# Enforcing Signal Temporal Logic Specifications in Multi-Agent Adversarial Environments: A Deep Q-Learning Approach

Devaprakash Muniraj, Kyriakos G. Vamvoudakis, and Mazen Farhood

Abstract—This work addresses the problem of learning optimal control policies for a multi-agent system in an adversarial environment. Specifically, we focus on multi-agent systems where the mission objectives are expressed as signal temporal logic (STL) specifications. The agents are classified as either defensive or adversarial. The defensive agents are maximizers, namely, they maximize an objective function that enforces the STL specification; the adversarial agents, on the other hand, are minimizers. The interaction among the agents is modeled as a finite-state team stochastic game with an unknown transition probability function. The synthesis objective is to determine optimal control policies for the defensive agents that implement the STL specification against the best responses of the adversarial agents. A multi-agent deep Q-learning algorithm, which is an extension of the minimax Q-learning algorithm, is then proposed to learn the optimal policies. The effectiveness of the proposed approach is illustrated through a simulation case study.

Index Terms—multi-agent system, signal temporal logic, deep Q-learning.

#### I. Introduction

Safety-critical autonomous systems typically operate in uncertain, dynamic, and adversarial environments. Additionally, many autonomous vehicle applications involve executing time-constrained tasks such as reaching a particular region within a specified time or visiting a particular region every few minutes. The problem setup considered in this work is that of a group of autonomous agents with unknown stochastic dynamics tasked to perform a time-constrained mission, which is expressed as a temporal logic specification, in an uncertain, adversarial environment. The agents in the multi-agent system are classified as either *defensive* agents or *adversarial* agents. The objective of the defensive agents is to complete the specified mission, while the objective of the adversarial agents is to act in a manner so as to prevent the defensive agents from achieving the mission.

The interaction among the agents of the multi-agent system is modeled as a team stochastic game, which is an extension of one-shot matrix games to environments modeled as Markov decision processes [1]. Temporal logic enables one to specify the mission objectives without any ambiguity (see [2] and the references therein). Signal temporal logic

Devaprakash Muniraj and Mazen Farhood are with the Kevin T. Crofton Department of Aerospace and Ocean Engineering, Virginia Tech, Blacksburg, VA 24061, USA. Kyriakos G. Vamvoudakis is with the School of Aerospace Engineering, Georgia Institute of Technology, Atlanta, GA 30332-0150, USA. Email: devapm@vt.edu,kyriakos@gatech.edu, farhood@vt.edu.

This work was supported in part by the Center for Unmanned Aircraft Systems (C-UAS), a National Science Foundation (NSF) sponsored industry/university cooperative research center (I/UCRC) under NSF Grant Nos. IIP-1539975 and CNS-1650465, along with significant contributions from C-UAS industry members, NSF CAREER under Grant No. CPS-1750789, and NATO under Grant No. SPS G5176.

(STL) [3], which can be used to specify temporal properties of real-valued signals over finite time intervals, is used as the temporal logic framework to specify the mission objective. The STL specification is enforced by maximizing the expected robustness degree, which is a real number that denotes how well a given signal satisfies or violates the STL specification. The control synthesis objective for the defensive agents is to find optimal policies that maximize the expected robustness degree of the STL specification against the best responses of the adversarial agents.

In this work, a model-free reinforcement learning (RL) approach is adopted to learn the optimal policies of the defensive agents. The objective function that enforces the STL specification is not in the form of a sum of discounted rewards, and therefore RL methods are not directly applicable. Following an approach similar to [4], the objective function is approximated using an extended state space and is written as a sum of discounted rewards. Since the number of states of the team stochastic game grows exponentially with the number of agents and the horizon length of the STL formula, conventional multi-agent RL methods [5]–[7] lose their viability as they use a tabular representation to store data pertaining to the action-value function. Recent advances in deep reinforcement learning reported in [8], [9] overcome many limitations of conventional RL algorithms by representing the action-value function as a deep neural network. In this work, we propose a minimax deep Qlearning algorithm that combines ideas from deep Q-learning [9] and minimax Q-learning [5] to learn the optimal policies. Related work: The problem of finding optimal policies for multi-agent systems with mission objectives specified as temporal logic specifications has been studied in [10]–[12]. Most of the existing approaches require precise knowledge of the system dynamics. In our proposed approach, the system model and the environment are assumed to be unknown and the agents learn the optimal policies by interacting with the environment through a sequence of observations, actions, and rewards. Minimax Q-learning is a model-free RL technique proposed in [5] to learn the optimal policies for two-player zero-sum stochastic games. In this work, an extension of the minimax Q-learning algorithm is proposed, wherein the Qfunction is represented using a convolutional neural network. Although extensions of deep RL algorithms to multi-agent systems have been proposed in [13]-[15], to the best of the authors' knowledge, the problem of finding optimal control policies that enforce an STL specification in an adversarial environment has not been studied before.

