

## 886 Broader Impacts

887 Our work aims to improve the generalization of reinforcement learning to satisfy arbitrary LTL  
 888 specifications. This framework enables flexible and expressive task definitions, but it also places  
 889 responsibility on the user to ensure that the specifications are well designed. The behavior of the  
 890 agent is tightly coupled with the given specification, so poorly constructed or harmful specifications  
 891 may lead to unsafe or unintended behaviors. Thus, careful design and verification of specifications  
 892 are essential when applying such methods in real-world settings.

## 893 A Experimental Settings

### 894 A.1 Environments

895 We use the LetterWorld [49] and ZoneEnv [49, 23] for evaluation. The LetterWorld is a  $7 \times 7$   
 896 grid world that contains 12 letters corresponding to atomic propositions  $AP = \{a, b, \dots, l\}$ . Each  
 897 letter appears twice and is randomly placed on the grid. The observation consists of the entire  
 898 grid in an egocentric view. The agent can move in four directions: up, down, left, and right. The  
 899 grid wraps around, i.e., if the agent moves out of bounds, it reappears on the opposite side. The  
 900 maximum episode length is  $T = 75$ . The ZoneEnv contains different colored regions corresponding  
 901 to atomic propositions  $AP = \{\text{blue, green, magenta, yellow}\}$  and walls acting as boundaries. Our  
 902 implementation is adapted from [23]. We use the point robot that has a continuous action space for  
 903 acceleration and steering. Observations consist of Lidar measurements of the colored regions and ego-  
 904 centric state information from other sensors. The maximum episode length is  $T = 1000$ . We further  
 905 construct several variants of the ZoneEnv to evaluate the generalization performance of our method  
 906 and baselines by modifying: (1) the region size, with  $r = \{0.2, 0.3, 0.4(\text{default}), 0.5, 0.6\}$ ; (2) the  
 907 number of regions per atomic proposition, with count  $= \{1, 2(\text{default}), 3, 4, 5\}$ ; (3) and the number  
 908 of atomic propositions, expanded to include new symbols  $\{\text{red, cyan, orange, purple, teal, lime}\}$ .  
 909 Illustrations of the environments are provided in Figure 7 and trajectories under different LTL  
 910 specifications are provided in Figure 11, 13, and 14.

911 We exclude the FlatWorld[23, 44] environment from our evaluation. It contains colored regions  
 912 similar to ZoneEnv. However, the region layout is fixed, meaning that certain Boolean formulas  
 913 over the atomic propositions cannot be evaluated since they are never true in this environment.  
 914 Additionally, the observation space is limited to the agent’s  $(x, y)$ -position. Due to this lack of  
 915 variability and expressiveness, we exclude FlatWorld from our evaluation.

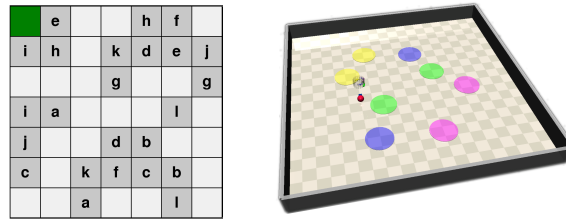


Figure 7: Visualisations of environments. Left: LetterWorld. Right: ZoneEnv

### 916 A.2 LTL Specifications

917 We evaluate zero-shot generalization across a range of specifications, including complex and nested  
 918 finite-horizon tasks (Table 3), complex infinite-horizon tasks (Table 4), varying specification complex-  
 919 ity with different sequence lengths of reach-only and reach-avoid specifications (Table 5), varying  
 920 environment complexity in reach-avoid tasks (Table 5), and different numbers of atopic propositions  
 921 (Table 6). The specifications  $\varphi_1 - \varphi_4$  and  $\varphi_9 - \varphi_{12}$  are from [23]. To better evaluate each method’s  
 922 generalization ability and its handling of the trade-off between task progression and safety constraints,  
 923 we construct more complex specifications,  $\varphi_5 - \varphi_8$  and  $\varphi_{13} - \varphi_{16}$ , which feature longer temporal  
 924 sequences and stricter safety requirements, i.e., more atomic propositions appear under negation.

Table 3: Complex finite-horizon specifications used in our evaluation. Results are shown in Table 1.

