



spec-driven-rl.github.io

QR Code  
for  
Code  
Videos

# Mission-driven Exploration for Accelerated Deep Reinforcement Learning with Temporal Logic Task Specifications

Jun Wang<sup>1</sup>, Hosein Hasanbeig<sup>2</sup>, Kaiyuan Tan<sup>3</sup>, Zihe Sun<sup>1</sup>, and Yiannis Kantaros<sup>1</sup>  
Preston M. Green Department of Electrical & Systems Engineering, Washington University in St Louis<sup>1</sup>  
Microsoft Research<sup>2</sup>

Department of Computer Science, Vanderbilt University<sup>3</sup>

Keywords: Reinforcement Learning, Temporal Logic Planning, Sample Efficiency

Washington  
University in St. Louis

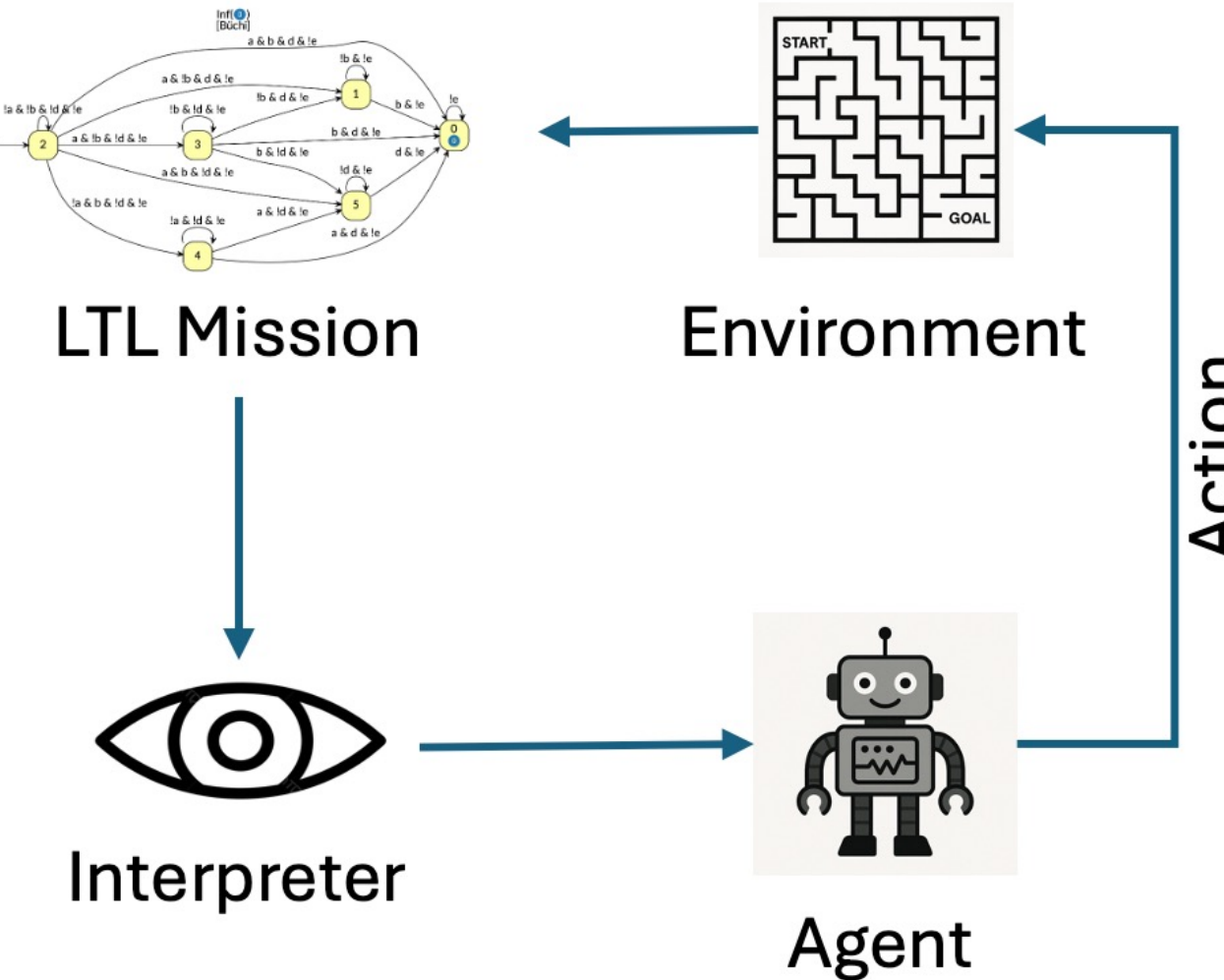
Microsoft

VANDERBILT  
UNIVERSITY

L4DC  
2025

## Motivation

- Linear Temporal Logic (LTL) has been used to encode complex tasks (e.g., navigation task with several ordered ROIs)



