



## Research article

## Reinforcement learning with soft temporal logic constraints using limit-deterministic generalized Büchi automaton

Mingyu Cai<sup>a</sup>, Zhangli Zhou<sup>b</sup>, Lin Li<sup>b</sup>, Shaoping Xiao<sup>c</sup>, Zhen Kan<sup>b,\*</sup><sup>a</sup> University of California, Riverside, CA, 92521, USA<sup>b</sup> University of Science and Technology of China, Hefei, 230026, China<sup>c</sup> The University of Iowa, IA, 52246, USA

## ARTICLE INFO

## Keywords:

Formal methods in robotics and automation

Linear temporal logic

Motion planning

Optimal control

## ABSTRACT

This paper investigates control synthesis for motion planning under conditions of uncertainty, specifically in robot motion and environmental properties, which are modeled using a probabilistic labeled Markov decision process (PL-MDP). To address this, a model-free reinforcement learning (RL) approach is designed to produce a finite-memory control policy that meets complex tasks specified by linear temporal logic (LTL) formulas. Recognizing the presence of uncertainties and potentially conflicting objectives, this study centers on addressing infeasible LTL specifications. A relaxed LTL constraint enables the agent to adapt its motion plan, allowing for partial satisfaction by accounting for necessary task violations. Additionally, a new automaton structure is introduced to increase the density of accepting rewards, facilitating deterministic policy outcomes. The proposed RL framework is rigorously analyzed and prioritizes two key objectives: (1) satisfying the acceptance condition of the relaxed product MDP, and (2) minimizing long-term violation costs. Simulation and experimental results are presented to demonstrate the framework's effectiveness and robustness.

## 1. Introduction

Formal logic enables the specification of complex, high-level tasks that go beyond traditional goal-directed navigation in robotic systems. As a formal language, linear temporal logic (LTL) has become increasingly popular for use in robotic motion planning [1–7]. Since robotic systems are often subject to a variety of uncertainties arising from the stochastic behaviors of the motion (e.g., an agent may not exactly follow the control inputs due to potential sensing noise or actuation failures) and uncertain environment properties (e.g., there exist mobile obstacles or time-varying areas of interest), Markov decision processes (MDPs) are often used to model the probabilistic motion of robotic systems [8,9]. Based on probabilistic model checking, control synthesis of MDP with LTL motion specifications has been widely investigated (cf. [10–15]). In particular, the topic of partial satisfaction of high-level tasks in deterministic and stochastic systems is investigated in [14,16–19]. Yet, new challenges arise when considering motion and environment uncertainties. Hence, learning to find a satisfying policy is paramount for the robot to operate in the presence of motion and environment uncertainties.

Reinforcement learning (RL) is a sequential decision-making process in which an agent continuously interacts with and learns from the

environment [20–23]. Model-based RL has been employed for motion planning with LTL specifications when full knowledge of MDP is available [24]. The work of [25] extends model-based RL to temporal logic constrained control of stochastic systems with unknown parameters by model approximation. In [26,27], transition probabilities are learned to facilitate the satisfaction of LTL specifications. However, these aforementioned works have to rely on the accuracy of transition probabilities for learning. On the other hand, model-free RL approaches with LTL-based rewards generate desired policies by directly optimizing the Q-values [28–31]. However, these approaches rely on the key assumption that at least one accepting maximum end component (AMEC) exists in a standard product MDP [1], which may not hold in practical scenarios. For example, certain areas of interest may be probabilistically inaccessible to the agent (e.g., areas that become unreachable due to flooding, preventing traversal by a ground robot). This can result in portions of the user-specified tasks being unachievable, leading to the absence of AMECs in the product MDP. Although minimal revision of motion plans in a potentially conflicting environment has been investigated in the works of [32–35], only deterministic transition systems are considered, and it is not yet clear how to address the above-mentioned issues in stochastic systems i.e., MDP.

Peer review under responsibility of Chongqing University.

\* Corresponding author.

E-mail addresses: [mingyu.cai@ucr.edu](mailto:mingyu.cai@ucr.edu) (M. Cai), [zhangli@ustc.edu.cn](mailto:zhangli@ustc.edu.cn) (Z. Zhou), [linlee1006@mail.ustc.edu.cn](mailto:linlee1006@mail.ustc.edu.cn) (L. Li), [shaoping-xiao@uiowa.edu](mailto:shaoping-xiao@uiowa.edu) (S. Xiao), [zkan@ustc.edu.cn](mailto:zkan@ustc.edu.cn) (Z. Kan).

<https://doi.org/10.1016/j.jai.2024.12.005>

Received 26 August 2024; Received in revised form 28 November 2024; Accepted 23 December 2024

Available online 31 December 2024

2949-8554/© 2024 The Authors. Published by Elsevier B.V. on behalf of KeAi Communications Co. Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

**Related Works:** Compared with most RL-based approaches [26,30,31,36,37], the reward functions are generally designed to enforce the convergence to AMECs so that the acceptance condition can be satisfied. However, these approaches are not applicable in this work due to the fact that, despite the existence of AMECs, they may not be reachable due to environmental uncertainties. As a result, the probability of achieving the specified tasks can become negligible. To address this, the approach in [38] introduces a reduced-variance deep Q-learning technique aimed at approximating optimal state-action values. Additionally, our prior work [39] constructs a relaxed product MDP by converting LTL formulas into a Deterministic Rabin Automaton (DRA) without relying on the existence of AMECs. However, [30] suggests that converting LTL into a DRA may not reliably yield the desired policies. Even though our recent work [40] shows how to leverage planning algorithms to guide minimally-violating policies for infeasible LTL tasks. But multi-objective reinforcement learning without the assistance of planning is still under investigation.

In [41], the LDGBA is applied, and a synchronous frontier function for RL reward is designed to synthesize control policies that mostly fulfill the given task by maximizing the visits of the accepting sets in LDGBA. However, it cannot record the visited or non-visited accepting sets in each round. Compared to the modular framework in [42], which designs deep deterministic policy gradients to handle continuous spaces but lacks explicit mechanisms to record unvisited automaton states, our method leverages E-LDGBA to ensure denser reward signals and more efficient policy optimization. Unlike the mission-driven exploration strategy in [43], which accelerates RL by prioritizing task-relevant state exploration, our approach provides a more robust mechanism to handle task infeasibility through a relaxed product MDP. Furthermore, in contrast to the certified RL algorithm in [44], which shapes rewards based on deterministic automata to ensure probabilistic satisfaction, we optimize multi-objective RL with explicit violation cost analysis, allowing partial satisfaction in stochastic environments. As an extension of [39], this work can handle the above issues by developing an E-LDGBA to record the non-visited accepting sets without adding extra computational complexity. The proposed relaxed product MDP and the designed utility function transforms the control synthesis problem into an expected utility optimization problem, in which a satisfactory policy is guaranteed to be found by leveraging RL to optimize the expected utility. LDGBA reduces the number of states through a simplified automata structure, avoiding the state-space overload problem common in traditional DRA transitions.

LDGBA formulates LTL formulas via limit deterministic structures and generalized Büchi acceptance conditions. Compared to other automaton, LDGBA typically generates a smaller state space, which is computationally more efficient. For many LTL formulas, LDGBA is an exponentially-sized automaton [45]. The size of the state space of LDGBA is typically similar to that of Büchi automaton and significantly smaller than that of Rabin automaton because it mitigates the state space expansion to some extent. The state space of Rabin automaton tends to be larger due to the complexity of their acceptance conditions, which typically need to handle LTL to Rabin automaton transitions on the double exponential level. Besides, LDGBA has more expressive power compared to Büchi automaton and is able to handle more complex LTL tasks. While LDGBA has a smaller state space than Rabin automaton and can express most LTL expressions. Moreover, the model-free RL-based approach is adopted and can synthesize optimal policies on-the-fly without explicitly memorizing structures of relaxed product MDP.

**Contributions:** Motivated by these challenges, this paper considers the transition probabilities of interactions between the environment and the mobile robot to be unknown. We study learning-based motion planning subject to uncertainties, where control objectives are defined as high-level LTL formulas to express complex tasks. The contributions of this work are multi-fold. (i). From the aspect of automaton theory, we leverage an embedded limit deterministic generalized Büchi automaton (E-LDGBA) that has several accepting sets to maintain a dense

reward and allows applying deterministic policies to achieve high-level objectives. (ii) Both motion and workspace uncertainties are considered to be unknown, leading to potentially conflicting tasks (i.e., the pre-specified LTL tasks cannot be fully satisfied). We design a relaxed product MDP from the PL-MDP and the novel automaton so that the RL agent can revise its motion plan without strictly following the desired LTL constraints. (iii) An expected return is defined, incorporating both violation rewards and acceptance rewards. The violation function is designed to measure the deviation between the adjusted and intended motion plans, while the acceptance rewards promote fulfillment of the acceptance condition in the relaxed product MDP. (iv) Rigorous analysis is provided to show how to properly design parameters of accepting and violating rewards for multi-objective RL (MORL). Based on that, the RL agent can find policies that fulfill pre-specified tasks as much as possible.

## 2. Preliminaries

### 2.1. Probabilistic labeled MDP

In real-world scenarios, uncertainties in robot motion and environmental characteristics can significantly impact a robotic system's performance and reliability. For instance, motion uncertainties may result in deviations from planned paths, necessitating robust corrective strategies, while environmental uncertainties can affect robots' perception and decision-making. Markov decision processes are often used to model these uncertain interactions between robots and the environment and are defined as follows:

**Definition 1.** A probabilistic labeled Markov decision process (PL-MDP) is defined as a tuple  $\mathcal{M} = (S, A, p_S, (s_0, l_0), \Pi, L, p_L, A)$ , where  $S$  represents a finite set of states, and  $A$  is the finite set of actions. The transition probability function  $p_S \times A \times S \rightarrow [0, 1]$  governs the likelihood of moving between states based on actions.  $\Pi$  denotes a set of atomic propositions, with  $L \rightarrow 2^\Pi$  serving as the labeling function. Let  $\xi$  be an action function, which can be either deterministic such that  $\xi : S \rightarrow A$  maps a state  $s \in S$  to an action in  $A(s)$ , or randomized such that  $\xi : S \times A \rightarrow [0, 1]$  represents the probability of taking an action in  $A(s)$  at  $s$ . The pair  $(s_0, l_0)$  denotes an initial state  $s_0 \in S$  with an initial label  $l_0 \in L(s_0)$ . The function  $p_L(s, l)$  denotes the probability of  $l \in L(s)$  associated with  $s \in S$  satisfying  $\sum_{l \in L(s)} p_L(s, l) = 1, \forall s \in S$ . The transition probability  $p_S$  captures the motion uncertainties of the agent while the labeling probability  $p_L$  captures the environment uncertainties.

It is assumed that the agent can fully observe its current state and the associated labels. The PL-MDP can be regarded as an advanced MDP, and we can use the MDP to denote it in the following sections. The MDP  $\mathcal{M}$  evolves by taking actions  $\xi_i$  at each stage  $i$ , where  $i \in \mathbb{N}_0$  with  $\mathbb{N}$  being the set of natural numbers.

**Definition 2.** A control policy  $\mu = \mu_0 \mu_1 \dots$  is a sequence of decision rules that generates a path  $s = s_0 s_1 s_2 \dots$  over  $\mathcal{M}$ , with  $p_S(s_i, a_i, s_{i+1}) > 0$  for all  $i$ . The policy  $\mu$  can be either deterministic, where  $\mu : S \rightarrow A$ , or stochastic, where  $\mu : S \times A \rightarrow [0, 1]$ . The control policy  $\mu$  is termed memoryless if each  $\mu_i$  depends only on the current state  $s_i$ . Conversely,  $\mu$  is called a finite-memory, or history-dependent, policy if each  $\mu_i$  considers past states.

This work shows how to apply the deterministic and memoryless policy that has more stable decision-making performance. Let  $\mu(s)$  denote the probability distribution of actions at state  $s$ , and  $\mu(s, a)$  represents the probability of generating action  $a$  at state  $s$  using the policy  $\mu$ . Let  $A : S \rightarrow \mathbb{R}$  denote a reward function over  $\mathcal{M}$ . Given a discount factor  $\gamma \in (0, 1)$ , the expected return under policy  $\xi$  starting from  $s \in S$  can be defined as

$$U^\xi(s) = \mathbb{E}^\xi \left[ \sum_{i=0}^{\infty} \gamma^i A(s_i) \mid s_0 = s \right] \quad (1)$$

The optimal policy  $\xi^*$  is a policy that maximizes the expected return for each state  $s \in S$  as  $\xi^* = \arg \max_{\xi} U^{\xi}(s)$ .

**Definition 3.** Given a PL-MDP  $\mathcal{M}$  under policy  $\pi$ , a Markov chain  $MC_{\mathcal{M}}^{\mu}$  of the PL-MDP  $\mathcal{M}$  induced by a policy  $\mu$  is a tuple  $(S, A, p_S^{\mu}, (s_0, l_0), L, p_L)$  where  $p_S^{\mu}(s, s') = p_S(s, a, s')$  with  $\mu(s, a) > 0$  for all  $s, s' \in S$ .

A sub-MDP  $\mathcal{M}_{(S', A')}$  of  $\mathcal{M}$  is a pair  $(S', A')$  where  $S' \subseteq S$  and  $A'$  is a finite action space of  $S'$  such that (i)  $S' \neq \emptyset$ , and  $A'(s) \neq \emptyset, \forall s \in S'$ ; (ii)  $\{s' \mid \forall s \in S' \text{ and } \forall a \in A'(s), p^P(s, a, s') > 0\} \subseteq S'$ . An induced graph of  $\mathcal{M}_{(S', A')}$  is denoted as  $\mathcal{G}_{(S', A')}$  that is a directed graph, where if  $p_S(s, a, s') > 0$  with  $a \in A'(s)$ , for any  $s, s' \in S'$ , there exists an edge between  $s$  and  $s'$  in  $\mathcal{G}_{(S', A')}$ . Note the evolution of a sub-MDP  $\mathcal{M}_{(S', A')}$  is restricted by the action space  $A'$ . A sub-MDP qualifies as a strongly connected component (SCC) if its induced graph is strongly connected, meaning that for any pair of nodes  $s, s' \in S'$ , a path exists from  $s$  to  $s'$ . A bottom strongly connected component (BSCC) is an SCC from which no external states can be reached under the restricted action space. For further details on MDP treatments, please refer to [1].

