# On Generating Explanations for Reinforcement Learning Policies: An Empirical Study

**Mikihisa (Miki) Yuasa[1], Huy T. Tran[1], and Ramavarapu S. Sreenivas[2]**

[1]Department of Aerospace Engineering / [2]Department of Industrial & Enterprise Systems Engineering, The Grainger College of Engineering, University of Illinois Urbana-Champaign

**LIRA@illinois**
Lab for Intellignet Robots and Agents

## INTRODUCTION

**Black-box nature of neural networks**: While learning-based methods have advanced robot decision-making and control, their lack of interpretability raises concerns for **safety-critical applications** like autonomous vehicles.

**Need for explainability**: Formal methods, such as **Linear Temporal Logic (LTL)**, offer a structured way to interpret robot policies by prioritizing constraints based on importance.

**Limitations of existing approaches**
Existing approaches typically learn from a fixed set of example trajectories. This makes them highly dependent on the provided data and unable to access the policy's internal logic, such as its preference or certainty for actions in unobserved states.

**Contribution**
- Introduce a novel method that overcomes these limitations by **assuming direct access** to the target policy itself.
- Automatically find an explanation in the form of a **Linear Temporal Logic (LTL)** formula.

**Key Idea**
Compare the **action distributions** of the target policy with policies explicitly optimized for candidate LTL formulas, allowing for a **more nuanced evaluation** that captures the agent's underlying intent.

## METHOD

We frame the problem as a search for the LTL formula that produces a policy most similar to the target policy ($\pi_{tar}$).

**Explanation Structure**
We define explanations as LTL formulas with a specific structure, capturing both a goal and a safety constraint:

$$\phi = \mathcal{F}(\phi_F) \wedge \mathcal{G}(\phi_G).$$

- $\mathcal{F}(\phi_F)$: "*Eventually, achieve some task $\phi_F$.*"
- $\mathcal{G}(\phi_G)$: "*Globally, always satisfy some safety constraint $\phi_G$.*"
- $\phi_F$ and $\phi_G$ are logical combinations of user-defined atomic predicates (e.g., distance_to_goal < 1).

**Evaluation Metric: Weighted KL Divergence**
How do we measure similarity? We compare the policies' action distributions ($\pi(a|s)$) in key states.

- **Utility Score** ($U^\phi$): We calculate the **weighted Kullback-Leibler (wKL) divergence** between distributions of the target policy and a candidate policy:

$$U^\phi = -\sum_s w_s D_{KL}\left(\pi_\phi(\cdot|s) \| \pi_{tar}(\cdot|s)\right).$$

- **Weighting** ($w_s$): The weights emphasize states where the target policy is most certain about its action, as these are the most informative states for understanding its intent.
- **Benefit:** This avoids trivial "catch-all" explanations (e.g., "*eventually do anything*") that might be satisfied by trajectories but do not capture the policy's specific logic.