- Model-free** DRL methods require a product state space that grows exponentially, result in *slow learning process*
- Model-based** RL methods rely on learned MDPs but are limited to *discrete state spaces*

## Goal

Design a **sample-efficient** DRL algorithm to learn control policies for agents with LTL-encoded tasks

## Contribution

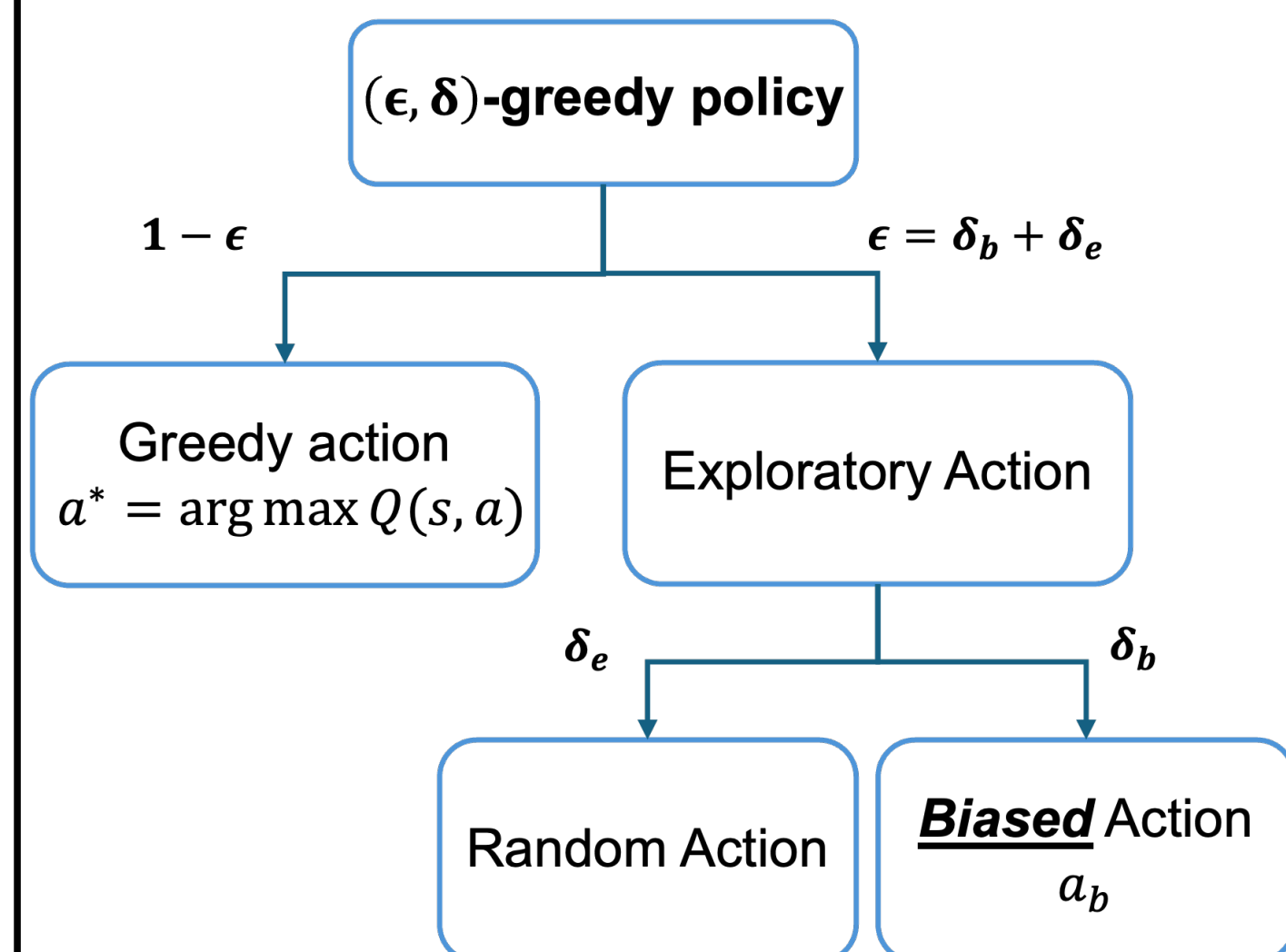
- Propose a new **DQN** algorithm for agents with *unknown MDPs*, *continuous state spaces* and LTL-encoded tasks
- Our policy is **complementary** with existing deep temporal-difference methods for LTL tasks to enhance sample efficiency
- Present comparative numerical and hardware experiments that demonstrate the **sample efficiency** of our method