LetterWorld	$\varphi_1$	$F(a \wedge (\neg b \vee c)) \wedge Fd$
	$\varphi_2$	$(Fd) \wedge (\neg f \vee U(d \wedge Fb))$
	$\varphi_3$	$\neg a \vee U(b \wedge (\neg c \vee U(d \wedge (\neg e \vee Uf))))$
	$\varphi_4$	$(a \vee b \vee c \vee d \Rightarrow F(e \wedge F(f \wedge Fg))) \vee U(h \wedge Fi)$
	$\varphi_5$	$F(d \wedge (\neg(a \vee b) \vee U(b \wedge (\neg e \vee Uc)))) \wedge F(\neg(f \vee g \vee h) \vee a)$
	$\varphi_6$	$F((k \wedge ((\neg b \vee c) \vee Uf)) \wedge (\neg(a \vee e \vee h) \vee Ug)) \wedge Fd$
	$\varphi_7$	$\neg(j \vee b \vee d) \vee U(a \wedge (\neg c \vee U(f \wedge F(g \wedge (\neg d \vee Ue))))$
	$\varphi_8$	$\neg(f \vee g) \vee U(a \wedge (\neg b \vee Uc) \wedge F(d \wedge (\neg e \vee Uf)))$
ZoneEnv	$\varphi_9$	$(F \text{ blue}) \wedge (\neg \text{blue} \vee U(\text{green} \wedge F \text{ yellow}))$
	$\varphi_{10}$	$\neg(\text{magenta} \vee \text{yellow}) \vee U(\text{blue} \wedge F \text{ green})$
	$\varphi_{11}$	$\neg \text{green} \vee U((\text{blue} \vee \text{magenta}) \wedge (\neg \text{green} \vee U \text{ yellow}))$
	$\varphi_{12}$	$(\text{green} \vee \text{blue} \Rightarrow (\neg \text{yellow} \vee U \text{ magenta})) \vee U \text{ yellow}$
	$\varphi_{13}$	$F(\text{green} \wedge (\neg(\text{blue} \vee \text{yellow}) \vee U(\text{yellow} \wedge (\neg \text{magenta} \vee U \text{ blue})))) \wedge F(\neg \text{green} \vee U \text{ yellow})$
	$\varphi_{14}$	$F((\text{blue} \vee \text{green}) \wedge (\neg \text{yellow} \vee U(\text{blue} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U \text{ magenta})))) \wedge F(\text{yellow} \wedge (\neg \text{blue} \vee U \text{ green}))$
	$\varphi_{15}$	$\neg(\text{magenta} \vee \text{yellow}) \vee U(\text{blue} \wedge (\neg \text{green} \vee U(\text{yellow} \wedge F(\text{green} \wedge (\neg \text{blue} \vee U \text{ magenta}))))$
	$\varphi_{16}$	$F(\text{blue} \wedge (\neg \text{yellow} \vee U(\text{green} \wedge F(\text{yellow} \wedge (\neg(\text{magenta} \vee \text{green}) \vee U \text{ blue}))))$

Table 4: Infinite-horizon specifications used in our evaluation. Results are shown in Table 2.

LetterWorld	$\psi_1$	$GF \text{ blue} \wedge GF \text{ green} \wedge G \neg(\text{yellow} \vee \text{magenta})$
	$\psi_2$	$GF \text{ blue} \wedge GF \text{ yellow} \wedge GF \text{ green} \wedge G \neg \text{magenta}$
	$\psi_3$	$F G \text{ yellow} \wedge G \neg(\text{green} \vee \text{blue} \vee \text{magenta})$
ZoneEnv	$\psi_4$	$GF(e \wedge (\neg a \vee Uf)) \wedge G \neg(c \vee d)$
	$\psi_5$	$GF a \wedge GF b \wedge GF c \wedge G \neg(e \vee f \vee i)$
	$\psi_6$	$GF c \wedge GF a \wedge GF(e \wedge (\neg f \vee Ug)) \wedge GF k \wedge G \neg(i \vee j)$

Table 5: Reach-only and reach-avoid specifications. Evaluation results are shown in Figure 3.

Reach-only	n=4	$F(\text{yellow} \wedge F(\text{blue} \wedge F(\text{green} \wedge F \text{ magenta})))$ $F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F \text{ yellow})))$ $F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F \text{ green})))$
	n=6	$F(\text{green} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{yellow} \wedge F(\text{blue} \wedge F \text{ green}))))$ $F(\text{yellow} \wedge F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F \text{ green}))))$ $F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{blue} \wedge F \text{ magenta}))))$
	n=8	$F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F \text{ magenta}))))))$ $F(\text{green} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F \text{ yellow}))))))$ $F(\text{magenta} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F \text{ yellow}))))))$
	n=10	$F(\text{magenta} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{blue} \wedge F(\text{green} \wedge F \text{ yellow}))))))$ $F(\text{blue} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F \text{ yellow}))))))$ $F(\text{yellow} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{blue} \wedge F \text{ magenta}))))))$
	n=12	$F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F \text{ magenta}))))))$ $F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{green} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F \text{ yellow}))))))$ $F(\text{yellow} \wedge F(\text{magenta} \wedge F(\text{blue} \wedge F(\text{green} \wedge F(\text{yellow} \wedge F(\text{blue} \wedge F(\text{magenta} \wedge F(\text{green} \wedge F(\text{blue} \wedge F(\text{yellow} \wedge F(\text{magenta} \wedge F \text{ green}))))))$
Reach-avoid	n=2	$(\neg \text{blue} \wedge \neg \text{green}) \vee U \text{ yellow}$ $(\neg \text{yellow} \wedge \neg \text{magenta}) \vee U \text{ green}$ $(\neg \text{green} \wedge \neg \text{magenta}) \vee U \text{ blue}$
	n=4	$\neg(\text{blue} \vee \text{yellow}) \vee U(\text{green} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U \text{ yellow}))$ $\neg(\text{green} \vee \text{yellow}) \vee U(\text{magenta} \wedge (\neg(\text{blue} \vee \text{yellow}) \vee U \text{ green}))$ $\neg(\text{yellow} \vee \text{magenta}) \vee U(\text{blue} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U \text{ yellow}))$
	n=6	$\neg(\text{blue} \vee \text{yellow}) \vee U(\text{green} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{blue}) \vee U \text{ magenta}))))$ $\neg(\text{yellow} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{yellow}) \vee U(\text{blue} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U \text{ yellow}))))$ $\neg(\text{green} \vee \text{yellow}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U \text{ green}))))$
	n=8	$\neg(\text{green} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{yellow} \vee \text{blue}) \vee U(\text{green} \wedge (\neg(\text{blue} \vee \text{yellow}) \vee U(\text{magenta} \wedge (\neg(\text{yellow} \vee \text{blue}) \vee U \text{ green}))))))$ $\neg(\text{green} \vee \text{yellow}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U(\text{blue} \wedge (\neg(\text{magenta} \vee \text{yellow}) \vee U \text{ green}))))))$ $\neg(\text{magenta} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U(\text{blue} \wedge (\neg(\text{magenta} \vee \text{green}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{blue}) \vee U \text{ magenta}))))))$
	n=10	$\neg(\text{magenta} \vee \text{yellow}) \vee U(\text{green} \wedge (\neg(\text{yellow} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{yellow} \vee \text{green}) \vee U(\text{blue} \wedge (\neg(\text{magenta} \vee \text{green}) \vee U(\text{yellow} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U \text{ green}))))))$ $\neg(\text{green} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U(\text{green} \wedge (\neg(\text{magenta} \vee \text{yellow}) \vee U(\text{blue} \wedge (\neg(\text{magenta} \vee \text{green}) \vee U \text{ yellow}))))))$ $\neg(\text{yellow} \vee \text{blue}) \vee U(\text{green} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U(\text{yellow} \wedge (\neg(\text{magenta} \vee \text{blue}) \vee U(\text{green} \wedge (\neg(\text{yellow} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{yellow}) \vee U \text{ blue}))))))$