**The Search Algorithm**
We employ a **greedy local-search algorithm** to navigate the space of possible LTL explanations.
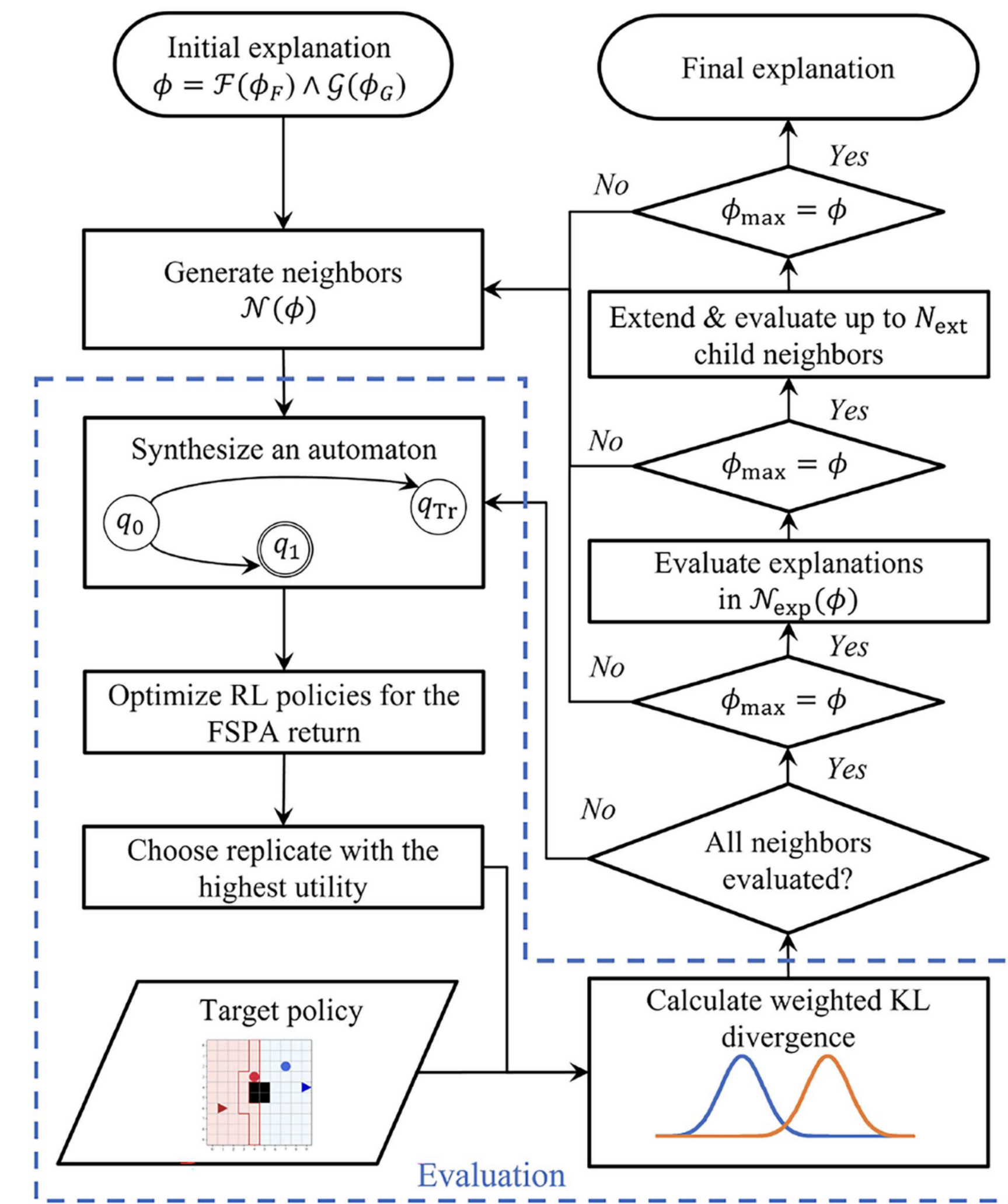
Fig 1. Overview of our proposed search algorithm.

## ANALYSIS

Our method successfully recovered the underlying objectives of the target policies.

**Ground-Truth Recovery:** In the **CtF**, **Parking**, and **Robot Navigation** environments, where the target policy was trained on a known LTL formula, our search **successfully identified the exact ground-truth formula** as the best explanation.

**Plausible Explanations:** For the **Robot Navigation** task with a non-LTL target policy, our method found a reasonable explanation: "*Eventually, the robot reaches the goal or does not hit the vase. Globally, the robot does not enter a hazard.*" This captures both the primary goal and the learned safety behavior.

**Ablation Study:** Removing any key component of our method (the search extension heuristics or the wKL weighting) caused the search to fail, confirming the importance of each part of the algorithm.

## CONCLUSIONS

We have developed a robust method for generating LTL explanations of RL policies by comparing action distributions. This approach provides deeper insight into an agent's learned behavior than existing methods.

**Limitations & Future Directions:**

**Scalability:** The search is computationally intensive. Future work could explore neural network representations of LTL to improve efficiency.