## Problem Formulation

- Environment  $\mathcal{W} \subseteq \mathbb{R}^d, d \in \{2, 3\}$
- LTL over set of  $\mathcal{AP}$ :  $\phi = \text{true} \mid \pi \mid \phi_1 \wedge \phi_2 \mid \neg \phi \mid \phi \mid \phi_1 \cup \phi_2$
- Fully observable MDP:  $\mathfrak{M} = (\mathcal{X}, \mathcal{A}, P, \mathcal{AP})$  with continuous state space  $x \in \mathcal{X}$  and a finite set of actions  $a \in \mathcal{A}$  (*transition probabilities unknown*)
- Problem 1:** Given a known LTL-encoded task specification  $\phi$ , develop a **sample-efficient** DRL method that can synthesize a finite memory control policy  $\xi^*$  for the unknown MDP that **maximizes satisfaction probability** of  $\phi$

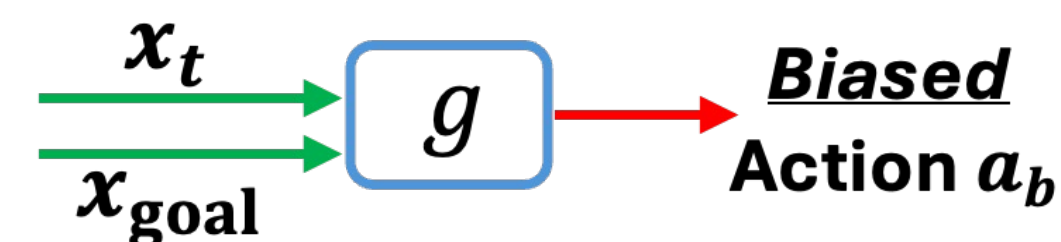
## Methods

- Translate  $\phi$  into DRA  $\mathfrak{D} = (\mathcal{Q}_{\mathfrak{D}}, q_{\mathfrak{D}}^0, \Sigma, \delta_{\mathfrak{D}}, \mathcal{F})$  with the set of accepting states  $\mathcal{F}$ , and compute product MDP  $\mathfrak{P} = \mathfrak{M} \times \mathfrak{D}$
- For generalization, we leverage features related to agent state (e.g., *distances to obstacles*)
- Our DQN algorithm ( **$(\epsilon, \delta)$ -greedy**) produces a policy  $\mu^*$  for the PMDP  $\mathfrak{P}$ ; Projection of  $\mu^*$  onto MDP  $\mathfrak{M}$  yields  $\xi^*$  (Problem 1)



