

---

# MODULAR DEEP REINFORCEMENT LEARNING FOR CONTINUOUS MOTION PLANNING WITH TEMPORAL LOGIC

---

Mingyu Cai<sup>1</sup>, Mohammadhosein Hasanbeig<sup>2</sup>, Shaoping Xiao<sup>1</sup>, Alessandro Abate<sup>2</sup> and Zhen Kan<sup>3</sup>

## ABSTRACT

This paper investigates the motion planning of autonomous dynamical systems modeled by Markov decision processes (MDP) with unknown transition probabilities over continuous state and action spaces. Linear temporal logic (LTL) is used to specify high-level tasks over infinite horizon, which can be converted into a limit deterministic generalized Büchi automaton (LDGBA) with several accepting sets. The novelty is to design an embedded product MDP (EP-MDP) between the LDGBA and the MDP by incorporating a synchronous tracking-frontier function to record unvisited accepting sets of the automaton, and to facilitate the satisfaction of the accepting conditions. The proposed LDGBA-based reward shaping and discounting schemes for the model-free reinforcement learning (RL) only depend on the EP-MDP states and can overcome the issues of sparse rewards. Rigorous analysis shows that any RL method that optimizes the expected discounted return is guaranteed to find an optimal policy whose traces maximize the satisfaction probability. A modular deep deterministic policy gradient (DDPG) is then developed to generate such policies over continuous state and action spaces. The performance of our framework is evaluated via an array of OpenAI gym environments.

**Keywords** Reinforcement Learning · Neural Networks · Formal Methods · Deep Learning

## 1 Introduction

The goal of motion planning is to generate valid configurations such that robotic systems can complete pre-specified tasks. In practice, however, dynamics uncertainties impose great challenges to motion planning. Markov decision processes (MDP) are often employed to model such uncertainties as transition probabilities. Growing research has been devoted to studying the motion planning modelled as an MDP when these transition probabilities are initially unknown. Reinforcement learning (RL) is a sequential decision-making process that learns optimal action policies for an unknown MDP via gathering experience samples from the MDP [1]. RL has achieved impressive results over the past few years, but often the learned solution is difficult to understand and examine by humans. Two main challenges existing in many RL applications are: (i) the design of an appropriate reward shaping mechanism to ensure correct mission specification, and (ii) the increasing sample complexity when considering continuous state and action spaces.

Temporal logics offer rich expressivity in describing complex tasks beyond traditional go-to-goal navigation for robotic systems [2]. Specifically, motion planning under linear temporal logic (LTL) constraints attracted growing research attention in the past few years [3–5]. Under the assumption of full knowledge of the MDP model, one common objective is to maximize the probability of accomplishing the given LTL task [6–8]. Once this assumption is relaxed, model-based RL is employed [9–11] by treating LTL specifications as reward shaping schemes to generate policies that satisfy LTL tasks by explicitly learning unknown transition probabilities of the MDP. This means that a model of the MDP is inferred over which an optimal policy is synthesized. However, scalability is a pressing issue for applying

<sup>1</sup>Department of Mechanical Engineering, The University of Iowa, Iowa City, IA, USA.

<sup>2</sup>Department of Computer Science, University of Oxford, Parks Road, Oxford, UK.

<sup>3</sup>Department of Automation, University of Science and Technology of China, Hefei, Anhui, China.

<sup>4</sup>Email: mingyu-cai@uiowa.edu, hosein.hasanbeig@icloud.com, shaoping-xiao@uiowa.edu, alessandro.abate@cs.ox.ac.uk, zkan@ustc.edu.cn.

model-based approaches due to the need to store the learned model. On the other hand, by relaxing the need to construct an MDP model, model-free RL is recently adopted where appropriate reward shaping schemes are proposed [12–21]. Signal temporal logic is also proposed in [22] and [23] where a task-guided reward function is proposed based on the robustness degrees. Despite the recent progresses, the aforementioned works can not handle many real-world applications that perform in high-dimensional continuous state and action spaces. In a pioneer work [24], Deep Q Network (DQN) addressed high-dimensional state space and is capable of human-level performance on many Atari video games. However, DQN can only handle discrete and low-dimensional action spaces. By leveraging actor-critic methods, deep networks, and the policy gradient methods, deep deterministic policy gradient (DDPG) was proposed to approximate optimal policies over a continuous action space to improve the learning performance [25].

In this work, we consider motion planning under LTL task specifications in continuous state and action spaces when the MDP is fully unknown. An unsupervised one-shot and on-the-fly DDPG-based motion planning framework is developed to learn the state of an underlying structure without explicitly constructing the model. The high-level LTL task is converted to a limit deterministic generalized Büchi automaton (LDGBA) [26] acting as task-guided reward shaping scheme and decomposing complex tasks into low-level and achievable modules. In particular, we consider LTL specifications over the infinite horizon, whose behaviors can be regarded as a repetitive pattern [27, 28], which consists of infinite rounds of visits of the accepting sets of LDGBA.

**Related works:** When considering deep RL with formal methods, deterministic finite automaton (DFA) was applied as reward machines in [29], [30] and [31], where DQN was employed to generate optimal policies. In [32], a truncated linear temporal logic was considered and its robustness degree was used as the reward signal to facilitate learning, based on which proximal policy optimization (PPO) was applied to obtain the optimal policy. However, [29–34] only consider tasks over finite horizons. In contrast, this work extends previous research to tasks over the infinite horizon, where finite horizon motion planning can be regarded as a special case of the infinite horizon setting. Along this line of research, the most relevant works include [35–41]. In [36], LTL constraints were translated to Deterministic Rabin Automata (DRA), which might fail to find policies that maximize LTL satisfaction probability [16]. The work in [39] proposed a binary vector to record the visited accepting sets of LDGBA and designed a varying reward function to improve the satisfaction of LTL specifications. However, the maximum probability of task satisfaction cannot be guaranteed and the standard DDPG algorithm cannot distinguish the sub-task module in [39], resulting in an unsatisfactory success rate. Modular DDPG that jointly optimizes LTL sub-policies was first introduced in [35, 38]. However, [35, 38] do not explicitly record visited or unvisited accepting sets of LDGBA in each round of repetitive pattern over the infinite horizon, which might be essential for synthesising a pure deterministic policy [19].

**Contributions:** The contributions of this work are multi-fold. In contrast to most existing works that consider either a discrete state space or a discrete action space, this paper proposes a modular DDPG architecture integrated with potential functions [42]. This allows our method to generate optimal policies that efficiently solve LTL motion planning of an unknown MDP over continuous state and action spaces. The novelty is to construct an embedded product MDP (EP-MDP) to record unvisited accepting sets of the automaton at each round of the repeated visiting pattern, by introducing a tracking-frontier function. Such a design ensures task completion by encouraging the satisfaction of LDGBA accepting conditions. To facilitate learning of optimal policies, the designed reward is enhanced with potential functions that effectively guide the agent toward task satisfaction without adding extra hyper-parameters to the algorithm. Unlike [19], rigorous analysis shows that the maximum probability of task satisfaction can be guaranteed. Compared to approaches based on limit deterministic Büchi automata (LDBA), e.g., [16, 17], LDGBA has several accepting sets while LDBA only has one accepting set which can result in sparse rewards during training. Consequently, our approach can maintain a higher density of training rewards by assigning positive rewards to the accepting states. Moreover, the proposed method does not require the construction of full EP-MDP, and its states can be obtained on-the-fly to generate optimal policies.

In summary, our approach can find the optimal policy in continuous state and action spaces to satisfy LTL specifications over infinite horizon with approximate maximum probability.

## 2 Preliminaries

### 2.1 Continuous Labeled MDP and Reinforcement Learning

A continuous labeled MDP is a tuple  $\mathcal{M} = (S, A, p_S, \Pi, L, A)$ , where  $S \subseteq \mathbb{R}^n$  is a continuous state space,  $A \subseteq \mathbb{R}^m$  is a continuous action space,  $\Pi$  is a set of atomic propositions,  $L : S \rightarrow 2^\Pi$  is a labeling function, and  $p_S : \mathfrak{B}(\mathbb{R}^n) \times A \times S \rightarrow [0, 1]$  is a Borel-measurable conditional transition kernel such that  $p_S(\cdot | s, a)$  is a probability measure of  $s \in S$  and  $a \in A$  over the Borel space  $(\mathbb{R}^n, \mathfrak{B}(\mathbb{R}^n))$ , where  $\mathfrak{B}(\mathbb{R}^n)$  is the set of all Borel sets on  $\mathbb{R}^n$ . The transition probability  $p_S$  captures the motion uncertainties of the agent. It is assumed that  $p_S$  is not known *a priori*, and the agent can only observe its current state and the associated label of the current state.

A deterministic policy  $\xi$  of continuous MDP is a function  $\xi : S \rightarrow A$  that maps each state to an action over the action space  $A$ . The MDP  $\mathcal{M}$  evolves by taking an action  $a_i$  based on the rule  $\xi_i$  at each stage  $i$ , and thus the control policy  $\xi = \xi_0 \xi_1 \dots$  is a sequence of rules, which yields a path  $s = s_0 s_1 s_2 \dots$  over  $\mathcal{M}$  with the transition  $s_i \xrightarrow{a_i} s_{i+1}$  that exists, i.e.,  $s_{i+1}$  belongs to the smallest Borel set  $B$  such that  $p_S(B | s_i, a_i) = 1$ . The control policy  $\xi$  is memoryless if each  $\xi_i$  only depends on its current state, and  $\xi$  is a finite memory policy if  $\xi_i$  depends on its past history.

Let  $\Lambda : S \times A \times S \rightarrow \mathbb{R}$  denote a reward function. Given a discounting function  $\gamma : S \times A \times S \rightarrow \mathbb{R}$ , the expected discounted return under policy  $\xi$  starting from  $s \in S$  is defined as

$$U^\xi(s) = \mathbb{E}^\xi \left[ \sum_{i=0}^{\infty} \gamma^i(s_i, a_i, s_{i+1}) \cdot \Lambda(s_i, a_i, s_{i+1}) \mid s_0 = s \right].$$

An optimal policy  $\xi^*$  maximizes the expected return for each state  $s \in S$ , i.e.,

$$\xi^* = \arg \max_{\xi} U^\xi(s).$$

The function  $U^\xi(s)$  is often referred to as the value function under policy  $\xi$ . If the MDP is not fully known, but the state and action spaces are countably finite, tabular approaches are usually employed [43]. However, traditional tabular RL methods are not applicable to MDPs with continuous state and action spaces. In this work, we propose a policy gradient method that relies on deep neural networks to parameterize the policy model.

## 2.2 LTL and Limit-Deterministic Generalized Büchi Automaton

Linear temporal logic (LTL) is a formal language that is widely used to describe complex mission tasks [2]. The semantics of an LTL formula are interpreted over a word, which is an infinite sequence  $o = o_0 o_1 \dots$  where  $o_i \in 2^\Pi$  for all  $i \geq 0$ , and  $2^\Pi$  represents the power set of  $\Pi$ . Denote by  $o \models \phi$  if the word  $o$  satisfies the LTL formula  $\phi$ . Given an LTL specification, its satisfaction can be evaluated by an LDGBA [26]. An LDGBA is a sub-class of generalized Büchi automata (GBA) that can express the set of words of an LTL formula. In this work, we restrict our attention to LTL formulas that exclude the *next* temporal operator, which is not meaningful for continuous time execution [44].

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

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 acceptance condition, i.e.,  $\inf(q) \cap F_i \neq \emptyset, \forall i \in \{1, \dots, f\}$ , where  $\inf(q)$  denotes the infinite part of  $q$ .

**Definition 2.** A GBA is an LDGBA if the transition function  $\delta$  is extended to  $Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$ , and the state set  $Q$  is partitioned into a deterministic set  $Q_D$  and a non-deterministic set  $Q_N$ , i.e.,  $Q_D \cup Q_N = Q$  and  $Q_D \cap Q_N = \emptyset$ , where

- the state transitions in  $Q_D$  are total and restricted within it, i.e.,  $|\delta(q, \alpha)| = 1$  and  $\delta(q, \alpha) \subseteq Q_D$  for every state  $q \in Q_D$  and  $\alpha \in \Sigma$ ,
- An  $\epsilon$ -transition is not allowed in the deterministic set, i.e., for any  $q \in Q_D, \delta(q, \epsilon) = \emptyset$ , and
- the accepting states are only in the deterministic set, i.e.,  $F_i \subseteq Q_D$  for every  $F_i \in F$ .

In Definition 2,  $\epsilon$ -transitions are only for state transitions from  $Q_N$  to  $Q_D$ , without reading an alphabet. We used OWL [45], an easily-accessible tool, to generate LDGBAs in this work. In the following analysis, we use  $\mathcal{A}_\phi$  to denote the LDGBA corresponding to an LTL formula  $\phi$ .

