

Funnel-based Reward Shaping for Signal Temporal Logic Tasks in Reinforcement Learning

Naman Saxena[†], Gorantla Sandeep[†], and Pushpak Jagtap

Abstract—Signal Temporal Logic (STL) is a powerful framework for describing the complex temporal and logical behaviour of the dynamical system. Numerous studies have attempted to employ reinforcement learning to learn a controller that enforces STL specifications; however, they have been unable to effectively tackle the challenges of ensuring robust satisfaction in continuous state space and maintaining tractability. In this paper, leveraging the concept of funnel functions, we propose a tractable reinforcement learning algorithm to learn a time-dependent policy for robust satisfaction of STL specification in continuous state space. We demonstrate the utility of our approach on several STL tasks using different environments.

I. INTRODUCTION

Temporal logic is an effective method of formally defining complex tasks involving spatial, temporal, and logical constraints [1]. The expressiveness of predicate logic combined with temporal dimension led to its widespread use of Linear Temporal Logic (LTL) [2], [3] in designing the specification of dynamical systems. Signal Temporal Logic (STL) [4] extends the applicability of LTL by allowing us to specify a behaviour for a fixed time interval. For example, a warehouse robot needs to reach particular locations in specific time intervals. The STL framework proves to be effective in efficiently capturing and modeling such requirements.

Several works in literature use the dynamics model to design controllers to enforce STL specifications (see [5] and references therein). However, they are restricted to a limited class of systems and a fragment of specifications. Recently, controller synthesis for STL specification without a mathematical model of systems using Reinforcement Learning (RL) has started gaining attention. Not only is RL beneficial in the context of STL specifications, but at the same time, STL facilitates the designing of rewards for RL in a structured manner, avoiding the loopholes in developing rewards in a heuristic manner [3]. [6] used Q-learning [7] to achieve tasks defined using STL by maximizing the robustness of STL satisfaction. [6] defined the τ -MDP framework to allow each state to store the history of states. Storing the history of states is required to check for the satisfaction of STL formulas and, at the same time, raises doubts about the tractability of the proposed method. Further, [8] tries to solve the tractability issue of Q-learning by proposing the use of flag variables. Flag variables are used to avoid

storing the history of states and check the satisfaction of STL formulas. One issue with this work is that it does not consider robustness values. [9] proposed a solution for multi-agent system using Deep Q-learning algorithm [10]. The authors used Deep Q-Network (DQN) to overcome the state-space explosion in the multi-agent setting. Their approach considers robustness value for the STL formula but again suffers from the drawback of storing the history of states.

On the other hand, a funnel-based control approach is employed to enforce the satisfaction of a fragment of STL specifications [11] by formulating exponentially decaying constraint functions referred to as ‘funnels’. The authors in [11] introduce a continuous and closed-form controller construction for enforcing those fragments of STL specifications. Achieving this objective necessitates making certain assumptions about the systems (such as requirements of control-affine systems, fully/overactuated systems, and unbounded inputs) and about specifications (such as STL tasks excluding ‘OR’ logical operators, convex predicates, and logical operations over temporal operators). Nevertheless, this approach holds great promise in effectively separating the temporal and logical components within an STL specification while capturing robustness on satisfaction.

Leveraging the advantages of the results in [11], this work proposes a funnel-based approach for reward shaping in reinforcement learning algorithms to enforce STL specification in a tractable manner. We demonstrated that the proposed approach resolves the issue of tractability by eliminating the requirement of storing state history ([6], [9], [12], [13]). This key enhancement enables the extension of our proposed approach to continuous-state environments. Furthermore, we present evidence that, for the first time, the proposed approach integrates robustness considerations while enforcing STL specifications in the RL framework in the context of continuous-state spaces. Moreover, by leveraging a learning framework, we can relax several assumptions made in [11] concerning systems and specifications. For instance, we observe that our approach can handle the logical operator ‘OR’, any type of predicate, and conjunction between temporal operators while allowing us to provide results for any general nonlinear systems with input constraints. It is important to note that the advantages we observe are based on empirical findings and lack formal guarantees, unlike the methodology presented in [11]. Finally, we demonstrated the effectiveness of the proposed approach with several simulation results on different environments and real-world sim-to-real transfer.

This work was supported in part by the Google Research Grant, the ATPARK, and the SERB Start-up Research Grant.

N. Saxena is with the Department of Computer Science and Automation, and G. Sandeep and P. Jagtap are with the Robert Bosch Center for Cyber-Physical Systems at the Indian Institute of Science, Bangalore, India. {namansaxena, sgorantla, pushpak}@iisc.ac.in

[†]Authors contributed equally

II. PRELIMINARIES

A. Deep Q-learning

Reinforcement learning (RL) is a learning paradigm based on the framework of Markov Decision Processes (MDP) [14]. An MDP is defined by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, r, \mathcal{P}, \pi, \gamma)$, where $\mathcal{S} \subset \mathbb{R}^m$ refers to the continuous state space, \mathcal{A} refers to the discrete action space and $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is the reward function. Further, $\mathcal{P}(\cdot|s, a)$ is the transition probability function defined as $\mathcal{P} : \mathcal{S} \times \mathcal{A} \mapsto \mu(\cdot)$. $\mu : \mathcal{B}(\mathcal{S}) \mapsto [0, 1]$ is a probability measure and $\mathcal{B}(\mathcal{S})$ is the Borel σ -algebra on state space \mathcal{S} . Here, $\pi : \mathcal{S} \mapsto \Delta(\mathcal{A})$ is a stochastic policy (controller) and $\Delta(\mathcal{A})$ is the probability simplex over action space \mathcal{A} . $\gamma \in (0, 1)$ is the discount factor. The policy π is obtained by optimizing the long-term discounted reward objective function $\eta(\pi)$ as defined below:

$$\eta(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right], \quad (1)$$

where s_t, a_t denotes the state and action taken at time t . To solve the above optimization problem, Q-learning [7] is one of the most widely used RL algorithms. It uses ϵ -greedy policy (3) based on the Q-value function (2) to explore and optimize the objective function in (1).

$$\begin{aligned} Q^\pi(s_t, a_t) &= E \left[\sum_{k=t}^{\infty} \gamma^{k-t} r(s_k, a_k) | s_t, a_t \right] \\ &= E \left[r(s_t, a_t) + \gamma Q^\pi(s_{t+1}, a_{t+1}) | s_t, a_t \right]. \quad (2) \\ \pi(a|s) &= \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}|} & a = \arg \max_{a'} Q^\pi(s, a') \\ \frac{\epsilon}{|\mathcal{A}|} & \text{otherwise.} \end{cases} \quad (3) \end{aligned}$$

Here, $Q^\pi(s, a)$ is the Q-value function (2) for $(s, a) \in \mathcal{S} \times \mathcal{A}$ pair that denotes the long-term discounted reward achieved after taking action a in state s and following the policy π after that. The Q-learning finds optimal policy π^* by finding solution to Bellman equation

$$Q^{\pi^*}(s, a) = E[r(s, a) + \gamma \max_{\bar{a} \in \mathcal{A}} Q^{\pi^*}(s', \bar{a})] \quad (4)$$

that ensures $Q^{\pi^*}(s, a) \geq Q^\pi(s, a) \forall \pi, s \in \mathcal{S}$, and $a \in \mathcal{A}$.