## Methods (Biased Action)

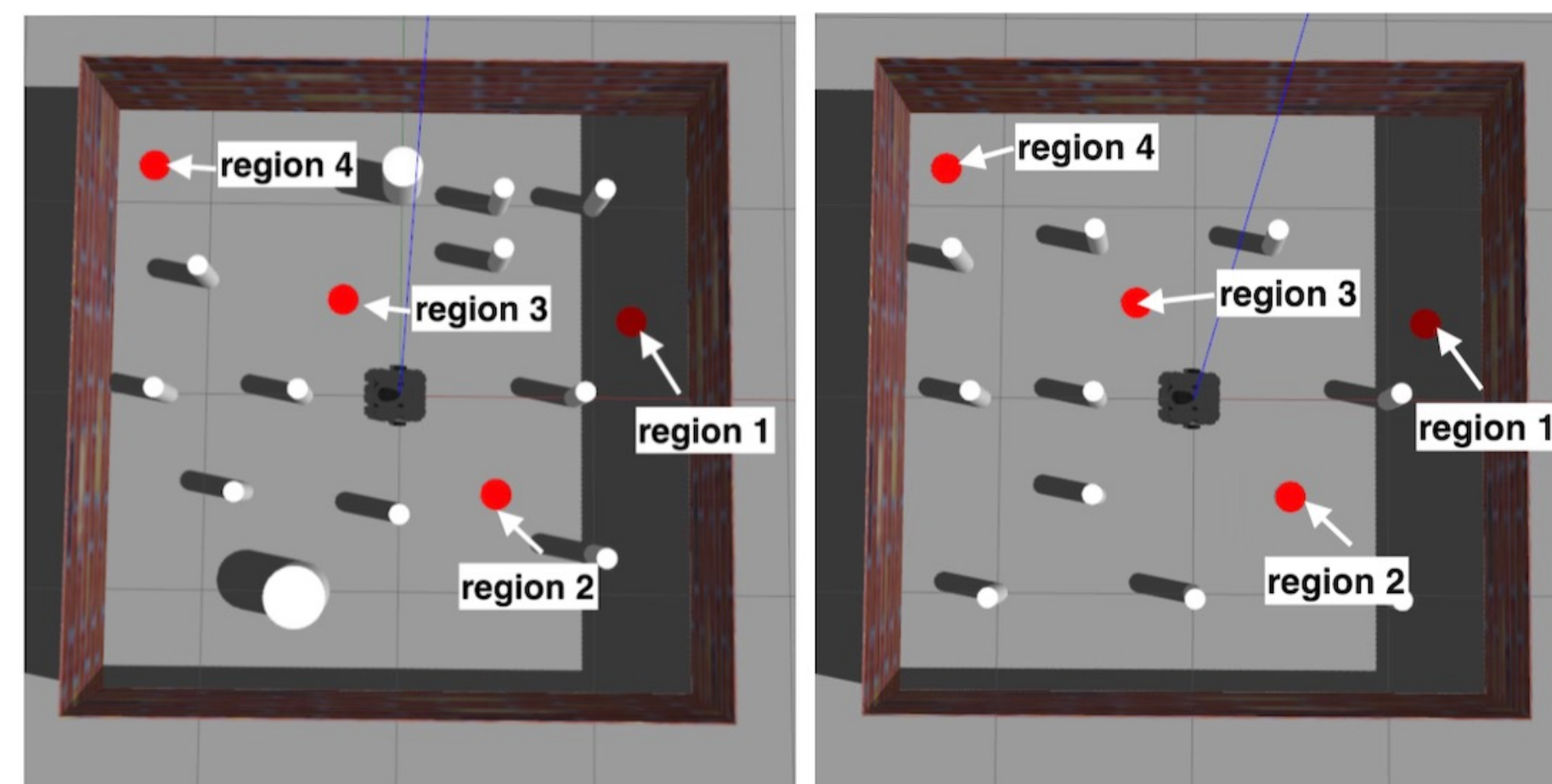
- Compute set  $\mathcal{Q}_{\text{goal}}(q_{\mathfrak{D}}^t)$  that is one-hop reachable from  $q_{\mathfrak{D}}^t$  and closer to accepting DRA states  $\mathcal{F}$ , and collect a set  $\mathcal{X}_{\text{goal}}(\mathcal{Q}_{\mathfrak{D}}^t)$  of 'goal' MDP states and randomly select one as  $x_{\text{goal}}$
- A **biased** action can drive **closer** to goal state
- Train a NN model  $g$  to compute the **biased action** prior to RL training**