**Definition 3.** A non-accepting sink component  $Q_{sink} \subseteq Q$  of an LDGBA is a strongly connected directed graph induced by a set of states, such that the accepting condition can not be satisfied if starting from any state in  $Q_{sink}$ .

## 3 Problem Statement

Consider a robot that performs a mission described by an LTL formula  $\phi$ . For instance, as shown in Fig. 1, the helicopter provides the map consisting of areas of interest, based on which the ground Mars rover is tasked to visit all regions marked with rectangles. Due to the complex terrain of Mars surface, there exist motion uncertainties during the rover movement. The interaction of the robot with the environment is modeled by a continuous MDP  $\mathcal{M}$ , which can be used

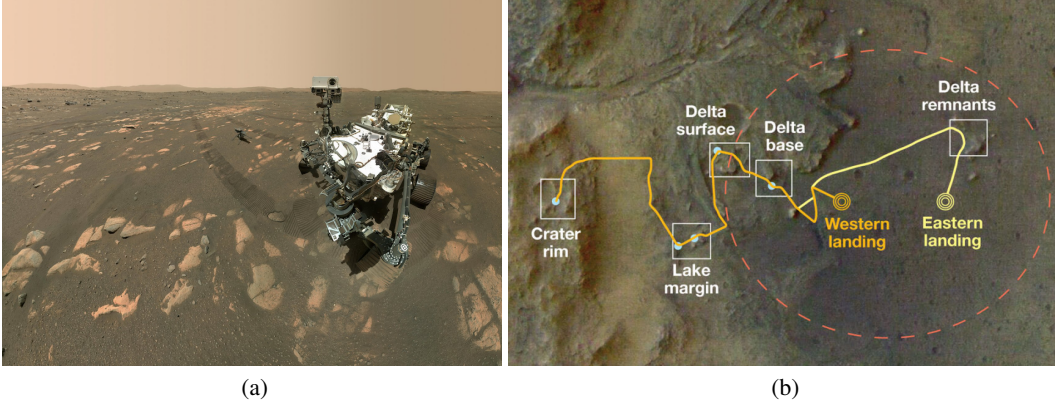


Figure 1: Example of Mars exploration (courtesy of NASA)

to model decision-making problems for general robotic systems. Under a policy  $\xi = \xi_0 \xi_1 \dots$ , the induced path over  $\mathcal{M}$  is  $s_\infty^\xi = s_0 \dots s_i s_{i+1} \dots$ . Let  $L(s_\infty^\xi) = l_0 l_1 \dots$  be the sequence of labels associated with  $s_\infty^\xi$  such that  $l_i \in L(s_i)$ . Denote the satisfaction relation by  $L(s_\infty^\xi) \models \phi$  if the induced trace satisfies  $\phi$ . The probabilistic satisfaction under the policy  $\xi$  from an initial state  $s_0$  can then be defined as

$$\Pr_M^\xi(\phi) = \Pr_M^\xi(L(s_\infty^\xi) \models \phi | s_\infty^\xi \in \mathcal{S}_\infty^\xi), \quad (1)$$

where  $\mathcal{S}_\infty^\xi$  is a set of admissible paths from the initial state under the policy  $\xi$ .

**Assumption 1.** It is assumed that there exists at least one policy whose induced traces satisfy the task  $\phi$  with non-zero probability.

Assumption 1 is mild and widely employed in the literature (cf. [9, 16, 17]), which indicates  $\phi$  can be completed. Based on Assumption 1, this work considers the following problem.

**Problem 1.** Given an LTL-specified task  $\phi$  and a continuous-state continuous-action labeled MDP  $\mathcal{M}$  with unknown transition probabilities, the goal is to learn a policy  $\xi^*$  that maximizes the satisfaction probability in the limit, i.e.,  $\xi^* = \arg \max_\xi \Pr_M^\xi(\phi)$ .

## 4 Automaton Design

To address Problem 1, Section 4.1 presents the construction of the EP-MDP between an MDP and an LDGBA. The advantages of incorporating EP-MDP are discussed in Section 4.2.

### 4.1 Embedded Product MDP

Given an LDGBA  $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ , a tracking-frontier set  $T$  is designed to keep track of unvisited accepting sets. Particularly,  $T$  is initialized as  $T_0 = F$  and  $\mathcal{B}$  is a Boolean variable, which is then updated according to the following rule:

$$(T, \mathcal{B}) = f_V(q, T) = \begin{cases} (T \setminus F_j, \text{False}), & \text{if } q \in F_j \text{ and } F_j \in T, \\ (F \setminus F_j, \text{True}) & \text{if } q \in F_j \text{ and } T = \emptyset, \\ (T, \text{False}), & \text{otherwise.} \end{cases} \quad (2)$$

Once a state  $q \in F_j$  is visited,  $F_j$  will be removed from  $T$  by rendering function  $f_V(q, T)$ . If  $T$  becomes empty before visiting set  $F_j$ , it will be reset as  $F \setminus F_j$ . Since the acceptance condition of LDGBA requires to infinitely visit all accepting sets, we call it one round if all accepting sets have been visited (i.e., a round ends if  $T$  becomes empty). The second output of  $f_V(q, T)$  is to indicate whether all accepting sets have been visited in current round, which is applied in Section 5.3 to design the potential function. Based on (2), the EP-MDP is constructed as follows.

**Definition 4.** Given an MDP  $\mathcal{M}$  and an LDGBA  $\mathcal{A}_\phi$ , the EP-MDP is defined as  $\mathcal{P} = \mathcal{M} \times \mathcal{A}_\phi = (X, U^\mathcal{P}, p^\mathcal{P}, x_0, F^\mathcal{P}, T, f_V, \mathcal{B})$ , where  $X = S \times Q \times 2^F$  is the set of product states and  $2^F$  denotes all subsets of  $F$  for  $\mathcal{A}_\phi$ , i.e.,  $x = (s, q, T) \in X$ ;  $U^\mathcal{P} = A \cup \{\epsilon\}$  is the set of actions, where the  $\epsilon$ -actions are only

---

**Algorithm 1** generating a random run of EP-MDP

---

```

1: procedure INPUT:  $(\mathcal{M}, \mathcal{A}_\phi, f_V, T$  and length  $L$ )
   Output: A valid run  $\mathbf{x}_P$  with length  $L$  in  $\mathcal{P}$ 
2: set  $x_0 = (s_0, q_0, T_0)$ ,  $s_{cur} = s_0$  and  $\mathbf{x}_P = (x_0)$ 
3: set  $T = F$  and  $count = 1$ 
4: while  $count \leq L$  do
5:   set  $x_{suc} = \emptyset$ 
6:   obtain  $s_{suc}$  from  $p_S(\cdot | s_{cur}, a)$  by random sampling action  $a$ 
7:   for each successor  $q_{next}$  of  $q_{cur}$  in  $\mathcal{A}_\phi$  do
8:     if  $\delta(q_{cur}, L(s_{cur})) = q_{next}$  or  $q_{next} \in \delta(q_{cur}, \epsilon)$  then
9:        $q_{suc} \leftarrow q_{next}$ 
10:       $x_{suc} \leftarrow (s_{suc}, q_{suc}, T)$ 
11:      check if  $x_{suc}$  is an accepting state
12:       $T = f_V(q_{suc}, T)$ 
13:      break
14:     end if
15:   end for
16:   if no successor  $q_{suc}$  found then
17:     fail to generate run with length  $L$ 
18:   end if
19:   add state  $x_{suc}$  to  $\mathbf{x}_P$ 
20:    $q_{cur} \leftarrow q_{suc}$  and  $s_{cur} \leftarrow s_{suc}$ 
21:    $count++$ ;
22: end while
23: end procedure

```

---

allowed for transitions from  $Q_N$  to  $Q_D$ ;  $x_0 = (s_0, q_0, T_0)$  is the initial state;  $F^P = \{F_1^P, F_2^P \dots F_f^P\}$  where  $F_j^P = \{(s, q, T) \in X | q \in F_j \wedge F_j \subseteq T\}$ ,  $j = 1, \dots, f$ , is a set of accepting states;  $p^P$  is the transition kernel for any transition  $p^P(x, u^P, x')$  with  $x = (s, q, T)$  and  $x' = (s', q', T)$  such that: (1)  $p^P(x, u^P, x') = p_S(s' | s, a)$  if  $s' \sim p_S(\cdot | s, a)$ ,  $\delta(q, L(s)) = q'$  where  $u^P = a \in A$  (2)  $p^P(x, u^P, x') = 1$  if  $u^P \in \{\epsilon\}$ ,  $q' \in \delta(q, \epsilon)$  and  $s' = s$ ; and (3)  $p^P(x, u^P, x') = 0$  otherwise. After completing each transition  $q' = \delta(q, \alpha)$  based on  $p^P$ ,  $T$  is synchronously updated as  $(T, B) = f_V(q', T)$  by (2).

The state-space is embedded with the tracking-frontier set  $T$  that can be practically represented via one-hot encoding based on the indices of the accepting set. Compared to the standard construction of product MDPs, any state of the accepting set in EP-MDP, e.g.,  $(s, q, T) \in F_j^P$ , requires the automaton state to satisfy  $q \in F_j \wedge F_j \subseteq T$ , and embedded tracking frontier set  $T$  is updated based on (2) after each transition. Consequently, to satisfy the accepting condition, the agent is encouraged to visit all accepting sets.

In this work, the EP-MDP is only used for theoretical analyses and it is not constructed in practice. The EP-MDP captures the intersections between all feasible paths over  $\mathcal{M}$  and all words accepted to  $\mathcal{A}_\phi$ , facilitating the identification of admissible agent actions that satisfy task  $\phi$ .

Algorithm 1 shows the procedure of obtaining a valid run  $\mathbf{x}_P$  on-the-fly within EP-MDP by randomly selecting an action. After each transition, the accepting state is determined based on the Definition 4 and the  $T$  is synchronously updated (line 7-15). Such property is the innovation of EP-MDP that encourages all accepting sets to be visited in each round. Since the action is selected randomly, there is a non-zero probability that  $\mathbf{x}_P$  violates the LTL task (line 14-16).

Let  $\pi$  denote a policy over  $\mathcal{P}$  and denote by  $\mathbf{x}_\infty^\pi = x_0 \dots x_i x_{i+1} \dots$  the infinite path generated by  $\pi$ . A path  $\mathbf{x}_\infty^\pi$  is accepted if  $\inf(\mathbf{x}_\infty^\pi) \cap F_i^P \neq \emptyset, \forall i \in \{1, \dots, f\}$ . We denote  $\Pr^\pi[x \models \text{Acc}_p]$  as the probability of satisfying the accepting condition of  $\mathcal{P}$  under policy  $\pi$ , and denote  $\Pr_{max}[x \models \text{Acc}_p] = \max_\pi \Pr_M^\pi(\text{Acc}_p)$  as the maximum probability of satisfying the accepting condition of  $\mathcal{P}$ . Let  $\pi^*$  denote an optimal policy that maximizes the expected discounted return over  $\mathcal{P}$ , i.e.,  $\pi^* = \arg \max_\pi U^\pi(s)$ . Note that the memory-less policy  $\pi^*$  over  $\mathcal{P}$  yields a finite-memory policy  $\xi^*$  over  $\mathcal{M}$ , allowing us to reformulate Problem 1 into the following:

**Problem 2.** Given a user-specified LTL task  $\phi$  and a general and unknown continuous-state continuous-action labeled MDP, the goal is to asymptotically find a policy  $\pi^*$  satisfying the acceptance condition of  $\mathcal{P}$  with maximum probability, i.e.,  $\Pr^{\pi^*}[x \models \text{Acc}_p] = \Pr_{max}[x \models \text{Acc}_p]$ .

## 4.2 Properties of EP-MDP

Consider a sub-EP-MDP  $\mathcal{P}'_{(X', U')}$ , where  $X' \subseteq X$  and  $U' \subseteq U^P$ . If  $\mathcal{P}'_{(X', U')}$  is a maximum end component (MEC) of  $\mathcal{P}$  and  $X' \cap F_i^P \neq \emptyset, \forall i \in \{1, \dots, f\}$ , then  $\mathcal{P}'_{(X', U')}$  is called an accepting maximum end component (AMEC) of

$\mathcal{P}$ . Once a path enters an AMEC, the subsequent path will stay within it by taking restricted actions from  $U'$ . There exist policies such that any state  $x \in X'$  can be visited infinitely often. As a result, satisfying task  $\phi$  is equivalent to reaching an AMEC. Moreover, a MEC that does not intersect with any accepting sets is called a rejecting accepting component (RMEC) and a MEC intersecting with only partial accepting sets is called a neutral maximum end component (NMEC) [2].