**Predicate Definition:** The method currently requires user-defined predicates. Automating predicate discovery is a key next step.

**Natural Language:** Translating the final LTL formula into natural language would further improve interpretability.

## ACKNOWLEDGEMENTS

## Experimental Setup

We tested our method in **three distinct simulated environments**.

**Target Policy Goals**

*CtF*: Capture the red flag while avoidng the red agent.

*Parking*: Park in a designated spot while avoiding walls and another car.

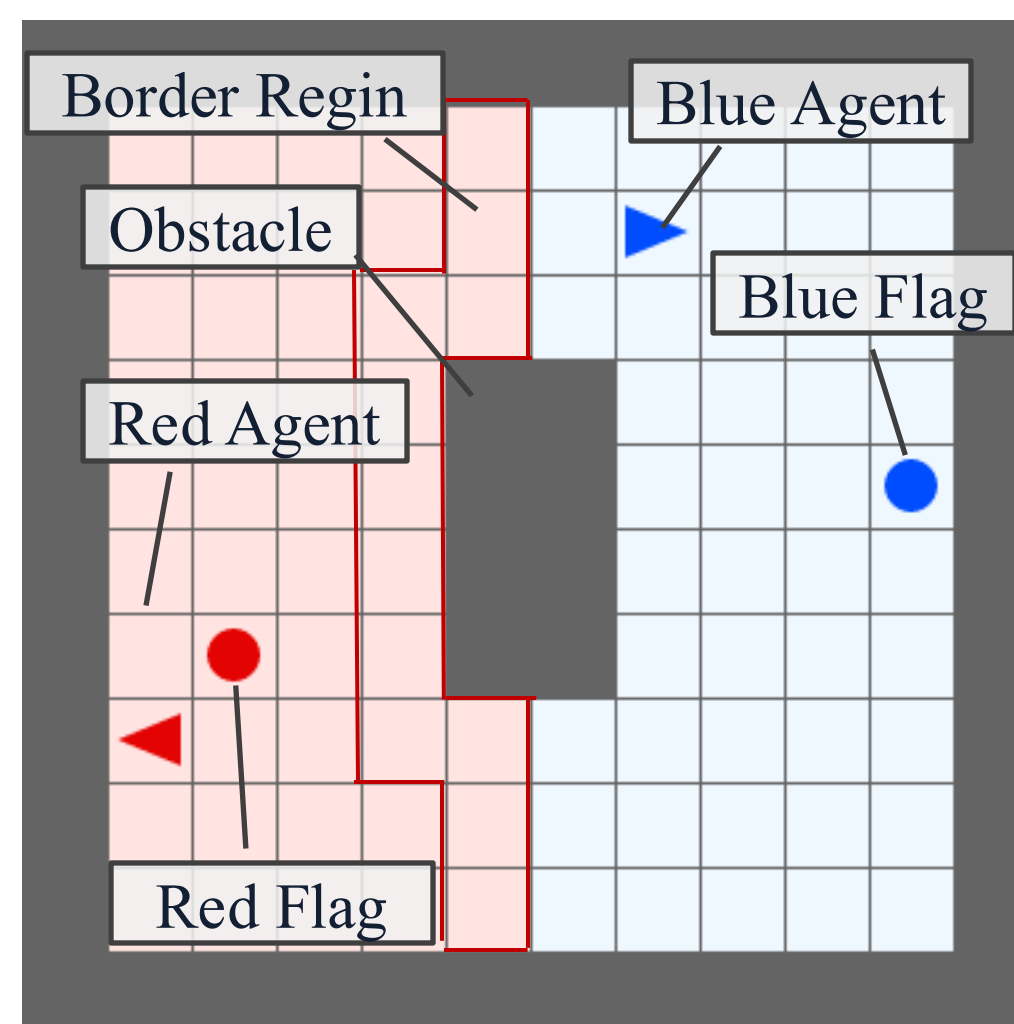*Robot Nav.*: Navigate to a goal while avoiding hazards.
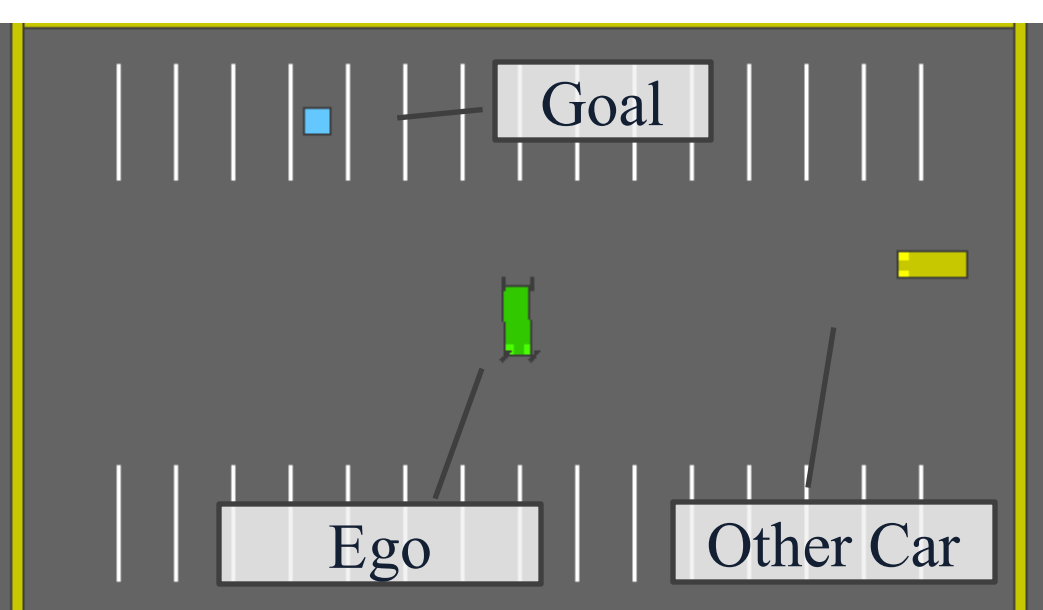
Fig 2. Capture-the-Flag (CtF)
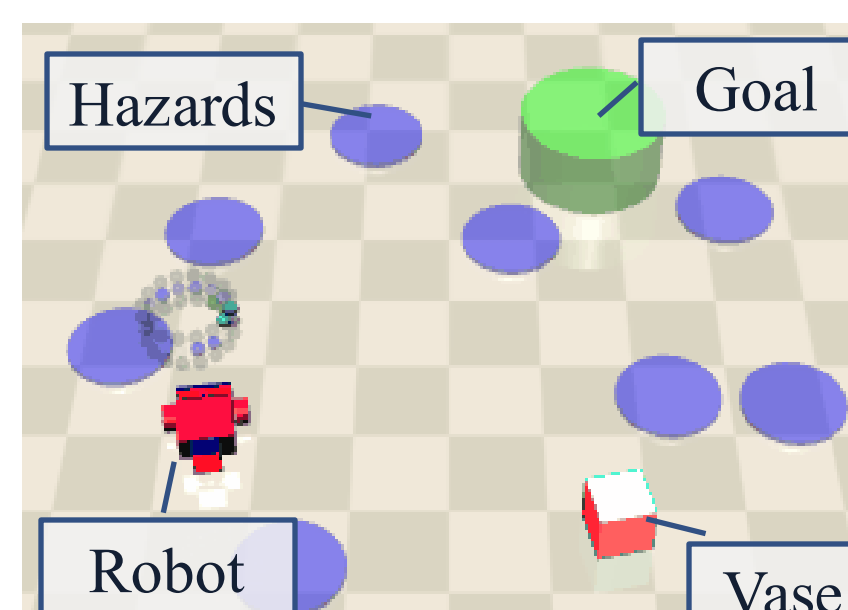
Fig 3. Adversarial Parking

Fig 4. Robot Navigation

## RESULTS

Table I. CtF & Parking Results.
Target policies were successfully found in Searches 1 & 2 (CtF) and 1, 2, & 3 (Parking).