- How to Collect Training Dataset:
- (i) Environment discretization (into grids) with goal states placed at center of each cell
- (ii) Randomly sample a set  $\mathcal{X}_{\text{start}}$ , assign each state to the nearest discrete cell  $i$
- (iii) Compute weighted Dijkstra distances to all goal cells, simulate  $Z$  times each for each action and compute the optimal action  $a_b$  to reach goal
- (iv) Collect a dataset of  $(x_{\text{start}}, x_{\text{goal}}, a_b)$

## Experiments

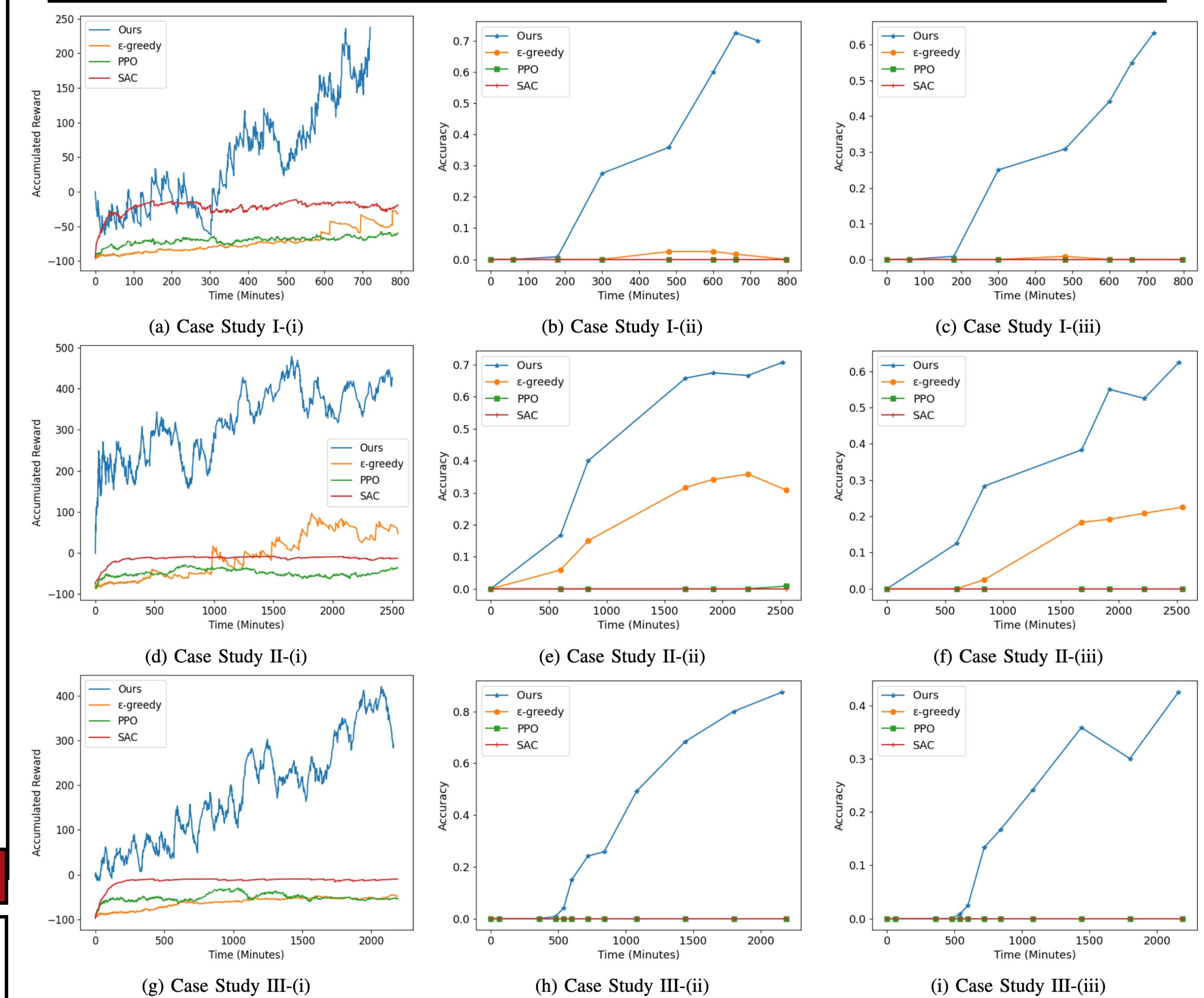
- Unknown robot dynamics with states  $x_t = [p_t^1, p_t^2, \theta_t]$  with  $|\mathcal{A}| = 23$  actions of  $(u, \omega)$ , and additive Gaussian actuation noise
- Baselines:**  $\epsilon$ -greedy DQN | PPO | SAC