**Definition 5.** Let  $MC_{\mathcal{P}}^{\pi}$  denote the Markov chain induced by a policy  $\pi$  on  $\mathcal{P}$ , whose states can be represented by a disjoint union of a transient class  $\mathcal{T}_{\pi}$  and  $n_R$  closed irreducible recurrent classes  $\mathcal{R}_{\pi}^j$ ,  $j \in \{1, \dots, n_R\}$  [46].

**Lemma 1.** Given an EP-MDP  $\mathcal{P} = \mathcal{M} \times \mathcal{A}_{\phi}$ , the recurrent class  $\mathcal{R}_{\pi}^j$  of  $MC_{\mathcal{P}}^{\pi}$ ,  $\forall j \in \{1, \dots, n_R\}$ , induced by  $\pi$  satisfies one of the following conditions:  $R_{\pi}^j \cap F_i^{\mathcal{P}} \neq \emptyset, \forall i \in \{1, \dots, f\}$ , or  $R_{\pi}^j \cap F_i^{\mathcal{P}} = \emptyset, \forall i \in \{1, \dots, f\}$ .

*Proof.* The strategy of the following proof is based on contradiction. Assume there exists a policy such that  $R_{\pi}^j \cap F_k^{\mathcal{P}} \neq \emptyset, \forall k \in K$ , where  $K$  is a subset of  $2^{\{1, \dots, f\}} \setminus \{\{1, \dots, f\}, \emptyset\}$ . As discussed in [46], for each state in recurrent class, it holds that  $\sum_{n=0}^{\infty} p^n(x, x) = \infty$ , where  $x \in R_{\pi}^j \cap F_k^{\mathcal{P}}$  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 EP-MDP, once  $x_k$  is visited, the corresponding  $F_k$  of  $F_k^{\mathcal{P}}$  is removed from  $T$ , and the tracking set  $T$  will not be reset until all accepting sets have been visited. As a result, neither  $q_k \in F_k$  nor  $x_k = (s, q_k, T) \in R_{\pi}^j \cap F_k^{\mathcal{P}}$  with  $s \in S$  will occur infinitely, which contradicts the property  $\sum_{n=0}^{\infty} p^n(x_k, x_k) = \infty$ .  $\square$

Lemma 1 indicates that, for any policy, all accepting sets will be placed either in the transient class or in the recurrent classes.

## 5 Learning-based Control Synthesis

In the following, we discuss a base reward design and present rigorous analysis to show how such a design can guide the RL-agent over the EP-MDP to find an optimal policy whose traces satisfies the LTL task with maximum probability. In order to improve the reward density, Section 5.2 proposes a reward shaping process via integrating a potential function under which the optimal policies remain invariant. Finally, Section 5.3 shows how to apply the shaped reward function with DDPG to construct a modular DDPG architecture and effectively solve Problem 2.

### 5.1 Base Reward

Let  $F_U^{\mathcal{P}}$  denote the union of accepting states, i.e.,  $F_U^{\mathcal{P}} = \{x \in X \mid x \in F_i^{\mathcal{P}}, \forall i \in \{1, \dots, f\}\}$ . For each transition  $(x, u^{\mathcal{P}}, x')$  in the EP-MDP, the reward and discounting function only depend on current state  $x$ , i.e.,  $R(x, u^{\mathcal{P}}, x') = R(x)$  and  $\gamma(x, u^{\mathcal{P}}, x') = \gamma(x)$ .

Inspired by [17], we propose a reward function as:

$$R(x) = \begin{cases} 1 - r_F, & \text{if } x \in F_U^{\mathcal{P}}, \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

and a discounting function as

$$\gamma(x) = \begin{cases} r_F, & \text{if } x \in F_U^{\mathcal{P}}, \\ \gamma_F, & \text{otherwise,} \end{cases} \quad (4)$$

where  $r_F(\gamma_F)$  is a function of  $\gamma_F$  satisfying  $\lim_{\gamma_F \rightarrow 1^-} r_F(\gamma_F) = 1$  and  $\lim_{\gamma_F \rightarrow 1^-} \frac{1 - \gamma_F}{1 - r_F(\gamma_F)} = 0$ .

Given a path  $\mathbf{x}_t = x_t x_{t+1} \dots$  starting from  $x_t$ , the return is denoted by

$$\mathcal{D}(\mathbf{x}_t) := \sum_{i=0}^{\infty} \left( \prod_{j=0}^{i-1} \gamma(\mathbf{x}_t[t+j]) \cdot R(\mathbf{x}_t[t+i]) \right) \quad (5)$$

where  $\prod_{j=0}^{-1} := 1$  and  $\mathbf{x}_t[t+i]$  denotes the  $(i+1)$ th state in  $\mathbf{x}_t$ . Based on (5), the expected return of any state  $x \in X$  under policy  $\pi$  can be defined as

$$U^{\pi}(x) = \mathbb{E}^{\pi}[\mathcal{D}(\mathbf{x}_t) \mid \mathbf{x}_t[t] = x]. \quad (6)$$

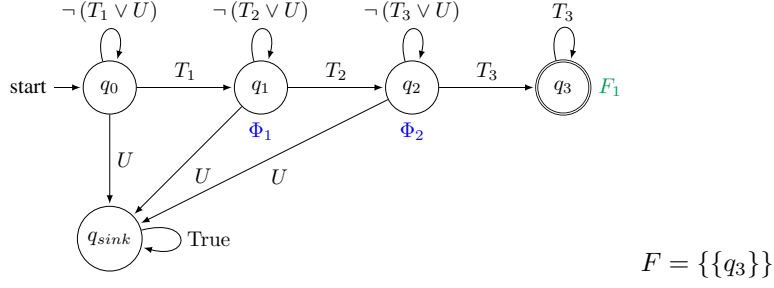


Figure 2: LDGBA  $\mathcal{A}_{\varphi_P}$  expressing  $\varphi_P = \diamond(T1 \wedge \diamond(T2 \wedge \diamond T3)) \wedge \neg \square U$

A bottom strongly connected component (BSCC) of the Markov chain  $MC_{\mathcal{P}}^{\pi}$  (Definition 5) is a strongly connected component with no outgoing transitions.

**Lemma 2.** *For any path  $\mathbf{x}_t$  and  $\mathcal{D}(\mathbf{x}_t)$  in (5), it holds that  $0 \leq \gamma_F \cdot \mathcal{D}(\mathbf{x}_t[t+1:]) \leq \mathcal{D}(\mathbf{x}_t) \leq 1 - r_F + r_F \cdot \mathcal{D}(\mathbf{x}_t[t+1:]) \leq 1$ , where  $\mathbf{x}_t[t+1:]$  denotes the suffix of  $\mathbf{x}_t$  starting from  $x_{t+1}$ . Let  $BSCC(MC_{\mathcal{P}}^{\pi})$  denote the set of all BSCCs of an induced Markov chain  $MC_{\mathcal{P}}^{\pi}$  and let  $X_{\mathcal{P}}^{\pi}$  denotes the set of accepting states that belongs to a BSCC of  $MC_{\mathcal{P}}^{\pi}$  s.t.  $X_{\mathcal{P}}^{\pi} := \{x \in X \mid x \in F_{\mathcal{U}}^{\mathcal{P}} \cap BSCC(MC_{\mathcal{P}}^{\pi})\}$ . Then, for any states  $x \in X_{\mathcal{P}}^{\pi}$ , it holds that  $\lim_{\gamma_F \rightarrow 1^-} U^{\pi}(x) = 1$ .*

The proof of Lemma 2 is omitted since it is a straightforward extension of Lemma 2 and Lemma 3 in [17], by replacing LDBA with LDGBA. Since we apply the LDGBA with several accepting sets which might result in more complicated situations, e.g., AMEC, NMEC and RMEC, we can not obtain the same results as in [17]. We then establish the following theorem which is one of the main contributions.

**Theorem 1.** *Given the EP-MDP  $\mathcal{P} = \mathcal{M} \times \mathcal{A}_{\phi}$ , for any state  $x \in X$ , the expected return under any policy  $\pi$  satisfies*

$$\exists i \in \{1, \dots, f\}, \lim_{\gamma_F \rightarrow 1^-} U^{\pi}(x) = \Pr^{\pi}[\diamond F_i^{\mathcal{P}}], \quad (7)$$

where  $\Pr^{\pi}[\diamond F_i^{\mathcal{P}}]$  is the probability that the paths starting from state  $x$  will eventually intersect a  $F_i^{\mathcal{P}} \in F^{\mathcal{P}}$ .

*Proof.* Proof can be found in Appendix 8.1. □

Next, we will show in the following sections how Lemma 1 and Theorem 1 can be leveraged to enforce the RL-agent satisfying the accepting condition of  $\mathcal{P}$ .

**Theorem 2.** *Consider an MDP  $\mathcal{M}$  and an LDGBA  $\mathcal{A}_{\phi}$  corresponding to an LTL formula  $\phi$ . Based on Assumption 1, there exists a discount factor  $\underline{\gamma}$ , with which any optimization method for (6) with  $\gamma_F > \underline{\gamma}$  and  $r_F > \underline{\gamma}$  can obtain a policy  $\bar{\pi}$ , such that the induced run  $r_{\bar{\mathcal{P}}}^{\bar{\pi}}$  satisfies the accepting condition  $\mathcal{P}$  with non-zero probability in the limit.*

*Proof.* Proof can be found in Appendix 8.2. □

Theorem 2 proves that by selecting  $\gamma_F > \underline{\gamma}$  and  $r_F > \underline{\gamma}$ , optimizing the expected return in (6) can find a policy satisfying the given task  $\phi$  with non-zero probability.

**Theorem 3.** *Given an MDP  $\mathcal{M}$  and an LDGBA  $\mathcal{A}_{\phi}$ , by selecting  $\gamma_F \rightarrow 1^-$ , the optimal policy in the limit  $\pi^*$  that maximizes the expected return (6) of the corresponding EP-MDP also maximizes the probability of satisfying  $\phi$ , i.e.,  $\Pr^{\pi^*}[x \models Acc_{\mathcal{P}}] = \Pr_{max}[x \models Acc_{\mathcal{P}}]$ .*

*Proof.* Since  $\gamma_F \rightarrow 1^-$ , we have  $\gamma_F > \underline{\gamma}$  and  $r_F > \underline{\gamma}$  from Theorem 2. There exists an induced run  $r_{\bar{\mathcal{P}}}^{\pi^*}$  satisfying the accepting condition of  $\mathcal{P}$ . According to Theorem 1,  $\lim_{\gamma_F \rightarrow 1^-} U^{\pi^*}(x)$  is exactly equal to the probability of visiting the accepting sets of an AMEC. Optimizing  $\lim_{\gamma_F \rightarrow 1^-} U^{\pi^*}(x)$  is equal to optimizing the probability of entering AMECs. □

## 5.2 Reward Shaping

Since the base reward function in Section 5.1 is always zero for  $x \notin F_U^{\mathcal{P}}$ , the reward signal might become sparse. To resolve this, we propose a potential function  $\Phi : X \rightarrow \mathbb{R}$ , and transform the reward as follows:

$$R'(x, u^{\mathcal{P}}, x') = R(x) + \gamma(x) \cdot \Phi(x') - \Phi(x) \quad (8)$$

As shown in [42], given any MDP model, e.g., EP-MDP  $\mathcal{P}$ , transforming the reward function using a potential function  $\Phi(x)$  as in  $R'(x, u^{\mathcal{P}}, x')$  will not change the set of optimal policies. Thus, a real-valued function  $\Phi(x)$  will improve the learning performance while guaranteeing that the resulted policies via  $R'(x, u^{\mathcal{P}}, x')$  are still optimal with respect to the base reward function  $R(x)$ .