Structure: The remainder of the paper is structured as follows. A brief overview of STL and stochastic games is provided in Section II. In Section III, the control synthesis problem is stated and then reformulated as a multi-agent RL problem. In Section IV, the proposed minimax deep Q-learning algorithm is presented. Section V presents a simulation case study that illustrates the effectiveness of the proposed approach. Finally, conclusions are presented in Section VI.

# II. PRELIMINARIES

## A. Signal Temporal Logic

In this work, the mission objective is formally specified in STL [3], which is an extension of linear temporal logic (LTL) [16] to real-valued signals. STL formulas are defined using the following syntax:

 $\varphi:=\chi^{\eta}\mid\neg\chi^{\eta}\mid\neg\varphi\mid\varphi\wedge\varphi\mid\varphi\vee\varphi\mid F_{[a,b]}\varphi\mid G_{[a,b]}\varphi,$  where  $\chi^{\eta}$  is an atomic predicate whose truth value depends on the value of the function  $\eta:\mathbb{R}^m\to\mathbb{R}$ , and  $\varphi$  is an STL formula. In this work,  $\eta$  is assumed to be a function of the form  $\eta:=d-f(\mathbf{x})$ , where  $\mathbf{x}:\mathbb{R}_+\to\mathbb{R}^m$  denotes the signal and d is a real-valued constant. The Boolean operators  $\neg$ ,  $\vee$ ,  $\wedge$  have their usual meanings, and a, b are non-negative scalars used to denote time bounds. The operators F and G are analogous to the LTL operators  $\diamondsuit$  (Eventually) and  $\square$  (Always), respectively.

The signal  ${\bf x}$  here is a continuous-time signal, however, in practical applications, measurements of the signal might be available only at discrete time instances. In this work, we assume that it is always possible to obtain a continuous-time signal from discrete measurements. For instance, given a sequence of discrete measurements  $x_k$  of the signal  ${\bf x}$  that correspond to time instances  $t_k$ , where k takes values in the set of natural numbers, a possible reconstruction of the continuous-time signal  ${\bf x}$  is  ${\bf x}(t){=}x_k$ , if  $t\in[t_k,t_{k+1})$ .  $({\bf x},t)$  denotes the trace of the signal starting at time t and going forward. For given STL formulas  $\phi$  and  $\psi$ , the Boolean semantics of STL are defined as follows:

$$\begin{split} (\mathbf{x},t) &\models \chi^{\eta} &\Leftrightarrow \eta(x_t) \geqslant 0, \\ (\mathbf{x},t) &\models \phi \lor \psi &\Leftrightarrow (\mathbf{x},t) \models \phi \text{ or } (\mathbf{x},t) \models \psi, \\ (\mathbf{x},t) &\models \phi \land \psi &\Leftrightarrow (\mathbf{x},t) \models \phi \text{ and } (\mathbf{x},t) \models \psi, \\ (\mathbf{x},t) &\models G_{[a,b]}\phi \Leftrightarrow (\mathbf{x},\hat{t}) \models \phi \ \forall \hat{t} \in [t+a,t+b], \text{ and} \\ (\mathbf{x},t) &\models F_{[a,b]}\phi \Leftrightarrow \exists \hat{t} \in [t+a,t+b] \quad \text{s.t. } (\mathbf{x},\hat{t}) \models \phi. \end{split}$$