Table 6: Reach-avoid specifications with more APs. Evaluation results are shown in Table 7.

ZoneEnv	I API = 4	$\neg(\text{blue} \vee \text{yellow}) \vee U(\text{green} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{blue}) \vee U \text{ magenta}))))$
		$\neg(\text{yellow} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{yellow}) \vee U(\text{blue} \wedge (\neg(\text{green} \vee \text{magenta}) \vee U \text{ yellow}))))$
		$\neg(\text{green} \vee \text{yellow}) \vee U(\text{magenta} \wedge (\neg(\text{green} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U \text{ green}))))$
	I API = 6	$\neg(\text{red} \vee \text{cyan}) \vee U(\text{blue} \wedge (\neg(\text{yellow} \vee \text{green}) \vee U(\text{magenta} \wedge (\neg(\text{cyan} \vee \text{yellow}) \vee U \text{ green}))))$
		$\neg(\text{green} \vee \text{red}) \vee U(\text{cyan} \wedge (\neg(\text{blue} \vee \text{magenta}) \vee U(\text{yellow} \wedge (\neg(\text{cyan} \vee \text{red}) \vee U \text{ blue}))))$
		$\neg(\text{yellow} \vee \text{cyan}) \vee U(\text{red} \wedge (\neg(\text{green} \vee \text{blue}) \vee U(\text{magenta} \wedge (\neg(\text{red} \vee \text{yellow}) \vee U \text{ green}))))$
	I API = 8	$\neg(\text{orange} \vee \text{red}) \vee U(\text{cyan} \wedge (\neg(\text{blue} \vee \text{green}) \vee U(\text{yellow} \wedge (\neg(\text{purple} \vee \text{blue}) \vee U \text{ magenta}))))$
		$\neg(\text{purple} \vee \text{cyan}) \vee U(\text{orange} \wedge (\neg(\text{red} \vee \text{yellow}) \vee U(\text{green} \wedge (\neg(\text{magenta} \vee \text{red}) \vee U \text{ blue}))))$
		$\neg(\text{yellow} \vee \text{purple}) \vee U(\text{green} \wedge (\neg(\text{cyan} \vee \text{orange}) \vee U(\text{blue} \wedge (\neg(\text{magenta} \vee \text{green}) \vee U \text{ orange}))))$
	I API = 10	$\neg(\text{teal} \vee \text{magenta}) \vee U(\text{orange} \wedge (\neg(\text{lime} \vee \text{blue}) \vee U(\text{yellow} \wedge (\neg(\text{green} \vee \text{purple}) \vee U \text{ red}))))$ $\neg(\text{cyan} \vee \text{purple}) \vee U(\text{magenta} \wedge (\neg(\text{red} \vee \text{teal}) \vee U(\text{lime} \wedge (\neg(\text{orange} \vee \text{yellow}) \vee U \text{ green}))))$ $\neg(\text{blue} \vee \text{green}) \vee U(\text{yellow} \wedge (\neg(\text{teal} \vee \text{lime}) \vee U(\text{purple} \wedge (\neg(\text{magenta} \vee \text{red}) \vee U \text{ orange}))))$

## 925 B Implementation Details

926 We use the official codebases of the baselines: LTL2Action<sup>1</sup>, GCRL-LTL<sup>2</sup>, and DeepLTL<sup>3</sup>. Our im-  
 927 plementation is adapted from DeepLTL by incorporating a safety value function and a state-dependent  
 928 Lagrangian multiplier and removing DeepLTL-specific modules. For DeepLTL, LTL2Action, and our  
 929 method, we use a fully connected actor network with three hidden layers of sizes [64, 64, 64]. The  
 930 value function network has two hidden layers of sizes [64, 64]. For our method, both the safety value  
 931 function and the Lagrangian multiplier network share the same architecture as the value function.  
 932 GCRL-LTL employs a larger actor-critic architecture with [512, 1024, 256] units due to its use of a  
 933 graph neural network. For fair comparison, we use the Adam optimizer with a learning rate of  $3e - 4$   
 934 and train all methods for 15M environment interactions. The discount factor is set to  $\gamma = 0.94$  for  
 935 LetterWorld and  $\gamma = 0.998$  for ZoneEnv. Additional hyperparameter details are provided in the  
 936 code in the supplementary materials. We train all the methods using a 16-core AMD CPU and an  
 937 NVIDIA GeForce RTX 4090 GPU, and it takes roughly 3 hours to train our model.