Deep Q-learning [10] is a function approximation-based Q-learning algorithm that uses a neural network to learn optimal policy online using a replay buffer. The algorithm uses a neural network with parameters θ to approximate Q-value function $Q_\theta(s, a)$ that satisfy the Bellman equation (4).

B. Signal Temporal Logic

Signal temporal logic (STL) [4] provides a formal framework to capture high-level specifications containing spatial, temporal, and logical constraints. It consists of a set of predicates φ that are evaluated using predicate functions $h : \mathcal{S} \rightarrow \mathbb{R}$ as $\varphi := \begin{cases} \text{True,} & \text{if } h(s) \geq 0 \\ \text{False,} & \text{if } h(s) < 0 \end{cases}$. The syntax for an STL formula ϕ is given by:

$$\phi ::= \text{True} \mid \varphi \mid \neg\phi \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid F_{[a,b]}\phi \mid G_{[a,b]}\phi,$$

where $a, b \in \mathbb{R}_0^+$ with $a \leq b$, ϕ_1 and ϕ_2 are STL formulas, \neg, \wedge and \vee are logical negation, conjunction and disjunction operator, respectively; and F and G are temporal *eventually* and *always* operators, respectively. The relation $s_t \models \phi$ indicates that the signal $s : \mathbb{R}_{\geq 0} \mapsto \mathcal{S}$ satisfies the STL formula ϕ at time t . The STL semantics for a signal s is recursively defined as follows:

$$\begin{aligned} s_t \models \varphi &\iff \varphi \text{ is True} \\ s_t \models \neg\phi &\iff \neg(s_t \models \phi) \\ s_t \models \phi_1 \wedge \phi_2 &\iff s_t \models \phi_1 \wedge s_t \models \phi_2 \\ s_t \models \phi_1 \vee \phi_2 &\iff s_t \models \phi_1 \vee s_t \models \phi_2 \\ s_t \models F_{[a,b]}\phi &\iff \exists t' \in [t+a, t+b] \text{ s.t. } s_{t'} \models \phi \\ s_t \models G_{[a,b]}\phi &\iff \forall t' \in [t+a, t+b] \text{ s.t. } s_{t'} \models \phi. \quad (5) \end{aligned}$$

Next, we recall the robust semantics for STL formulas introduced by [15], which will later be used to construct rewards.

$$\begin{aligned} \rho_\varphi(s_t) &= h(s_t) \\ \rho_{\neg\phi}(s_t) &= -\rho_\phi(s_t) \\ \rho_{\phi_1 \wedge \phi_2}(s_t) &= \min(\rho_{\phi_1}(s_t), \rho_{\phi_2}(s_t)) \\ \rho_{\phi_1 \vee \phi_2}(s_t) &= \max(\rho_{\phi_1}(s_t), \rho_{\phi_2}(s_t)) \\ \rho_{F_{[a,b]}\phi}(s_t) &= \max_{t' \in [t+a, t+b]} \rho_\phi(s_{t'}) \\ \rho_{G_{[a,b]}\phi}(s_t) &= \min_{t' \in [t+a, t+b]} \rho_\phi(s_{t'}). \quad (6) \end{aligned}$$

In this paper, we consider the following fragment of STL:

$$\begin{aligned} \psi &:= \varphi \mid \neg\psi \mid \psi_1 \wedge \psi_2 \mid \psi_1 \vee \psi_2 \\ \phi_{[a,b]} &:= F_{[a,b]}\psi \mid G_{[a,b]}\psi \mid F_{[a,c_1]}G_{[c_2,b]}\psi \\ &\quad \bigwedge_{i=1}^k \phi_{[a_i, b_i]} \mid \bigwedge_{i=1}^k \phi_{[\alpha_i, \beta_i]}, \quad (7) \end{aligned}$$

where $0 \leq a \leq c_1, c_2 \leq b, b_i < a_{i+1}, \forall i \in \{1, \dots, k-1\}$, ψ and ϕ denote non-temporal and temporal formulas, respectively. Further, there may exist temporal formula with overlapping time intervals such that for some $i, j \in \{1, \dots, k-1\}$, $\alpha_i < \beta_j$ and $\alpha_j < \beta_i$.

In the next section, we discuss a funnel-based construction of rewards for deep Q-learning to learn a control policy enforcing the fragment of STL specifications given in (7).

III. PROPOSED APPROACH

In this section, we describe the utilization of funnel-based control concepts to construct time-varying rewards capturing robust satisfaction of STL specifications.

A. Construction of Rewards using Funnel Functions

Funnel based approach was first used by [11] to develop controllers that satisfy a fragment of STL specifications for known control systems. Several works in the literature now build upon this direction [16], [17]. [11] proposed to find a controller that satisfies the following relation:

$$\forall t \geq 0, \quad -\gamma(t) + \rho_{max} < \rho_\psi(s_t) < \rho_{max}, \quad (8)$$

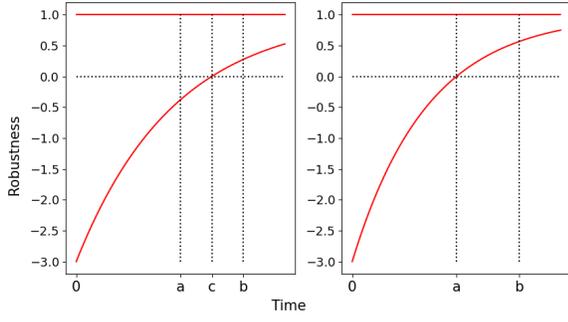


Fig. 1. The funnel for eventually operator with $t^* = c \in [a, b]$ (left) and funnel for always operator (right).

where $\gamma(t)$ is a non-increasing and continuously differentiable positive function referred to as *funnel* and defined as $\gamma(t) = (\gamma_0 - \gamma_\infty)e^{-lt} + \gamma_\infty$, where γ_0, γ_∞ , and l are positive constants with $\gamma_0 \geq \gamma_\infty$, and ρ_{max} is the maximum robustness defined for the system for corresponding non-temporal specification ψ and obtained as $\rho_{max} = \max_{s \in \mathcal{S}} \rho_\psi(s)$.