**Definition 4 ([1]).** A sub-MDP  $\mathcal{M}_{(S', A')}$  is an end component (EC) of  $\mathcal{M}$  if it is a BSCC. An EC  $\mathcal{M}_{(S', A')}$  is called a maximal end component (MEC) if there is no other EC  $\mathcal{M}_{(S'', A'')}$  such that  $S' \subseteq S''$  and  $A'(s) \subseteq A''(s), \forall s \in S$ .

## 2.2. LTL and limit-deterministic generalized Büchi automaton

Linear temporal logic (LTL) is a formal language used to specify high-level requirements for a system. An LTL formula is constructed from a set of atomic propositions, e.g.,  $a \in AP$ , standard Boolean operators like  $\wedge$  (and),  $\neg$  (not), along with temporal operators:  $\mathcal{U}$  (until),  $\bigcirc$  (next),  $\diamond$  (eventually), and  $\square$  (always). The syntax of an LTL formula is defined inductively as follows:

$$\phi := \text{True} \mid a \mid \phi_1 \wedge \phi_2 \mid \neg \phi \mid \phi_1 \mathcal{U} \phi_2.$$

The semantics of an Linear Temporal Logic formula are evaluated over words, which are infinite sequences  $o = o_0 o_1 \dots$ , where  $o_i \in 2^{AP}$  for all  $i \geq 0$ . Here,  $2^{AP}$  denotes the power set of  $AP$ , defined as follows:

$$\begin{aligned} o &\models \text{true} \\ o &\models a &\Leftrightarrow a \in L(o[0]) \\ o &\models \phi_1 \wedge \phi_2 &\Leftrightarrow o \models \phi_1 \text{ and } o \models \phi_2 \\ o &\models \neg \phi &\Leftrightarrow o \not\models \phi \\ o &\models \phi &\Leftrightarrow o[1:] \models \phi \\ o &\models \phi_1 \mathcal{U} \phi_2 &\Leftrightarrow \exists t \text{ s.t. } o[t:] \models \phi_2, \forall t' \in [0, t), o[t':] \models \phi_1 \end{aligned}$$

Denote by  $o \models \phi$  if the word  $o$  satisfies the LTL formula  $\phi$ . More expressions can be achieved by combining temporal and Boolean operators. Detailed descriptions of the syntax and semantics of LTL can be found in [1]. Given an LTL formula that specifies the missions, the satisfaction of the LTL formula can be evaluated by an LDGBA [45].

**Definition 5.** A GBA is a tuple  $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ , where  $Q$  is a finite set of states,  $\Sigma = 2^I$  is a finite alphabet;  $\delta: Q \times \Sigma \rightarrow 2^Q$  is a transition function,  $q_0 \in Q$  is an initial state, and  $F = \{F_1, F_2, \dots, F_f\}$  is a set of acceptance conditions with  $F_i \subseteq Q, \forall i \in \{1, \dots, f\}$ .

Denote by  $q = q_0 q_1 \dots$  a run of a GBA, where  $q_i \in Q, i = 0, 1, \dots$ . The run  $q$  is accepted by the GBA, if it satisfies the generalized Büchi accepting sets, i.e.,  $\inf(q) \cap F_i \neq \emptyset, \forall i \in \{1, \dots, f\}$ , where  $\inf(q)$  denotes the set of states that is visited infinitely often.

**Definition 6.** A Generalized Büchi Automaton (GBA) is called a Limit-deterministic Generalized Büchi Automaton (LDGBA) if the transition function  $\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$  is defined, and the state set  $Q$  can be partitioned into a deterministic set  $Q_D$  and a non-deterministic set  $Q_N$ , such that  $Q = Q_D \cup Q_N$  and  $Q_D \cap Q_N = \emptyset$ , where:

- State transitions within  $Q_D$  are total and restricted to  $Q_D$ , i.e.,  $|\delta(q, \alpha)| = 1$  and  $\delta(q, \alpha) \subseteq Q_D$  for every state  $q \in Q_D$  and input  $\alpha \in \Sigma$ ,
- $\epsilon$ -transitions are only allowed from states in  $Q_N$  to states in  $Q_D$  and are not permitted within  $Q_D$ , i.e., for any  $q \in Q_D, \delta(q, \epsilon) = \emptyset$ ,
- All accepting states are contained within  $Q_D$ , i.e.,  $F_i \subseteq Q_D$  for every  $F_i \in F$ .

In Definition 6, the  $\epsilon$ -transitions are only defined for state transitions from  $Q_N$  to  $Q_D$  that do not consume the input atomic proposition. Readers are referred to [46] for algorithms with free implementations to convert an LTL formula to an LDGBA. Note that the state-based LDGBA is used in this work for demonstration purposes. Transition-based LDGBA can be constructed based on basic graph transformations. More discussions about state-based and transition-based LDGBA in HOA format can be found in Owl [47].

**Remark 1.** Unlike the widely used deterministic Rabin Automaton (DRA), LDGBA has the Generalized Büchi Accepting condition that purely involves reachability problems. On the other hand, compared with Limit-Deterministic Büchi Automaton (LDBA) applied in [30,31], LDGBA has more accepting sets shown in Fig. 1(a) to increase the density of rewards since the positive rewards are always assigned to the accepting states to enforce the acceptance condition.

## 3. Problem statement and challenge

The task specification to be performed by the agent is described by an LTL formula  $\phi$  over  $\Pi$ . Given  $\phi, \mathcal{M}$ , and  $\xi = \xi_0 \xi_1 \dots$ , the induced infinite path is denoted by  $s_{\infty}^{\xi} = s_0 \dots s_i s_{i+1} \dots$  that satisfies  $s_{i+1} \in \{s \in S \mid p_S(s_i, a_i, s) > 0\}$ . Let  $L(s_{\infty}^{\xi}) = l_0 l_1 \dots$  be the sequence of labels associated with  $s_{\infty}^{\xi}$  such that  $l_i \in L(s_i)$  and  $p_L(s_i, l_i) > 0$ . Denote by  $L(s_{\infty}^{\xi}) \models \phi$  if  $s_{\infty}^{\xi}$  satisfies  $\phi$ . The satisfaction probability under  $\mu$  from an initial state  $s_0$  can be computed as

$$\Pr_{\mathcal{M}}^{\mu}(\phi) = \Pr_{\mathcal{M}}^{\mu}(s_{\infty}^{\mu} \in S_{\infty}^{\mu} \mid L(s_{\infty}^{\mu}) \models \phi), \quad (2)$$

where  $S_{\infty}^{\mu}$  is a set of all admissible paths under policy  $\mu$ , and the computation of  $\Pr_{\mathcal{M}}^{\mu}(s_{\infty}^{\mu})$  can be found in [1].

**Definition 7.** Given a PL-MDP  $\mathcal{M}$ , an LTL task  $\phi$  is fully feasible if and only if  $\Pr_{\mathcal{M}}^{\mu}(\phi) > 0$  s.t. there exists a path  $s_{\infty}^{\mu}$  over the infinite horizons under the policy  $\mu$  satisfying  $\phi$ .

Based on Definition 7, an infeasible case means there does not exist any policy  $\mu$  to satisfy the task i.e.,  $\Pr_{\mathcal{M}}^{\mu}(\phi) = 0$ .

**Example 1.** To illustrate an infeasible case, the environment, its abstracted PL-MDP, and the LDGBA of an LTL formula  $\phi = \square \Diamond a \wedge \square \Diamond b$  are shown in Fig. 1(a) and (b), respectively. Since obstacles surround region  $s_3$ , the given LTL task is infeasible in this case.

We define the expected discount violation cost of task satisfaction as follows, considering feasible and infeasible tasks.

**Definition 8.** Given a PL-MDP  $\mathcal{M}$  and an LTL task  $\phi$ , the discount expected violation under the policy  $\mu$  is defined as

$$J_V(\mathcal{M}^{\mu}, \phi) = \mathbb{E}_{\mathcal{M}}^{\mu} \left[ \sum_{i=0}^{\infty} \gamma^i c_V(s_i, a_i, s_{i+1}, \phi) \right], \quad (3)$$

where  $c_V(s, a, s', \phi)$  is defined as the violation cost of a transition  $(s, a, s')$  with respect to  $\phi$ , and  $a_i$  is the action generated based on the policy  $\mu(s_i)$ .

Most existing results (cf. [26,30,31]) assume the existence of at least one policy  $\mu$  s.t.  $\Pr_{\mathcal{M}}^{\mu}(\phi) > 0$ , which may not be true in practice if the task is only partially feasible. In addition, previous work [18] addresses the infeasible cases via an optimization approach that assumes the

transitions probabilities of  $\mathcal{M}$  are known. We remove this assumption and apply the RL approach to learn the desired policies in this work. We consider the following problem to account for challenges of both infeasible cases and unknown transitions.

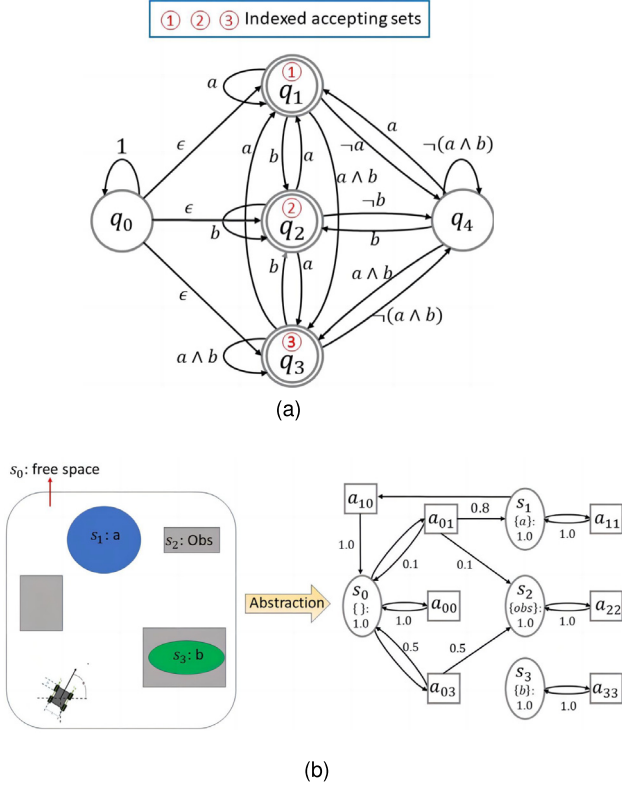


Fig. 1. A running example of an infeasible LTL task with respect to the PL-MDP. (a) The LDGBA of an LTL formula  $\phi = \square \Diamond a \wedge \square \Diamond b$ . (b) The abstraction process from a robotic system into a PL-MDP  $\mathcal{M}$ , where the region labeled with  $b$  is surrounded by obstacles and cannot be accessed. As a result,  $\phi$  is infeasible in this case.

**Problem 1.** Given an LTL-specified task  $\phi$  and a PL-MDP with unknown transition probabilities (representing motion uncertainties) and unknown probabilistic label functions (representing workspace uncertainties), the objective is to determine a policy that satisfies the following conditions: (1) if feasible,  $\Pr_{\mathcal{M}}^{\mu}(\phi) > 0$ ; (2) if  $\Pr_{\mathcal{M}}^{\mu}(\phi) = 0$  for all policies  $\mu$ , identify the policy  $\mu$  that maximally fulfills the task  $\phi$  (i.e., under infeasible constraints) over an infinite horizon by minimizing  $J_V(\mathcal{M}^{\mu}, \phi)$ .

#### 4. Automaton analysis

To solve Problem 1, Section 4.1 first presents how the LDGBA in Definition 6 can be extended to E-LDGBA to keep tracking the non-visited accepting sets. Section 4.2 shows the traditional approach and the corresponding challenges for infeasible tasks. Section 4.3 presents the construction of a relaxed product MDP to handle soft LTL constraints. The benefits of incorporating E-LDGBA with relaxed product MDP are discussed in Section 4.4.

##### 4.1. E-LDGBA

To identify the desired policy in PL-MDP  $\mathcal{M}$  that satisfies the LTL formula  $\phi$  specified by the user, one can construct the standard product MDP between  $\mathcal{M}$  and the LDGBA of  $\phi$  as outlined in [1]. This approach reformulates the problem as finding a policy that meets the acceptance

condition of the resulting product MDP. However, using the LDGBA directly may lead to unsatisfied LTL specifications when applying deterministic policies, as discussed in our conference paper [36]. To address these limitations, we introduce the E-LDGBA as follows.

Given an LDGBA  $\mathcal{A} = (Q, \Sigma, \delta \cup \{\epsilon\}, q_0, F)$ , and following [41], a tracking-frontier set  $T$  is introduced to monitor unvisited accepting sets. Initially,  $T$  is set to  $F$  and is updated according to

$$f_V(q, T) = \begin{cases} T \setminus F_j, & \text{if } q \in F_j \text{ and } F_j \in T, \\ F \setminus F_j, & \text{if } q \in F_j \text{ and } T = \emptyset, \\ T, & \text{otherwise.} \end{cases} \quad (4)$$

Each time an accepting set  $F_j$  is visited, it is removed from  $T$ . When  $T$  is empty, it resets to  $F \setminus F_j$ . Since the LDGBA acceptance condition requires infinite visits to all accepting sets, we define a “round” as a complete visit to all sets (i.e., a round concludes when  $T$  becomes empty). If a state  $q$  belongs to multiple sets in  $T$ , all such sets are removed from  $T$  upon visiting  $q$ .