938 Our method is summarized in Algorithm 1. We begin by constructing a subgoal set that includes all  
 939 feasible subgoals: we first enumerate each assignment  $\alpha \in 2^{AP}$  as a candidate  $\alpha^+$ . For each such  
 940  $\alpha^+$ , we filter out the remaining assignments that conflict with it. We then enumerate all possible  
 941 combinations of the filtered assignments to form  $A^-$ , resulting in the full set of feasible reach-avoid  
 942 subgoals  $\xi = \{(\alpha^+, A^-)_i\}_{i=1}^M$ , from which subgoals are sampled uniformly. The agent receives a  
 943 reward  $r = 1$  if the subgoal is satisfied and  $r = 0$  otherwise. At each timestep, it receives a safety  
 944 signal  $h = -1$  if the subgoal is not violated, and  $h = 1$  otherwise. The reward and safety functions  
 945 are defined as  $r = \mathbb{1}[L(s^\sigma) \in \alpha^+]$  and  $h = 2 \cdot \mathbb{1}[L(s^\sigma) \in A^-] - 1$ , respectively. If the current  
 946 subgoal is satisfied, we sample a new subgoal whose avoid set does not contain the current state, i.e.,  
 947  $L(s_t) \notin A^-$ , and whose reach set also excludes the current state, i.e.,  $L(s_t) \notin \alpha^+$ , to prevent the new  
 948 subgoal from being immediately violated or trivially satisfied. If a subgoal is violated, the episode  
 949 terminates. We use Generalized Advantage Estimation (GAE) [41] to estimate value functions:

$$A_r^\pi = \sum_{l=0}^{T-t-1} (\gamma \lambda_{\text{GAE}})^l \delta_{r,t+l} \quad \delta_{r,t} = r_t + \gamma V(s_{t+1}^\sigma) - V(s_t^\sigma) \quad (4)$$

$$A_h^\pi = \sum_{l=0}^{T-t-1} (\gamma \lambda_{\text{GAE}})^l \delta_{h,t+l} \quad \delta_{h,t} = (1 - \gamma) h_t + \max(h_t, V(s_{t+1}^\sigma)) \quad (5)$$

950 where  $\lambda_{\text{GAE}} = 0.95$  is the discount factor of GAE. At test time, given a target LTL specification, we  
 951 convert it into a Büchi automaton and apply the DFS to enumerate all possible paths to accepting  
 952 states. For each path, we evaluate the current subgoal using the learned value functions and select the  
 953 optimal subgoal to execute. This process is applied iteratively until an accepting state is reached.

---

### Algorithm 1 GenZ-LTL

---

**Require:** Initial policy  $\pi$ , value and safety value function  $V_r, V_h$ , state-dependent Lagrangian multiplier  $\lambda$ , subgoal set  $\xi$ , maximum training iterations  $K$ , interactions per iteration  $N$

- 1: **for**  $k = 0, 1, 2, \dots, K$  **do**
  - 2:   **for**  $n = 0, 1, 2, \dots, N$  **do**
  - 3:     **if** the current subgoal is not specified **or** the current subgoal is satisfied/violated **then**
  - 4:       Sample a subgoal  $\sigma \sim \text{Unif}(\xi)$
  - 5:     **end if**
  - 6:     Collect trajectory  $\tau = \{s_t^\sigma, a_t, r_t, h_t\}$
  - 7:   **end for**
  - 8:   Compute reward-to-go  $\hat{R}_t \doteq \sum_{i=t}^T \gamma^i r_i$  and cost-to-go  $\hat{H}_t \doteq \max_t h_t$  for each  $\tau$
  - 9:   Compute advantage functions  $A_r^\pi, A_h^\pi$ , based on  $V_r$  and  $V_h$  using (4) and (5).
  - 10:   Compute the Lagrangian multiplier  $\lambda$ .
  - 11:   Fit value function, safety value function by regression on mean-square error.
  - 12:   Update the policy parameters by maximizing (3).
  - 13:   Update the multiplier parameters by minimizing (3).
  - 14: **end for**
- 

<sup>1</sup><https://github.com/LTL2Action/LTL2Action>

<sup>2</sup><https://github.com/RU-Automated-Reasoning-Group/GCRL-LTL>

<sup>3</sup><https://github.com/mathiasj33/deep-ltl>

## 954 B.1 Subgoal-induced observation reduction.

955 For Lidar measurements in the ZoneEnv,  $f(\cdot)$  is implemented as the element-wise min to obtain the  
 956 closest distance along each Lidar beam, such that the resulting  $s^\sigma$  For grid maps in LetterWorld, we  
 957 assign distinct values to each cell associated with  $\sigma$  at location  $(i, j)$ :  $f(s_{i,j}) = v_{\text{reach}}$  if  $L(s_{i,j}) = \alpha^+$ ;  
 958  $v_{\text{avoid}}$  if  $L(s_{i,j}) \in A^-$ ;  $v_{\text{neutral}}$  otherwise. This reduction allows the agent to focus on "reach" and  
 959 "avoid" rather than specific set of atomic propositions and eliminates the need of encoding subgoals,  
 960 enabling efficient policy learning and generalization to new atomic propositions that satisfies the  
 961 conditioned discussed in Section 4.2 as shown later in Appendix C. If multiple observation functions  
 962 are available, each capable of observing an assignment  $\alpha \in 2^{AP}$ , we can first fuse the observations  
 963 within each observation type based on the current subgoal  $\sigma$  and then concatenate the results to form  
 964 the final state  $s^\sigma$ . In cases where each atomic proposition is tied to a distinct observation type, our  
 965 method is not applicable. However, learning policies in such settings is inherently intractable due to  
 966 the exponential number of subgoal-state combinations.