Let us consider the STL fragment defined in (7). The STL formula Φ can consist of a single eventually (F) or always (G) operator, or it could contain these operators combined using a conjunction operator. The parameter l of funnel function $\gamma(t)$ for $F_{[a,b]} \psi$, $G_{[a,b]} \psi$, and $F_{[a,c_1]}G_{[c_2,b]} \psi$ is chosen as given in Table I, while $\gamma_0 = \rho_{max} - \min_{s \in \mathcal{S}} \rho_\psi(s)$ and $\gamma_\infty \in (0, \min(\gamma_0, \rho_{max}))$ for all the temporal operators.

	t^*	l
$G_{[a,b]} \psi$	$t^* = a$	$\frac{1}{t^*} \ln \frac{\gamma_0 - \gamma_\infty}{\rho_{max} - \gamma_\infty}$
$F_{[a,b]} \psi$	$t^* \in [a, b]$	$\frac{1}{t^*} \ln \frac{\gamma_0 - \gamma_\infty}{\rho_{max} - \gamma_\infty}$
$F_{[a,c_1]}G_{[c_2,b]} \psi$	$t^* \in [a + c_2, c_1 + c_2]$	$\frac{1}{t^*} \ln \frac{\gamma_0 - \gamma_\infty}{\rho_{max} - \gamma_\infty}$

TABLE I
SELECTION OF FUNNEL FUNCTION PARAMETER l .

The illustration of the funnel for eventually and always operators is shown in Figure 1. The value for l is chosen according to the interval for temporal operators. For $F_{[a,b]}$ operator l is $\frac{\ln((\gamma_0 - \gamma_\infty)/(\rho_{max} - \gamma_\infty))}{t^*}$, where $t^* \in [a, b]$, so that $\gamma(t^*) = 0$. Note that t^* lies in $[a, b]$ because for the *eventually* operator, we want the robustness to be positive at least once in the interval $[a, b]$. For $G_{[a,b]}$ operator l is $\frac{\ln((\gamma_0 - \gamma_\infty)/(\rho_{max} - \gamma_\infty))}{t^*}$, where $t^* = a$, so that $\gamma(a) = 0$ and the robustness is positive throughout the interval $[a, b]$. Similar reasoning follows for $F_{[a,c_1]}G_{[c_2,b]}$ operator.

Now we will discuss cases where temporal operators appear in conjunction. Let us take the STL formula $\Phi = F_{[a_1,b_1]} \psi_1 \wedge G_{[a_2,b_2]} \psi_2$ with $b_1 < a_2$, where ψ_1 and ψ_2 are as defined in (7). The funnel function for Φ is given as

$$\gamma(t) = \begin{cases} (\gamma_0 - \gamma_\infty)e^{-\ln(\frac{\gamma_0 - \gamma_\infty}{\rho_{max} - \gamma_\infty})\frac{t}{c}} + \gamma_\infty, & \text{for } 0 \leq t \leq b_1, \\ (\gamma_0 - \gamma_\infty)e^{-\ln(\frac{\gamma_0 - \gamma_\infty}{\rho_{max} - \gamma_\infty})\frac{t-b_1}{a_2-b_1}} + \gamma_\infty, & \text{for } t > b_1, \end{cases} \quad (9)$$

and the plot is shown in Figure 2. In (9), $t^* = c$ lies in $[a_1, b_1]$. For $t \leq b_1$, the funnel function $\gamma(t)$ is defined according to $F_{[a_1,b_1]}$ operator and for $t > b_1$, the funnel function is defined according to $G_{[a_2,b_2]}$ operator. In this

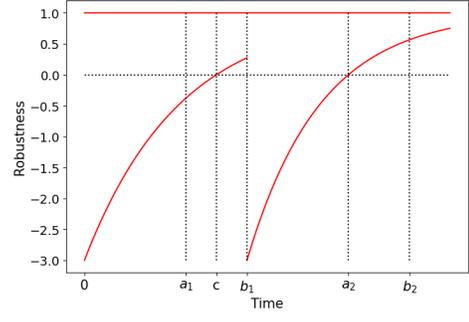


Fig. 2. The funnel function for conjunction of $F_{[a_1,b_1]}$ and $G_{[a_2,b_2]}$ operator.

case, the F operator is considered with G , but the funnel function can also be designed similarly to handle two F operators and/or two G operators. Further, this method of designing funnel function is not limited to two operators but can be extended to several operators in conjunction.

Given the construction of funnel function $\gamma(t)$ for the temporal part and the robustness measure $\rho_\psi(s_t)$ for the non-temporal part of the STL formula, we can now define the reward function for the deep Q-learning algorithm as

$$r'(s_t, a_t, t) = \rho_\psi(s_t) + \gamma(t) - \rho_{max}. \quad (10)$$

The reward function in (10) is positive at time t if the current state of the system (agent) follows the bounds given in (8); otherwise, it is negative. Further, the reward is more positive if $\rho_\psi(s_t)$ is close to ρ_{max} and more negative as it is farther away from the lower bound in (8) in the negative direction.

In general, one needs a history of states to check whether a predicate is satisfied for the time interval of the temporal operator. However, in the case of a funnel-based reward, the temporal satisfaction is captured using time-varying funnel constraints as discussed above. Thus, being inside the funnel constraints at each time instance indicates that the predicate is satisfied for the given time interval. Hence, we do not need to explicitly store the history of the states.