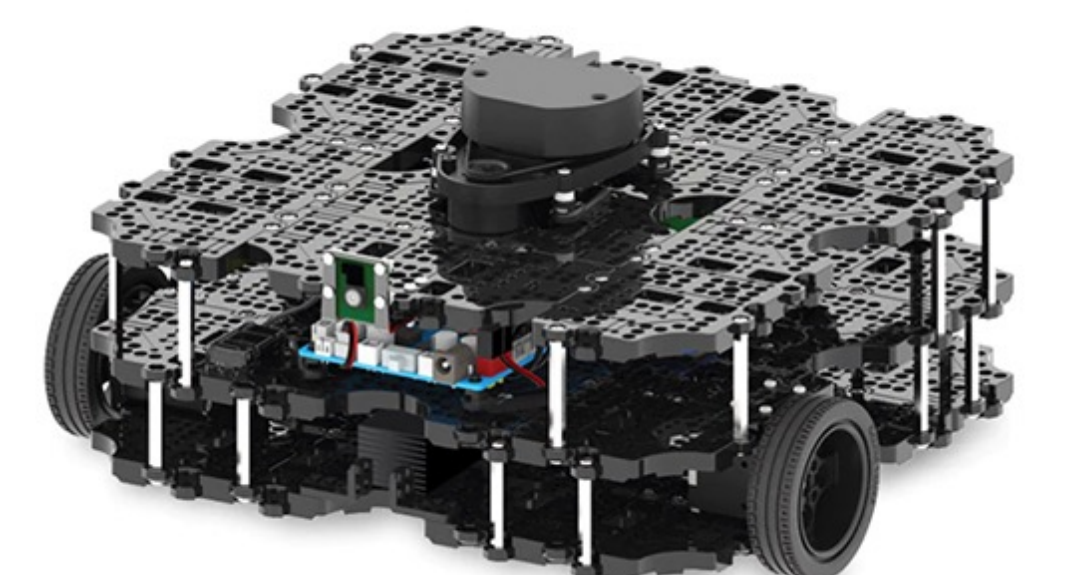
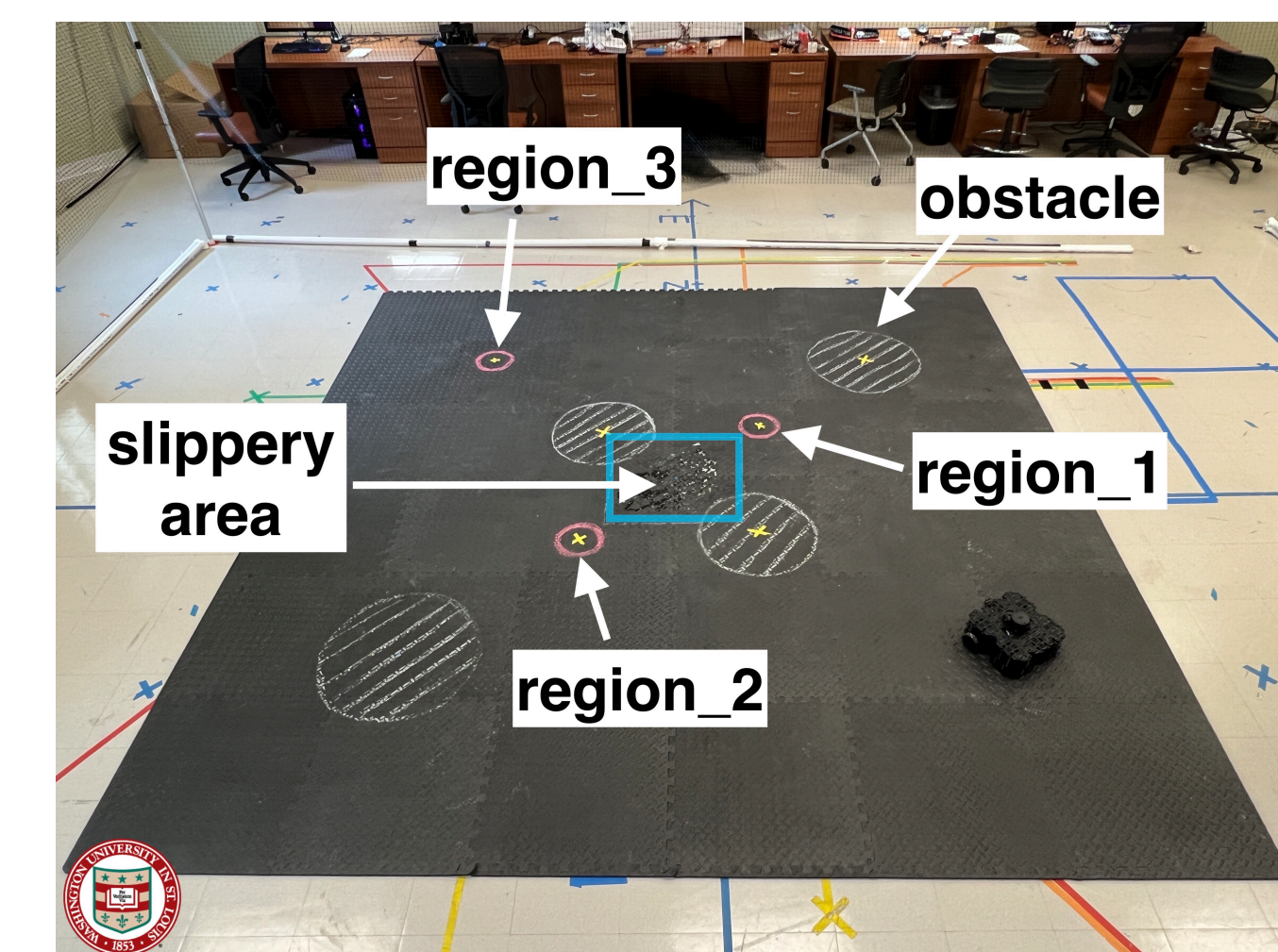
Training (Left) and Test (Right) Gazebo Envs

## Experiments and Results

- Eval Metric: (i) Average return per episode; (ii) Accuracy on **training** environments; (iii) Accuracy on **test** environments
- We outperformed baselines in sample efficiency, especially when task or environment complexity increases.



Ours (Blue), DQN (Orange), PPO (Green), SAC (Red)  
Columns 1 ~ 3 plot metrics (i), (ii), and (iii), respectively.



Hardware Environment (Left) and Turtlebot Waffle Pi (Right)