## 967 B.2 Timeout & Subgoal-switching Mechanism

968 At the test time, given a target LTL specification, we convert it into a Büchi automaton and apply a  
 969 depth-first search (DFS) to enumerate all possible paths to accepting states. However, not all of these  
 970 paths, or in some cases, any, are feasible in the physical environment, as some may include subgoals  
 971 that are unsatisfiable. Although the labeling function is known, meaning we can determine which  
 972 propositions are true for any given state, it does not necessarily imply that we can exhaustively iterate  
 973 over the entire state space to identify unsatisfiable assignments.

974 For example, consider the specification  $\neg \text{yellow} \cup ((\text{blue} \wedge \text{green}) \vee \text{magenta})$ . To satisfy this  
 975 formula, the agent must eventually reach either a magenta region or a region where blue and green  
 976 overlap. However, the latter becomes unsatisfiable if the blue and green regions are disjoint. In such  
 977 cases, the agent must be able to adapt by pursuing an alternative subgoal, a capability that is essential  
 978 for goal-conditioned policies, including our approach as well as existing baselines such as DeepLTL  
 979 and GCRL-LTL. However, these baselines do not address this issue.

980 To address this, we implement a timeout mechanism that allows the agent to revisit the automaton  
 981 and search for alternative subgoals when the current one cannot be satisfied. We determine subgoal  
 982 satisfiability by enforcing a timeout threshold  $(1 + \epsilon_{\text{scale}}) \cdot \mu_{\text{subgoal}}$  as discussed in Section 4.1: if the  
 983 agent fails to achieve this subgoal within that timestep, it is considered unsatisfiable, and the agent  
 984 switches its subgoal. If no feasible subgoals remain, the episode terminates.

---

### Algorithm 2 Timeout-based subgoal-switching mechanism

---

**Require:** Trained agent, Büchi automaton  $\mathcal{B}_\varphi := (\mathcal{Q}, \Sigma, \delta, \mathcal{F}, q_0)$ , timeout value  $t_{\text{max}}$

- 1:  $\mathcal{Q}_c \leftarrow \{q_0\}$  ▷ Current set of states.
- 2:  $\mathcal{U} \leftarrow \emptyset$  ▷ Set of unsatisfiable transitions.
- 3: **while** FindSubgoals( $\mathcal{B}_\varphi, \mathcal{Q}_c, \mathcal{U}$ )  $\neq \emptyset$  **do**
- 4:   Select subgoal  $\sigma = (\alpha^+, A^-)$  from FindSubgoals( $\mathcal{B}_\varphi, \mathcal{Q}_c, \mathcal{U}$ ) with maximum value.
- 5:   **for**  $t = 1, 2, \dots, t_{\text{max}}$  **do**
- 6:     Get an action from the agent, apply it, and observe  $\alpha \in \Sigma$ .
- 7:      $\mathcal{Q}_o \leftarrow \mathcal{Q}_c$  ▷ Remember the old states.
- 8:      $\mathcal{Q}_c \leftarrow \{q' \mid \exists q \in \mathcal{Q}_o, (q, \alpha, q') \in \delta\}$  ▷ Update the current set of states.
- 9:     **if**  $\mathcal{Q}_c \neq \mathcal{Q}_o$  **then**
- 10:       **break** ▷ A new state set is reached; the subgoal should be updated.
- 11:     **else if**  $t = t_{\text{max}}$  **then**
- 12:        $\mathcal{U} \leftarrow \mathcal{U} \cup \{(q, \alpha^+) \mid q \in (\mathcal{Q}_c - \mathcal{F})\}$  ▷ If  $q$  is not accepting, mark this subgoal as unsatisfiable.
- 13:     **end if**
- 14:   **end for**
- 15: **end while**
- 16: All subgoals are marked as unsatisfiable. Terminate.

---

985 This timeout & subgoal-switching mechanism is explained in Algorithm 2, which outlines the general  
 986 mode of operation during the test time. We use a set  $\mathcal{U} : \mathcal{Q} \times \Sigma$  to store the transitions that are marked  
 987 as unsatisfiable after a timeout. These unsatisfiable transitions are ignored by the FindSubgoals

function. Apart from this, the only difference between FindSubgoals and the DFS algorithm used in DeepLTL is that we separate the subgoals after identifying the paths with accepting cycles.