Given  $\mathcal{P} = \mathcal{M} \times \mathcal{A}_\phi = (X, U^{\mathcal{P}}, p^{\mathcal{P}}, x_0, F^{\mathcal{P}}, T, f_V, \mathcal{B})$  with  $\mathcal{A}_\phi = (Q, \Sigma, \delta, q_0, F)$ , let  $F_U = \{q \in Q \mid q \in F_i, \forall i \in \{1, \dots, f\}\}$  denote the union of automaton accepting states. For the states of  $\mathcal{P}$  whose automaton states belong to  $Q \setminus (F_U \cup q_0 \cup Q_{sink})$ , it is desirable to assign positive rewards when the agent first visits them and assign large value of reward to the accepting states to enhance the convergence of neural network. This is because starting from the automaton initial state, any automaton state that can reach any of the accepting sets has to be explored. To this end, an automaton tracking-frontier set  $T_\Phi$  is designed to keep track of unvisited automaton components  $Q \setminus (q_0 \cup Q_{sink})$ , and  $T_\Phi$  is initialized as  $T_{\Phi_0} = Q \setminus (q_0 \cup Q_{sink})$ . The set  $T_{\Phi_0}$  is then updated after each transition  $((s, q, T), u^{\mathcal{P}}, (s', q', T))$  of  $\mathcal{P}$  as:

$$f_\Phi(q', T_\Phi) = \begin{cases} T_\Phi \setminus q', & \text{if } q \in T_\Phi, \\ T_{\Phi_0} \setminus q' & \text{if } \mathcal{B} = \text{True}, \\ T_\Phi, & \text{otherwise.} \end{cases} \quad (9)$$

The set  $T_\Phi$  will only be reset when  $\mathcal{B}$  in  $f_V$  becomes True, indicating that all accepting sets in the current round have been visited. Based on (9), the potential function  $\Phi(x)$  for  $x = (s, q, T)$  is constructed as:

$$\Phi(x) = \begin{cases} \eta_\Phi \cdot (1 - r_F), & \text{if } q \in T_\Phi, \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

where  $\eta_\Phi > 0$  is the shaing parameter. Intuitively, the value of potential function for unvisited and visited states in  $T_{\Phi_0}$  is equal to  $\eta_\Phi \cdot (1 - r_F)$  and 0 respectively, in each round, which will guide the system finally reach any of the accepting sets.

To illustrate the idea, consider an LTL formula  $\varphi_{\mathcal{P}} = \diamond(T1 \wedge \diamond(T2 \wedge \diamond T3)) \wedge \neg \square U$ , e.g., the agent should always avoid obstacles and visit T1, T2, T3 sequentially. The corresponding LDGBA  $\mathcal{A}_{\varphi_{\mathcal{P}}}$  is shown in Fig. 2. Let's denote any product state with the same automaton component and an arbitrary MDP state as  $x(\llbracket s \rrbracket, q, T)$ , where the MDP component can be different. Note that only state  $q_3$  belongs to the accepting set and the original reward function  $R(x) = 0$ ,  $\forall x = (\llbracket s \rrbracket, q_{12}, T)$  with  $q_{12} \neq q_3$ . However, for  $x_1 = (\llbracket s \rrbracket, q_1, T)$  and  $x_2 = (\llbracket s \rrbracket, q_2, T)$ , we have  $\Phi(x_1) \neq 0$  and  $\Phi(x_2) \neq 0$  if the corresponding automaton component has not been visited yet (i.e.,  $q_1, q_2$  still in  $T_\Phi$ ) by (9) and (10). For instance, given a run  $\mathbf{x} = (\llbracket s \rrbracket, q_0, T) u_0^{\mathcal{P}}(\llbracket s \rrbracket, q_1, T) u_1^{\mathcal{P}}(\llbracket s \rrbracket, q_2, T) u_2^{\mathcal{P}}(\llbracket s \rrbracket, q_3, T)$ , the associated shaped reward for each transition is  $R'((\llbracket s \rrbracket, q_0, T), u_0^{\mathcal{P}}, (\llbracket s \rrbracket, q_1, T)) = R'((\llbracket s \rrbracket, q_1, T), u_1^{\mathcal{P}}, (\llbracket s \rrbracket, q_2, T)) = \gamma_F \cdot (1 - r_F)$ , where  $\gamma_F$  is obtained based on (4). Note that every time, if the  $q_i$  has been visited, it'll be removed from  $T_\Phi$ , resulting in  $\Phi(x) = 0$  in future transitions until  $T_\Phi$  is reset.

## 5.3 Modular Deep Deterministic Policy Gradient

To deal with continuous-state and continuous-action MDPs, a deep deterministic policy gradient (DDPG) [25] is adopted in this work to approximate the current deterministic policy via a parameterized function  $\pi(x | \theta^u)$  called actor. The actor is a deep neural network whose set of weights are  $\theta^u$ . The critic function also applies a deep neural network with parameters  $\theta^Q$  to approximate action-value function  $Q(x, u^{\mathcal{P}} | \theta^Q)$ , which is updated by minimizing the following loss function:

$$L(\theta^Q) = \mathbb{E}_{s \sim \rho^\beta} \left[ \left( y - Q(x, \pi(x | \theta^u) | \theta^Q) \right)^2 \right], \quad (12)$$



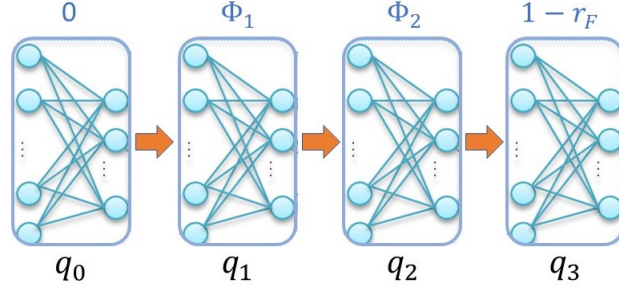


Figure 3: Modular DDPG for LDGBA  $\mathcal{A}_{\varphi_P}$  in Fig. 2

---

### Algorithm 2 Modular DDPG

---

1: **procedure** INPUT: (MDP  $\mathcal{M}$ , LDGBA  $\mathcal{A}_\phi$ )

Output: modular DDPG for optimal policy  $\pi^*$

Initialization:  $|Q|$  actor  $\pi_{q_i}(x|\theta^{u_{q_i}})$  and critic networks  $Q_{q_i}(x, u^P|\theta^{Q_{q_i}})$  with arbitrary weights  $\theta^{u_{q_i}}$  and  $\theta^{Q_{q_i}}$  for all  $q_i \in Q$ ;  $|Q|$  corresponding target networks  $\pi'_{q_i}(x|\theta'^{u_{q_i}})$  and  $Q'_{q_i}(x, u^P|\theta'^{Q_{q_i}})$  with weights  $\theta'^{u_{q_i}}$  and  $\theta'^{Q_{q_i}}$  for each  $q_i \in Q$ , respectively;  $|Q|$  replay buffers  $B_{q_i}$ ;  $|Q|$  random processes noise  $N_{q_i}$

2: set  $r_F = 0.99$  and  $\gamma_F = 0.9999$  to determine  $R(x)$  and  $\gamma(x)$

3: set maximum episodes  $E$  and iteration number  $\tau$

4: **for** each episode in  $E$  **do**

5: set  $t = 0$

6: sample an initial state  $s_0$  of  $\mathcal{M}$  and  $q_0$  of  $\mathcal{A}_\phi$  as  $s_t, q_t$

7: set  $t = 0$  and construct an initial product state  $x_t = (s_t, q_t, T)$ ,

where  $T = T_0$

8: **while**  $t \leq \tau$  **do**

9: select action  $u_t^P = \pi_{q_t}(x|\theta^{u_{q_t}}) + R_{q_t}$  based on exploitation versus exploration noise

10: execute  $u_t^P, \delta$  and observe  $x_{t+1} = (s_{t+1}, q_{t+1}, T), R(x_{t+1}), \Phi(x_{t+1}), \gamma(x_{t+1})$

11: obtain  $\Phi(x_t)$  based on current  $T_\Phi$  and calculate  $R'(x_t, u_t^P, x_{t+1})$

12: execute the updates via  $f_V(q_{t+1}, T)$  and  $f_\Phi(q_{t+1}, T_\Phi)$

13: store the sample  $(x_t, u_t^P, R'(x_t, u_t^P, x_{t+1}), \gamma(x_{t+1}), x_{t+1})$  in replay buffers  $B_{q_t}$

14: mini-batch sampling  $N$  data from the replay buffers  $B_{q_t}$

15: calculate target values for each  $i \in N$  as:

$$y_i = R'(x_i, u_i^P, x_{i+1}) + \gamma(x_i) \cdot Q'_{q_{i+1}}(x_{i+1}, u_{i+1}^P|\theta'^{Q_{q_{i+1}}})$$

16: update weights  $\theta^{Q_{q_t}}$  of critic neural network  $Q_{q_t}(x, u^P|\theta^{Q_{q_t}})$  by minimizing the loss function:

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - Q_{q_t}(x_i, u_i^P|\theta^{Q_{q_t}}))^2$$

17: update weights  $\theta^{u_{q_t}}$  of actor neural network  $\pi_{q_t}(x|\theta^{u_{q_t}})$  by maximizing the policy gradient:

$$\begin{aligned} \nabla_{\theta^{u_{q_t}}} U^{q_t} &\approx \frac{1}{N} \sum_{i=1}^N (\nabla_{u^P} Q_{q_t}(x_i, u^P|\theta^{Q_{q_t}})|_{u^P=\pi_{q_t}(x_i|\theta^{u_{q_t}})}) \\ &\quad \cdot \nabla_{\theta^{u_{q_t}}} \pi_{q_t}(x_i|\theta^{u_{q_t}}) \end{aligned}$$

18: soft update of target networks:

$$\begin{aligned} \theta^{u'_{q_t}} &\leftarrow \tau \theta^{u_{q_t}} + (1 - \tau) \theta'^{u'_{q_t}} \\ \theta^{Q'_{q_t}} &\leftarrow \tau \theta^{Q_{q_t}} + (1 - \tau) \theta'^{Q'_{q_t}} \end{aligned} \tag{11}$$

19:  $x_t \leftarrow x_{t+1}$  and  $t++$

20: **end while**

21: **end for**

22: **end procedure**

---

where  $\rho^\beta$  is the state distribution under any arbitrary policy  $\beta$ , and  $y = R'(x, u^{\mathcal{P}}, x') + \gamma(x) Q(x', u^{\mathcal{P}'} | \theta^Q)$  with  $u^{\mathcal{P}'} = \pi(x' | \theta^u)$ . The actor can be updated by applying the chain rule to the expected return with respect to actor parameters  $\theta^u$  as the following policy gradient theorem [25]:

$$\begin{aligned} \nabla_{\theta^u} U^u(x) &\approx \mathbb{E}_{s \sim \rho^\beta}^{\pi} [\nabla_{\theta^u} Q(x, \pi(x | \theta^u) | \theta^Q)] \\ &= \mathbb{E}_{s \sim \rho^\beta}^{\pi} [\nabla_{u^{\mathcal{P}}} Q(x, u^{\mathcal{P}} | \theta^Q) \Big|_{u^{\mathcal{P}} = \pi(x | \theta^u)} \nabla_{\theta^u} \pi(x | \theta^u)]. \end{aligned} \quad (13)$$

Inspired by [35] and [38], the complex LTL task  $\phi$  can be divided into simple composable modules. Each state of the automaton in the LDGBA is module and each transition between these automaton states is a ‘‘task divider’’. In particular, given  $\phi$  and its LDGBA  $\mathcal{A}_\phi$ , we propose a modular architecture of  $|Q|$  DDPG respectively, i.e.,  $\pi_{q_i}(x | \theta^u)$  and  $Q_{q_i}(x, u^{\mathcal{P}} | \theta^Q)$  with  $q_i \in Q$ , along with their own replay buffer. Experience samples are stored in each replay modular buffer  $B_{q_i}$  in the form of  $(x, u^{\mathcal{P}}, R(x), \gamma(x), x')$ . By dividing the LTL task into sub-stages, the set of neural nets acts in a global modular DDPG architecture, which allows the agent to jump from one module to another by switching between the set of neural nets based on transition relations of  $\mathcal{A}_\phi$ .

Fig. 3 shows the modular DDPG architecture corresponding to the LDGBA  $\mathcal{A}_{\varphi_{\mathcal{P}}}$  shown in Fig. 2 without the sink node, where each network represents the standard DDPG structure along with an automaton state, and the transitions between each DDPG are consistent with the edges in  $\mathcal{A}_{\varphi_{\mathcal{P}}}$ . The shaped reward function consisting of  $R(x)$  and  $\Phi(x)$  is capable of guiding the transitions among the modular neural networks.

The proposed method to solve a continuous MDP with LTL specifications is summarized in Alg. 2, and the product states of EP-MDP are synchronized on-the-fly (line 9-12). We assign each DDPG an individual replay buffer  $B_{q_i}$  and a random process noise  $N_{q_i}$ . The corresponding weights of modular networks, i.e.,  $Q_{q_i}(x, u^{\mathcal{P}} | \theta^{Q_{q_i}})$  and  $\pi_{q_i}(x | \theta^{u_{q_i}})$ , are also updated at each iteration (line 15-20). All neural networks are trained using their own replay buffer, which is a finite-sized cache that stores transitions sampled from exploring the environment. Since the direct implementation of updating target networks can be unstable and divergent [24], the soft update (11) is employed, where target networks are slowly updated via relaxation (line 20). Note that for each iteration we first observe the output of the shaped reward function  $R'$ , then execute the update process via  $f_V$  and  $f_\Phi$  (line 10-12). Note that since DDPG leverages a deep neural network to approximate the action-value function, and that in practice we have to stop the training after finite number of steps the synthesised policy might be sub-optimal with respect to the true  $\pi^*$  as in Theorem 3. This is due to the nature of DDPG algorithm and non-linear function approximation in deep architectures whose analysis is out of the scope of this paper.