**Definition 9 (Embedded LDGBA).** Given an LDGBA  $\mathcal{A} = (Q, \Sigma, \delta \cup \{\epsilon\}, q_0, F)$ , the corresponding E-LDGBA is denoted by  $\bar{\mathcal{A}} = (\bar{Q}, \Sigma, \bar{\delta} \cup \{\epsilon\}, \bar{q}_0, \bar{F}, f_V, T)$ , where:

- $T$  is initially set as  $T = F$ .
- $\bar{Q} = Q \times 2^F$  is the set of augmented states, e.g.,  $\bar{q} = (q, T)$ . The finite alphabet  $\Sigma$  is the same as that in the LDGBA.
- The transition  $\bar{\delta}: \bar{Q} \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^{\bar{Q}}$  is defined by  $\bar{q}' = \bar{\delta}(\bar{q}, \sigma)$  with  $\bar{\sigma} \in (\Sigma \cup \{\epsilon\})$ , where  $\bar{q} = (q, T)$  and  $\bar{q}' = (q', T')$ , subject to two conditions: (1)  $q' = \delta(q, \bar{\sigma})$ , and (2)  $T$  is updated as  $T' = f_V(q', T)$  following the transition  $q' = \delta(q, \bar{\sigma})$ .
- $\bar{F} = \{\bar{F}_1, \bar{F}_2, \dots, \bar{F}_f\}$  is the set of accepting states, where  $\bar{F}_j = \{(q, T) \in \bar{Q} \mid q \in F_j \wedge F_j \subseteq T\}$  for  $j = 1, \dots, f$ .

In Definition 9, we adjust the tuple structure such that the frontier set  $T$  is synchronously updated after each transition. The state space is augmented with the tracking frontier set  $T$ , which can be efficiently represented using one-hot encoding based on the indices of the accepting sets. The accepting state is determined by both the current automaton state and the frontier set  $T$ . This feature is a key innovation of the E-LDGBA, ensuring that all accepting sets are visited in each round. Alg. 1 demonstrates the procedure of obtaining a valid run  $q_{\bar{\mathcal{A}}}$  over E-LDGBA  $\bar{\mathcal{A}}_{\phi}$ . Given an input alphabet  $\alpha$  at each time step, lines 6–10 show how to generate the next state of  $\bar{\mathcal{A}}_{\phi}$  and update the tracking frontier set  $T$  synchronously. Also, lines 6–8 indicate the difference and relationship between  $\bar{\mathcal{A}}_{\phi}$  and its corresponding LDGBA  $\mathcal{A}_{\phi}$ .

Given  $\bar{\mathcal{A}}_{\phi}$  and  $\mathcal{A}_{\phi}$  for the same LTL formula, the E-LDGBA  $\bar{\mathcal{A}}_{\phi}$  tracks unvisited accepting sets of  $\mathcal{A}_{\phi}$  through the inclusion of  $f_V$  and  $T$ , allowing it to distinguish and enforce acceptance satisfaction related to different accepting sets. The set  $T$  is reset once all accepting sets of  $\mathcal{A}_{\phi}$  have been visited. Let  $\mathcal{L}(\mathcal{A}_{\phi}) \subseteq \Sigma^{\omega}$  and  $\mathcal{L}(\bar{\mathcal{A}}_{\phi}) \subseteq \Sigma^{\omega}$  represent the languages accepted by  $\mathcal{A}_{\phi}$  and  $\bar{\mathcal{A}}_{\phi}$ , respectively, over the same alphabet  $\Sigma$ . According to [1],  $\mathcal{L}(\mathcal{A}_{\phi}) \subseteq \Sigma^{\omega}$  is the set of all infinite words accepted by  $\mathcal{A}_{\phi}$  that satisfy the LTL formula  $\phi$ .

**Lemma 1.** For any LTL formula  $\phi$ , we can construct LDGBA  $\mathcal{A}_{\phi} = (Q, \Sigma, \delta, q_0, F)$  and E-LDGBA  $\bar{\mathcal{A}}_{\phi} = (\bar{Q}, \Sigma, \bar{\delta}, \bar{q}_0, \bar{F}, f_V, T)$ . Then it holds that

$$\mathcal{L}(\bar{\mathcal{A}}_{\phi}) = \mathcal{L}(\mathcal{A}_{\phi}). \quad (5)$$

**Proof.** Details of the proof can be found in our work [36] that has shown  $\mathcal{L}(\bar{\mathcal{A}}_{\phi}) \supseteq \mathcal{L}(\mathcal{A}_{\phi})$  and  $\mathcal{L}(\bar{\mathcal{A}}_{\phi}) \subseteq \mathcal{L}(\mathcal{A}_{\phi})$ . This indicates that the E-LDGBA accepts the same language as the LDGBA and can be used to verify LTL satisfaction.  $\square$



**Algorithm 1** Synthesizing a valid run of E-LDGBA

---

```

1: procedure INPUT: (LDGBA  $\mathcal{A}$ ,  $f_V, T$  and length  $L$ )
   Output: A valid run  $\bar{q}_A$  with length  $L$  in  $\mathcal{A}_\phi$ 
2: set  $T = F$  and  $count = 1$ 
3: set  $\bar{q}_{cur} = (q_0, T)$  and  $\bar{q}_A = (q_{cur})$ 
4: set  $q_{cur} = q$ 
5: while  $count \leq L$  do
6:   given an input  $\alpha$ ,  $q_{next} = \delta(q_{cur}, \alpha)$ 
7:   obtain  $\bar{q}_{cur} \leftarrow (q_{next}, T)$ 
8:   check if  $\bar{q}_{cur}$  is an accepting state
9:   add state  $\bar{q}_{cur}$  to  $\bar{q}_A$ 
10:  update  $T = f_V(q_{next}, T)$ 
11:   $count++$  and  $q_{cur} \leftarrow q_{next}$ 
12: end while
13: end procedure

```

---

**4.2. Product MDP and challenges**

To find a policy  $\mu$  satisfying  $\phi$  i.e.,  $\Pr_{\mathcal{M}}^{\mu}(\phi) > 0$ , one can construct the standard product MDP as follows.

**Definition 10.** Given the labeled MDP  $\mathcal{M}$  and the E-LDGBA  $\bar{\mathcal{A}}_\phi$  corresponding to  $\phi$ , the product MDP is defined as  $\bar{\mathcal{P}} = (Y, U^{\bar{\mathcal{P}}}, p^{\bar{\mathcal{P}}}, x_0, F^{\bar{\mathcal{P}}})$ , where

- $Y = S \times 2^I \times \bar{Q}$  is the set of labeled states, i.e.,  $y = (s, l, \bar{q}) \in Y$  with  $l \in L(s)$  satisfying  $p_L(s, l) > 0$ ;
- $x_0 = (s_0, l_0, \bar{q}_0)$  is the initial state;
- $U^{\bar{\mathcal{P}}} = A \cup \{\epsilon\}$  is the set of actions, where the  $\epsilon$ -actions are only allowed for transitions of E-LDGBA components from  $\bar{Q}_N$  to  $\bar{Q}_D$ ;
- $F^{\bar{\mathcal{P}}} = \{\bar{F}_1^{\bar{\mathcal{P}}}, \bar{F}_2^{\bar{\mathcal{P}}}, \dots, \bar{F}_f^{\bar{\mathcal{P}}}\}$  is the set of accepting states. where  $\bar{F}_j^{\bar{\mathcal{P}}} = \{(s, l, \bar{q}) \in X \mid \bar{q} \in \bar{F}_j\}$ ,  $j = 1, \dots, f$ ;
- $p^{\bar{\mathcal{P}}} : Y \times U^{\bar{\mathcal{P}}} \times Y \rightarrow [0, 1]$  is transition probability defined as: (1)  $p^{\bar{\mathcal{P}}}(y, u^{\bar{\mathcal{P}}}, y') = p_L(s', l') \cdot p_S(s, a, s')$  if  $\delta(\bar{q}, l) = \bar{q}'$  and  $u^{\bar{\mathcal{P}}} = a \in A(s)$ ; (2)  $p^{\bar{\mathcal{P}}}(y, u^{\bar{\mathcal{P}}}, y') = 1$  if  $u^{\bar{\mathcal{P}}} \in \{\epsilon\}$ ,  $\bar{q}' \in \delta(\bar{q}, \epsilon)$ , and  $(s', l') = (s, l)$ ; and (3)  $p^{\bar{\mathcal{P}}}(y, u^{\bar{\mathcal{P}}}, y') = 0$  otherwise.

The product MDP  $\bar{\mathcal{P}}$  represents the intersections between all feasible paths in  $\mathcal{M}$  and all words accepted by  $\mathcal{A}_\phi$ , thereby enabling the identification of admissible agent motions that fulfill the task  $\phi$ . Let  $\pi$  denote a policy over  $\bar{\mathcal{P}}$ , and let  $y_\infty^\pi = y_0 y_1 y_2 \dots$  represent the infinite path generated by  $\pi$ . A path  $y_\infty^\pi$  is accepted if  $\inf(y_\infty^\pi) \cap \bar{F}_i^{\bar{\mathcal{P}}} \neq \emptyset$  for all  $i \in \{1, \dots, f\}$ . If  $y_\infty^\pi$  is an accepting run, then there exists a policy  $\xi$  in  $\mathcal{M}$  that satisfies  $\phi$ .

Consider a sub-product MDP  $\bar{\mathcal{P}}'_{(Y', U')}$ , where  $Y' \subseteq X$  and  $U' \subseteq U^{\bar{\mathcal{P}}}$ . If  $\bar{\mathcal{P}}'_{(Y', U')}$  is a MEC of  $\bar{\mathcal{P}}$  and  $Y' \cap \bar{F}_i^{\bar{\mathcal{P}}} \neq \emptyset$ ,  $\forall i \in \{1, \dots, f\}$ , then  $\bar{\mathcal{P}}'_{(Y', U')}$  is called an accepting maximum end component (AMEC) of  $\bar{\mathcal{P}}$ . Once a path enters an AMEC, The subsequent path remains within the component by taking restricted actions from  $U'$ . There exist policies allowing any state  $x \in X'$  to be visited infinitely often. Consequently, satisfying the task  $\phi$  is equivalent to reaching an AMEC. Furthermore, a MEC that lacks any accepting sets is termed a rejecting maximum end component (RMEC), while a MEC containing only some accepting sets is referred to as a neutral maximum end component (NMEC).

**4.3. Relaxed product MDP**

For the product MDP  $\bar{\mathcal{P}}$  or  $\bar{\mathcal{P}}$  introduced above, the satisfaction of  $\phi$  is based on the assumption that at least one AMEC exists in the product MDP, i.e., at least one policy satisfies the given LTL formula with respect to the PL-MDP. Otherwise, the task is infeasible with respect to the PL-MDP. We treat the LTL task as soft constraints to address the cases where tasks can be potentially infeasible. The relaxed product

MDP is designed to allow the agent to revise its motion plan without strictly following the desired LTL constraints.

**Definition 11.** The relaxed product MDP is constructed from  $\bar{\mathcal{P}}$  as a tuple  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi = (X, U^{\bar{\mathcal{R}}}, p^{\bar{\mathcal{R}}}, x_0, F^{\bar{\mathcal{R}}}, c_V^{\bar{\mathcal{R}}}, f_V, T)$ , where

- $X, x_0$ , and  $F^{\bar{\mathcal{R}}} = F^{\bar{\mathcal{P}}}$  are the same as in  $\bar{\mathcal{P}}$ ;
- $U^{\bar{\mathcal{R}}}$  is the set of extended actions that jointly consider the actions of  $\mathcal{M}$  and the input atomic proposition of  $\mathcal{A}_\phi$ . Specifically, given a state  $x = (s, l, q) \in X$ , the available actions are  $u^{\bar{\mathcal{R}}}(x) = \{(a, i) \mid a \in A(s), i \in (2^{AP} \cup \{\epsilon\})\}$ . Given an action  $u^{\bar{\mathcal{R}}} = (a, i) \in u^{\bar{\mathcal{R}}}(x)$ , the projections of  $u^{\bar{\mathcal{R}}}$  to  $A(s)$  in  $\mathcal{M}$  and to  $2^{AP} \cup \{\epsilon\}$  in  $\mathcal{A}_\phi$  are denoted by  $u|_{\mathcal{M}}^{\bar{\mathcal{R}}}$  and  $u|_{\mathcal{A}}^{\bar{\mathcal{R}}}$ , respectively;
- $p^{\bar{\mathcal{R}}} : X \times u^{\bar{\mathcal{R}}} \times X \rightarrow [0, 1]$  is the transition function. The transition probability  $p^{\bar{\mathcal{R}}}$  from a state  $x = (s, l, \bar{q})$  to a state  $x' = (s', l', \bar{q}')$  is defined as: (1)  $p^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x') = p_L(s', l') \cdot p_S(s, a, s')$  with  $a = u|_{\mathcal{M}}^{\bar{\mathcal{R}}}$ , if  $\bar{q}$  can be transitioned to  $\bar{q}'$  and  $\bar{q}' = \delta(\bar{q}, u|_{\mathcal{A}}^{\bar{\mathcal{R}}})$  s.t.  $u|_{\mathcal{A}}^{\bar{\mathcal{R}}} \neq \epsilon$ ; (2)  $p^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x') = 1$ , if  $u|_{\mathcal{A}}^{\bar{\mathcal{R}}} = \epsilon$ ,  $\bar{q}' \in \delta(\bar{q}, \epsilon)$ , and  $(s', l') = (s, l)$ ; (3)  $p^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x') = 0$  otherwise. Under an action  $u^{\bar{\mathcal{R}}} \in U^{\bar{\mathcal{R}}}(x)$ , it holds that  $\sum_{x' \in X} p^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x') = 1$ ;
- $c_V^{\bar{\mathcal{R}}} : X \times u^{\bar{\mathcal{R}}} \times X \rightarrow \mathbb{R}$  is the violation cost. To define the violation cost  $c_V^{\bar{\mathcal{R}}}$ , consider  $\Pi = \{i_1, i_2, \dots, i_M\}$  and an evaluation function  $\text{Eval} : 2^I \rightarrow \{0, 1\}^M$ , defined as  $\text{Eval}(l) = [v_i]^M$  with  $v_i = 1$  if  $i_i \in l$  and  $v_i = 0$  otherwise, where  $i = 1, 2, \dots, M$  and  $l \in 2^I$ . To quantify the difference between two elements in  $2^I$ , we introduce  $\rho(l, l') = \|v - v'\|_1 = \sum_{i=1}^M |v_i - v'_i|$ , where  $v = \text{Eval}(l)$ ,  $v' = \text{Eval}(l')$ ,  $l, l' \in 2^I$ , and  $\|\cdot\|_1$  is the  $\ell_1$  norm. The distance from  $l \in 2^I$  to a set  $\mathcal{X} \subseteq 2^I$  is then defined as  $\text{Dist}(l, \mathcal{X}) = 0$  if  $l \in \mathcal{X}$ , and  $\text{Dist}(l, \mathcal{X}) = \min_{l' \in \mathcal{X}} \rho(l, l')$  if  $l \notin \mathcal{X}$ . The violation cost of the transition from  $x = (s, l, q)$  to  $x' = (s', l', q')$  under an action  $u^{\bar{\mathcal{R}}}$  is defined as

$$c_V^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x') = \begin{cases} p_L(s', l') \cdot w_V(x, x') & \text{if } u|_{\mathcal{A}}^{\bar{\mathcal{R}}} \neq \epsilon, \\ 0 & \text{otherwise,} \end{cases}$$

where  $w_V(x, x') = \text{Dist}(L(s), \mathcal{X}(q, q'))$  with

$\mathcal{X}(q, q') = \{l \in 2^I \mid q \xrightarrow{l} q'\}$  being the set of input alphabets that enables the transition from  $q$  to  $q'$ . Borrowed from [35], the function  $\text{Dist}(L(s), \mathcal{X}(q, q'))$  measures the distance from  $L(s)$  to the set  $\mathcal{X}(q, q')$ .