In addition to the Boolean semantics, STL formulas can also be interpreted using quantitative semantics, which are defined using the robustness degree. For a given signal  ${\bf x}$  and an STL specification  $\phi$ , the robustness degree denoted by  $\rho^\phi({\bf x},t)$  is defined as follows:

$$\rho^{\mathbf{x}^{\eta}}(\mathbf{x},t) = \eta(x_t), \ \rho^{-\mathbf{x}^{\eta}}(\mathbf{x},t) = -\eta(x_t),$$

$$\rho^{\phi \vee \psi}(\mathbf{x},t) = \max \ (\rho^{\phi}(\mathbf{x},t),\rho^{\psi}(\mathbf{x},t)),$$

$$\rho^{\phi \wedge \psi}(\mathbf{x},t) = \min \ (\rho^{\phi}(\mathbf{x},t),\rho^{\psi}(\mathbf{x},t)),$$

$$\rho^{G_{[a,b]}\phi}(\mathbf{x},t) = \min_{\hat{t} \in [t+a,t+b]} \rho^{\phi}(\mathbf{x},\hat{t}), \text{ and }$$

$$\rho^{F_{[a,b]}\phi}(\mathbf{x},t) = \max_{\hat{t} \in [t+a,t+b]} \rho^{\phi}(\mathbf{x},\hat{t}).$$

The signal  ${\bf x}$  satisfies the STL formula  $\phi$  at time t if and only if  $\rho^\phi({\bf x},t){\geqslant}0$ . The robustness degree provides a quantitative measure of the extent to which a signal trace satisfies an STL

specification. Another definition which will be used in the sequel is the horizon length of an STL formula  $\phi$ , denoted by  $H(\phi)$ , which is the duration of time required to verify whether a signal satisfies the formula  $\phi$ . For instance, the horizon length of the STL formula  $G_{[2,5]}\phi$  is 5. The horizon length of any other STL formula is computed as in [17].

#### B. Stochastic Games

Stochastic games, also known as Markov games, are an extension of Markov decision processes to multiple agents [1]. An n-player discounted stochastic game is defined by the tuple  $G = \langle S, A_1, \ldots, A_n, P, R_1, \ldots, R_n, \gamma \rangle$ , where S is the set of states of the multi-agent system and  $A_i$  is the finite set of actions available to agent  $i.\ P:\ S \times A_1 \times \cdots \times A_n \times S \rightarrow [0,1]$  is the transition probability function that gives a probability distribution for the next state of the system given the current state and the joint action of the agents. For instance,  $P(s,a_1,a_2,s')$  is the probability of ending in state s', given that the agents start in state s and each agent s is the reward function of agent s, and s is the reward function of agent s, and s is the discount factor.

A (stationary) control policy  $\pi_i$  for agent i is a mapping  $\pi_i \colon S \to PD(A_i)$  from the set of states to the set of discrete probability distributions over the set of control actions of agent i. As in [6], given the joint policy  $\pi = (\pi_1, \dots, \pi_n)$ , the action-value function  $Q_i^{\pi}(s, a_1, \dots, a_n)$  can be defined as

$$Q_i^{\pi}(s, a_1, \dots, a_n) = \sum_{s' \in S} P(s, a_1, \dots, a_n, s') \times [R_i(s, a_1, \dots, a_n, s') + \gamma V_i^{\pi}(s')].$$

 $V_i^\pi(\cdot)$  in the above equation is the state-value function under  $\pi$  and is defined as

$$V_i^{\pi}(s) = \sum_{j=0}^{\infty} \gamma^j E^{\pi}(R_i(s_j, a_1, \dots, a_n, s_{j+1})),$$

where  $E^{\pi}(\cdot)$  is the expectation operator over  $\pi$ . The Q-functions are used to relate the equilibrium in stage games to the equilibrium of a stochastic game [18], [19]. Specifically, if the Q-function at each state is considered as the payoff function for the stage game, then the stage policies are in equilibrium if and only if the multistage policies are in equilibrium.

