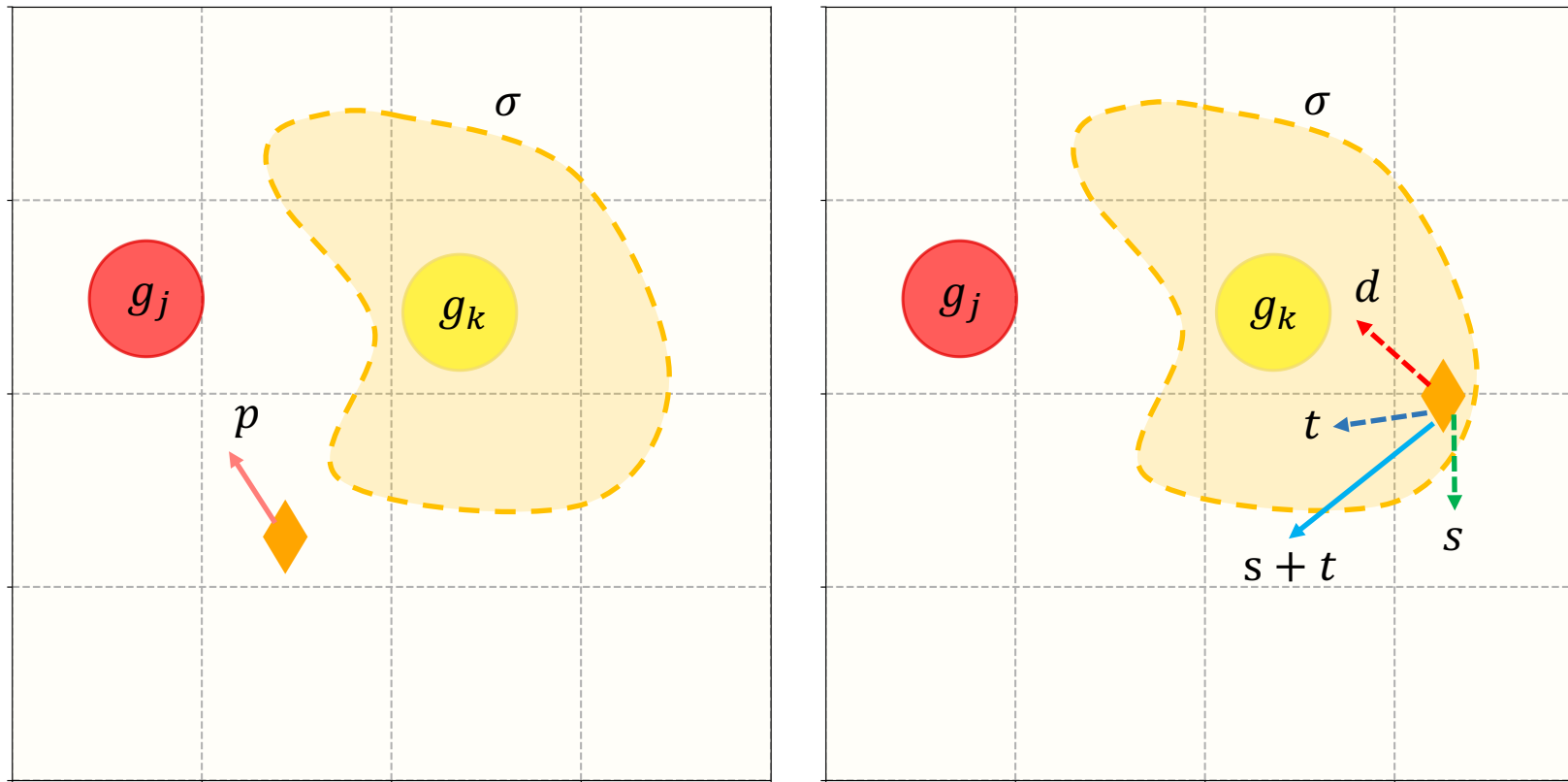
$$W_{\text{edge}} = -\log \mathcal{V}(s_0, \text{room}_{(3,1)}, \text{room}_{(3,0)})$$