## C Additional Experiments

**Training curves.** We first evaluate simple reach-only and reach-avoid specifications with sequence length of  $n = 3$  every specific training interval, and the success and violation rates are shown in Figure 8. We can observe that our method outperforms the baselines in terms of achieving a lower violation rate (close to zero constraint violation) while maintaining a high success rate. The violation rates of the baselines also decrease over training because task progression and safety constraints are jointly encoded in the LTL specification, so optimizing for satisfaction of the specification implicitly improves safety performance. However, due to the inherent trade-off between these requirements, the baselines struggle to further reduce violation rates. In contrast, by formulating the problem using safe RL framework with reachability constraints, our method achieves stricter constraint satisfaction, thereby improving the overall probability of satisfying the full LTL specifications.

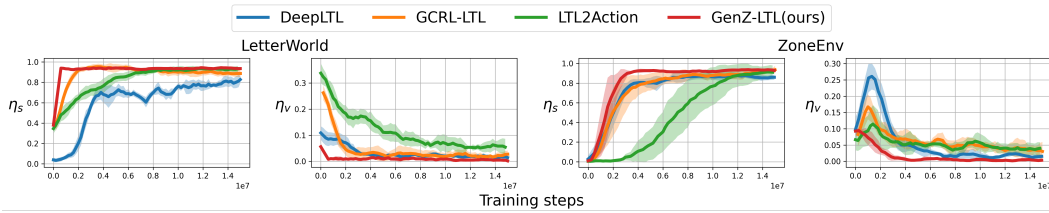


Figure 8: Evaluation curves on simple reach-only and reach-avoid specifications with  $n = 3$ . Each value is averaged over 5 seeds with 50 trajectories with randomly sampled tasks.

**Partial observability.** To satisfy arbitrary LTL specifications, we leverage the structure of the Büchi automaton to decompose each formula into a sequence of subgoals and handle one subgoal at a time, progressing through the sequence incrementally. While our method could be extended to condition on the full subgoal sequence, we argue that this introduces significant drawbacks that outweigh the potential benefits. Among the baselines, DeepLTL conditions its policy on the full subgoal sequence to achieve non-myopic behavior. Intuitively, this "look-ahead" capability may help the agent satisfy the specification in fewer steps. However, subgoal sequences come with several disadvantages. First, policies trained on full sequences are more likely to encounter OOD issues at test time, as the sequence length can vary across specifications. Although one might limit the policy to observe at most  $k$  subgoals, where  $k$  is the maximum length seen during training, choosing  $k$  is nontrivial. A large  $k$  can lead to exponential sample complexity due to the combinatorial explosion of subgoal sequences, e.g., each subgoal is drawn from  $2^{AP}$ , while a small  $k$  limits the look-ahead benefit. Moreover, in practice, the agent may not perceive all relevant subgoals at once, rendering the subgoal sequence ineffective. To illustrate this, we conduct experiments where the agent's observability is restricted to half of the environment. We train both DeepLTL and our method in LetterWorld and ZoneEnv, and then evaluate them using specifications  $\varphi_1$ – $\varphi_3$  and reach-avoid tasks with  $n = 4$ . We report the average success rate  $\eta_s$  and violation rate  $\eta_v$ . As shown in the left and middle plots of Figure 9, our single-subgoal-at-a-time method outperforms DeepLTL in terms of both the success and violation rates, despite the latter using full subgoal sequences.

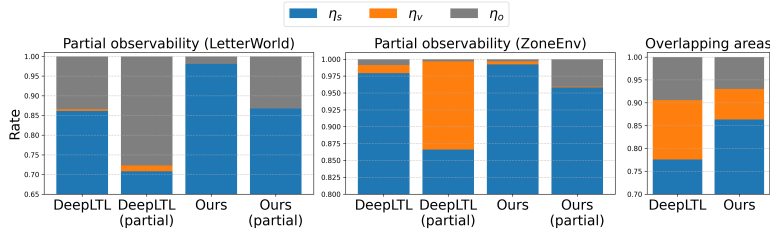


Figure 9: Evaluation results of DeepLTL and our method in two settings: (1) partial observability (left two plots) and (2) overlapping areas (the rightmost plot).

**Overlapping areas.** We construct a variant of ZoneEnv to introduce a more challenging setting where multiple atomic propositions can be true simultaneously, i.e., different colored regions may overlap in arbitrary configurations. Note that FlatWorld is an overly simplified case where the layout is fixed. In ZoneEnv, overlapping areas can be handled in two ways: (1) by adding separate observations for each assignment involving multiple true propositions, or (2) by using observations only for individual atomic propositions (default setting for ZoneEnv) and learning policies through interaction with the environment, during which the agent implicitly discovers the relationships among propositions and observations. The first approach does not offer meaningful advantages over the second, as our observation reduction technique can extract the relevant parts of the observation for any target assignment and fuse them as needed. Thus, we adopt the second approach, which presents a more challenging generalization problem. To accommodate this scenario, we modify DeepLTL’s curriculum training to allow multiple atomic propositions to be true at the same time. The positions of both the agent and the regions are randomly sampled. We evaluate both methods using reach-avoid specifications with  $n = 2$ :  $\phi_2 : \neg(\text{magenta} \vee \text{yellow}) \cup (\text{blue} \wedge \text{green})$ ,  $\phi_3 : \neg(\text{blue} \vee \text{magenta}) \cup (\text{yellow} \wedge \text{green})$ ,  $\phi_4 : \neg(\text{yellow} \vee \text{blue}) \cup (\text{green} \wedge \text{magenta})$ , and  $\phi_5 : \neg(\text{green} \vee \text{magenta}) \cup (\text{yellow} \wedge \text{blue})$ . Trajectory illustrations are shown in Figure 10. We report the average rates and the results in the rightmost plot in Figure 9 show that our method achieves a higher success rate and a lower violation rate. Our method can better satisfy safety constraints due to the integration of reachability constraints.