*Remark 1.* Standard embeddings, such as the one-hot and integer encoding, have been applied with a single DDPG network to estimate the Q-function [39]. However, it is observed that single DDPG network exhibits undesirable success rate. The modular DDPG with the EP-MDP in this work is capable of recording the unvisited accepting sets for each round, which outperforms the result of [39].

## 6 Case Studies

The developed modular DDPG-based control synthesis is implemented in Python, and all simulation videos and source codes can be found in our Github repository<sup>1</sup>. Owl [45] is used to convert LTL specifications to LDGBA. Simulations are carried out on a machine with 2.80 GHz quad-core CPU, 16 GB RAM, and an external Nvidia RTX 1080 GPU. We test the algorithm in 4 different environments in OpenAI gym and with 6 LTL tasks. In particular, we first assign complex tasks, formulated via LTL, to Ball-Pass and CartPole problems. To show the advantage of EP-MDP with modular DDPG, we compare the method with (i) the standard product MDP using modular DDPG and (ii) the EP-MDP using standard DDPG. Then, we test the scalability of our algorithm in pixel-based Mars exploration environments. Finally, we analyze and compare the performances of the modular and standard DDPG algorithms via success rates of task accomplishments.

### 6.1 Ball-Pass and Cart-Pole

**(1). Ball-Pass:** We first demonstrate the developed control synthesis in a Ball-Pass environment ( $600m \times 600m$ ). Consider a red ball moving according to the following dynamics:  $\ddot{x} = a_x, \ddot{y} = a_y + g$  in Fig. 4 (a), where  $(x, y)$  is the planar position of the ball,  $a_x, a_y$  represent accelerations along  $x$  and  $y$  induced by an external force (i.e., the control input), respectively, and  $g$  is the gravitational acceleration. The simulation step size is  $\Delta t = 0.05s$ , and the

<sup>1</sup>[https://github.com/mingyucui/Modular\\_Deep\\_RL\\_E-LDGBA](https://github.com/mingyucui/Modular_Deep_RL_E-LDGBA)

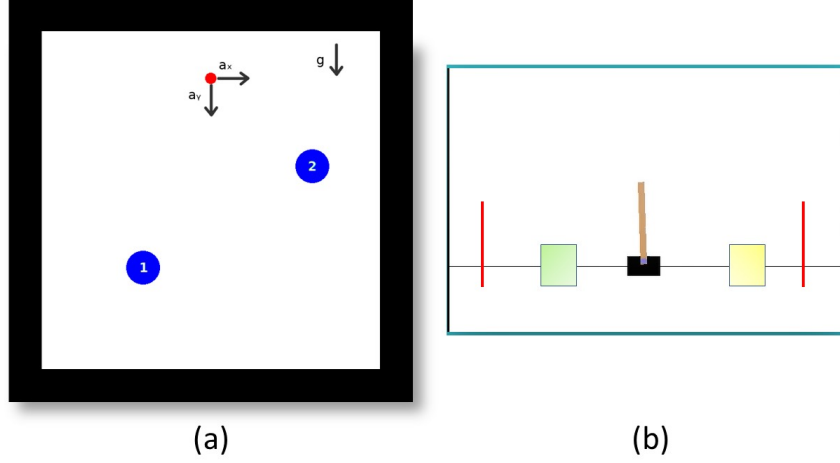


Figure 4: (a) Ball-Pass and (b) Cart-Pole OpenAI environment.

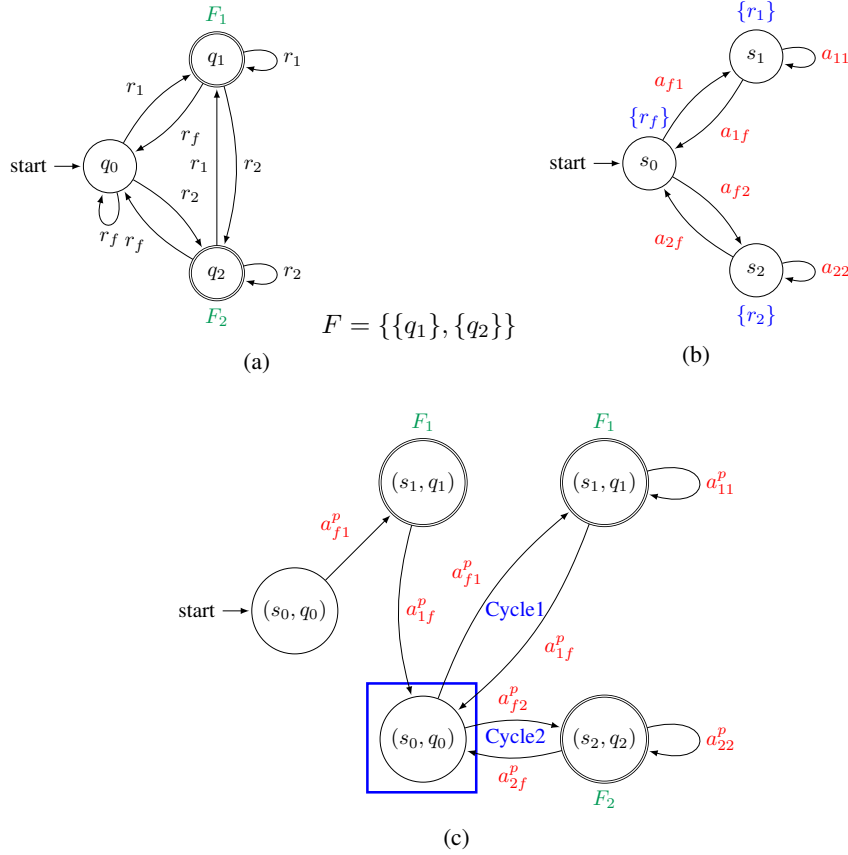


Figure 5: (a) LDGBA of LTL formula  $\varphi_{B1}$ . (b) Generalized MDP model for Ball-pass case. (c) The standard product MDP.

acceleration range is  $a_x, a_y \in [0, 1]$  m/s. We consider two LTL tasks, e.g.,  $\varphi_{B1} = (\Box\Diamond\text{Region1}) \wedge (\Box\Diamond\text{Region2})$  and  $\varphi_{B2} = \Diamond(\text{Region1} \wedge \Diamond\text{Region2})$ .  $\varphi_{B1}$  is a task over the infinite horizon and requires the agent repetitively visits Region1 and Region2, and  $\varphi_{B2}$  is a task over the finite horizon and requires the agent first to visit Region1 and then Region2.

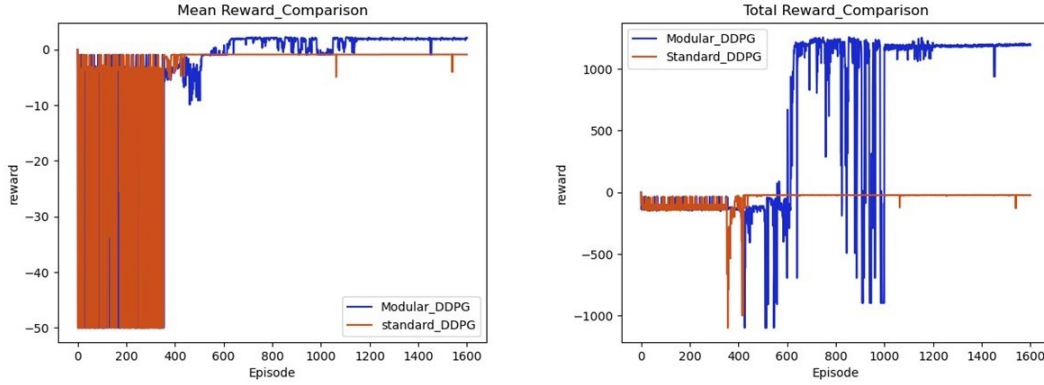


Figure 6: The evolution of reward for  $\varphi_{B1}$ . (a) Average reward. (b) Total reward.

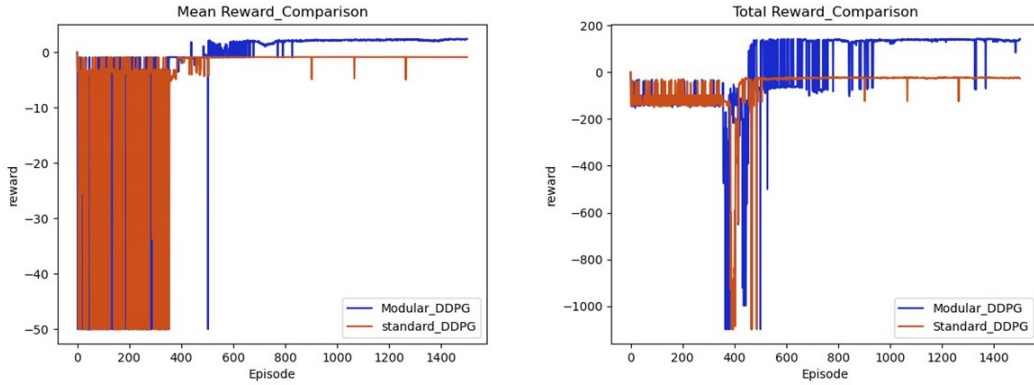


Figure 7: The evolution of reward for  $\varphi_{B2}$ . (a) Average reward. (b) Total reward.

For comparison of method (i), Fig. 5 (a) shows the corresponding LDGBA of  $\varphi_{B1}$  with two accepting sets  $F = \{\{q_1\}, \{q_2\}\}$ , and  $r_1, r_2$  and  $r_f$  represent Region 1, Region 2 and free space being visited, respectively. Fig. 5 (b) shows the Ball-pass MDP model, where  $a_{fi}$  (or  $a_{if}$ ) represent a sequence of continuous actions that drive the ball from free space to Region  $i$  (or inverse). In addition, Fig. 5 (c) shows the resulting standard product MDP. By Definition 2, the policy that satisfies  $\varphi_{B1}$  should enforce the repetitive trajectories, i.e., cycle 1 and cycle 2 in Fig. 5 (c). However, there exists no deterministic policy that can periodically select two actions  $a_{f1}^P$  and  $a_{f2}^P$  at state  $(s_f, q_0)$  (marked with a blue rectangle) in Fig. 5 (c). As a result, applying standard product MDP cannot generate a pure deterministic optimal policy to complete task  $\varphi_{B1}$ . Such a scenario also happens for task  $\varphi_{C1}$ . In contrast, the tracking-frontier set of EP-MDP developed in this work can resolve this issue by recording unvisited acceptance sets and being embedded with each state at every time-step via one-hot encoding. We simulate 10 runs for tasks  $\varphi_{B1}$  and  $\varphi_{C1}$  and select the worst case of applying standard product MDP as comparison. For comparison of method (ii), we conduct 1600 episodes for each task ( $\varphi_{B1}$  and  $\varphi_{B2}$ ) of Ball-pass, and the reward collections are shown in Fig. 6 and 7..