$$\mathcal{V} = \min_{\mathcal{V}} (V(s_t, g_k) - \mathcal{V}(s_j, g_t, g_k))^2 \quad \text{where } g_k \in L(s_k) \wedge g_t \in L(s_t)$$

with $\tau \sim B$ (a replay buffer), $t \sim \{0 \dots t_{max}\}$, $s_t, a_t, s_{t+1} \sim \tau$, $j \sim \{0, t\}$, $k \sim \{t+1 \dots t_{max}\}$, and L as the state labeling function. We train \mathcal{V} together with a goal-conditioned learning algorithm.

Goal Reaching and Obstacle Avoidance

- For example, reach-avoid task on one edge : $\neg \bullet U \bullet$

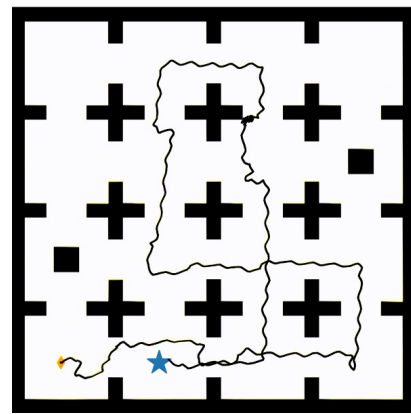
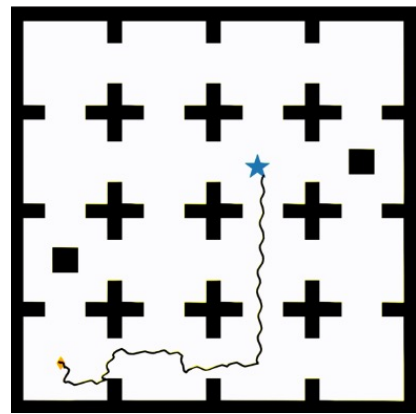
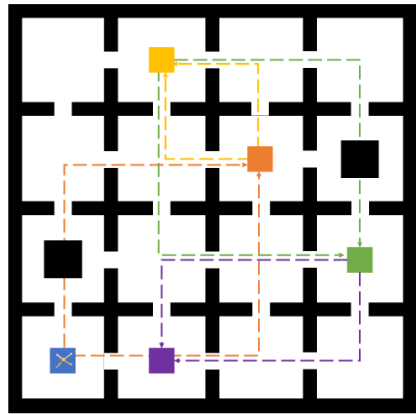


Goal ● Obstacle ●

$$\tilde{\pi}(\cdot | s, \bigwedge_j g_j, \bigwedge_k \neg g_k) \equiv$$

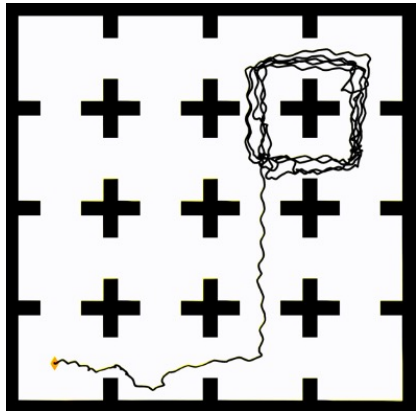
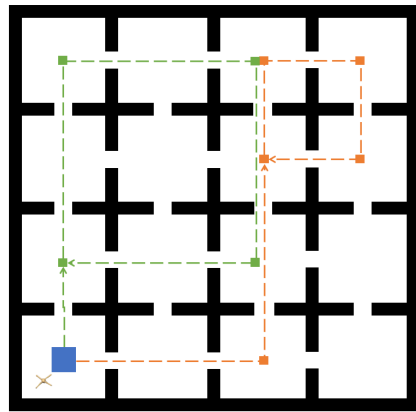
let $k = \arg \max_k V(s, g_k)$ **in**
if $V(s, g_k) < \sigma$ *reaching action* p
then $\arg \max_a \min_j Q^\pi(s, g_j, a)$
else
let $d = \arg \max_a Q(s, g_k, a)$
let $s = \arg \min_a Q(s, g_k, a)$
let $t = \arg \max_{a \neq d} \min_j Q^\pi(s, g_j, a)$
 $s + t$