Alg. 2 illustrates executing a given policy and generating a valid run over the relaxed product MDP  $\bar{\mathcal{R}}$  on-the-fly. “On-the-fly” means the procedure does not need to construct  $\bar{\mathcal{R}}$  completely. In contrast, it tracks the states based on the evolution in Definition 11. In line 8 of Alg. 2, the next automaton state of  $\bar{\mathcal{A}}_\phi$  is reached according to lines 6 – 7 of Alg. 1 by assigning the input alphabet as  $u|_{\mathcal{A}}^{\bar{\mathcal{R}}}$ . In lines 5–9, the next state is jointly determined by the transitions of PL-MDP and E-LDGBA (Alg. 1). Line 14 of Alg. 2 also updates the tracking frontier set synchronously as Alg. 1, and line 11 computes its violation cost after generating the valid transition.

The weighted violation function  $w_V(x, x')$  quantifies how much the transition from  $x$  to  $x'$  in a product MDP violates the constraints imposed by  $\phi$ . The transition probability  $p^{\bar{\mathcal{R}}}$  jointly considers the probability of an event  $p_L(s', l')$  and the transition  $p_S$ . Then the measurement of violation in (3) can be reformulated over relaxed product MDP as the expected return of violation cost:

**Definition 12.** Given a relaxed product MDP  $\bar{\mathcal{R}}$  generated from a PL-MDP  $\mathcal{M}$  and an E-LDGBA  $\bar{\mathcal{A}}_\phi$ , the measurement of expected violation in (3) can be reformulated over relaxed product MDP under policy  $\pi$

**Algorithm 2** Policy Execution over  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi$ 


---

```

1: procedure INPUT: (policy  $\pi$ ,  $\mathcal{M}$ ,  $\bar{\mathcal{A}}_\phi$ ,  $N$ )
   Output: a valid run  $\mathbf{x}_N$  with length of  $N$  under policy  $\pi$  over  $\bar{\mathcal{R}}$ 
2: set  $t = 0$  and  $x_t = x_0$ 
3: set  $\mathbf{x}_N = (x_t)$ 
4: while  $t < N$  do
5:   obtain  $u_t^{\bar{\mathcal{R}}}$  based on  $\pi(x_t)$ 
6:   execute  $u_t^{\bar{\mathcal{R}}}$  for  $\mathcal{M}$  and obtain  $s_{t+1}$ 
7:   observe label  $l_{t+1}$  according to  $L(s_t)$ 
8:   obtain  $q_{t+1} = \delta(q_t, u_t^{\bar{\mathcal{R}}})$  s.t.  $\bar{q}_{t+1} = (q_{t+1}, T)$ 
9:   obtain  $x_{t+1} = (s_{t+1}, l_{t+1} q_{t+1}, T)$ 
10:  check if the  $x_{t+1}$  is an accepting state
11:  compute  $c_V^{\bar{\mathcal{R}}}(x, u^{\bar{\mathcal{R}}}, x')$  based on Definition 11
12:  add  $x_{t+1}$  to  $\mathbf{x}_N$ 
13:   $x_t = x_{t+1}$  and  $t++$ 
14:  update  $T \leftarrow f_V(q', T)$  for  $\bar{\mathcal{A}}_\phi$ 
15: end while
16: end procedure

```

---

as:

$$J_V(\bar{\mathcal{R}}^\pi) = \mathbb{E}_{\bar{\mathcal{R}}^\pi} \left[ \sum_{i=0}^{\infty} \gamma^i c_V^{\bar{\mathcal{R}}}(x_i, u_i^{\bar{\mathcal{R}}}, x_{i+1}) \right] \quad (6)$$

Finding a policy to minimize  $J_V(\bar{\mathcal{R}}^\pi)$  will enforce the planned path towards more fulfilling the LTL task  $\phi$  by penalizing  $w_V$ . A run  $r_{\bar{\mathcal{R}}^\pi}^\pi$  induced by a policy  $\pi$  that satisfies the accepting condition of the relaxed product MDP  $\bar{\mathcal{R}}$  completes the corresponding LTL mission  $\phi$  exactly if and only if the violation return (6) is equal to zero.

#### 4.4. Properties of relaxed product MDP

This section discusses the properties of the relaxed product MDP and standard product MDP. Applying the same process in Definition 10, we can construct the product MDP between LDGBA and PL-MDP denoted as  $\mathcal{P} = \mathcal{M} \otimes \mathcal{A}_\phi$ . The detailed procedure can be found in [1]. The relaxed product MDP between MDP and LDGBA can also be constructed as the procedure Definition 11 by replacing E-LDGBA with LDGBA denoted as  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi$ . Note, the product and relaxed product are the procedures of constructing the interaction between the automaton structure and the MDP model, which can be applied to any automata.

Similarly,  $\mathcal{R}$  and  $\mathcal{P}$  for the same  $\mathcal{A}_\phi$  and PL-MDP  $\mathcal{M}$  share identical states. Thus,  $\mathcal{R}$  and  $\mathcal{P}$  can be viewed as two distinct directed graphs with the same set of nodes. From a graph-theoretic perspective, a MEC can be considered a bottom strongly connected component (BSCC). Similarly, let ABSCC denote the BSCC that intersects with all accepting sets in  $\mathcal{R}$  or  $\mathcal{P}$ . We then arrive at the following conclusion.

**Theorem 1.** Given an MDP  $\mathcal{M}$  and an LDGBA  $\mathcal{A}_\phi$  for the LTL formula  $\phi$ , let  $\mathcal{R} = \mathcal{M} \otimes \mathcal{A}_\phi$  denote the relaxed product MDP and  $\mathcal{P}$  the corresponding standard product MDP. The following properties hold:

1. The directed graph of the standard product MDP  $\mathcal{P}$  is a subgraph of the directed graph of  $\mathcal{R}$ .
2. There always exists at least one AMEC in  $\mathcal{R}$ .
3. If the LTL formula  $\phi$  is feasible over  $\mathcal{M}$ , then any directed graph of an AMEC in  $\mathcal{P}$  is a subgraph of a directed graph of an AMEC in  $\mathcal{R}$ .

**Proof.** The proof is similar as our previous work [18] by replacing the LDGBA with LDGBA.  $\square$

Since the E-LDGBA is an extension of LDGBA to enable the recording ability regarding accepting sets, which accept the same LTL formulas as in Lemma 1. We can also have the same conclusion for  $\bar{\mathcal{R}}$  and  $\bar{\mathcal{P}}$ .

**Lemma 2.** Given an MDP  $\mathcal{M}$  and an E-LDGBA  $\bar{\mathcal{A}}_\phi$  of  $\phi$ , the relaxed product MDP  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi$  and its corresponding  $\bar{\mathcal{P}}$  have the following properties

1. the directed graph of standard product  $\bar{\mathcal{P}}$  is a sub-graph of the directed graph of  $\bar{\mathcal{R}}$ ,
2. there always exists at least one AMEC in  $\bar{\mathcal{R}}$ ,
3. if the LTL formula  $\phi$  is feasible over  $\mathcal{M}$ , any AMEC of  $\bar{\mathcal{P}}$  is the sub-graph of an AMEC of  $\bar{\mathcal{R}}$ .

Lemma 2 can be verified by employing the proof of Theorem 1. Both of them demonstrate the advantages of applying the relaxed product MDP, which allows us to handle infeasible situations.

**Example 2.** Fig. 2(a) and (b) provide examples of a standard product MDP and the corresponding relaxed one, respectively, to illustrate the benefits of designing the relaxed product MDP. Both product MDPs are constructed between the PL-MDP and the LDGBA of LTL formula  $\phi = \square \Diamond a \wedge \square \Diamond b$ , as shown in Fig. 1. Obviously, the product MDP in Fig. 2(a) has no AMEC, while the relaxed product MDP in Fig. 2(b) has one AMEC (blue rectangular) intersecting with all accepting sets. Moreover, the transitions with nonzero violation costs are gray dashed lines. Note that due to the complicated graph structure, we use the LDGBA in this example to illustrate the novelty of the proposed relaxed product MDP. Its advantages for infeasible cases are also applicable to the E-LDGBA.

Given a relaxed product MDP  $\bar{\mathcal{R}}$ , let  $MC_{\bar{\mathcal{R}}}^\pi$  represent the Markov chain induced by policy  $\pi$  on  $\bar{\mathcal{R}}$ . The states of this Markov chain can be expressed as a disjoint union of a transient class  $T_\pi$  and  $n$  closed irreducible recurrent sets  $R_\pi^j$  for  $j \in \{1, \dots, n\}$ , i.e.,  $MC_{\bar{\mathcal{R}}}^\pi = T_\pi \sqcup R_\pi^1 \sqcup R_\pi^2 \sqcup \dots \sqcup R_\pi^n$  [48].

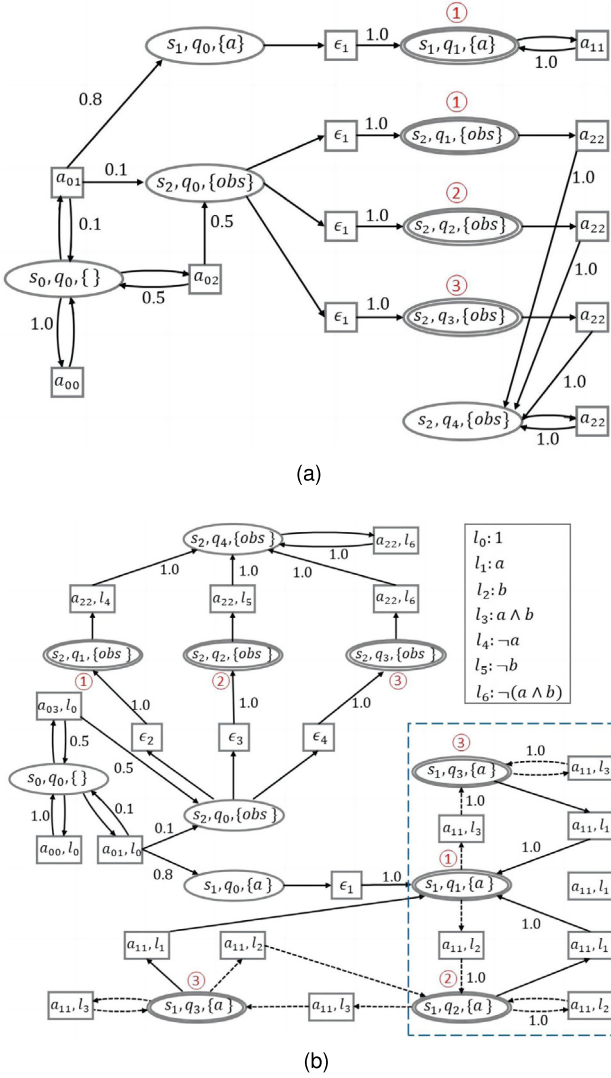
**Lemma 3.** Given a relaxed product MDP  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi$ , the recurrent class  $R_\pi^j$  of  $MC_{\bar{\mathcal{R}}}^\pi$   $\forall j \in \{1, \dots, n\}$ , induced by  $\pi$  satisfies one of the following conditions:

1.  $R_\pi^j \cap \bar{F}_i^{\bar{\mathcal{R}}} \neq \emptyset, \forall i \in \{1, \dots, f\}$ , or
2.  $R_\pi^j \cap \bar{F}_i^{\bar{\mathcal{R}}} = \emptyset, \forall i \in \{1, \dots, f\}$ .

**Proof.** The following proof is based on contradiction. Assume there exists a policy such that  $R_\pi^j \cap \bar{F}_k^{\bar{\mathcal{R}}} \neq \emptyset, \forall k \in K$ , where  $K$  is a subset of  $2^{\{1, \dots, f\}} \setminus \{\{1, \dots, f\}, \emptyset\}$ . As discussed in [48], for each state in recurrent class, it holds that  $\sum_{n=0}^{\infty} p^n(x, x) = \infty$ , where  $x \in R_\pi^j \cap \bar{F}_k^{\bar{\mathcal{R}}}$  and  $p^n(x, x)$  denotes the probability of returning from a transient state  $x$  to itself in  $n$  steps. This means that each state in the recurrent class occurs infinitely often. However, based on the embedded tracking frontier function of E-LDGBA in Definition 9, the tracking set  $T$  will not be reset until all accepting sets have been visited. As a result, neither  $\bar{q}_k \in \bar{F}_k$  nor  $x_k = (s, \bar{q}_k) \in R_\pi^j \cap \bar{F}_k^{\bar{\mathcal{R}}}$  with  $s \in S$  will occur infinitely, which contradicts the property  $\sum_{n=0}^{\infty} p^n(x_k, x_k) = \infty$ .  $\square$

Lemma 3 indicates that, for any policy, all accepting sets will be placed either in the transient class or the recurrent class. As a result, the issue of NMEC, as in many existing methods, can be avoided. Based on Theorem 1 and Lemma 3, Problem 1 can be reformulated as follows.