B. Temporal Operator with Overlapping time intervals

Let us consider the STL specification given in (11) using two temporal operator with overlapping time intervals. Here, in the specification $\alpha_1 < \beta_2$ and $\alpha_2 < \beta_1$.

$$\Phi = \underbrace{G_{[\alpha_1,\beta_1]} \phi_1}_{\psi_1} \wedge \underbrace{F_{[\alpha_2,\beta_2]} \phi_2}_{\psi_2} \quad (11)$$

Let $r'(s_t, a_t, t, \psi_1) = \rho_{\psi_1}(s_t) + \gamma_1(t) - \rho_{max,1}$ and $r'(s_t, a_t, t, \psi_2) = \rho_{\psi_2}(s_t) + \gamma_2(t) - \rho_{max,2}$. Note that $r'(s_t, a_t, t, \psi)$ denotes the reward for predicate ψ . We design the reward by defining the reward as earlier in the time intervals with no overlap and taking the minimum of the reward for different predicates with overlapping time intervals. One can write a reward for the STL formula (11) as:

$$r'(s_t, a_t, t) = \begin{cases} r'(s_t, a_t, t, \psi_1), & t \in [\alpha_1, \alpha_2) \\ \min(r'(s_t, a_t, t, \psi_1), r'(s_t, a_t, t, \psi_2)), & t \in [\alpha_2, \min(\beta_1, \beta_2)) \\ r'(s_t, a_t, t, \psi_2), & t \in [\min(\beta_1, \beta_2), \max(\beta_1, \beta_2)] \end{cases}$$

This approach is not limited to handling simultaneous overlapping between time intervals of two temporal operators but can also be used for an arbitrary number of simultaneously overlapping temporal operators (Algorithm 1). Further, in Section IV-B, we show experimentally that the reward structure proposed for overlapping intervals satisfies the STL specification robustly.

Algorithm 1 Reward Calculation for Overlapping Interval

```

1: function REWARD( $s_t, a_t, t$ )
2:   intervals =  $\{[a_i, b_i]\}_{i=0}^n \triangleright [a_i, b_i]$  is the time interval
   of  $i^{\text{th}}$  temporal operator
3:   reward = 0
4:   for  $i \in \{0, \dots, n\}$  do
5:     if  $t \in \text{interval}[i]$  then
6:       reward = min(reward,  $r'(s_t, a_t, t, \psi_i)$ )
   return reward

```

C. Time-aware Deep Q-learning

Our time-aware Deep Q-learning algorithm uses the funnel-based time-dependent reward function in (10), which is not only a function of state and action but also a function of time t . The modified MDP is defined as $\mathcal{M}' = \{\mathcal{S}, \mathcal{A}, r', \mathcal{P}, \pi', \gamma\}$, where $r' : \mathcal{S} \times \mathcal{A} \times \mathbb{N} \cup \{0\} \mapsto \mathbb{R}$ is the new reward function. The stochastic policy is now defined as $\pi' : \mathcal{S} \times \mathbb{N} \cup \{0\} \mapsto \Delta(\mathcal{A})$. The Markov property of the transition probability function is intact because the transition to s_{t+1} still depends on (s_t, a_t) and a_t depends on s_t and the current time t . Now, since the reward depends on time and consequently Q-value function also depends on time and is defined as follows:

$$\begin{aligned}
Q^\pi(s_t, a_t, t) &= E \left[\sum_{k=t}^{\infty} \gamma^{k-t} r(s_k, a_k, k) \mid s_t, a_t, t \right] \\
&= E \left[r(s_t, a_t, t) + \gamma Q^\pi(s_{t+1}, a_{t+1}, t+1) \mid s_t, a_t, t \right].
\end{aligned}$$

Further, the policy will depend on time because we use the ϵ -greedy policy, which depends on the Q-value function (3). The proposed method is summarized in Algorithm 2. The next section discusses the results obtained using our time-aware Deep Q-learning algorithm.

IV. EXPERIMENTAL RESULTS

We performed experiments on three case studies with various systems and STL properties to demonstrate the merits of the proposed approach. We trained the RL agent using a time-aware Q-learning algorithm (Algorithm 2).

A. Pendulum System

Consider an inverted pendulum model defined as:

$$\begin{aligned}
\theta_{t+1} &= \theta_t + \tau \omega_t, \\
\omega_{t+1} &= \omega_t + \tau \left(\frac{g}{l} \sin \theta_t - \frac{\mu}{ml^2} \omega_t + \frac{1}{ml^2} a_t \right),
\end{aligned}$$

where θ , ω , and $a \in \{-3, -2.9, \dots, 2.9, 3\}$ are the angle of the pendulum, angular velocity, and actions representing torque applied, respectively. $\tau = 0.01$ is the sampling time.

Algorithm 2 Time-aware Deep Q-learning

Initialize Q-value function parameter θ .

Initialize target Q-value function parameter $\bar{\theta} \leftarrow \theta$

α is the step size for parameter update.

```

1:  $k = 0, s_0 = \text{env.reset()}, t = 0$ 
2: while  $k \leq \text{total steps do}$ 
3:    $a_t \sim \pi(\cdot \mid s_t, t) \triangleright \pi$  is  $\epsilon$ -greedy policy
4:    $s_{t+1} \sim P(\cdot \mid s_t, a_t)$  and  $r_t = r'(s_t, a_t, t)$ 
5:   Store  $\{s_t, a_t, r_t, s_{t+1}, t\}$  in Replay Buffer  $\triangleright s_{t+1} = s'_t$ 
6:   if  $k \% \text{eval\_freq} == 0$  then
7:     Evaluate(agent)
8:   Sample  $\mathbb{B}_k = \{s_i, a_i, r_i, s'_i, t_i\}_{i=0}^{M-1}$  from the Replay
   Buffer
9:   Update  $\theta \leftarrow \theta + \alpha \nabla_{\theta} \left( \frac{1}{M} \sum_{i=0}^{M-1} (r'(s_i, a_i, t_i) + \right.$ 
    $\left. \gamma \max_{\bar{a}} Q_{\bar{\theta}}(s'_i, \bar{a}, t_i + 1) - Q_{\theta}(s_i, a_i, t_i) \right)^2$ 
10:  if  $k \% \text{target\_update\_freq} == 0$  then
11:    Update  $\bar{\theta} \leftarrow \theta$ 
12:   $k = k + 1$ 
13:  if  $s_{t+1}$  is terminal then
14:     $s_t = \text{env.reset()}, t = 0$ 
15:  else
16:     $s_t = s_{t+1}, t = t + 1$ 

```

The constants $g = 9.8m/s^2$, $m = 0.15$, $l = 0.5m$, and $\mu = 0.05$ represent acceleration due to gravity, mass, length of pendulum and friction coefficient, respectively. We consider the following STL specification:

$$\begin{aligned}
\Phi &= G_{[400,700]} \underbrace{(|\theta| \leq 0.05 \wedge |\omega| \leq 0.05)}_{\psi_1} \\
&\wedge G_{[1000,1200]} \underbrace{(|1.57 - \theta| \leq 0.05 \wedge |\omega| \leq 0.05)}_{\psi_2} \\
&\wedge G_{[1700,2000]} \underbrace{(|-1.57 - \theta| \leq 0.05 \wedge |\omega| \leq 0.05)}_{\psi_3}.
\end{aligned} \tag{12}$$