Experiments on Ant16rooms



$$\phi_2 = \text{blue} \rightarrow \text{orange}$$

$$\phi_5 = \text{blue} \rightarrow \text{orange} \rightarrow \text{yellow} \rightarrow \text{green} \rightarrow \text{purple}$$



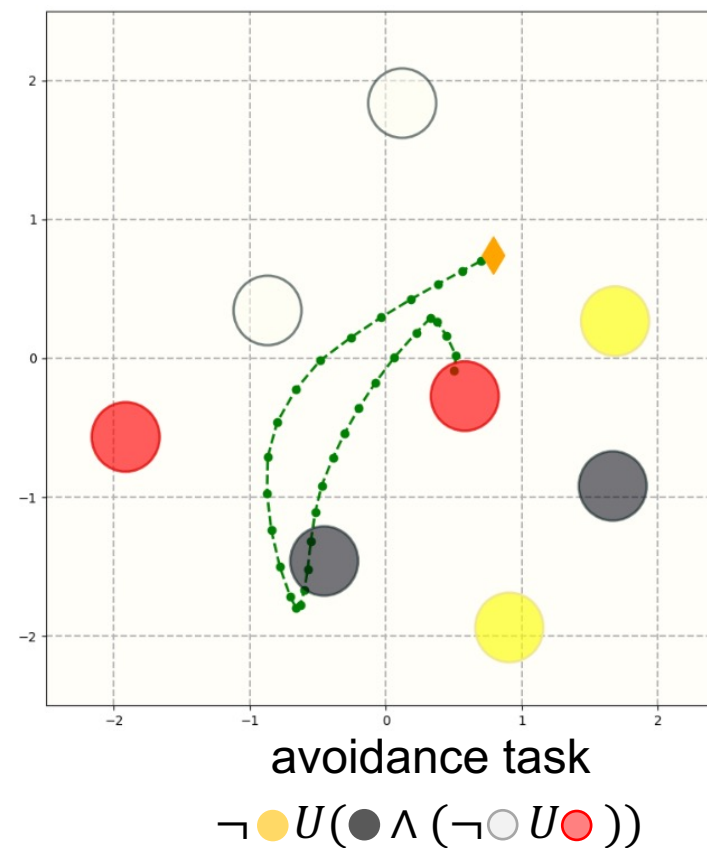
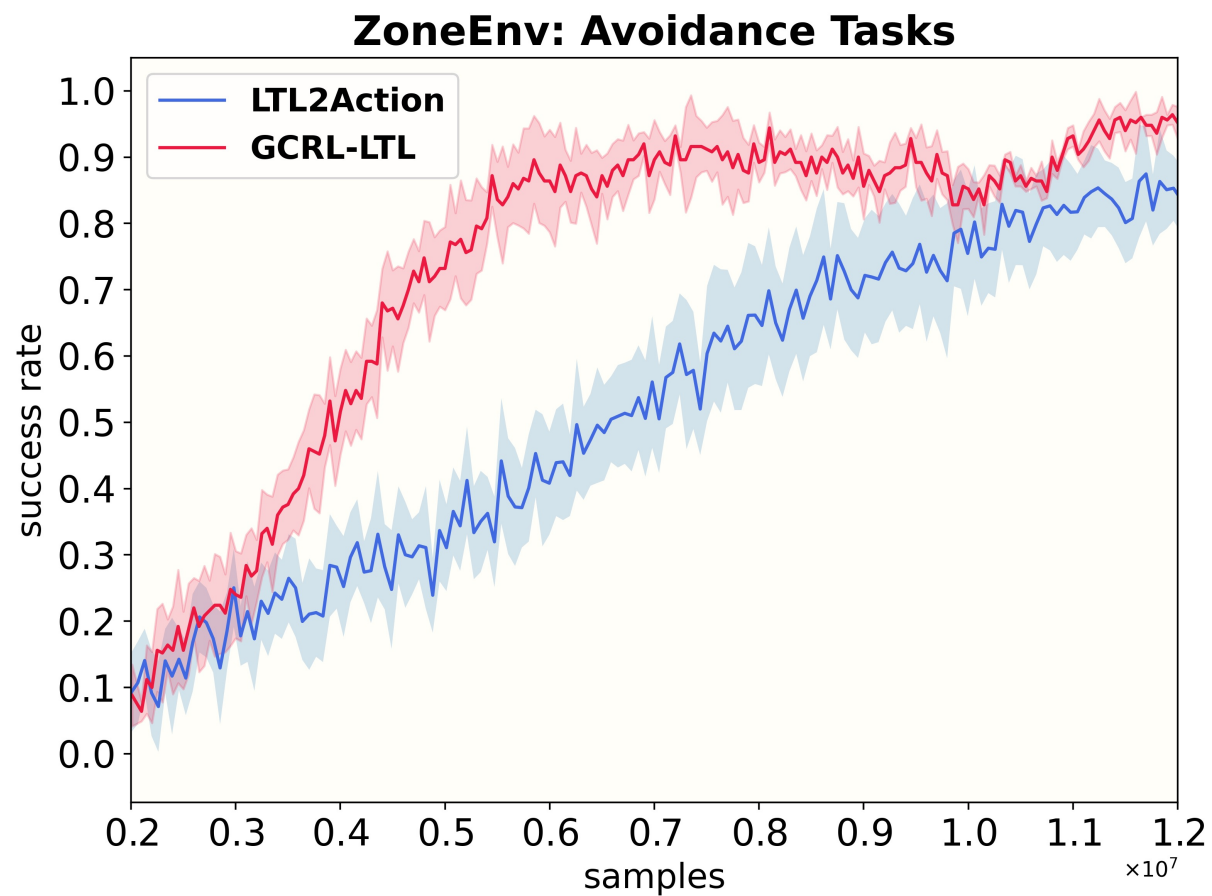
$$\phi_8 = \text{blue} \rightarrow (\text{orange})^\omega \vee \text{blue} \rightarrow (\text{green})^\omega$$