Definition 1: (Nash equilibrium) A Nash equilibrium (NE) in a stochastic game G is an n-tuple of strategies  $(\pi_1^*, \ldots, \pi_n^*)$  such that  $\forall s \in S$  and  $i \in \{1, \ldots, n\}$ , the following inequalities hold:

$$V_i^{(\pi_1^*,\dots,\pi_i^*,\dots,\pi_n^*)}(s) \geqslant V_i^{(\pi_1^*,\dots,\pi_i,\dots,\pi_n^*)}(s) \quad \forall \pi_i \in \Pi_i,$$

where  $\Pi_i$  is the set of strategies available to player i. If the reward functions of the agents are in conflict with each other, then a special case of NE called the adversarial

each other, then a special case of NE called the adversarial equilibrium can be defined.

Definition 2: (Adversarial equilibrium) An adversarial

equilibrium in a stochastic game G is an n-tuple of strategies  $(\pi_1^*,\ldots,\pi_n^*)$  such that  $\forall s{\in}S$  and  $i{\in}\{1,\ldots,n\}$ , the following are true:  $V_i^{(\pi_1^*,\ldots,\pi_i^*,\ldots,\pi_n^*)}(s)\geqslant V_i^{(\pi_1^*,\ldots,\pi_i,\ldots,\pi_n^*)}(s)\quad\text{and}$ 

$$\begin{split} V_{i}^{(\pi_{1},...,\pi_{i},...,\pi_{n})}(s) &\geqslant V_{i}^{(\pi_{1},...,\pi_{i},...,\pi_{n})}(s) \quad \text{and} \\ V_{i}^{(\pi_{1}^{*},...,\pi_{i}^{*},...,\pi_{n}^{*})}(s) &\leqslant V_{i}^{(\pi_{1},...,\pi_{i}^{*},...,\pi_{n})}(s) \quad \forall \pi_{i} \in \Pi_{i}. \end{split}$$

If an agent unilaterally deviates from the adversarial equi-

librium, it not only decreases its expected payoff, but also increases the expected payoffs of the other agents. If a stochastic game has only two players and their reward functions are such that  $R_1(s,a_1,a_2,s')=-R_2(s,a_1,a_2,s')$   $\forall s{\in}S,\ s'{\in}S,\ a_1{\in}A_1,$  and  $a_2{\in}A_2,$  then the stochastic game is called a zero-sum stochastic game. All Nash equilibria in a zero-sum stochastic game are adversarial equilibria [20]. The concepts of Nash equilibrium and adversarial equilibrium are important in multi-agent reinforcement learning, as they are used as learning goals by most of the RL algorithms.

## III. PROBLEM FORMULATION

### A. System Model and Control Objective

The multi-agent system that we consider in this work consists of two teams of agents, namely the team of defensive agents and the team of adversarial agents. The interaction among the agents is modeled as a team stochastic game, which is an n-player stochastic game consisting of a team of p maximizers and (n-p) minimizers [21]. During the game, the defensive agents (maximizers) and the adversarial agents (minimizers) choose actions that respectively maximize and minimize their expected discounted return.

Let the team stochastic game be denoted by the tuple  $\bar{G} = \langle S, \bar{A}_1, \bar{A}_2, P, R, \gamma \rangle$ , where  $\bar{A}_1 = A_1 \times \cdots \times A_p$  and  $\bar{A}_2 = A_{p+1} \times \cdots \times A_n$  are the joint action sets of the defensive agents and the adversarial agents, respectively. The state space S is composed of discrete partitions of the continuous state space. For example, if the continuous state space is a two-dimensional grid, then  $s_k \in \mathbb{R}^{2n}$  and is of the form  $s_k = (x_1, y_1, \dots, x_n, y_n)^T$ , where  $(x_i, y_i)$  denotes the centroid of the cell that contains the  $i^{th}$  agent. All the agents share the same reward function R. A state trajectory of the system is of the form  $(s_0s_1 \dots s_N)$ , where  $s_k$  corresponds to the state of the system at time  $t_k$ , and for  $k = 1, \dots, N$ ,  $t_k - t_{k-1} = \Delta t$  with  $\Delta t$  being the time step. The state trajectory of the system during the time interval [0, T] is denoted by  $s_{0:T}$  and is given by  $(s_0s_1 \dots s_N)$ , where  $T = N\Delta t$ .