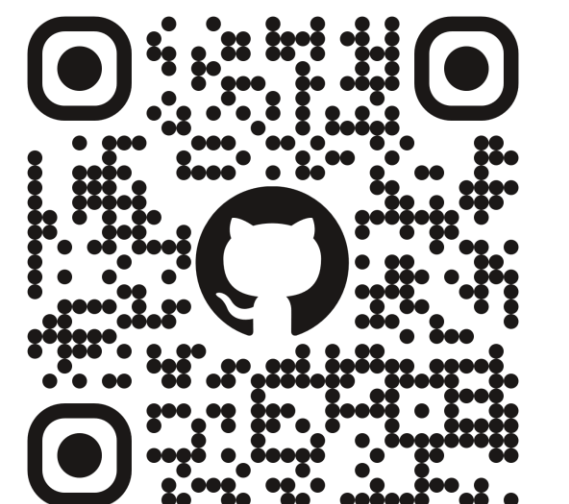| Search | CtF explanations | wKL div. [-] | Searched specs [%] | Parking explanations | wKL div. [-] | Searched specs [%] |
|---|---|---|---|---|---|---|
| 1 | $\mathcal{F}(\psi_{ba,rf} \wedge \neg\psi_{ra,bf}) \wedge \mathcal{G}(\neg\psi_{ba,ra} \vee \psi_{ba,bt})$ | $8.00 \times 10^{-8}$ | 8.13 | $\mathcal{F}(\psi_g) \wedge \mathcal{G}(\neg\psi_o \wedge \neg\psi_w)$ | 0.00 | 37.5 |
| 2 | $\mathcal{F}(\psi_{ba,rf} \wedge \neg\psi_{ra,bf}) \wedge \mathcal{G}(\neg\psi_{ba,ra} \vee \psi_{ba,bt})$ | $8.00 \times 10^{-8}$ | 10.2 | $\mathcal{F}(\psi_g) \wedge \mathcal{G}(\neg\psi_o \wedge \neg\psi_w)$ | 0.00 | 29.2 |
| 3 | $\mathcal{F}(\neg\psi_{ba,bt}) \wedge \mathcal{G}((\neg\psi_{ba,ra} \vee \psi_{ba,rf}) \vee (\neg\psi_{ra,bf}))$ | $1.46 \times 10^{-7}$ | 6.56 | $\mathcal{F}(\psi_g) \wedge \mathcal{G}(\neg\psi_o \wedge \neg\psi_w)$ | 0.00 | 37.5 |
| 4 | $\mathcal{F}(\psi_{ba,rf}) \wedge \mathcal{G}((\neg\psi_{ba,bt} \wedge \neg\psi_{ra,bf}) \vee (\neg\psi_{ba,ra}))$ | $7.42 \times 10^{-7}$ | 8.44 | $\mathcal{F}(\psi_o) \wedge \mathcal{G}(\neg\psi_g \vee \neg\psi_w)$ | $4.61 \times 10^{-4}$ | 44.8 |
| 5 | $\mathcal{F}(\neg\psi_{ba,bt} \wedge \neg\psi_{ra,bf}) \wedge \mathcal{G}((\neg\psi_{ba,rf}) \vee (\neg\psi_{ba,ra}))$ | $1.57 \times 10^{-6}$ | 6.56 | $\mathcal{F}(\psi_o) \wedge \mathcal{G}(\neg\psi_g \vee \neg\psi_w)$ | $4.61 \times 10^{-4}$ | 33.3 |
| 6 | $\mathcal{F}((\neg\psi_{ra,bf}) \vee (\psi_{ba,rf})) \wedge \mathcal{G}(\neg\psi_{ba,ra} \vee \psi_{ba,bt})$ | $5.95 \times 10^{-6}$ | 9.06 | $\mathcal{F}(\psi_w) \wedge \mathcal{G}(\neg\psi_g \vee \psi_o)$ | $6.40 \times 10^{-4}$ | 30.2 |
| 7 | $\mathcal{F}(\psi_{ba,ra}) \wedge \mathcal{G}((\neg\psi_{ba,bt}) \vee (\neg\psi_{ba,rf} \wedge \neg\psi_{ra,bf}))$ | $1.18 \times 10^{-5}$ | 8.59 | $\mathcal{F}(\neg\psi_o \wedge \psi_w) \wedge \mathcal{G}(\neg\psi_g)$ | $7.35 \times 10^{-4}$ | 28.1 |
| 8 | $\mathcal{F}((\psi_{ba,ra}) \wedge (\neg\psi_{ba,rf})) \wedge \mathcal{G}((\psi_{ba,bt}) \vee (\neg\psi_{ra,bf}))$ | $1.18 \times 10^{-5}$ | 11.4 | $\mathcal{F}(\neg\psi_o \wedge \psi_w) \wedge \mathcal{G}(\neg\psi_g)$ | $7.35 \times 10^{-4}$ | 27.1 |
| 9 | $\mathcal{F}((\psi_{ra,bf}) \vee (\neg\psi_{ba,bt})) \wedge \mathcal{G}((\neg\psi_{ba,ra}) \vee (\psi_{ba,rf}))$ | $1.62 \times 10^{-5}$ | 7.19 | - | - | - |
| 10 | $\mathcal{F}((\neg\psi_{ba,ra}) \vee (\neg\psi_{ra,bf})) \wedge \mathcal{G}((\psi_{ba,bt}) \vee (\psi_{ba,ra}))$ | $8.52 \times 10^{-5}$ | 6.56 | - | - | - |