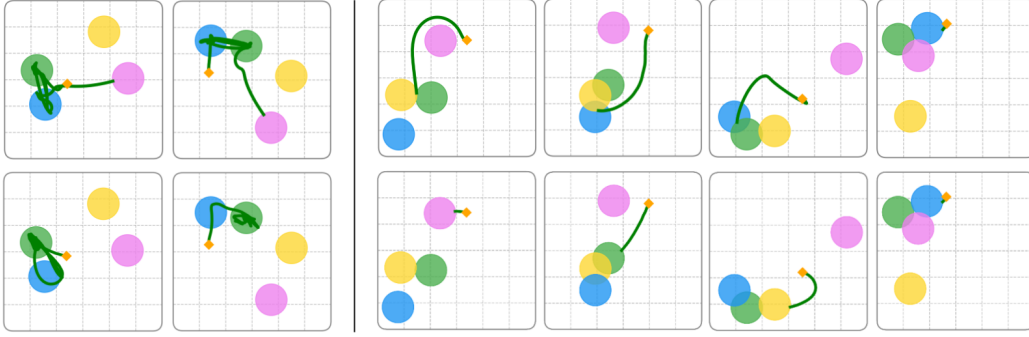


Figure 10: Illustration of trajectories in ZoneEnv: (1) unsatisfiable subgoals (left two columns,  $\phi_1$ ) and (2) overlapping areas (right four columns,  $\phi_2$ – $\phi_5$ ). The top row shows trajectories generated by our method, while the bottom row shows those from DeepLTL.

**Number of atomic propositions.** We also evaluate performance as the number of atomic propositions increases. In this experiment, we emphasize that the observations corresponding to the additional atomic propositions satisfy the conditions outlined in Section 4.2. It is important to note that when new atomic propositions are observed through modalities not encountered during training — such as those representing the agent’s speed or heading — zero-shot adaptation remains an open challenge. The results are presented in Table 7, and example trajectories are shown in Figure 11. Among the evaluated methods, only our method can generalize in a zero-shot manner to new atomic propositions whose observations are produced by identical observation functions under output transformation, as the subgoal-induced observation reduction technique fuses these propositions into "reach" and "avoid", allowing the agent to make decisions without relying on their specific labels  $L(s)$ .

Table 7: Performance in ZoneEnv with increasing number of atomic propositions. We evaluate the reach-avoid specifications with sequence length of  $n = 6$ , using an increasing number of atomic propositions as listed in Table 6. We report success rate  $\eta_s$ , violation rate  $\eta_v$ , other rate  $\eta_o$ , and average steps to satisfy the specifications  $\mu$ .

Metrics	$ \text{API}  = 4$	$ \text{API}  = 6$	$ \text{API}  = 8$	$ \text{API}  = 10$
$\eta_s \uparrow$	$0.96 \pm 0.02$	$0.94 \pm 0.02$	$0.91 \pm 0.03$	$0.93 \pm 0.03$
$\eta_v \downarrow$	$0.02 \pm 0.01$	$0.02 \pm 0.01$	$0.03 \pm 0.01$	$0.02 \pm 0.01$
$\eta_o \downarrow$	$0.02 \pm 0.02$	$0.04 \pm 0.02$	$0.06 \pm 0.03$	$0.05 \pm 0.02$
$\mu \downarrow$	$306.58 \pm 16.88$	$314.88 \pm 17.02$	$313.10 \pm 10.61$	$323.30 \pm 12.57$

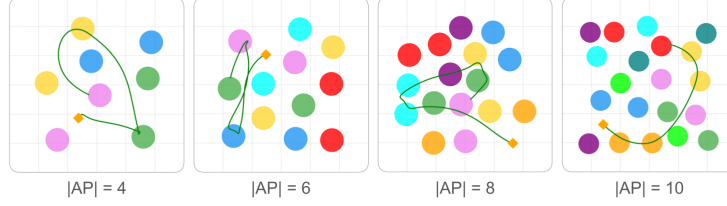


Figure 11: Illustration of ZoneEnv with an increasing number of atomic propositions. The trajectories shown correspond to the first specification for  $|AP| = 4, 6, 8, 10$  in Table 6.

**Deterministic v.s. stochastic policies.** As discussed in Section 5, we follow the standard practice of using a deterministic policy during evaluation. We also note that DeepLTL is evaluated using a stochastic policy in LetterWorld in its original paper. To provide a comprehensive comparison, we report results for both stochastic and deterministic policies in LetterWorld, as shown in Figure 12. We observe that DeepLTL exhibits inconsistent behavior: while the success rate improves under a stochastic policy, the violation rate also increases. This trade-off is unacceptable in safety-critical applications and indicates that DeepLTL struggles to maintain safety. In contrast, our method delivers consistent performance without compromising safety for task progression, since we explicitly model the violation of subgoals as a reachability constant in a safe RL formulation. In our view, stochastic policies are not suitable for evaluation. For reach-only specifications, they contain no safety constraints; a stochastic policy may satisfy them simply by randomly exploring the environment, thereby overstating performance. For reach-avoid specifications, stochasticity further introduces inconsistency: the mean action may be safe, while a sampled action may lead to a violation, or vice versa. This randomness complicates evaluation and prevents a fair comparison between methods.