Choose to execute the loop with a smaller path cost!

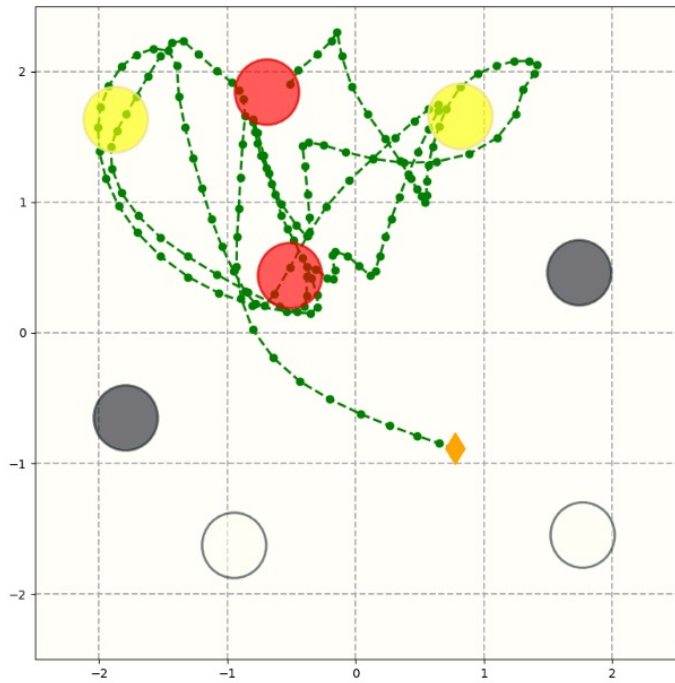
ω

Methods	DiRL	DiRL+ GCSL	Ours
ϕ_1	0.910 (0.022)	0.923 (0.082)	0.967 (0.006)
ϕ_2	0.770 (0.083)	0.953 (0.037)	0.925 (0.049)
ϕ_3	0.367 (0.147)	0.967 (0.017)	0.935 (0.028)
ϕ_4	0.183 (0.046)	0.937 (0.031)	0.875 (0.041)
ϕ_5	0.043 (0.061)	0.913 (0.017)	0.868 (0.038)
ϕ_6	/	/	0.857 (0.004)
ϕ_7	/	/	0.882 (0.018)
ϕ_8	/	/	0.903 (0.045)