Table II. Robot Navigation Results for LTL (left) & Non-LTL (right) Target Policies.

| Search | Robot explanations | wKL div. [-] | Searched specs [%] | Robot explanations | wKL div. [-] | Searched specs [%] |
|---|---|---|---|---|---|---|
| 1 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 21.8 | $\mathcal{F}(\psi_{gl} \vee \neg\psi_{vs}) \wedge \mathcal{G}(\neg\psi_{hz})$ | $3.0501 \times 10^{-4}$ | 39.6 |
| 2 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 30.2 | $\mathcal{F}(\neg\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \vee \neg\psi_{vs})$ | $3.0504 \times 10^{-4}$ | 29.2 |
| 3 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 26.0 | $\mathcal{F}(\neg\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \vee \neg\psi_{vs})$ | $3.0504 \times 10^{-4}$ | 27.1 |
| 4 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 29.2 | $\mathcal{F}(\neg\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | $3.0504 \times 10^{-4}$ | 27.1 |
| 5 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 21.9 | $\mathcal{F}(\neg\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \vee \neg\psi_{vs})$ | $3.0504 \times 10^{-4}$ | 29.2 |
| 6 | $\mathcal{F}(\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \wedge \neg\psi_{vs})$ | 0.00 | 24.0 | $\mathcal{F}(\neg\psi_{gl}) \wedge \mathcal{G}(\neg\psi_{hz} \vee \neg\psi_{vs})$ | $3.0504 \times 10^{-4}$ | 33.3 |
| 7 | $\mathcal{F}(\neg\psi_{vs}) \wedge \mathcal{G}(\psi_{gl} \vee \neg\psi_{hz})$ | $2.64 \times 10^{-4}$ | 25.0 | $\mathcal{F}(\psi_{gl} \vee \psi_{hz}) \wedge \mathcal{G}(\neg\psi_{vs})$ | $3.0508 \times 10^{-4}$ | 29.2 |
| 8 | $\mathcal{F}(\neg\psi_{vs}) \wedge \mathcal{G}(\psi_{gl} \vee \neg\psi_{hz})$ | $2.64 \times 10^{-4}$ | 27.1 | $\mathcal{F}(\psi_{gl} \vee \psi_{hz}) \wedge \mathcal{G}(\neg\psi_{vs})$ | $3.0508 \times 10^{-4}$ | 21.9 |

## LINKS

Paper (IEEE L-CSS)    GitHub

**ILLINOIS**