In simple words, the specification in (12) says that "the pendulum should maintain $\theta = 0$ and $\omega = 0$ in the interval of 400 to 700 timesteps (which is equivalent to 4 to 7 seconds as τ is 0.01 seconds) with a tolerance value of 0.05. Subsequently, in the interval of 1000 to 1200 time steps, the pendulum should be balanced at $\theta = 1.57$ and $\omega = 0$ followed by $\theta = -1.57$ and $\omega = 0$ from 1700 to 2000 time steps with a tolerance value of 0.05". The funnel-based reward function for (12) with $s_t = [\theta_t, \omega_t]$ is computed as discussed in Section III-A and given as follows:

$$\begin{aligned}
r'(s_t, a_t, t) &= \begin{cases} \rho_{\psi_1}(s_t) + (\gamma_{0,1} - \gamma_{\infty,1}) e^{-l_1 t} + \gamma_{\infty,1} - \rho_{\max,1} & t \in [0, 700] \\ \rho_{\psi_2}(s_t) + (\gamma_{0,2} - \gamma_{\infty,2}) e^{-l_2(t-700)} + \gamma_{\infty,2} - \rho_{\max,2} & t \in [700, 1200] \\ \rho_{\psi_3}(s_t) + (\gamma_{0,3} - \gamma_{\infty,3}) e^{-l_3(t-1200)} + \gamma_{\infty,3} - \rho_{\max,3} & t \in [1700, 2000], \end{cases} \\
&\tag{13}
\end{aligned}$$

where $\rho_{\psi_1}(s_t) = 0.05 - \min(|\theta_t|, |\omega_t|)$, $\rho_{\psi_2}(s_t) = 0.05 - \min(|1.57 - \theta_t|, |\omega_t|)$, $\rho_{\psi_3}(s_t) = 0.05 - \min(|-1.57 - \theta_t|, |\omega_t|)$, $l_1 = 0.0103$ $l_2 = 0.0138$ $l_3 = 0.0083$, $\gamma_{0,i} = \pi$,

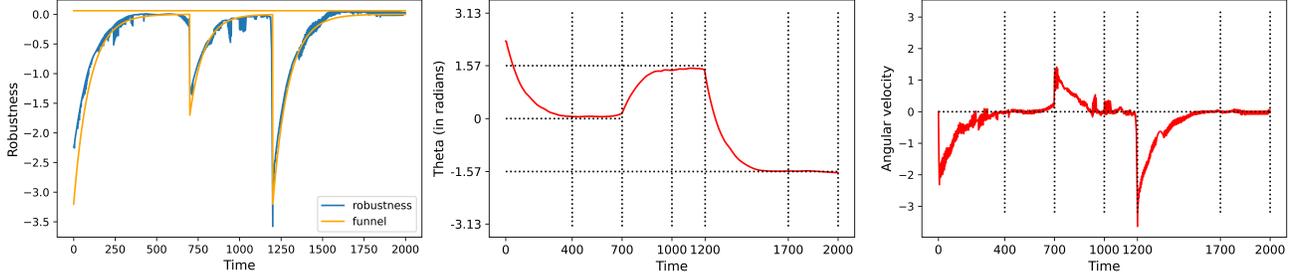


Fig. 3. The evolution of robustness values (left), angle (middle), and angular velocity (right) of the pendulum.

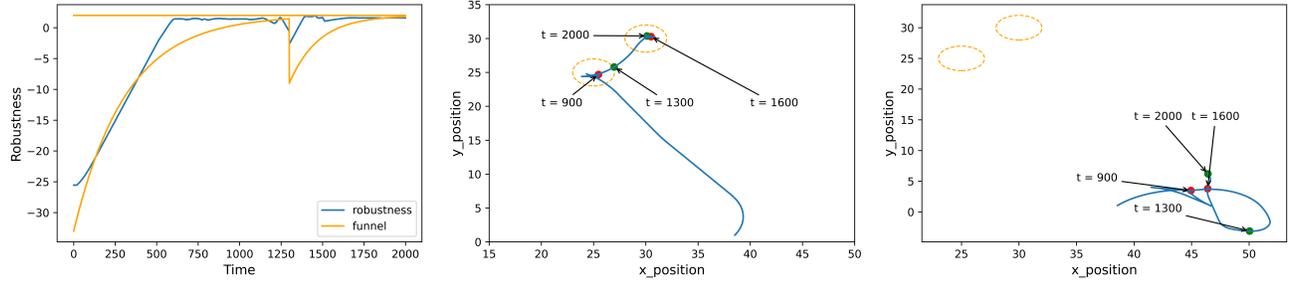


Fig. 4. The evolution of robustness values (left), the trajectory in the 2-d plane with controller learned using reward function with (middle) and without (right) funnel for a mobile robot.

$\gamma_{\infty,i} = 0.01$, $\rho_{max,i} = 0.05$, for all $i \in \{1, 2, 3\}$. Then, we trained the RL agent using Algorithm 2 to obtain policy enforcing desired STL specifications. Figure 3 shows the evolution of robustness values, angle and angular velocity of the pendulum over time. One can readily observe that the robustness values are always inside the constructed funnels, and the property is satisfied. Notice that the dynamics is under-actuated, so one can not use results in [11].

B. Mobile Robot Navigation

For the second case study, we consider the differential drive mobile robot described by: $x_{t+1} = x_t + \tau v_t \cos \theta_t$, $y_{t+1} = y_t + \tau v_t \sin \theta_t$, $\theta_{t+1} = \theta_t + \tau \omega_t$, where, x and y represent the location of robot in $x - y$ plane, θ represents orientation, and $\tau = 0.01$ is the sampling time. The actions $v_t \in \{-5, -4.5, \dots, 4.5, 5\}$ and $\omega_t \in \{-3, -2.5, \dots, 2.5, 3\}$ represent forward and angular velocity, respectively. The STL specification considered is

$$\Phi = G_{[900,1300]} (\|(x, y) - (25, 25)\|_2 \leq 2) \wedge G_{[1600,2000]} (\|(x, y) - (30, 30)\|_2 \leq 2).$$

To satisfy the above STL specification, the agent has to reach inside a circle of radius two centered at (25,25) in the interval from 900 to 1300 timesteps. Further, it should move inside the circle of radius two centered at (30,30). By constructing funnel-based reward, we learn the time-dependent policy using Algorithm 2 to enforce the STL specification. Figure 4 shows the plot of robustness values following the constructed funnel (left plot) and trajectory followed by the trained RL agent in the 2-d plane (middle plot).