**(2). CartPole:** We also test our control framework for the Cart-pole<sup>2</sup> in Fig. 4 (b). The pendulum starts upright with an initial angle between  $-0.05$  and  $0.05$  rads. The horizontal force exerted on the cart is defined over a continuous space (action-space) with a range  $(-10N, 10N)$ . The green and yellow regions range from  $-1.44$  to  $-0.96$  m and from  $0.96$  to  $1.44$  m, respectively. The objective is to prevent the pendulum from falling over while moving the cart between the yellow and green regions. Similarly, we consider two tasks  $\varphi_{C1} = (\Box \Diamond \text{Green}) \wedge (\Box \Diamond \text{Yellow}) \wedge \neg \Box \text{Unsafe}$  and  $\varphi_{C2} = \Diamond (\text{Green} \wedge \Diamond \text{Yellow}) \wedge \neg \Box \text{Unsafe}$ , where  $\text{Unsafe}$  represents the condition that the pendulum falls over or exceeds the red line, and  $\text{Green}$ ,  $\text{Yellow}$  represent colored areas. Method (i) has the similar issues as shown in Fig. for task  $\varphi_{C1}$  over infinite horizon. To compare the method (ii), we conduct mean reward collections over 1500 episodes for each task ( $\varphi_{C1}$  and  $\varphi_{C2}$ ) of CartPole shown in Fig. 8. We can conclude the modular DDPG has a better performance of

<sup>2</sup><https://gym.openai.com/envs/CartPole-v0/#bart083>

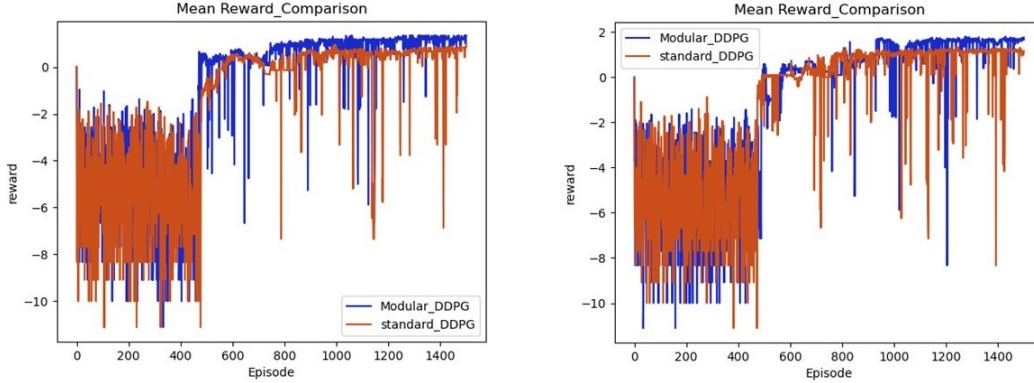


Figure 8: The evolution of average reward. (a) Task  $\varphi_{C1}$ . (b) Task  $\varphi_{C2}$ .

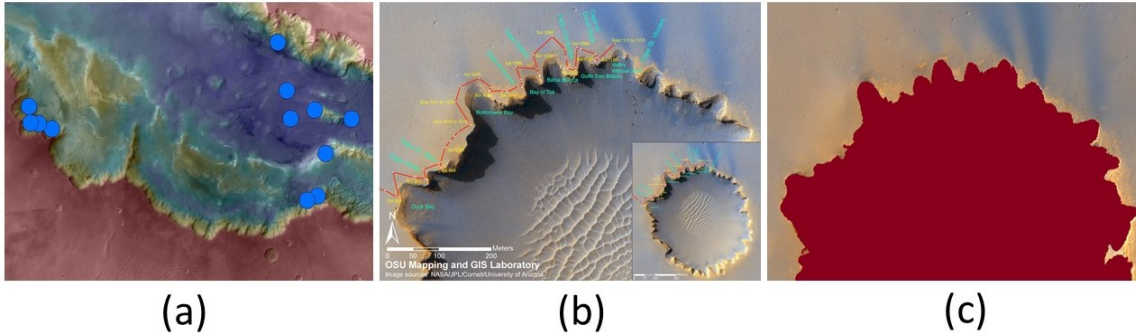


Figure 9: Mars exploration mission. (a) Melas Chasma with possible location of water (blue dots). (b) Victoria Crater with exploration path (red line). (c) Processed image of Victoria Crater.

collecting positive rewards during learning. The figures also illustrate that the modular DDPG has a better performance than the standard DDPG for tasks over both finite and infinite horizons. We simulate 10 runs and select the worst case of applying standard product MDP as comparison that can be found in our Github repository.

## 6.2 Image-based Environment

In this section, we test our algorithm in large scale continuous environments, and conduct motion planning to complete two Mars exploration missions using satellite images as shown in Fig. 9. The missions are to explore areas around the Melas Chasma [47] and the Victoria Crater [48]. In the Melas Chasma of Fig. 9 (a), there are a number of spots with the potential presence of water, possible river valleys and lakes. According to NASA, the blue spots are possible locations of water. The LTL task, in this case, is to visit both clusters by visiting any blue dot in each cluster while avoiding unsafe (red land) regions in Fig. 9 (a). The Victoria Crater in Fig. 9 (b) is an impact crater located near the equator of Mars. Layered sedimentary rocks are exposed along the wall of crater, providing information about the ancient surface condition of Mars. The mission is related to visiting all spots along with the path of the well-known Mars Rover Opportunity that are given in Fig. 9 (b), and avoiding the unsafe areas (red regions in Fig. 9 (c)). We employ LTL to specify such missions. The LTL specifications for Melas Chasma and Victoria Crater are expressed as following separately:

$$\varphi_{Melas} = \square \diamond M_1 \wedge \square \diamond M_2 \wedge \neg \square M_{unsafe},$$

$$\varphi_{Victoria} = \square \diamond V_1 \wedge \square \diamond V_2 \dots \square \diamond V_{12} \wedge \neg \square V_{unsafe},$$

where  $M_i$  represents  $i$ -th target and  $M_{unsafe}$  indicates unsafe areas in Melas Chasma. The task  $\varphi_{Victoria}$  in Victoria Crater has similar settings. At each stage, the rover agent has a continuous action-range  $[0, 2\pi)$  and the resulting outcome of each action is to move toward the direction of the action within a range drawn from  $(0, D]$ . The dimensions and action ranges of Fig. 9 (a) and Fig. 9 (b) are  $456km \times 322km$ ,  $D = 2km$  and  $746m \times 530m$ ,  $D = 10m$ , respectively.

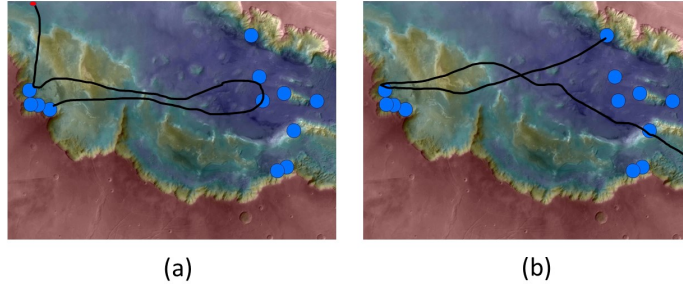


Figure 10: The trajectories with different initial locations (image coordinates). (a) Top-side with initial coordinates (2,18). (b) Right-side with initial coordinate (105,199).

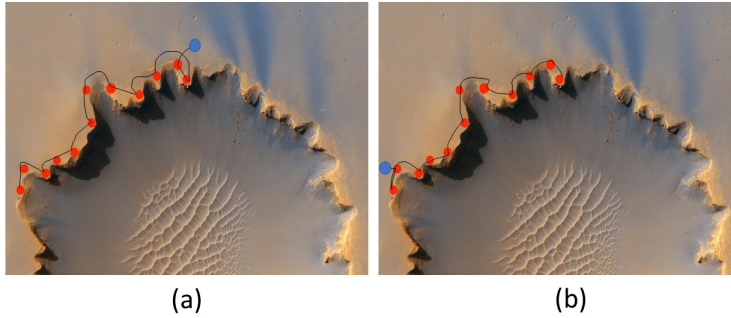


Figure 11: The trajectories with different initial locations (marked as blue cycles). (a) top-right. (b) bottom-left

The LDGBA for  $\varphi_{Melas}$  has 3 states and 2 accepting sets, and the one for  $\varphi_{Victoria}$  has 13 states and 12 accepting sets. The scenario is to train a modular DDPG with the tuning reward design described in Section 5.1, which can autonomously satisfy the complex tasks by accessing the images.

The surveillance task  $\varphi_{Melas}$  requires visiting right and left clusters in Fig. 9 (b), and the similar task  $\varphi_{Victoria}$  requires visiting all the spots associated with the path in Fig. 9 (b). After training, we show the trajectories of one round satisfaction with different initial locations for the two tasks in Fig. 10 and 11. The training time is influenced by the number of actor-critic neural network pairs.

### 6.3 Results Summary

Based on the experimental results, we can summarize our contributions compared with two other methods: (i) the modular DDPG with standard product MDP, (ii) the standard DDPG with EP-MDP.

As analyzed in Fig. 5, the modular DDPG with standard product MDP can not be guaranteed to synthesise pure deterministic policies for a number of LTL properties, e.g.,  $\varphi_{B1}$ ,  $\varphi_{C1}$ ,  $\varphi_{Melas}$  and  $\varphi_{Victoria}$ . The reason is that the corresponding LDGBA for each of them has multiple accepting sets, and there exist no deterministic policies for one state of  $\mathcal{P}$  to select different actions.

Then, we take 200 runs for all tasks above and analyze the success rate for all aforementioned tasks compared with method (ii) in Table 1. In practice, since the training process of a deep neural network has limited steps for each episode and finite number of total episodes, and dimensions of automaton structure grows for more complex tasks, it becomes difficult for standard DDPG to explore the whole tasks especially for tasks over infinite horizon. We can conclude that the standard DDPG with a recording method as [39] might yield poor performance results for tasks with repetitive pattern (infinite horizon), whereas the modular DDPG has better performance from the perspective of training and

Table 1: Comparison of standard and modular DDPG with EP-MDP. Statistics are taken over 200 runs.

LTL Task	DDPG	Success rate
$\varphi_{B1}$	Modular DDPG	100%
	Standard DDPG	0%
$\varphi_{B2}$	Modular DDPG	100%
	Standard DDPG	77.5%
$\varphi_{C1}$	Modular DDPG	100%
	Standard DDPG	23.5%
$\varphi_{C2}$	Modular DDPG	100%
	Standard DDPG	61%
$\varphi_{Melas}$	Modular DDPG	100%
	Standard DDPG	0%
$\varphi_{Victoria}$	Modular DDPG	100%
	Standard DDPG	0%

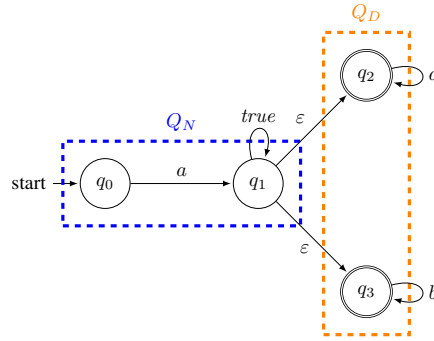


Figure 12: LDGBA for the formula  $a \wedge \bigcirc(\bigtriangleleft \square a \vee \bigtriangleleft \square b)$ .

success rate. It also shall be noted that the surveillance tasks were mainly considered in this paper. However, the temporal operators in LTL, such as "next" and "until", enable the modular DDPG algorithm to solve problems with other types of complex tasks, e.g., fairness and safety [15, 18]. For instance, consider the following LTL formula

$$\varphi = a \wedge \bigcirc(\bigtriangleleft \square a \vee \bigtriangleleft \square b)$$

In those cases, LTL formulas can be converted to LDGBAs [26], as in Fig. 12, and the same framework is applied.

At last, although several actor-critic neural network pairs are adopted in the modular architecture, they are synchronously trained online, and each of them is only responsible for a sub-task. This setup is more effective to complete LTL complex tasks. The training time for each task with two different algorithms (the standard DDPG and the modular DDPG) is shown in Table 2. It can be seen that the runtime complexity is not increased in the modular DDPG.

Table 2: Training time analysis of standard and modular DDPG with EP-MDP.

Tasks and Training Parameters			Training Time (minute, hour, day)	
LTL Task	Maximum steps	Episode	Standard DDPG	Modular DDPG
$\varphi_{B1}$	200	1600	11.3 min	12.1 min
$\varphi_{B2}$	500	1500	17.0 min	20.4 min
$\varphi_{C1}$	200	1500	14.1 min	13.6 min
$\varphi_{C2}$	500	1500	20.2 min	21.0 min
$\varphi_{Melas}$	2000	10000	5.5 hr	5.4 hr
$\varphi_{Victoria}$	8000	50000	48.0 hr	50.0 hr

## 7 Conclusions

In this paper, a model-free deep RL learning is developed to synthesize control policies in continuous-state and continuous-action MDPs. The designed EP-MDP can enforce the task satisfaction, and a modular DDPG framework is proposed to decompose the temporal automaton into interrelated sub-tasks such that the resulted optimal policies are shown to satisfy the LTL specifications with a higher success rate compared to the standard DDPG. Based on the comparisons, we demonstrate the benefits of applying LDGBA and modular DDPG to improve the learning performance and enforce the task satisfaction with high success rates.

## 8 Appendix

### 8.1 Proof of Theorem 1

Based on whether or not the path  $\mathbf{x}_t$  intersects with accepting states of  $F_i^{\mathcal{P}}$ , the expected return in (6) can be rewritten as

$$U^\pi(x) = \mathbb{E}^\pi [\mathcal{D}(\mathbf{x}_t) | \mathbf{x}_t \models \diamond F_i^{\mathcal{P}}] \cdot \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}] + \mathbb{E}^\pi [\mathcal{D}(\mathbf{x}_t) | \mathbf{x}_t \not\models \diamond F_i^{\mathcal{P}}] \cdot \Pr^\pi [x \not\models \diamond F_i^{\mathcal{P}}] \quad (14)$$

where  $\Pr^\pi [x \models \diamond F_i^{\mathcal{P}}]$  and  $\Pr^\pi [x \not\models \diamond F_i^{\mathcal{P}}]$  represent the probability of eventually reaching and not reaching  $F_i^{\mathcal{P}}$  eventually under policy  $\pi$  starting from state  $x$ , respectively.