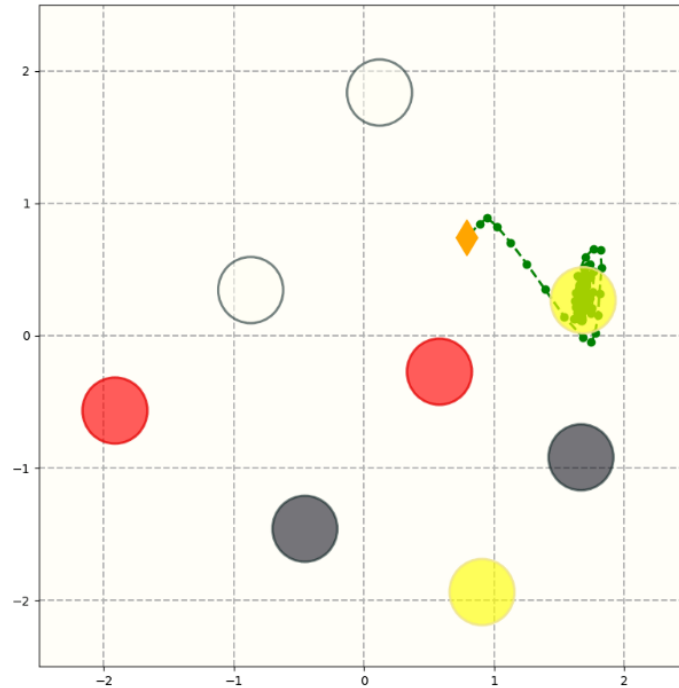
Experiments on ZoneEnv



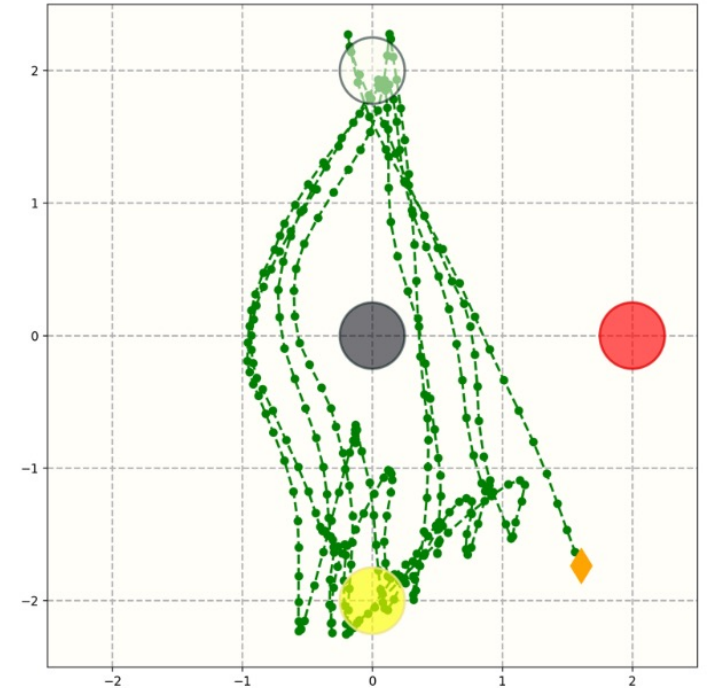
More complex experiments on ZoneEnv



$GF(\text{red} \wedge XF\text{yellow}) \wedge G(\neg \text{white})$



$FG\text{yellow}$



$GF\text{white} \wedge GF\text{yellow} \wedge G(\neg \text{grey})$