**Problem 2.** Given a user-specified LTL task  $\phi$  and an MDP with unknown transition probabilities (representing motion uncertainties) and unknown labeling probabilities (representing environment uncertainties), the objective is to derive a policy that prioritizes, in decreasing order: (1) satisfying the acceptance condition of the relaxed product MDP, and (2) minimizing the violation cost within the expected return.



**Fig. 2.** Two different product MDPs between the PL-MDP and the LDGBA provided in Fig. 1. (a) Standard product MDP. There are no criteria to evaluate a policy for task satisfaction, and we cannot find a path satisfying the acceptance condition. (b) Partial structure of the relaxed product MDP.  $l_i, \forall i \in \{0, 1, \dots, 7\}$  denotes the input alphabet of each automaton transition as shown in the gray rectangular. The gray dashed lines represent the transitions with nonzero violation costs, and the graph component in blue rectangular contains the path that satisfies the acceptance condition.

## 5. Learning-based control synthesis

In this section, RL is leveraged to identify policies for Problem 2. Specifically, a model-free multi-objective RL (MORL) is designed.

### 5.1. Reward design

Referring to [41], the accepting reward function  $\Lambda : X \rightarrow \mathbb{R}$  is designed as

$$\Lambda(x) = \begin{cases} r_A & \text{if } \exists i \in \{1, \dots, f\} \text{ such that } x \in F_i^{\bar{R}}, \\ 0 & \text{otherwise,} \end{cases} \quad (7)$$

where  $r_A > 0$ . The violation function  $V : X \times U^{\bar{R}} \times X \rightarrow \mathbb{R}$  is designed as

$$V(x, u^{\bar{R}}, x') = -c_V^{\bar{R}}(x, u^{\bar{R}}, x'). \quad (8)$$

The non-negative  $\Lambda(x)$  enforces the accepting condition of  $\bar{R}$ , while the non-positive function  $V(x, u^{\bar{R}}, x')$  indicates the penalty of violations. Let

$U_\pi = [U_\pi(x_0) \ U_\pi(x_1) \ \dots]^T \in \mathbb{R}^N$  denote the stacked expected return induced by  $\pi$  over  $\bar{R}$  with  $N = |X|$ , the expected return is designed as

$$U = \sum_{n=0}^{\infty} \gamma^n P_\pi^n (\Lambda_\pi + \beta V_\pi \cdot \mathbf{1}_N), \quad (9)$$

where  $0 < \gamma < 1$  is a discount factor,  $P_\pi \in \mathbb{R}^{N \times N}$  is a matrix with entries representing the probabilities  $p_\pi^{\bar{R}}(x, \pi(x), x')$  under  $\pi$  for all  $x, x' \in X$ ,  $\Lambda_\pi = [\Lambda(x_0) \ \Lambda(x_1) \ \dots]^T \in \mathbb{R}^N$  is the stacked state rewards,  $\beta \in \mathbb{R}^+$  is a weight indicating the relative importance,  $\mathbf{1}_N$  is an  $N$ -dimensional vector of ones, and  $V_\pi = P_\pi \circ V \in \mathbb{R}^{N \times N}$  is the Hadamard product of  $P_\pi$  and  $V$ , i.e.,  $V_\pi = [p_\pi^{\bar{R}}(x, a_\pi(x), x') \cdot V(x, u^{\bar{R}}, x')]_{N \times N}$  and  $a_\pi(x)$  represents taking action  $x$  from policy  $\pi$ .

The objective is to identify a stationary policy  $\pi^*$  that maximizes the expected return

$$\pi^* = \arg \max_{\pi} \sum_{n=0}^{\infty} \gamma^n P_\pi^n (\Lambda_\pi + \beta V_\pi \cdot \mathbf{1}_N) \quad (10)$$

$U_\pi(x) \leq U_{\pi^*}(x)$  for all  $x \in X$  if  $\pi^*$  in (10) is optimal.

**Theorem 2.** Consider a relaxed MDP product  $\bar{R} = \mathcal{M} \otimes \bar{A}_\phi$ . If there exists a policy  $\bar{\pi}$  such that an induced run  $r_{\bar{\pi}}^{\bar{R}}$  satisfies the acceptance condition of  $\bar{R}$ , any optimization method that solves (10) can find the policy  $\bar{\pi}$ .

**Proof.** Detailed rigorous analysis can be found in Appendix A.1  $\square$

To solve Problem 2 by optimizing the expected return of (9),  $r_A$ ,  $\beta$ ,  $\gamma$  can be determined as follows. Firstly, we can choose a fixed  $r_A$ . Then,  $\gamma$  can be obtained by solving (21) and (22). Finally, the range of  $\beta$  can be determined by solving (18) and (22). In order to minimize the violation cost, a great value of  $\beta$  is preferred.

**Remark 2.** Problem 2 is a MORL problem, and Theorem 2 shows how it can be addressed to ensure acceptance satisfaction while trying to minimize the long-term violation. In literature, MORL based approaches often seek to identify Pareto fronts [49]. However, our problem considers a trade-off between two possibly conflicting objectives. There are no existing MORL methods to provide formal guarantee of both acceptance satisfaction and violation depreciation, which might result in a safety issue during the execution of optimal policies. To overcome this challenge, we can always divide the task into two parts: hard and soft constraints, with two parallel automatons. The hard task should always be satisfied, and it can be applied with the traditional product MDP as Definition 10. The soft part can be partially infeasible, and it can be relaxed as Definition 11. Such an idea can be found in our previous results [17].

It is worth pointing out that  $V_\pi$  is a sparse matrix because most transitions have zero violation cost, and Theorem 2 provides a performance guarantee for the worst cases of  $\beta$  for acceptance satisfaction. By applying E-LDGBA, the sparse reward issue can be improved compared with using LDGBA.

### 5.2. Model-free reinforcement learning

Q-learning is a model-free RL method [20], which can be used to find the optimal policy for a finite MDP. In particular, the agent updates its Q-value from  $x$  to  $x'$  according to

$$Q(x, u^{\bar{R}}) \leftarrow (1 - \alpha) Q(x, u^{\bar{R}}) + \alpha \left[ R(x, u^{\bar{R}}, x') + \gamma \max_{u_*^{\bar{R}} \in U^{\bar{R}}(x)} Q(x, u_*^{\bar{R}}) \right], \quad (11)$$

where  $Q(x, u^{\bar{R}})$  is the Q-value of the state-action pair  $(x, u^{\bar{R}})$ ,  $\alpha \in (0, 1)$  is the learning rate,  $0 < \gamma \leq 1$  is the discount factor, and

$$R(x, u^{\bar{R}}, x') = \Lambda(x) + \beta V(x, u^{\bar{R}}, x') \quad (12)$$

**Algorithm 3** Model-free RL-based control on MDPs under potentially soft LTL specifications.

```

1: procedure INPUT:  $(\mathcal{M}, \phi, \Lambda)$ 
   Output: optimal policy  $\pi^*$ 
   Initialization: Set  $episode = 0$ ,  $iteration = 0$  and
                   $\tau$  as maximum allowed learning steps
2: Initialize  $r_A, \beta, \gamma$ ;
3: for all  $x \in X$  do
4:    $U(x) = 0$  and  $Q(x, u^{\bar{R}}) = 0 \forall x, u^{\bar{R}} \in U^{\bar{R}}(x)$ 
5:    $Count(x, u^{\bar{R}}) = 0$  for all  $u^{\bar{R}} \in U^{\bar{R}}(x)$ 
6: end for
7:  $x = x_0$ ;
8: while  $U$  are not converged do
9:    $episode++$ ;
10:   $\epsilon = 1/episode$ ;
11:  while  $iteration < \tau$  do
12:     $iteration++$ 
13:     $u = \arg \max_{u^{\bar{R}} \in U^{\bar{R}}} Q(x, u^{\bar{R}})$  with probability
         $1 - \epsilon$ , or select  $u$  randomly from  $U^{\bar{R}}(x)$ 
14:    Obtain and execute  $a$  of  $\mathcal{M}$  from  $u$ 
15:    Observe  $x', \Lambda(x), c_V^{\bar{R}}(x, u, x')$ 
16:     $r \leftarrow \Lambda(x) - \beta c_V^{\bar{R}}(x, u, x')$ 
17:     $Count(x, u)++$ 
18:     $\alpha = 1/Count(x, u)$ 
19:     $Q(x, u) \leftarrow (1 - \alpha)Q(x, u) +$ 
         $\alpha \left[ r + \gamma \cdot \max_{u^{\bar{R}} \in U^{\bar{R}}} Q(x', u^{\bar{R}}) \right]$ 
20:     $x = x'$ 
21:  end while
22: end while
23: for all  $x \in X$  do
24:    $U(x) = \max_{u^{\bar{R}} \in U^{\bar{R}}} Q(x, u^{\bar{R}})$ 
25:    $\pi^*(x) = \arg \max_{u^{\bar{R}} \in U^{\bar{R}}} U(x)$ 
26: end for
27: end procedure

```

denotes the immediate reward from  $x$  to  $x'$  under  $u^{\bar{R}}$ . The learning strategy is outlined in Alg. 3 that shows the steps of applying Q-learning into our framework. By applying the off-the-shelf RL, we have the following theorem.

**Theorem 3.** Given a finite MDP, i.e.,  $\bar{\mathcal{R}} = \mathcal{M} \otimes \bar{\mathcal{A}}_\phi = (X, U^{\bar{R}}, p^{\bar{R}}, x_0, F^{\bar{R}}, c_V^{\bar{R}}, f_V, T)$ , let  $Q^*(x, u^{\bar{R}})$  be the optimal Q-function for every pair of state-action  $(x, u^{\bar{R}})$ . Consider the RL with the updating rule

$$Q_{k+1}(x, u^{\bar{R}}) \leftarrow (1 - \alpha_k)Q_k(x, u^{\bar{R}}) + \alpha_k \left[ R(x, u^{\bar{R}}, x') + \gamma \max_{u^{\bar{R}}_* \in U^{\bar{R}}(x)} Q_k(x, u^{\bar{R}}_*) \right],$$

where  $k$  is the updating step,  $\gamma \in (0, 1]$  is the discount factor, and the learning rate  $\alpha_k$  satisfies  $\sum_k \alpha_k = \infty$  and  $\sum_k \alpha_k^2 < \infty$ . Then,  $Q_k(x, u^{\bar{R}})$  converges to  $Q^*(x, u^{\bar{R}})$  with a probability of 1 as  $k \rightarrow \infty$ .

**Theorem 3** is an immediate result of [50]. With standard learning rate  $\alpha$  and discount factor  $\gamma$  as Alg. 3, Q-value will converge to a unique limit  $Q^*$ . Therefore, the optimal expected utility and policy can be obtained as  $U_{\pi^*}(x) = \max_{u^{\bar{R}} \in U^{\bar{R}}(x)} Q^*(x, u^{\bar{R}})$  and  $\pi^*(x) = \arg \max_{u^{\bar{R}} \in U^{\bar{R}}(x)} U_{\pi^*}(x)$ . In (11), the discount  $\gamma$  is tuned to improve the trade-off between immediate and future rewards. Note that our novel design can be also applied with modern advanced RL algorithms.

**Complexity Analysis** The number of states is  $|S| \times |Q|$ , where  $|Q|$  is determined by the original LDGBA  $\mathcal{A}_\phi$  because the construction of E-LDGBA  $\bar{\mathcal{A}}_\phi$  will not increase the size, and  $|S|$  is the size of the labeled MDP model. Due to the consideration of a relaxed product MDP and extended actions in Definition 11, the maximum complexity of actions

available at  $x = (s, l, q)$  is  $O(|A(s)| \times (|\Sigma| + 1))$ , since  $U^{\bar{R}}(x)$  are created from  $A(s)$  and  $\Sigma \cup \{\epsilon\}$ .

## 6. Case studies

The developed RL-based control synthesis is implemented in Python. Owl [47] is used to convert LTL specifications into LDGBA, and P\_MDP package [13] is used to construct state transition models. All simulations are carried out on a laptop with 2.60 GHz quad-core CPU and 8 GB of RAM. For Q-learning, the optimal policies of each case are generated using  $10^5$  episodes with random initial states. The learning rate  $\alpha$  is determined by Alg. 3 with  $\gamma = 0.999$  and  $r_A = 10$ . To validate the effectiveness of our approach, we first carry out simulations over grid environments and then validate the approach in a more realistic office scenario with a TIAGo robot.

### 6.1. Simulation results

Consider a partitioned workspace of  $10\text{m} \times 10\text{m}$ , as illustrated in Figs. 3 and 5, where each cell represents an area of  $2\text{m} \times 2\text{m}$ . Cells are color-coded to designate specific areas of interest, including Base1, Base2, Base3, Obs, and Sply, where Obs and Sply stand for obstacle and supply zones, respectively. To incorporate environmental uncertainties, each cell is assigned a probability indicating the likelihood of a particular property appearing there. For instance, Obs : 0.1 denotes a 0.1 probability that the cell contains an obstacle. The robot operates based on a unicycle model, with dynamics given by  $\dot{x} = v \cdot \sin(\theta)$ ,  $\dot{y} = v \cdot \cos(\theta)$ , and  $\dot{\theta} = \omega$ , where  $x, y$ , and  $\theta$  represent the robot's position and orientation. The control inputs are the linear and angular velocities, denoted as  $u = (v, \omega)$ .