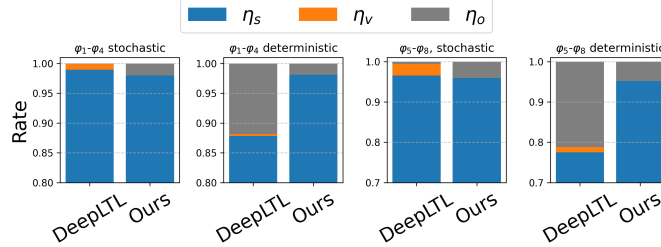


Figure 12: Evaluation results of DeepLTL and our method in LetterWorld using stochastic and deterministic policies. Each value is averaged over 5 seeds, with 100 trajectories per seed.

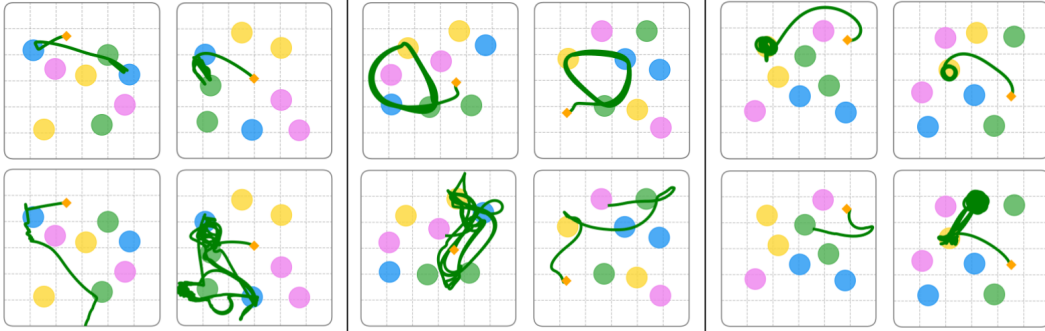


Figure 13: Illustration of infinite-horizon tasks in ZoneEnv. The trajectories correspond to specifications  $\psi_1 - \psi_3$ , shown from left to right (two columns per specification). The top row shows trajectories generated by our method, while the bottom row shows those from DeepLTL.

## D Visualizations

We provide visualizations for the experiments in Section 5.1 and Section 5.2, including infinite-horizon tasks (Figure 13), increasing environment complexity (Figure 14) and increasing numbers

1066 of atomic propositions (Figure 11). For infinite-horizon tasks, our method exhibits a clear recurring  
 1067 pattern that visits the specified regions in the correct order, while the baseline struggles to do so  
 1068 and suffers from safety constraint violations. For different environment complexities, we vary the  
 1069 zone size and zone count per atomic proposition. Importantly, we increase the zone count only for  
 1070 propositions the agent must avoid, since doing so for propositions the agent need to satisfy would  
 1071 simplify the task. Larger zones or more zones per proposition force the agent to navigate around  
 1072 constrained regions to reach the target, and the resulting trajectories demonstrate that our method  
 1073 satisfies the safety constraints. For experiments with an increasing number of atomic propositions,  
 1074 we incrementally introduce new propositions, up to 10 in total, each corresponding to a new colored  
 1075 region with associated lidar observations. This makes zero-shot generalization challenging for  
 1076 baseline methods, as they explicitly encode atomic propositions and cannot generalize to unseen  
 1077 ones without retraining. Even substituting one proposition for another, for example, replacing blue  
 1078 with red, causes these baselines to fail, due to their dependence on fixed encodings. In contrast, our  
 1079 method uses subgoal-induced observation reduction, which removes this dependency when atomic  
 1080 propositions share the same property. This allows our method to focus on the semantics of the  
 1081 propositions, i.e., "reach" or "avoid", rather than the specific proposition symbols, enabling zero-shot  
 1082 generalization to unseen atomic propositions.

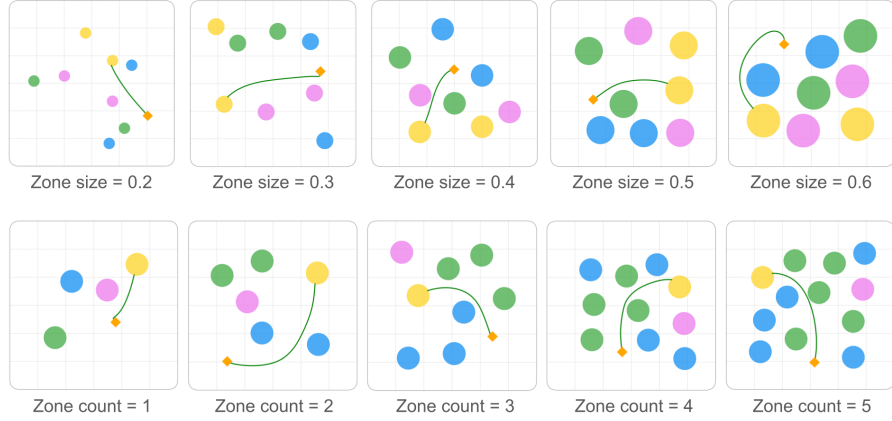


Figure 14: Illustration of ZoneEnv under increasing environment complexity, including different zone sizes and different numbers of zones per atomic proposition. The default zone size is 0.4, and the default number of zones per atomic proposition is 2. The trajectories are generated by our method and correspond to the evaluated specification  $(\neg \text{blue} \wedge \neg \text{green}) \cup \text{yellow}$ .