In this work, it is assumed that the defensive agents and the adversarial agents are fully cooperative among themselves, by which we mean that the agents within a team either maximize or minimize the reward function together. Therefore, the agents within a team also share the same Q-function. We do not make any assumptions on the transition probability function P, instead, it is considered as an unknown. Thus, the proposed approach obviates the need for a system model.

As mentioned earlier, the mission objective is specified as an STL formula and is denoted as  $\Phi$ . The control synthesis problem is then to find the optimal policies for the defensive agents such that the STL specification  $\Phi$  is enforced by maximizing the expected robustness degree against the best responses of the adversarial agents. Since a team stochastic game can be considered as a two-player zero-sum stochastic game, therefore, without any loss of generality, we assume that  $p{=}1$  and  $n{=}2$  and model the interaction between the agents as a two-player zero-sum stochastic game. The control objective for the case of two agents is now formally stated. Control objective: Given an STL specification  $\Phi$  with  $H(\Phi){=}T$ , a zero-sum stochastic game denoted by the tuple  $G{=}\langle S, A_1, A_2, P, R, \gamma \rangle$  with unknown P, and an initial

partial state trajectory  $\tilde{s}_{0:t'}$  for some  $0 \le t' < T$ , the control objective is to find a control policy  $\pi_1^*$  such that

$$\pi_1^* = \arg \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} E^{(\pi_1, \pi_2)} [\rho^{\Phi}(s_{0:T})], \tag{1}$$

where  $\Pi_1$  and  $\Pi_2$  denote the set of policies available to the maximizing agent and the minimizing agent, respectively.  $s_{0:T}$  denotes a state trajectory of the system during the interval [0,T] with  $s_{0:t'} = \tilde{s}_{0:t'}$ .  $E^{(\pi_1,\pi_2)}[\rho^{\Phi}(s_{0:T})]$  is the expected robustness degree of  $s_{0:T}$  with respect to  $\Phi$  under the policies  $\pi_1$  and  $\pi_2$ . In order to ensure that the minimum is achieved in (1), we assume that the policy of the adversarial agent,  $\pi_2$ , is deterministic [6].

## B. Reformulation of the Control Synthesis Problem

The objective function in the control synthesis problem posed in the previous subsection is not in the form of a sum of discounted rewards, and so RL methods are not directly applicable. Moreover, if the STL specification  $\Phi$  contains a nested operator such as  $\Phi = F_{[0,a]}G_{[0,b]}\phi$ , the control policy of the agent should take into account not only the current state of the system but also a sufficient length of the state history. To overcome these two issues, we adopt an approach similar to the one proposed in [4] to redefine the state space based on  $\Phi$ . We assume that the STL specification  $\Phi$  is of the form  $\Phi = G_{[0,T]}\phi$  or  $\Phi = F_{[0,T]}\phi$ , where  $\phi$  is any STL formula as defined in Section II. Let  $\tau$  be a non-negative integer defined as  $\tau = [H(\phi)/\Delta t] + 1$ , where  $[\cdot]$  is the ceiling function and  $\Delta t$  is the time step. The state space of the zero-sum stochastic game is redefined as  $S^{\tau}$  and an extended zero-sum stochastic game is formed.

Definition 3: (Extended zero-sum stochastic game) Given a zero-sum stochastic game denoted by the tuple  $G = \langle S, A_1, A_2, P, R, \gamma \rangle$  and a non-negative integer  $\tau$ , an extended zero-sum stochastic game is defined as the tuple  $G^{\tau} = \langle S^{\tau}, A_1, A_2, P^{\tau}, R^{\tau}, \gamma \rangle$ , where