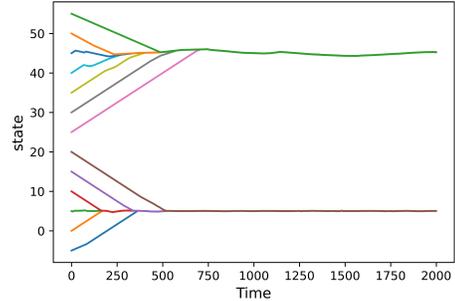


Fig. 5. The evolution of trajectories from different initial locations for the discrete-time integrator.

To show the importance of a funnel-based reward structure, we modified the reward function by eliminating the funnel part as described below:

$$r'(s_t, a_t, t) = \begin{cases} \rho_{\psi_1} & \text{for } 0 \leq t \leq 1300 \\ \rho_{\psi_2} & \text{for } 1300 \leq t \leq 2000. \end{cases} \quad (14)$$

Our ablation study reveals that using robustness values without funnel function $\gamma(t)$ does not help the agent learn the task. The trajectory in the 2-d plane obtained after using the reward function given in (14) is shown in Figure 4 (right).

Specification with convex predicates: To showcase the applicability of the results for the convex predicate (which is one of the limitations in [11]), we consider the following

STL specification:

$$\Phi = G_{[300,2000]} \left(\underbrace{\|(x, y) - (5, 5)\|_2 \geq 2}_{\psi_1: (\text{convex predicate})} \right. \\ \left. \wedge \underbrace{\|(x, y) - (5, 5)\|_2 \leq 5}_{\psi_2: (\text{concave predicate})} \right). \quad (15)$$

According to specification (15), the robot has to stay inside an annular region centered at (5,5) with the inner radius two and the outer radius 5 for the interval of 300 to 2000 time steps. The agent is always reset at a randomly chosen point in a $[0,15] \times [0,15]$ grid. Figure 7 shows the plot of robustness values and the trajectory followed by the robot starting from some random point under the policy learned using our proposed funnel-based STL satisfaction approach in the 2-d plane. One can readily observe that the trajectory achieves the best possible robustness by staying in the centre of the annular region and satisfying the convex predicate ψ_1 . Hence our method is capable of handling convex predicate.

Specification with overlapping time intervals: We tried STL specification with overlapping time intervals of temporal operator for the differential drive mobile robot and found that our method is able to satisfy specifications with overlapping time intervals as well. The STL specification $\Phi = G_{[0,100]} \left((\|(x, y) - (2, 2)\|_\infty \leq 2) \wedge (\|(x, y) - (2, 2)\|_\infty \geq 0.5) \wedge (\|(x, y) - (2, 0.5)\|_\infty \geq 0.5) \wedge (\|(x, y) - (3, 3)\|_\infty \geq 0.5) \right) \wedge F_{[0,50]} (\|(x, y) - (3, 1)\|_2 \leq 0.3) \wedge F_{[50,90]} (\|(x, y) - (1, 3)\|_2 \leq 0.3)$, and the time intervals are provided in seconds. The robustness plots obtained are given in Figure 6. We implemented the trained RL agent on hardware, and a video demonstration* is also available for the same.

C. Discrete-time Integrator for specification with Disjunction inside Temporal Operator

In this section, we show the applicability of our approach to learning policy for STL specification with the disjunction between predicates inside the temporal operator (which is one of the limitations in [11]). Consider the system (discrete-time integrator) with dynamics $x_{t+1} = x_t + \tau v_t$, where x_t is the location at time t and $v_t \in \{-3, -2.5, \dots, 2.5, 3\}$ is the velocity given as input to the system. We train the RL agent for the STL specification $G_{[0,2000]}(\varphi_1 \vee \varphi_2)$. Here, $\varphi_1 := |x - 5| \leq 5$ and $\varphi_2 := |x - 45| \leq 5$. We plotted trajectories of the system in Figure 5 for different initial locations and found that the agent has learnt to reach either $x = 5$ or $x = 45$.

D. 7-DoF Fetch Mobile Manipulator

Here we consider a manipulator arm (Figure 8) environment [18] of "Gymnasium-Robotics" suite of tasks. The task of the manipulator is to move to different points in 3D space according to different time intervals. The actual task is defined using the STL specification $\Phi = G_{[50,100]} (\|(x, y, z) - (1.5, 0.43, 0.47)\|_2 \leq 0.1) \wedge$

* The video can be found at this YouTube link: <https://www.youtube.com/watch?v=f60-LhD-8PM>

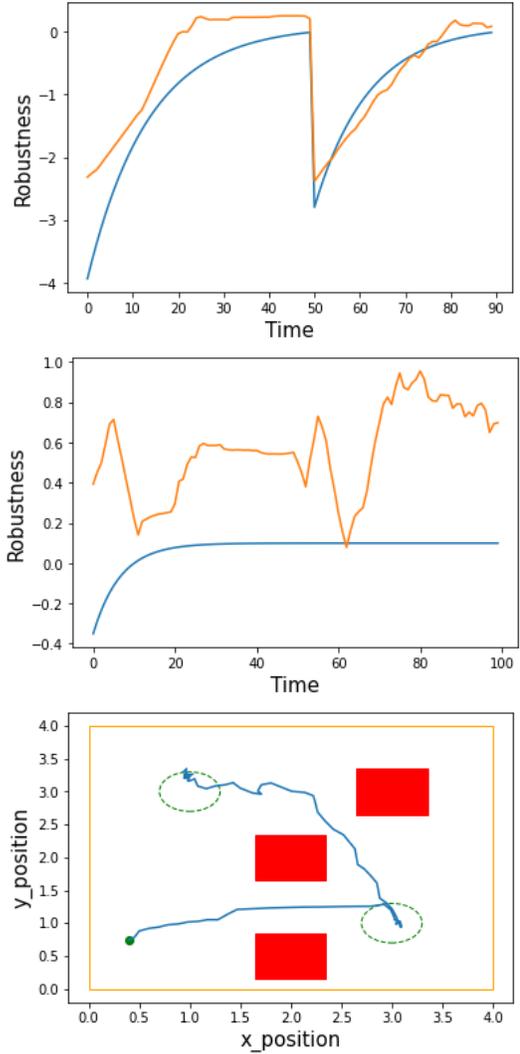


Fig. 6. Evolution of robustness value for eventually reaching goals (top), robustness for always avoiding obstacles (middle) and trajectory followed by differential drive robot hardware (bottom)

$G_{[150,200]} (\|(x, y, z) - (1.5, 1.05, 0.47)\|_2 \leq 0.1)$. Because of the continuous action space, we used the TD3 algorithm [19] to train the RL agent. The robustness plot is given as Figure 9 and a video demonstration* is also available for the manipulator. The robustness plot shows that our method is capable of handling both continuous state and action space. To substantiate our claim we will provide robustness results on two tasks with continuous action space.