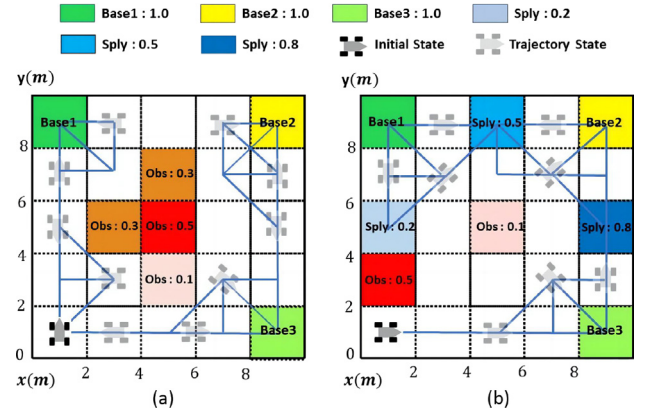


Fig. 3. Simulated trajectories of 50 time steps under corresponding optimal policies in feasible workspaces with (a) a simple task, i.e.,  $\varphi_{case1}$ , and (b) a relatively complex task, i.e.,  $\varphi_{case2}$ .

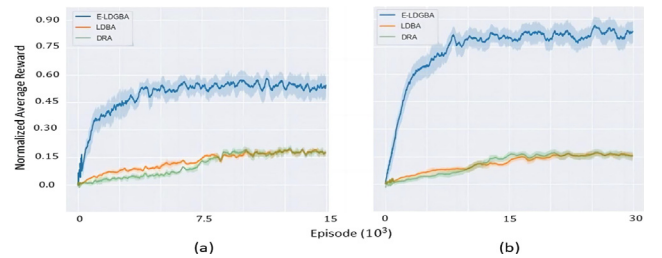


Fig. 4. Normalized average rewards during training using E-LDGBA, LDGBA, and DRA for (a) Task  $\varphi_{case1}$  in Fig. 3 and (b) Task  $\varphi_{case2}$  in Fig. 3.



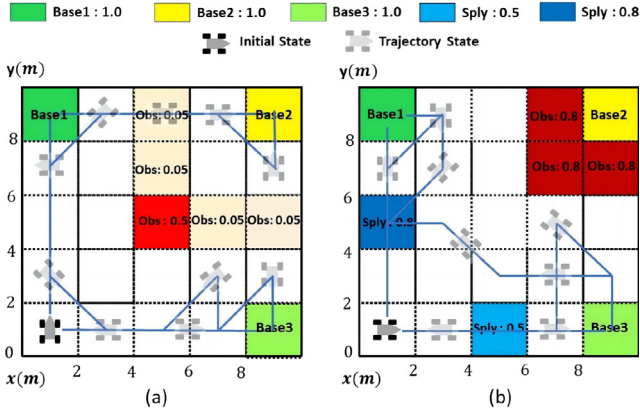


Fig. 5. Simulated trajectories of 50 time steps under corresponding optimal policies for partially infeasible tasks on workspaces where the agent visits the grid labeled as Base2 with (a) a low risk or (b) a high risk.

In addition, we assume the robot cannot consistently successfully execute the action primitives to model motion uncertainties. For instance, action primitives “FR” and “BK” mean the robot can successfully move forward and backward 2 m with a probability of 0.9, respectively, and fail with a probability of 0.1. On the other hand, action primitives “TR” and “TL” mean the robot can successfully turn right and left for an angle of  $\frac{\pi}{2}$  exactly with a probability of 0.9, respectively, and fail by an angle of  $\frac{\pi}{4}$  (undershoot) with a probability of 0.05 and an angle of  $\frac{3\pi}{4}$  (overshoot) with another probability of 0.05. Finally, action primitive “ST” means the robot remains at its current cell. The resulting MDP model has 25 states.

**Case 1:** As shown in Fig. 3(a), we first consider a case in which user-specified tasks can be successfully executed. The desired surveillance task to be performed is formulated as

$$\varphi_{case1} = (\Box\Diamond\text{Base1}) \wedge (\Box\Diamond\text{Base2}) \wedge (\Box\Diamond\text{Base3}) \wedge \Box\neg\text{Obs}, \quad (13)$$

which requires the mobile robot to visit all base stations infinitely while avoiding obstacles. Its corresponding LDGBA has 2 states with 3 accepting sets, and the relaxed product MDP has 50 states. In this case, each episode terminates after  $\tau = 100$  steps and  $\beta = 8$ . Fig. 3(a) shows the generated optimal trajectory, which indicates  $\varphi_{case1}$  is completed.

**Case 2:** We validate our approach with more complex task specifications in Fig. 3(b). The task is expressed as

$$\varphi_{case2} = \varphi_{case1} \wedge \Box\Diamond(\text{Sply} \rightarrow \bigcirc((\neg\text{Sply})\mathcal{U}\varphi_{one1})), \quad (14)$$

where  $\varphi_{one1} = \text{Base1} \vee \text{Base2} \vee \text{Base3}$ .  $\varphi_{case2}$  requires the robot to visit the supply station and then go to one of the base stations while avoiding obstacles. It also requires all base stations to be visited. Its corresponding LDGBA has 24 states with 4 accepting sets, and the relaxed product MDP has 600 states. The generated optimal trajectory is shown in Fig. 3(b).

**Case 3:** We consider more challenging workspaces in Fig. 5 (a) and (b), where user-specified tasks might not be fully feasible due to potential obstacles (i.e., environment uncertainties). The task specification is

$$\varphi_{case3} = \varphi_{case1} \wedge \Box\Diamond(\varphi_{one1} \rightarrow \bigcirc((\neg\varphi_{one1})\mathcal{U}\text{Sply})), \quad (15)$$

where  $\varphi_{case3}$  requires the robot to visit one of the base stations and then go to one of the supply stations while avoiding obstacles. Its corresponding LDGBA has 24 states with 4 accepting sets, and the product-MDP has 600 states.

For  $\varphi_{case2}$  and  $\varphi_{case3}$ , each episode terminates after  $\tau = 1000$  steps and  $\beta = 5$ . Note that in the case of Fig. 5 (a) and (b), AMECs might not exist since Base2 is surrounded by probabilistic obstacles and may not be visited. The simulated optimal trajectories indicate that the robot

Table 1

Simulation results of large scale workspaces.

Workspace size[cell]	MDP States	$\bar{R}$ States	Episode Steps
15 × 15	225	450	30000
25 × 25	625	1250	50000
40 × 40	1600	3200	100000

takes a risk to accomplish the task in the case of Fig. 5(a). In contrast, the robot decides not to visit Base2 to avoid the high risk of running into obstacles in the case of Fig. 5(b).

**Scalability Analysis:** RL-based policy synthesis is conducted for  $\varphi_{case1}$  across workspaces of varying sizes (with each grid further subdivided) to illustrate computational complexity. The simulation results are presented in Table 1. The number of episodes in Table 1 indicates the time used to converge to optimal satisfaction planning. It is also verified that the given task  $\varphi_{case1}$  can be successfully carried out in large workspaces.

## 6.2. Comparison

We compare our algorithm with the existing learning-based works [26,30,31,38,39,41]. From the automaton perspective, DRA has been leveraged in [26,38,39] to follow the Rabin acceptance condition with rejecting sets and accepting sets. In addition, the Limit-Deterministic Büchi Automaton (LDDBA) used in [30,31] only has one accepting set. Figs. 7(a) and 1(b) illustrate that (E-)LDGBA and LDDBA have different numbers of accepting sets for the same LTL formula  $\phi = \Box\Diamond a \wedge \Box\Diamond b$ . On the other hand, the reward is typically designed based on the acceptance condition after integrating LTL with the RL algorithm. Consequently, applying the LDDBA and DRA results in more sparse rewards than E-LDGBA (with several accepting sets). Figs. 4 and 6 show the convergence and comparison of learning results.

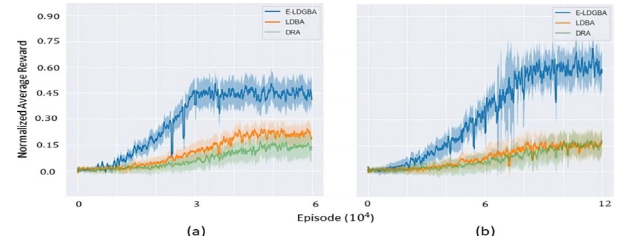
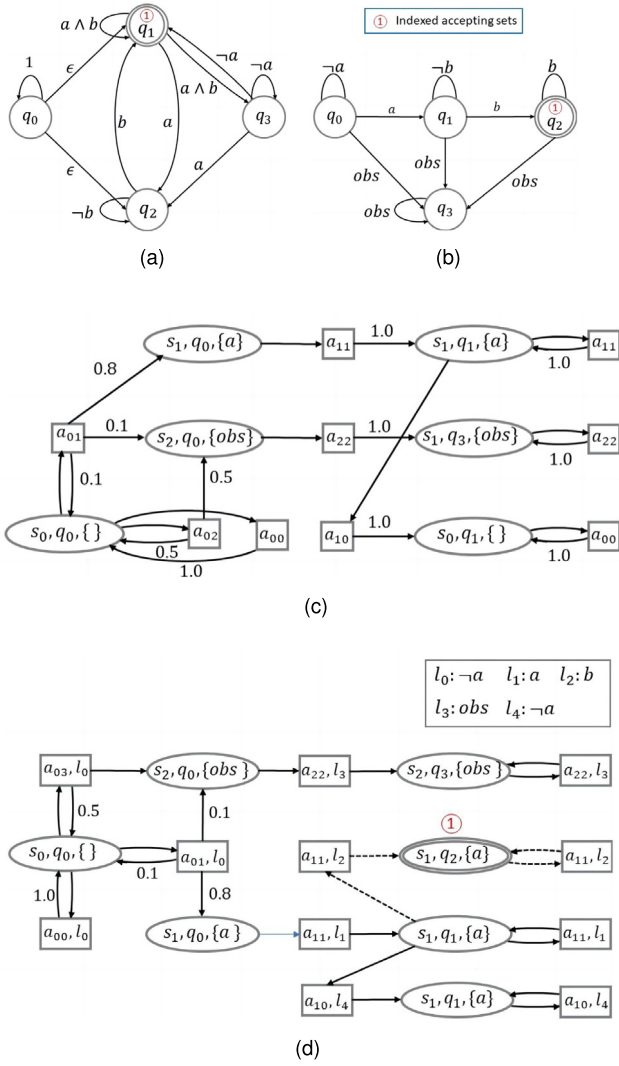


Fig. 6. Normalized average rewards during training using E-LDGBA, LDDBA, and DRA for Task  $\varphi_{case3}$  on two environments in Fig. 5, respectively.

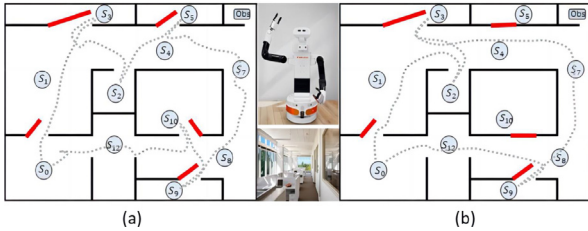
Regarding infeasible LTL tasks with respect to environment uncertainties, the most recent work [41] designed the reward to maximize the satisfaction of infeasible tasks by visiting accepting sets in a standard product MDP as much as possible. However, for some specific LTL formulas, e.g.,  $\phi = \Diamond(a \wedge \Diamond b)$ , the corresponding LDGBA only has one accepting set, as shown in Fig. 7(b). After taking the product of this LDGBA and the PL-MDP in Fig. 1(a), it can be observed in Fig. 7(c) that the standard product MDP has no accepting set. Consequently, the work [41] cannot handle this infeasible case. In contrast, the relaxed product MDP in Fig. 7(d) from our work allows finding the policy with the least expected violation costs.

## 6.3. Experimental results

This section validates our algorithm for high-level decision-making tasks within a real-world environment, demonstrating that the framework can integrate with low-level noisy controllers to establish a hierarchical architecture. Consider an office environment constructed in



**Fig. 7.** Examples of the baseline comparisons regarding automaton structures and infeasible cases. (a) The LDBA of the LTL formula  $\phi = \square\Diamond a \wedge \square\Diamond b$  only has one accepting set. (b) The LDGBA of the LTL formula  $\phi = \Diamond(a \wedge b)$ . (c) The standard product MDP between the PL-MDP in Fig. 1(a) and the LDGBA of the LTL formula  $\phi = \Diamond(a \wedge b)$  in Fig. 7(b). There are no accepting states, and this is an infeasible case. (d) The relaxed product MDP (Definition 11) corresponding to the standard product MDP in (c). An accepting state exists, and the dashed arrows represent the transitions with nonzero violation costs.



**Fig. 8.** Experimental trajectories for the task  $\varphi_{case4}$ .

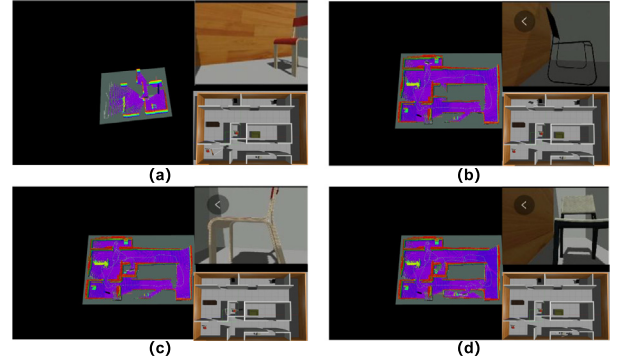
ROS Gazebo, illustrated in Fig. 8, which consists of 7 rooms denoted by  $S_0, S_2, S_3, S_5, S_9, S_{10}, Obs$  and 5 corridors denoted by  $S_1, S_4, S_7, S_8, S_{12}$ . The TIAGo robot navigates collision-free from the center of one region to another, avoiding obstacles and ensuring no unintended crossings.