To find the lower bound of  $U^\pi(x)$ , for any  $\mathbf{x}_t$  with  $\mathbf{x}_t[t] = x$ , let  $t + N_t$  be the index that  $\mathbf{x}_t$  first intersects a state in  $X_{\mathcal{P}}^\pi$ , i.e.,  $N_t = \min [i | \mathbf{x}_t[t+i] \in X_{\mathcal{P}}^\pi]$ . The following holds

$$\begin{aligned} & \mathbb{E}^\pi [\mathcal{D}(\mathbf{x}_t) | \mathbf{x}_t \models \diamond F_i^{\mathcal{P}}] \\ & \stackrel{(1)}{\geq} \mathbb{E}^\pi [\mathcal{D}(\mathbf{x}_t) | \mathbf{x}_t \cap X_{\mathcal{P}}^\pi \neq \emptyset] \\ & \stackrel{(2)}{\geq} \mathbb{E}^\pi [\gamma_F^{N_t} \cdot \mathcal{D}(\mathbf{x}_t[t+N_t]) | \mathbf{x}_t[t+N_t] = x | \mathbf{x}_t \cap X_{\mathcal{P}}^\pi \neq \emptyset] \\ & \stackrel{(3)}{\geq} \mathbb{E}^\pi [\gamma_F^{N_t} | \mathbf{x}_t \cap X_{\mathcal{P}}^\pi \neq \emptyset] \cdot U_{\min}^\pi(\mathbf{x}_t[t+N_t]) \\ & \stackrel{(4)}{\geq} \gamma_F^{\mathbb{E}^\pi[N_t | \mathbf{x}_t[t]=x | \mathbf{x}_t \cap X_{\mathcal{P}}^\pi \neq \emptyset]} \cdot U_{\min}^\pi(x_{Acc}) \\ & = \gamma_F^{n_t} \cdot U_{\min}^\pi(x_{Acc}), \end{aligned} \quad (15)$$

where  $x_{Acc} \in X_{\mathcal{P}}^\pi$ ,  $U_{\min}^\pi(x_{Acc}) = \min_{x \in X_{\mathcal{P}}^\pi} U^\pi(x)$ , and  $n_t$  is a constant. By Lemma 2, one has  $\lim_{\gamma_F \rightarrow 1^-} U_{\min}^\pi(x_{Acc}) = 1$ . In (15), the first inequality (1) holds because visiting  $X_{\mathcal{P}}^\pi$  is one of the cases for  $\diamond F_i^{\mathcal{P}}$  that satisfy  $\mathbf{x}_t \models \diamond F_i^{\mathcal{P}}$ , e.g.,  $F_i^{\mathcal{P}}$  can be placed outside of all BSCCs; the second inequality (2) holds due to Lemma 2; the third inequality (3) holds due to the Markov properties of (5) and (6); the fourth inequality (4) holds due to Jensen's inequality. Based on (15), the lower bound of (14) is  $U^\pi(x) \geq \gamma_F^{n_t} \cdot U_{\min}^\pi(x_{Acc}) \cdot \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}]$  from which one has

$$\lim_{\gamma_F \rightarrow 1^-} U^\pi(x) \geq \gamma_F^{n_t} \cdot \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}]. \quad (16)$$

Similarly, let  $t + M_t$  denote the index that  $\mathbf{x}_t$  first enters the BSCC that contains no accepting states. We have

$$\begin{aligned} \mathbb{E}^\pi [\mathcal{D}(\mathbf{x}_t) | \mathbf{x}_t \not\models \diamond F_i^{\mathcal{P}}] & \stackrel{(1)}{\leq} \mathbb{E}^\pi [1 - r_F^{M_t} | \mathbf{x}_t \not\models \diamond F_i^{\mathcal{P}}] \\ & \stackrel{(2)}{\leq} 1 - r_F^{\mathbb{E}^\pi[M_t | \mathbf{x}_t[t]=x, \mathbf{x}_t \not\models \diamond F_i^{\mathcal{P}}]} = 1 - r_F^{m_t} \end{aligned} \quad (17)$$

where  $m_t$  is a constant and (17) holds due to Lemma 2 and Markov properties.

Hence, the upper bound of (14) is obtained as

$$\lim_{\gamma_F \rightarrow 1^-} U^\pi(x) \leq \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}] + (1 - r_F^{m_t}) \Pr^\pi [x \not\models \diamond F_i^{\mathcal{P}}]. \quad (18)$$

By (16) and (18), we can conclude

$$\begin{aligned} \gamma_F^{n_t} \cdot \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}] & \leq \lim_{\gamma_F \rightarrow 1^-} U^\pi(x) \\ & \leq \Pr^\pi [x \models \diamond F_i^{\mathcal{P}}] + (1 - r_F^{m_t}) \cdot \Pr^\pi [x \not\models \diamond F_i^{\mathcal{P}}] \end{aligned}$$

According to  $\lim_{\gamma_F \rightarrow 1^-} r_F(\gamma_F) = 1$  in the reward function, (7) can be concluded.



## 8.2 Proof of Theorem 2

For any policy  $\pi$ ,  $MC_{\mathcal{P}}^{\pi} = \mathcal{T}_{\pi} \sqcup \mathcal{R}_{\pi}^1 \sqcup \mathcal{R}_{\pi}^2 \dots \mathcal{R}_{\pi}^{n_R}$ . Let  $\mathbf{U}_{\pi} = [U^{\pi}(x_0) \ U^{\pi}(x_1) \ \dots]^T \in \mathbb{R}^{|\mathcal{X}|}$  denote the stacked expected return under policy  $\pi$ , which can be reorganized as

$$\begin{aligned} \begin{bmatrix} \mathbf{U}_{\pi}^{tr} \\ \mathbf{U}_{\pi}^{rec} \end{bmatrix} &= \sum_{n=0}^{\infty} \left( \prod_{j=0}^{n-1} \begin{bmatrix} \gamma_{\pi}^{\mathcal{T}} & \gamma_{\pi}^{tr} \\ \mathbf{0}_{\sum_{i=1}^m N_i \times r} & \gamma_{\pi}^{rec} \end{bmatrix} \right) \\ &\cdot \begin{bmatrix} \mathbf{P}_{\pi}(\mathcal{T}, \mathcal{T}) & \mathbf{P}_{\pi}^{tr}(\mathcal{R}, \mathcal{R}) \\ \mathbf{0}_{\sum_{i=1}^m N_i \times r} & \mathbf{P}_{\pi}(\mathcal{R}, \mathcal{R}) \end{bmatrix}^n \begin{bmatrix} \mathbf{R}_{\pi}^{tr} \\ \mathbf{R}_{\pi}^{rec} \end{bmatrix}, \end{aligned} \quad (19)$$

where  $\mathbf{U}_{\pi}^{tr}$  and  $\mathbf{U}_{\pi}^{rec}$  are the expected return of states in transient and recurrent classes under policy  $\pi$ , respectively. In (19),  $\mathbf{P}_{\pi}(\mathcal{T}, \mathcal{T}) \in \mathbb{R}^{r \times r}$  is the probability transition matrix between states in  $\mathcal{T}_{\pi}$ , and  $\mathbf{P}_{\pi}^{tr} = [P_{\pi}^{tr_1} \dots P_{\pi}^{tr_m}] \in \mathbb{R}^{r \times \sum_{i=1}^m N_i}$  is the probability transition matrix where  $P_{\pi}^{tr_i} \in \mathbb{R}^{r \times N_i}$  represents the transition probability from a transient state in  $\mathcal{T}_{\pi}$  to a state of  $\mathcal{R}_{\pi}^i$ . The  $\mathbf{P}_{\pi}(\mathcal{R}, \mathcal{R})$  is a diagonal block matrix, where the  $i$ th block is a  $N_i \times N_i$  matrix containing transition probabilities between states within  $\mathcal{R}_{\pi}^i$ . Note that  $\mathbf{P}_{\pi}(\mathcal{R}, \mathcal{R})$  is a stochastic matrix since each block matrix is a stochastic matrix [46]. Similarly, the rewards  $\mathbf{R}_{\pi}$  can also be partitioned into  $\mathbf{R}_{\pi}^{tr}$  and  $\mathbf{R}_{\pi}^{rec}$ .

The following proof is based on contradiction. Suppose there exists a policy  $\pi^*$  that optimizes the expected return, but does not satisfy the accepting condition of  $\mathcal{P}$  with non-zero probability. Based on Lemma 1, the following is true:  $F_k^{\mathcal{P}} \subseteq \mathcal{T}_{\pi^*}, \forall k \in \{1, \dots, f\}$ , where  $\mathcal{T}_{\pi^*}$  denotes the transient class of Markov chain induced by  $\pi^*$  on  $\mathcal{P}$ . First, consider a state  $x_R \in \mathcal{R}_{\pi^*}^j$  and let  $\mathbf{P}_{\pi^*}^{x_R R_j}$  denote a row vector of  $\mathbf{P}_{\pi^*}^n(\mathcal{R}, \mathcal{R})$  that contains the transition probabilities from  $x_R$  to the states in the same recurrent class  $\mathcal{R}_{\pi^*}^j$  after  $n$  steps. The expected return of  $x_R$  under  $\pi^*$  is then obtained from (19) as

$$U_{\pi^*}^{rec}(x_R) = \sum_{n=0}^{\infty} \gamma^n \left[ \mathbf{0}_{k_1}^T \mathbf{P}_{\pi^*}^{x_R R_j} \mathbf{0}_{k_2}^T \right] \mathbf{R}_{\pi^*}^{rec},$$

where  $k_1 = \sum_{i=1}^{j-1} N_i, k_2 = \sum_{i=j+1}^n N_i$ . Since  $\mathcal{R}_{\pi^*}^j \cap F_i^{\mathcal{P}} = \emptyset, \forall i \in \{1, \dots, f\}$ , by the designed reward function, all entries of  $\mathbf{R}_{\pi^*}^{rec}$  are zero. We can conclude  $U_{\pi^*}^{rec}(x_R) = 0$ . To show contradiction, the following analysis will show that  $U_{\bar{\pi}}^{rec}(x_R) > U_{\pi^*}^{rec}(x_R)$  for any policy  $\bar{\pi}$  that satisfies the accepting condition of  $\mathcal{P}$ . Thus, it's true that there exists  $\mathcal{R}_{\bar{\pi}}^j$  such that  $\mathcal{R}_{\bar{\pi}}^j \cap F_k^{\mathcal{P}} \neq \emptyset, \forall k \in \{1, \dots, f\}$ . We use  $\underline{\gamma}$  and  $\bar{\gamma}$  to denote the lower and upper bound of  $\gamma$ .

**Case 1:** If  $x_R \in \mathcal{R}_{\bar{\pi}}^j$ , there exist states such that  $x_A \in \mathcal{R}_{\bar{\pi}}^j \cap F_i^{\mathcal{P}}$ . From Lemma 1, the entries in  $\mathbf{R}_{\bar{\pi}}^{rec}$  corresponding to the recurrent states in  $\mathcal{R}_{\bar{\pi}}^j$  have non-negative rewards and at least there exist  $f$  states in  $\mathcal{R}_{\bar{\pi}}^j$  from different accepting sets  $F_i^{\mathcal{R}}$  with positive reward  $1 - r_F$ . From (19),  $U_{\bar{\pi}}^{rec}(x_R)$  can be lower bounded as

$$U_{\bar{\pi}}^{rec}(x_R) \geq \sum_{n=0}^{\infty} \underline{\gamma}^n (P_{\bar{\pi}}^{x_R x_A} r_F) > 0,$$

where  $P_{\bar{\pi}}^{x_R x_A}$  is the transition probability from  $x_R$  to  $x_A$  in  $n$  steps. We can conclude in this case  $U_{\bar{\pi}}^{rec}(x_R) > U_{\pi^*}^{rec}(x_R)$ .