E. DeepMind Cartpole Balance

We test our proposed method using the CartPole-Balance benchmarking task from the DeepMind control suite ([20]). In the task considered the RL controller has to balance the pole at one position on the x -axis during a particular time interval, and subsequently, the pole has to be balanced at different positions. The temporal constraints-based task is described using the STL specification: $\Phi = G_{[400,500]} (|x - 0.5| \leq 0.2 \wedge |\theta| \leq 0.1 \wedge |\omega| \leq 0.5) \wedge G_{[900,1000]} (|x + 0.5| \leq$

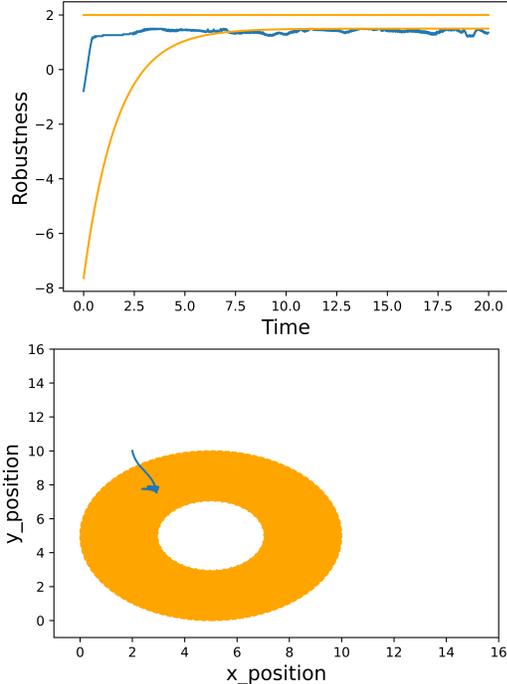


Fig. 7. The evolution of robustness values (top), the trajectory of the agent in the 2-d plane (bottom), colored area represents the safe region.

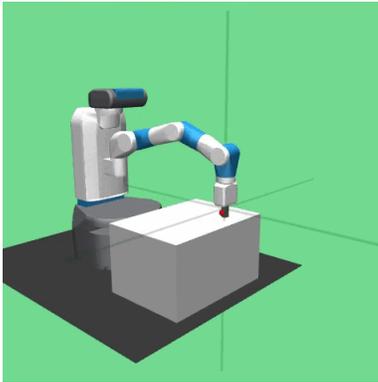


Fig. 8. Manipulator arm from the Gymnasium-Robotics suite of tasks

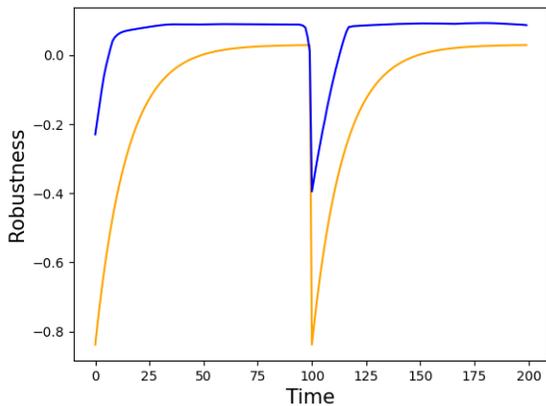


Fig. 9. Evolution of robustness values for manipulator arm

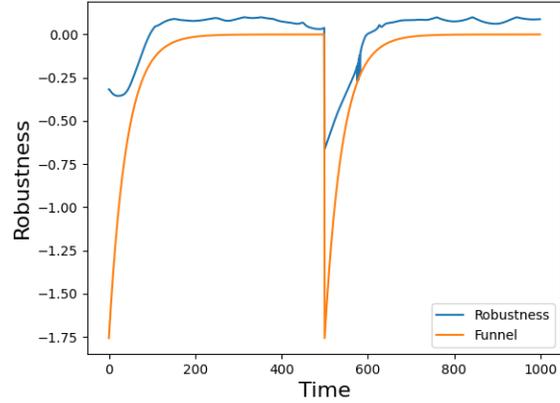
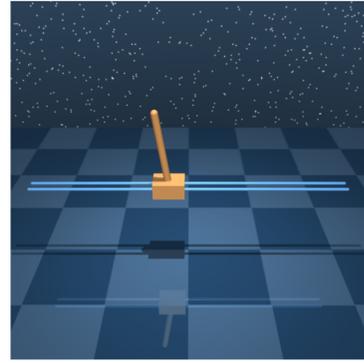


Fig. 10. DeepMind control suite CartPole-Balance task(top). Evolution of robustness values for cartpole (bottom).

$0.2 \wedge |\theta| \leq 0.1 \wedge |\omega| \leq 0.5$). Here x is the position of the cart on the x -axis, θ is the angle made by the pole with position y -axis, and ω is the angular velocity. The robustness plot is shown in Figure 10. One can readily observe that the robustness of learned policy satisfies the funnel constraints which ensures satisfaction of specifications.

F. DeepMind Ball-in-Cup

We further evaluated our method on the ball-in-cup benchmarking task from the DeepMind control suite ([20]). In this task, the RL agent has to move the cup to capture the ball first, then release the ball, and finally capture the ball again with specific time constraints. The exact task is described by the STL specification as: $\Phi = G_{[200,300]}(\|(x, z)_{ball} - (x, z)_{cup}\|_2 \leq 0.1) \wedge G_{[500,600]}(\|(x, z)_{ball} - (x, z)_{cup}\|_2 \geq 0.2) \wedge G_{[800,1000]}(\|(x, z)_{ball} - (x, z)_{cup}\|_2 \leq 0.1)$. Here, $(x, z)_{ball}$ and $(x, z)_{cup}$ are the location of the ball and the cup respectively in a 2-D plane. The robustness plot is available in Figure 11 which implies satisfaction of specification.