To account for motion uncertainties, we assume the robot reaches its target region with a probability of 0.9, while there is a 0.1 probability it will move to an adjacent region instead. The resulting MDP comprises 12 states. The service to be performed by TIAGo is expressed as

$$\varphi_{case4} = \varphi_{all} \wedge \square\neg Obs \quad (16)$$

where  $\varphi_{all} = \square\Diamond S_0 \wedge \square\Diamond S_2 \wedge \square\Diamond S_3 \wedge \square\Diamond S_5 \wedge \square\Diamond S_9 \wedge \square\Diamond S_{10}$ . In (16),  $\varphi_{all}$  requires the robot to always service all rooms (e.g., pick trash) and return to  $S_0$  (e.g., release trash) while avoiding  $Obs$ . Its corresponding LDGBA has 6 states with 6 accepting sets, and the relaxed product MDP has 72 states.

All room doors are open, except the doors of room  $S_5, S_{10}$  in Fig. 8(b). As a result,  $\varphi_{case4}$  cannot be fully completed in the case of Fig. 8(b), and the task is infeasible. It is also worth pointing out there do not exist AMECs in the corresponding product automaton  $\bar{P}$  or  $\bar{P}$  in Fig. 8 (a) because the robot has the nonzero probability of entering  $Obs$  at state  $S_7$  as the motion uncertainties arise. The optimal policies for the two cases are generated, and each episode terminates after  $\tau = 150$  steps with  $\beta = 4$ . The generated satisfying trajectories (without collision) of one round are shown in Fig. 8 (a) and (b), and the robot tries to complete the feasible part of task  $\varphi_{case4}$  in Fig. 8(b). To further illustrate the effectiveness of our approach, we consider the semantics of objects as atomic propositions in temporal logic formulas, following the method presented in [51,52]. Assuming that the environment depicted in Fig. 8 is unknown, the robot must explore the environment while executing the task. The experimental results are shown in Fig. 9. A video of the experiment is available at the following webpage.<sup>1</sup>



**Fig. 9.** The figure illustrates a quadrupedal robot navigating the environment while performing a task using the methods outlined in this article. In each sub-figure, the top right corner displays the robot's initial viewpoint, while the bottom right provides a synchronized image from Gazebo.

## 7. Conclusions

This paper presents a learning-based control synthesis of motion planning subject to motion and environment uncertainties. The LTL formulas are applied to express complex tasks via the automaton theory. Differently, in this work, an LTL formula is converted into a designed E-LDGBA to improve the performance of mission accomplishment and resolve the drawbacks of DRA and LDBA. Furthermore, the innovative relaxed product MDP and utility schemes consisting of the accepting reward and violation reward are proposed to accomplish the satisfaction of soft tasks. In order to provide formal guarantees of achieving the goals of multi-objective optimizations, future research will consider

<sup>1</sup> [https://sites.google.com/d/1Gx374qPGITeinU\\_8IS6vd6IVnVY-1kJp/p/13rdD73dGargssrmlQyfXue6-nsiw3qc/edit](https://sites.google.com/d/1Gx374qPGITeinU_8IS6vd6IVnVY-1kJp/p/13rdD73dGargssrmlQyfXue6-nsiw3qc/edit).

more advanced multi-objective learning methods. Also problems over continuous state and/or action spaces will be studied by incorporating deep neural networks.

### CRedit authorship contribution statement

**Mingyu Cai:** Formal analysis. **Zhangli Zhou:** Formal analysis, Conceptualization. **Lin Li:** Formal analysis, Conceptualization. **Shaoping Xiao:** Formal analysis, Conceptualization. **Zhen Kan:** Formal analysis, Conceptualization.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

This work was supported by the National Natural Science Foundation of China under Grant 62173314.

### Data availability

No data was used for the research described in the article.

### Appendix

#### A.1. Proof of Theorem 2

With reference to [41], we provide a proof for Theorem 2. For any policy  $\pi$ , let  $MC_{\bar{R}}^{\pi}$  denote the Markov chain induced by  $\pi$  on  $\bar{R}$ . Since  $MC_{\bar{R}}^{\pi}$  can be written as  $MC_{\bar{R}}^{\pi} = T_{\pi} \sqcup R_{\pi}^1 \sqcup R_{\pi}^2 \dots R_{\pi}^n$ , (9) can be reformulated as

$$\begin{bmatrix} U_{\pi}^{tr} \\ U_{\pi}^{rec} \end{bmatrix} = \sum_{n=0}^{\infty} \gamma^n \begin{bmatrix} P_{\pi}(\mathcal{T}, \mathcal{T}) & P_{\pi}^{tr} \\ 0_{\sum_{i=1}^m N_i \times r} & P_{\pi}(R, R) \end{bmatrix}^n \cdot \left( \begin{bmatrix} A_{\pi}^{tr} \\ A_{\pi}^{rec} \end{bmatrix} + \beta \begin{bmatrix} O_{\pi}^{tr} \\ O_{\pi}^{rec} \end{bmatrix} \right), \quad (17)$$

where  $U_{\pi}^{tr}$  and  $U_{\pi}^{rec}$  are the utilities of states in the transient and recurrent classes, respectively. In (17),  $P_{\pi}(\mathcal{T}, \mathcal{T}) \in \mathbb{R}^{r \times r}$  denotes the probability transition matrix between states in  $\mathcal{T}^{\pi}$ .  $P_{\pi}^{tr} = [P_{\pi}^{tr_1} \dots P_{\pi}^{tr_m}] \in \mathbb{R}^{r \times \sum_{i=1}^m N_i}$  is a probability transition matrix where  $P_{\pi}^{tr_i} \in \mathbb{R}^{r \times N_i}$  represents the probability of transiting from a transient state in  $\mathcal{T}^{\pi}$  to the states of  $R_{\pi}^i$ . The  $P_{\pi}(R, R)$  is a diagonal block matrix, where the  $i$ th block is an  $N_i \times N_i$  matrix containing transition probabilities between states within  $R_{\pi}^i$ .  $P_{\pi}(\bar{R}, \bar{R})$  is a stochastic matrix since each block matrix is a stochastic matrix [48]. The rewards vector  $A_{\pi}$  can also be partitioned to  $A_{\pi}^{tr}$  and  $A_{\pi}^{rec}$ . Similarly,  $O_{\pi} = [O_{\pi}(x_0) \ O_{\pi}(x_1) \ \dots]^T = V_{\pi} \cdot \mathbf{1}_N$  can be divided into transient class  $O_{\pi}^{tr}$  and recurrent class  $O_{\pi}^{rec}$ .

We prove this theorem by contradiction. Suppose there exists an optimal policy  $\pi^*$  not satisfying the acceptance condition of  $\bar{R}$ . Based on Lemma 3, the following is true:  $F_k^{\bar{R}} \subseteq T_{\pi^*}, \forall i \in \{1, \dots, f\}$ . As a result, for any  $j \in \{1, \dots, n\}$ , we have  $R_{\pi^*}^j \cap F_i^{\bar{R}} \neq \emptyset, \forall i \in \{1, \dots, f\}$ .

The strategy is to show that there always exists a policy  $\bar{\pi}$  with greater utility than  $\pi^*$ , which contradicts to the optimality of  $\pi^*$ . Let us consider a state  $x_R \in R_{\pi^*}^j$  and let  $P_{\pi^*}^{x_R R_{\pi^*}^j}$  denote a row vector of  $P_{\pi^*}^n(R, R)$  that contains the transition probabilities from  $x_R$  to the states in the same recurrent class  $R_{\pi^*}^j$  in  $n$  steps. The expected return  $U_{\pi^*}(x_R)$  of  $x_R$  under  $\pi^*$  is then obtained from (17) as:  $U_{\pi^*}^{rec}(x_R) = \sum_{n=0}^{\infty} \gamma^n \left[ 0_{k_1}^T P_{\pi^*}^{x_R R_{\pi^*}^j} 0_{k_2}^T \right] (A_{\pi^*}^{rec} + \beta O_{\pi^*}^{rec})$  where  $k_1 = \sum_{i=1}^{j-1} N_i, k_2 = \sum_{i=j+1}^n N_i$ . Since  $\bar{R}_{\pi^*}^j \cap F_i^{\bar{R}} = \emptyset, \forall i \in \{1, \dots, f\}$ , all the

elements of  $A_{\pi^*}^{rec}$  are equal to zero and each entry of  $O_{\pi^*}^{rec}$  is non-positive. We can conclude  $U_{\pi^*}^{rec}(x_R) \leq 0$ . To prove that optimal policy  $\pi^*$  raises a contradiction, the following analysis will show that  $U_{\bar{\pi}}^{rec}(x_R) \geq U_{\pi^*}^{rec}(x_R)$  for some policies  $\bar{\pi}$  that satisfy the acceptance condition of  $\bar{R}$ . Thus we have  $R_{\bar{\pi}}^j \cap F_i^{\bar{R}} \neq \emptyset, \forall i \in \{1, \dots, f\}$ .

**Case 1:** If  $x_R \in R_{\bar{\pi}}^j$ , by Lemma 3 and (7), there exist a minimum of  $f$  accepting states such that  $X_A = \{x_A \mid x_A \in R_{\bar{\pi}}^j \cap F_i^{\bar{R}}, \forall i \in \{1, \dots, f\}\}$  in  $R_{\bar{\pi}}^j$  with positive rewards  $r_A$ . From (17),  $U_{\bar{\pi}}^{rec}(x_R)$  can be lower bounded as  $U_{\bar{\pi}}^{rec}(x_R) \geq \sum_{n=0}^{\infty} \gamma^n \left( P_{\bar{\pi}}^{x_R X_A} r_A + \beta P_{\bar{\pi}}^{x_R R_{\bar{\pi}}^j} V_{\bar{\pi}}^{rec} \mathbf{1}_{N_j} \right)$ , where  $P_{\bar{\pi}}^{x_R X_A}$  is the transition probability from  $x_R$  to  $x_A$  in  $n$  steps, and  $V_{\bar{\pi}}^{rec} \in \mathbb{R}^{N_j \times N_j}$  represents the violation cost of states in  $\bar{R}_{\bar{\pi}}^j$ . Since  $x_R$  and  $x_A$  are recurrent states, there always exists a lower bound  $\underline{P}_{\bar{\pi}}^{x_R X_A} \in \mathbb{R}^+$  of the transition probability  $P_{\bar{\pi}}^{x_R X_A}$ . We can select a positive reward  $r_A$  such that

$$U_{\bar{\pi}}^{rec}(x_R) \geq \underline{P}_{\bar{\pi}}^{x_R X_A} r_A + \beta N_j^2 V_{\bar{\pi}}^{rec} \geq 0 \quad (18)$$

where  $\underline{V}_{\bar{\pi}}^{rec} \in \mathbb{R}^-$  represents the minimal entry in  $V_{\bar{\pi}}^{rec}$ . By selecting  $r_A$  to satisfy (18), we can conclude in this case  $U_{\bar{\pi}}^{rec}(x_R) > U_{\pi^*}^{rec}(x_R)$ .

**Case 2:** If  $x_R \in \mathcal{T}_{\bar{\pi}}$ , we know  $\mathcal{T}_{\bar{\pi}} \cap F_i^{\bar{R}} = \emptyset, \forall i \in \{1, \dots, f\}$ . As demonstrated in [48], for a transient state  $x_{tr} \in \mathcal{T}_{\bar{\pi}}$ , there always exists an upper bound  $\Delta < \infty$  such that  $\sum_{n=0}^{\infty} p^n(x_{tr}, x_{tr}) < \Delta$ , where  $p^n(x_{tr}, x_{tr})$  denotes the probability of returning from a transient state  $x_{tr}$  to itself in  $n$  time steps. For a recurrent state  $x_{rec} \in \bar{R}_{\bar{\pi}}^j$ , it is always true that

$$\sum_{n=0}^{\infty} \gamma^n p^n(x_{rec}, x_{rec}) > \frac{1}{1 - \gamma^n} \bar{p} \quad (19)$$

where there exists  $\bar{n}$  such that  $p^{\bar{n}}(x_{rec}, x_{rec})$  is nonzero and can be lower bounded by  $\bar{p}$  [48]. From (17), one has

$$\begin{aligned} U_{\bar{\pi}}^{tr} &> \sum_{n=0}^{\infty} \gamma^n P_{\bar{\pi}}^n(\mathcal{T}, \mathcal{T}) (A_{\bar{\pi}}^{tr} + \beta O_{\bar{\pi}}^{tr}) \\ &+ \sum_{n=0}^{\infty} \gamma^n P_{\bar{\pi}}^{tr} P_{\bar{\pi}}^n(\bar{R}, \bar{R}) (A_{\bar{\pi}}^{rec} + \beta O_{\bar{\pi}}^{rec}) \end{aligned} \quad (20)$$

Let  $\max(\cdot)$  and  $\min(\cdot)$  represent the maximum and minimum entry of an input vector, respectively. The upper bound  $\bar{m} = \left\{ \max(\bar{M}) \mid \bar{M} < P_{\bar{\pi}}^{tr} \bar{P} (A_{\bar{\pi}}^{rec} + \beta V_{\bar{\pi}}^{rec} \mathbf{1}_{\bar{N}}) \right\}$ , where  $\bar{N} = \sum_{j=1}^m N_j$  and  $\bar{P}$  is a block matrix whose nonzero entries are derived similarly to the  $\bar{p}$  in (19). Using the fact  $\sum_{n=0}^{\infty} \gamma^n P_{\bar{\pi}}^n(\mathcal{T}, \mathcal{T}) \leq \Delta \mathbf{1}_{r \times r}$  [48], where  $\mathbf{1}_{r \times r}$  is a  $r \times r$  matrix of all ones, the utility  $U_{\bar{\pi}}^{tr}(x_R)$  can be lower bounded from (19) and (20) as

$$U_{\bar{\pi}}^{tr}(x_R) > \Delta \cdot r \cdot \beta V_{\bar{\pi}}^{tr} + \frac{1}{1 - \gamma^n} \bar{m} \quad (21)$$

Since  $U_{\pi^*}^{rec}(x_R) = 0$ , the contradiction  $U_{\bar{\pi}}^{tr}(x_R) > 0$  in this case will be achieved if  $\Delta \cdot r \cdot \beta V_{\bar{\pi}}^{tr} + \frac{1}{1 - \gamma^n} \bar{m} \geq 0$ . Because  $\Delta r (w_L + \beta V_{\bar{\pi}}^{tr}) < 0$ , it needs  $\bar{m} > 0$  as

$$\bar{m} > r_A + \beta \bar{N} V_{\bar{\pi}}^{rec} > 0. \quad (22)$$

Thus, by choosing  $\gamma$  to be sufficiently close to 1 with  $\bar{m} > 0$ , we have  $U_{\bar{\pi}}^{tr}(x_R) > 0 \geq U_{\pi^*}^{rec}(x_R)$ . The above procedure shows the contradiction of the assumption that  $\pi_V^*$  is optimal.