**Case 2:** If  $x_R \in \mathcal{T}_{\bar{\pi}}$ , there are no states of any accepting set  $F_i^{\mathcal{P}}$  in  $\mathcal{T}_{\bar{\pi}}$ . As demonstrated in [46], 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  $x_T$  to itself in  $n$  time steps. In addition, for a recurrent state  $x_{rec}$  of  $\mathcal{R}_{\bar{\pi}}^j$ , it is always true that

$$\sum_{n=0}^{\infty} \gamma^n p^n(x_{rec}, x_{rec}) > \frac{1}{1 - \gamma \bar{p}}, \quad (20)$$

where there exists  $\bar{p}$  such that  $p^{\bar{n}}(x_{rec}, x_{rec})$  is nonzero and can be lower bounded by  $\bar{p}$  [46]. From (19), one has

$$\begin{aligned} \mathbf{U}_{\bar{\pi}}^{tr} &> \sum_{n=0}^{\infty} \left( \prod_{j=0}^{n-1} \gamma_{\bar{\pi}}^{tr} \right) \cdot \mathbf{P}_{\bar{\pi}}^{tr} \mathbf{P}_{\bar{\pi}}^n(\mathcal{R}, \mathcal{R}) \mathbf{R}_{\bar{\pi}}^{rec} \\ &> \underline{\gamma}^n \cdot \mathbf{P}_{\bar{\pi}}^{tr} \mathbf{P}_{\bar{\pi}}^n(\mathcal{R}, \mathcal{R}) \mathbf{R}_{\bar{\pi}}^{rec}. \end{aligned} \quad (21)$$

Let  $\max(\cdot)$  and  $\min(\cdot)$  represent the maximum and minimum entry of an input vector, respectively. The upper bound  $\bar{m} = \{ \max(\bar{M}) \mid \bar{M} < \mathbf{P}_{\bar{\pi}}^{tr} \bar{\mathbf{P}} \mathbf{R}_{\bar{\pi}}^{rec} \}$  and  $\bar{m} \geq 0$ , where  $\bar{\mathbf{P}}$  is a block matrix whose nonzero entries are derived

similarly to  $\bar{p}$  in (20). The utility  $U_{\pi}^{tr}(x_R)$  can be lower bounded from (20) and (21) as  $U_{\pi}^{tr}(x_R) > \frac{1}{1-\gamma^n} \bar{m}$ . Since  $U_{\pi^*}^{rec}(x_R) = 0$ , the contradiction  $U_{\pi}^{tr}(x_R) > 0$  is achieved if  $\frac{1}{1-\gamma^n} \bar{m}$ . Thus, there exist  $0 < \underline{\gamma} < 1$  such that  $\gamma_F > \underline{\gamma}$  and  $r_F > \underline{\gamma}$ , which implies  $U_{\pi}^{tr}(x_R) > \frac{1}{1-\gamma^n} \bar{m} \geq 0$ . The procedure shows the contradiction of the assumption that  $\pi^*$  that does not satisfy the acceptance condition of  $\mathcal{P}$  with non-zero probability is optimal, and Theorem 2 is proved.

## References

- [1] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [2] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT press, 2008.
- [3] M. Guo and D. V. Dimarogonas, “Multi-agent plan reconfiguration under local LTL specifications,” *Int. J. Robotics Res.*, vol. 34, no. 2, pp. 218–235, 2015.
- [4] C. I. Vasile, X. Li, and C. Belta, “Reactive sampling-based path planning with temporal logic specifications,” *Int. J. Robot. Res.*, p. 0278364920918919, 2020.
- [5] M. Cai, H. Peng, Z. Li, H. Gao, and Z. Kan, “Receding horizon control based motion planning with partially infeasible LTL constrains,” *IEEE Control Syst. Lett.*, vol. 5, no. 4, pp. 1279–1284, 2020.
- [6] X. Ding, S. L. Smith, C. Belta, and D. Rus, “Optimal control of Markov decision processes with linear temporal logic constraints,” *IEEE Trans. Autom. Control*, vol. 59, no. 5, pp. 1244–1257, 2014.
- [7] B. Lacerda, F. Faruq, D. Parker, and N. Hawes, “Probabilistic planning with formal performance guarantees for mobile service robots,” *Int. J. Robot. Res.*, vol. 38, no. 9, pp. 1098–1123, 2019.
- [8] M. Kloetzer and C. Mahulea, “LTL-based planning in environments with probabilistic observations,” *IEEE Trans. Autom. Sci. Eng.*, vol. 12, no. 4, pp. 1407–1420, 2015.
- [9] D. Sadigh, E. S. Kim, S. Coogan, S. S. Sastry, and 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.
- [10] J. Fu and U. Topcu, “Probably approximately correct mdp learning and control with temporal logic constraints,” *Robotics: Science and Systems*, 2014.
- [11] J. Wang, X. Ding, M. Lahijanian, I. C. Paschalidis, and C. A. Belta, “Temporal logic motion control using actor–critic methods,” *Int. J. Robotics Res.*, vol. 34, no. 10, pp. 1329–1344, 2015.
- [12] M. Hasanbeig, A. Abate, and D. Kroening, “Logically-Constrained Reinforcement Learning,” *arXiv preprint arXiv:1801.08099*, 2018.
- [13] M. Cai, H. Peng, Z. Li, and Z. Kan, “Learning-based probabilistic LTL motion planning with environment and motion uncertainties,” *IEEE Trans. Autom. Control*, 2020, to appear.
- [14] M. Hasanbeig, Y. Kantaros, A. Abate, D. Kroening, G. J. Pappas, and I. Lee, “Reinforcement learning for temporal logic control synthesis with probabilistic satisfaction guarantees,” in *Proc. IEEE CDC*. IEEE, 2019, pp. 5338–5343.
- [15] M. Hasanbeig, A. Abate, and D. Kroening, “Certified reinforcement learning with logic guidance,” *arXiv preprint arXiv:1902.00778*, 2019.
- [16] E. M. Hahn, M. Perez, S. Schewe, F. Somenzi, A. Trivedi, and D. Wojtczak, “Omega-regular objectives in model-free reinforcement learning,” in *Int. Conf. Tools Alg. Constr. Anal. Syst.* Springer, 2019, pp. 395–412.
- [17] A. K. Bozkurt, Y. Wang, M. M. Zavlanos, and M. Pajic, “Control synthesis from linear temporal logic specifications using model-free reinforcement learning,” in *Int. Conf. Robot. Autom.* IEEE, 2020, pp. 10 349–10 355.
- [18] M. Hasanbeig, A. Abate, and D. Kroening, “Cautious reinforcement learning with logical constraints,” in *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2020, pp. 483–491.
- [19] R. Oura, A. Sakakibara, and T. Ushio, “Reinforcement learning of control policy for linear temporal logic specifications using limit-deterministic generalized büchi automata,” *IEEE Control Syst. Lett.*, vol. 4, no. 3, pp. 761–766, 2020.
- [20] M. Cai, S. Xiao, B. Li, Z. Li, and Z. Kan, “Reinforcement learning based temporal logic control with maximum probabilistic satisfaction,” in *Int. Conf. Robot. Autom.* IEEE, 2021, pp. 806–812.
- [21] M. Cai, S. Xiao, and Z. Kan, “Reinforcement learning based temporal logic control with soft constraints using limit-deterministic büchi automata,” *arXiv preprint arXiv:2101.10284*, 2021.

- 
- [22] D. Aksaray, A. Jones, Z. Kong, M. Schwager, and C. Belta, “Q-learning for robust satisfaction of signal temporal logic specifications,” in *Proc. IEEE Conf. Decis. Control*, 2016, pp. 6565–6570.
- [23] H. Venkataraman, D. Aksaray, and P. Seiler, “Tractable Reinforcement Learning of Signal Temporal Logic Objectives,” *Learning for Dynamics and Control*, 308-317, 2020.
- [24] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [25] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
- [26] S. Sickert, J. Esparza, S. Jaax, and J. Křetínský, “Limit-deterministic Büchi automata for linear temporal logic,” in *Int. Conf. Comput. Aided Verif.* Springer, 2016, pp. 312–332.
- [27] S. L. Smith, J. Tumova, C. Belta, and D. Rus, “Optimal path planning for surveillance with temporal-logic constraints,” *Int. J. Robotics Res.*, vol. 30, no. 14, pp. 1695–1708, 2011.
- [28] M. Cai, Z. Li, H. Gao, S. Xiao, and Z. Kan, “Optimal probabilistic motion planning with partially infeasible LTL constraints,” *arXiv preprint arXiv:2007.14325*, 2020.
- [29] M. Hasanbeig, N. Yogananda Jeppu, A. Abate, T. Melham, and D. Kroening, “DeepSynth: Program synthesis for automatic task segmentation in deep reinforcement learning,” in *AAAI Conference on Artificial Intelligence*. Association for the Advancement of Artificial Intelligence, 2021.
- [30] R. T. Icarte, T. Klassen, R. Valenzano, and S. McIlraith, “Using reward machines for high-level task specification and decomposition in reinforcement learning,” in *Int. Conf. Mach. Learn.*, 2018, pp. 2107–2116.
- [31] A. Camacho, R. T. Icarte, T. Q. Klassen, R. A. Valenzano, and S. A. McIlraith, “LTL and beyond: Formal languages for reward function specification in reinforcement learning,” in *IJCAI*, vol. 19, 2019, pp. 6065–6073.
- [32] X. Li, Z. Serlin, G. Yang, and C. Belta, “A formal methods approach to interpretable reinforcement learning for robotic planning,” *Sci. Robot.*, vol. 4, no. 37, 2019.
- [33] G. Rens and J.-F. Raskin, “Learning non-Markovian reward models in MDPs,” *arXiv preprint arXiv:2001.09293*, 2020.
- [34] F. Memarian, Z. Xu, B. Wu, M. Wen, and U. Topcu, “Active task-inference-guided deep inverse reinforcement learning,” in *2020 59th IEEE Conference on Decision and Control (CDC)*. IEEE, 2020, pp. 1932–1938.
- [35] L. Z. Yuan, M. Hasanbeig, A. Abate, and D. Kroening, “Modular deep reinforcement learning with temporal logic specifications,” *arXiv preprint arXiv:1909.11591*, 2019.
- [36] Q. Gao, D. Hajinezhad, Y. Zhang, Y. Kantaros, and 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.
- [37] M. Hasanbeig, A. Abate, and D. Kroening, “Logically-Constrained Neural Fitted Q-Iteration,” in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2019, pp. 2012–2014.
- [38] M. Hasanbeig, D. Kroening, and A. Abate, “Deep reinforcement learning with temporal logics,” in *International Conference on Formal Modeling and Analysis of Timed Systems*. Springer, 2020, pp. 1–22.
- [39] C. Wang, Y. Li, S. L. Smith, and J. Liu, “Continuous motion planning with temporal logic specifications using deep neural networks,” *arXiv preprint arXiv:2004.02610*, 2020.
- [40] M. Kazemi and S. Soudjani, “Formal policy synthesis for continuous-state systems via reinforcement learning,” in *International Conference on Integrated Formal Methods*. Springer, 2020, pp. 3–21.
- [41] A. Lavaei, F. Somenzi, S. Soudjani, A. Trivedi, and M. Zamani, “Formal controller synthesis for continuous-space MDPs via model-free reinforcement learning,” in *2020 ACM/IEEE 11th International Conference on Cyber-Physical Systems (ICCP)*. IEEE, 2020, pp. 98–107.
- [42] A. Y. Ng, D. Harada, and S. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *ICML*, vol. 99, 1999, pp. 278–287.
- [43] C. J. Watkins and P. Dayan, “Q-learning,” *Mach. Learn.*, vol. 8, no. 3-4, pp. 279–292, 1992.
- [44] M. Kloetzer and C. Belta, “A fully automated framework for control of linear systems from temporal logic specifications,” *IEEE Transactions on Automatic Control*, vol. 53, no. 1, pp. 287–297, 2008.

- 
- [45] J. Kretínský, T. Meggendorfer, and S. Sickert, “Owl: A library for  $\omega$ -words, automata, and LTL,” in *Autom. Tech. Verif. Anal.* Springer, 2018, pp. 543–550. [Online]. Available: [https://doi.org/10.1007/978-3-030-01090-4\\_34](https://doi.org/10.1007/978-3-030-01090-4_34)
- [46] R. Durrett and R. Durrett, *Essentials of stochastic processes*. Springer, 1999, vol. 1.
- [47] A. S. McEwen, C. M. Dundas, S. S. Mattson, A. D. Toigo, L. Ojha, J. J. Wray, M. Chojnacki, S. Byrne, S. L. Murchie, and N. Thomas, “Recurring slope lineae in equatorial regions of mars,” *Nature geoscience*, vol. 7, no. 1, pp. 53–58, 2014.
- [48] S. W. Squyres, A. H. Knoll, R. E. Arvidson, J. W. Ashley, J. Bell, W. M. Calvin, P. R. Christensen, B. C. Clark, B. A. Cohen, P. De Souza *et al.*, “Exploration of victoria crater by the mars rover opportunity,” *Science*, vol. 324, no. 5930, pp. 1058–1061, 2009.