G. Comparative Study

We compare our proposed method for STL satisfaction with the flag-based method proposed in [8] by training the RL agent for the differential drive mobile robot (Section IV-B) considering the STL specification $\Phi = G_{[0,200]}(\|(x, y) - (2, 2)\|_2 \leq 1)$. We calculated the robustness value by taking

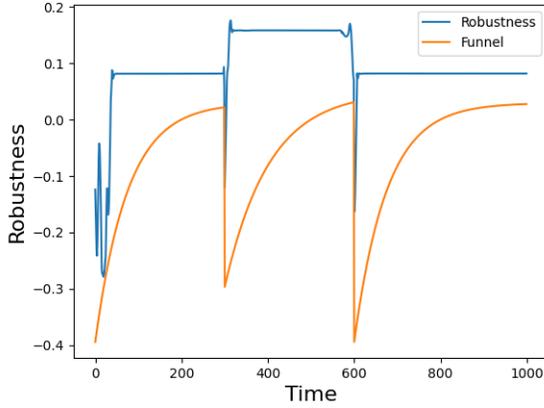


Fig. 11. DeepMind control suite Ball-in-cup task(top). Evolution of robustness values for Ball-in-cup task (bottom).

a minimum of robustness for all time steps and obtained a robustness value of 0.938 for our method and a robustness value of 0.102 for the method suggested in [8]. Also, from Figure 12, it is clear that the trajectory generated by using our approach appears more robust. Further, we obtain an execution time of 237.72 minutes for our proposed method as compared to 236.84 minutes for [8] that shows our method achieves better robustness while consuming almost the same amount of time as the state-of-the-art method in [8]. We could not compare our method on infinite state space tasks with the method proposed in [6] because the method is proposed for finite state space only.

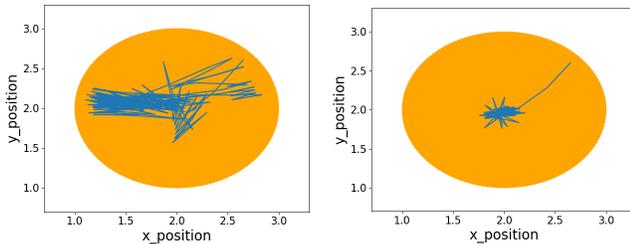


Fig. 12. Trajectories generated using the flag-based method [8] (left) and the proposed method (right).

V. CONCLUSION

In this paper, we proposed a tractable method to learn a controller for robust satisfaction of STL specification using time-aware deep Q-learning. We described how funnel functions could be used to design reward functions for reinforcement learning algorithms that allow learning controllers for STL specifications. We showed the performance of our method on various environments, such as the pendulum and mobile robot, in accomplishing time-constraint sequential goals. One of the significant advantages of the proposed approach is the satisfaction of the STL formulas with convex predicates. Further, we demonstrated, using a simple environment, that we can learn the controller for STL formula with disjunction operator inside temporal operator.

REFERENCES

- [1] C. Baier and J.-P. Katoen, *Principles of model checking*. MIT press, 2008.
- [2] X. Li, Y. Ma, and C. Belta, “A policy search method for temporal logic specified reinforcement learning tasks,” in *American Control Conference (ACC)*. IEEE, 2018, pp. 240–245.
- [3] X. Li, C.-I. Vasile, and C. Belta, “Reinforcement learning with temporal logic rewards,” in *IEEE/RSS International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 3834–3839.
- [4] O. Maler and D. Nickovic, “Monitoring temporal properties of continuous signals,” in *Formal Techniques, Modelling and Analysis of Timed and Fault-Tolerant Systems*. Springer, 2004, pp. 152–166.
- [5] M. Lahijanian, S. B. Andersson, and C. Belta, “Formal verification and synthesis for discrete-time stochastic systems,” *IEEE Transactions on Automatic Control*, vol. 60, no. 8, pp. 2031–2045, 2015.
- [6] D. Aksaray, A. Jones, Z. Kong, M. Schwager, and C. Belta, “Q-learning for robust satisfaction of signal temporal logic specifications,” in *IEEE 55th Conference on Decision and Control (CDC)*. IEEE, 2016, pp. 6565–6570.
- [7] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3, pp. 279–292, 1992.
- [8] H. Venkataraman, D. Aksaray, and P. Seiler, “Tractable reinforcement learning of signal temporal logic objectives,” in *Learning for Dynamics and Control*. PMLR, 2020, pp. 308–317.
- [9] D. Muniraj, K. G. Vamvoudakis, and M. Farhood, “Enforcing signal temporal logic specifications in multi-agent adversarial environments: A deep q-learning approach,” in *IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 4141–4146.
- [10] 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.
- [11] L. Lindemann, C. K. Verginis, and D. V. Dimarogonas, “Prescribed performance control for signal temporal logic specifications,” in *IEEE 56th Annual Conference on Decision and Control (CDC)*. IEEE, 2017, pp. 2997–3002.
- [12] J. Wang, S. Yang, Z. An, S. Han, Z. Zhang, R. Mangharam, M. Ma, and F. Miao, “Multi-agent reinforcement learning guided by signal temporal logic specifications,” *arXiv preprint arXiv:2306.06808*, 2023.
- [13] J. Ikemoto and T. Ushio, “Deep reinforcement learning under signal temporal logic constraints using lagrangian relaxation,” *IEEE Access*, vol. 10, pp. 114 814–114 828, 2022.
- [14] M. L. Puterman, “Markov decision processes,” *Handbooks in operations research and management science*, vol. 2, pp. 331–434, 1990.
- [15] A. Donzé and O. Maler, “Robust satisfaction of temporal logic over real-valued signals,” in *International Conference on Formal Modeling and Analysis of Timed Systems*. Springer, 2010, pp. 92–106.
- [16] L. Lindemann and D. V. Dimarogonas, “Feedback control strategies for multi-agent systems under a fragment of signal temporal logic tasks,” *Automatica*, vol. 106, pp. 284–293, 2019.
- [17] P. Varnai and D. V. Dimarogonas, “Prescribed performance control guided policy improvement for satisfying signal temporal logic tasks,” in *American Control Conference (ACC)*. IEEE, 2019, pp. 286–291.

- [18] M. Plappert, M. Andrychowicz, A. Ray, B. McGrew, B. Baker, G. Powell, J. Schneider, J. Tobin, M. Chociej, P. Welinder, V. Kumar, and W. Zaremba, "Multi-goal reinforcement learning: Challenging robotics environments and request for research," 2018.
- [19] S. Fujimoto, H. Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *International conference on machine learning*. PMLR, 2018, pp. 1587–1596.
- [20] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, *et al.*, "Deepmind control suite," *arXiv preprint arXiv:1801.00690*, 2018.