### References

- [1] C. Baier, J.-P. Katoen, Principles of Model Checking, MIT Press, 2008.
- [2] M. Kloetzer, C. Mahulea, LTL-based planning in environments with probabilistic observations, IEEE Trans. Autom. Sci. Eng. 12 (4) (2015) 1407–1420.
- [3] Y. Kantaros, M.M. Zavlanos, STyLuS\*: A temporal logic optimal control synthesis algorithm for large-scale multi-robot systems, Int. J. Robot. Res. 39 (7) (2020) 812–836.
- [4] A. Pacheco, H. Kress-Gazit, Physically feasible repair of reactive, linear temporal logic-based, high-level tasks, IEEE Trans. Robot. (2023).
- [5] X. Sun, Y. Shoukry, Neurosymbolic motion and task planning for linear temporal logic tasks, IEEE Trans. Robot. (2024).

- [6] Z. Chen, M. Cai, Z. Zhou, L. Li, Z. Kan, Fast motion planning in dynamic environments with extended predicate-based temporal logic, *IEEE Trans. Autom. Sci. Eng.* (2024).
- [7] Z. Chen, Z. Kan, Real-time reactive task allocation and planning of large heterogeneous multi-robot systems with temporal logic specifications, *The International Journal of Robotics Research* (2024).
- [8] S. Thrun, W. Burgard, D. Fox, *Probabilistic Robotics*, MIT Press, 2005.
- [9] M. Guo, T. Liao, J. Wang, Z. Li, Hierarchical motion planning under probabilistic temporal tasks and safe-return constraints, *IEEE Trans. Autom. Control* 68 (11) (2023) 6727–6742.
- [10] X.C. Ding, S.L. Smith, C. Belta, D. Rus, MDP optimal control under temporal logic constraints, in: *Proc. IEEE Conf. Decis. Control.*, 2011, pp. 532–538.
- [11] X. Ding, S.L. Smith, C. Belta, D. Rus, Optimal control of Markov decision processes with linear temporal logic constraints, *IEEE Trans. Autom. Control* 59 (5) (2014) 1244–1257.
- [12] A. Ulusoy, T. Wongpiromsarn, C. Belta, Incremental controller synthesis in probabilistic environments with temporal logic constraints, *Int. J. Robot. Res.* 33 (8) (2014) 1130–1144.
- [13] M. Guo, M.M. Zavlanos, Probabilistic motion planning under temporal tasks and soft constraints, *IEEE Trans. Autom. Control* 63 (12) (2018) 4051–4066.
- [14] B. Lacerda, F. Faruq, D. Parker, N. Hawes, Probabilistic planning with formal performance guarantees for mobile service robots, *Int. J. Robot. Res.* 38 (9) (2019) 1098–1123.
- [15] C.I. Vasile, X. Li, C. Belta, Reactive sampling-based path planning with temporal logic specifications, *Int. J. Robot. Res.* (2020) 0278364920918919.
- [16] C.-I. Vasile, J. Tumova, S. Karaman, C. Belta, D. Rus, Minimum-violation sctl motion planning for mobility-on-demand, in: *2017 IEEE International Conference on Robotics and Automation, ICRA, IEEE*, 2017, pp. 1481–1488.
- [17] M. Cai, H. Peng, Z. Li, H. Gao, Z. Kan, Receding horizon control-based motion planning with partially infeasible LTL constraints, *IEEE Control Syst. Lett.* 5 (4) (2020) 1279–1284.
- [18] M. Cai, S. Xiao, Z. Li, Z. Kan, Optimal probabilistic motion planning with potential infeasible LTL constraints, *IEEE Trans. Autom. Control* (2021).
- [19] N. Mehdipour, M. Althoff, R.D. Tebbens, C. Belta, Formal methods to comply with rules of the road in autonomous driving: State of the art and grand challenges, *Automatica* 152 (2023) 110692.
- [20] C.J. Watkins, P. Dayan, Q-learning, *Mach. Learn.* 8 (3–4) (1992) 279–292.
- [21] C. Sun, X. Li, C. Belta, Automata guided semi-decentralized multi-agent reinforcement learning, in: *2020 American Control Conference, ACC, IEEE*, 2020, pp. 3900–3905.
- [22] H. Wang, H. Zhang, L. Li, Z. Kan, Y. Song, Task-driven reinforcement learning with action primitives for long-horizon manipulation skills, *IEEE Trans. Cybern.* (2023).
- [23] Y. Meng, C. Fan, Signal temporal logic neural predictive control, *IEEE Robot. Autom. Lett.* (2023).
- [24] T. Brázdil, K. Chatterjee, M. Chmelik, V. Forejt, J. Křetínský, M. Kwiatkowska, D. Parker, M. Ujma, Verification of Markov decision processes using learning algorithms, in: *Int. Symp. Autom. Tech. Verif. Anal.*, Springer, 2014, pp. 98–114.
- [25] J. Fu, U. Topcu, Probably approximately correct MDP learning and control with temporal logic constraints, 2014, *arXiv preprint arXiv:1404.7073*.
- [26] D. Sadigh, E.S. Kim, S. Coogan, S.S. Sastry, S.A. Seshia, A learning based approach to control synthesis of Markov decision processes for linear temporal logic specifications, in: *Proc. IEEE Conf. Decis. Control.*, 2014, pp. 1091–1096.
- [27] J. Wang, X. Ding, M. Lahijanian, I.C. Paschalidis, C.A. Belta, Temporal logic motion control using actor-critic methods, *Int. J. Robot. Res.* 34 (10) (2015) 1329–1344.
- [28] D. Aksaray, A. Jones, Z. Kong, M. Schwager, C. Belta, Q-learning for robust satisfaction of signal temporal logic specifications, in: *Proc. IEEE Conf. Decis. Control*, 2016, pp. 6565–6570.
- [29] X. Li, Z. Serlin, G. Yang, C. Belta, A formal methods approach to interpretable reinforcement learning for robotic planning, *Science Robotics* 4 (37) (2019).
- [30] E.M. Hahn, M. Perez, S. Schewe, F. Somenzi, A. Trivedi, D. Wojtczak, Omega-regular objectives in model-free reinforcement learning, in: *Int. Conf. Tools Alg. Constr. Anal. Syst.*, Springer, 2019, pp. 395–412.
- [31] A.K. Bozkurt, Y. Wang, M.M. Zavlanos, M. Pajic, Control synthesis from linear temporal logic specifications using model-free reinforcement learning, 2019, *arXiv preprint arXiv:1909.07299*.
- [32] K. Kim, G.E. Fainekos, Approximate solutions for the minimal revision problem of specification automata, in: *IEEE Int. Conf. Intell. Robot. and Syst.*, 2012, pp. 265–271.
- [33] K. Kim, G.E. Fainekos, S. Sankaranarayanan, On the revision problem of specification automata, in: *Proc. IEEE Int. Conf. Robot. Autom.*, 2012, pp. 5171–5176.
- [34] J. Tumova, G.C. Hall, S. Karaman, E. Frazzoli, D. Rus, Least-violating control strategy synthesis with safety rules, in: *Proc. Int. Conf. Hybrid Syst., Comput. Control*, 2013, pp. 1–10.
- [35] M. Guo, D.V. Dimarogonas, Multi-agent plan reconfiguration under local LTL specifications, *Int. J. Robot. Res.* 34 (2) (2015) 218–235.
- [36] M. Cai, S. Xiao, B. Li, Z. Li, Z. Kan, Reinforcement learning based temporal logic control with maximum probabilistic satisfaction, in: *Int. Conf. Robot. Autom.*, IEEE, 2021, pp. 806–812.
- [37] R. Alur, O. Bastani, K. Jothimurugan, M. Perez, F. Somenzi, A. Trivedi, Policy synthesis and reinforcement learning for discounted LTL, in: *International Conference on Computer Aided Verification*, Springer, 2023, pp. 415–435.
- [38] Q. Gao, D. Hajinezhad, Y. Zhang, Y. Kantaros, M.M. Zavlanos, Reduced variance deep reinforcement learning with temporal logic specifications, in: *Proc. ACM/IEEE Int. Conf. Cyber-Physical Syst.*, 2019, pp. 237–248.
- [39] M. Cai, H. Peng, Z. Li, Z. Kan, Learning-based probabilistic LTL motion planning with environment and motion uncertainties, *IEEE Trans. Autom. Control* (2020).
- [40] M. Cai, M. Mann, Z. Serlin, K. Leahy, C.-I. Vasile, Learning minimally-violating continuous control for infeasible linear temporal logic specifications, in: *American Control Conference*, 2023.
- [41] M. Hasanbeig, Y. Kantaros, A. Abate, D. Kroening, G.J. Pappas, I. Lee, Reinforcement learning for temporal logic control synthesis with probabilistic satisfaction guarantees, 2019, *arXiv preprint arXiv:1909.05304*.
- [42] M. Cai, M. Hasanbeig, S. Xiao, et al., Modular deep reinforcement learning for continuous motion planning with temporal logic, *IEEE Robot. Autom. Lett.* 6 (4) (2021) 7973–7980.
- [43] J. Wang, H. Hasanbeig, K. Tan, et al., Mission-driven exploration for accelerated deep reinforcement learning with temporal logic task specifications, 2023, *arXiv preprint arXiv:2311.17059*.
- [44] M. Hasanbeig, A. Abate, D. Kroening, Certified reinforcement learning with logic guidance, 2019, *arXiv preprint arXiv:1902.00778*.
- [45] S. Sickert, J. Esparza, S. Jaax, J. Křetínský, Limit-deterministic büchi automata for linear temporal logic, in: *Int. Conf. Comput. Aided Verif.*, Springer, 2016, pp. 312–332.
- [46] E.M. Hahn, G. Li, S. Schewe, A. Turrini, L. Zhang, Lazy probabilistic model checking without determinisation, 2013, *arXiv preprint arXiv:1311.2928*.
- [47] J. Křetínský, T. Meggendorfer, S. Sickert, Owl: A library for  $\omega$ -words, automata, and LTL, in: *Autom. Tech. Verif. Anal.*, Springer, 2018, pp. 543–550.
- [48] R. Durrett, second ed., *Essentials of Stochastic Processes*, vol. 1, Springer, 2012.
- [49] P. Vamplew, R. Dazeley, A. Berry, R. Issabekov, E. Dekker, Empirical evaluation methods for multiobjective reinforcement learning algorithms, *Mach. Learn.* 84 (1–2) (2011) 51–80.
- [50] J.N. Tsitsiklis, Asynchronous stochastic approximation and Q-learning, *Mach. Learn.* 16 (3) (1994) 185–202.
- [51] Z. Zhou, S. Wang, Z. Chen, M. Cai, H. Wang, Z. Li, Z. Kan, Local observation based reactive temporal logic planning of human-robot systems, *IEEE Trans. Autom. Sci. Eng.* (2023) 1–13.
- [52] Z. Zhou, Z. Chen, M. Cai, Z. Li, Z. Kan, C.-Y. Su, Vision-based reactive temporal logic motion planning for quadruped robots in unstructured dynamic environments, *IEEE Trans. Ind. Electron.* (2023) 1–10.





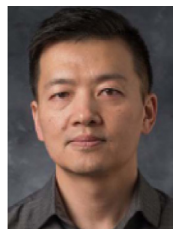
**Mingyu Cai** received the Ph.D. degree in mechanical engineering from the University of Iowa, Iowa City, IA, USA, in 2021. From 2021 to 2023, he was a Postdoctoral Associate with the Department of Mechanical Engineering at Lehigh University, Bethlehem, USA. In 2023, he worked at Honda Research Institute as a research scientist. From 2024, he is an assistant professor at the University of California, Riverside. His research interests encompass robotics, machine learning, control theory, and formal methods, with a focus on applications in motion planning, decision-making, nonlinear control, and autonomous driving.



**Zhangli Zhou** received his B.E. degree from Wuhan University of Technology, China, in 2016, and his Ph.D. degree from the University of Science and Technology of China, Hefei, China, in 2024. He is currently a Postdoctoral Fellow at the University of Science and Technology of China, with research interests in robotics and reactive task and motion planning.



**Lin Li** received the B.E. degree in mechanical engineering and automation from the Hefei University of Technology, Hefei, China, in 2021. She is currently pursuing the M.S. degree in control engineering with the University of Science and Technology of China, Hefei. Her current research interests include heterogeneous multirobot system and motion planning.



**Shaoping Xiao** received the Ph.D. degree in mechanical engineering from Northwestern University in 2003. He is currently an Associate Professor with the Department of Mechanical Engineering, University of Iowa, Iowa City, IA, USA. His primary expertise lies in computational mechanics and materials science, and one of his papers has been cited over 1000 times. In recent years, he has extended his research to artificial intelligence (AI) and its applications in engineering problem-solving. His research interests include machine learning-enhanced numerical modeling of composite materials, reinforcement learning and formal methods for

robotics control, AI-powered design of distributed reservoir systems for flood mitigation, intelligent traffic light systems, and quantum computing.



**Zhen Kan** (Senior Member, IEEE) received his Ph.D. degree in 2011 from the Department of Mechanical and Aerospace Engineering at the University of Florida. He was a Postdoctoral Research Fellow with the Air Force Research Laboratory (AFRL) at Eglin AFB and University of Florida REEF from 2012 to 2016, and was an Assistant Professor in the Department of Mechanical Engineering at the University of Iowa from 2016–2019. He is currently a Professor in the Department of Automation at the University of Science and Technology of China. His research interests include networked control systems, nonlinear control, formal methods,

and robotics. He currently serves on program committees of several internationally recognized scientific and engineering conferences and is an Associate Editor of IEEE Transactions on Automatic Control.