- each state  $s^{\tau} \in S^{\tau}$  corresponds to a state trajectory with  $\tau$  samples; in the case that the length of the trajectory k is less than  $\tau$ , then  $s^{\tau}$  is formed by appending the empty string  $\epsilon$  to the trajectory  $\tau k$  times;
- $\begin{array}{lll} \ P^\tau \colon & S^\tau \times A_1 \times A_2 \times S^\tau {\longrightarrow} [0,1] & \text{is the transition} \\ & \text{probability function and is related to} \ P & \text{as} \\ & \text{follows: given} & s_i^\tau {=} s_a s_b \dots s_c & \text{and} & s_j^\tau {=} s_d \dots s_f s_g, \\ & P^\tau (s_i^\tau, a_1, a_2, s_j^\tau) {>} 0 & \text{if and only if} \ P(s_c, a_1, a_2, s_g) \in \\ & [0,1] & \text{and} \ s_d \dots s_f = s_b \dots s_c; \end{array}$
- $A_i$  is the finite set of actions for agent i; and
- $R^{\tau}$ :  $S^{\tau} \times A_1 \times A_2 \times S^{\tau} \to \mathbb{R}$  is the reward function.

The optimal policy of the extended stochastic game is given by the map  $\pi_1^*: S^{\tau} \to PD(A_1)$ . Given a state trajectory  $s_{0:T}$ , the corresponding trajectory in the extended state space can be written as  $s_{(\tau-1)\Delta t:T}^{\tau} = s_{\tau-1}^{\tau} \dots s_N^{\tau}$ , where  $T = N\Delta t$  and  $s_k^{\tau} = s_{(k-\tau+1)\Delta t:k\Delta t}$  for  $(\tau-1) \leqslant k \leqslant N$ . Using this relation and the quantitative semantics of the robustness degree provided in Section II, the robustness degree of  $s_{0:T}$  with respect to  $\Phi$  can be written as

$$\rho^{\Phi}(s_{0:T}) = f(\rho^{\Phi}(s_{\tau-1}^{\tau}), \dots, \rho^{\Phi}(s_{T}^{\tau})), \tag{2}$$

where the function  $f(\cdot)$  depends on the temporal operator used in the STL specification  $\Phi$ . If  $\Phi=F_{[0,T]}\phi$ , then  $f(\cdot)$  is the maximum function, and if  $\Phi=G_{[0,T]}\phi$ ,  $f(\cdot)$ 

is the minimum function. The maximum function can be approximated by the log-sum-exp function as

$$\max(z_1, \dots, z_n) \approx \frac{1}{\beta} \log \sum_{i=1}^n e^{\beta z_i}, \tag{3}$$

where  $\beta$  is a positive constant. Using equations (2) and (3), the objective function in (1) can be written as

$$E^{(\pi_1, \pi_2)}[\rho^{\Phi}(s_{0:T})] \approx E^{(\pi_1, \pi_2)} \left[ \frac{1}{\beta} \log \sum_{k=\tau-1}^{N} e^{\beta \rho^{\Phi}(s_k^{\tau})} \right], \quad (4)$$

for  $\Phi = F_{[0,T]}\phi$ . Similarly, for  $\Phi = G_{[0,T]}\phi$ , we can write

$$E^{(\pi_1, \pi_2)}[\rho^{\Phi}(s_{0:T})] \approx E^{(\pi_1, \pi_2)} \left[ -\frac{1}{\beta} \log \sum_{k=\tau-1}^{N} e^{-\beta \rho^{\Phi}(s_k^{\tau})} \right]. \tag{5}$$

The objective function in the control synthesis problem can now be interpreted as a sum of rewards, where the agents receive an immediate reward at every time instant  $t_k$  after both agents take their corresponding actions. Using the approximations given in (4) and (5), the control synthesis problem given in (1) is reformulated as follows.

**Approximated control synthesis problem**: Given an STL specification  $\Phi$  with  $H(\Phi) = T$ , where  $T = N \Delta t$  for some positive integer N, and a zero-sum stochastic game denoted by the tuple  $G = \langle S, A_1, A_2, P, R, \gamma \rangle$ , let  $G^{\tau} = \langle S^{\tau}, A_1, A_2, P^{\tau}, R^{\tau}, \gamma \rangle$  be the extended zero-sum stochastic game defined as in Definition 3. For a given initial partial state trajectory  $s^{\tau}_{\tau-1} = s_{0:(\tau-1)\Delta t}$  and  $\beta > 0$ , the control objective is to find a control policy  $\pi_1^*: S^{\tau} \to PD(A_1)$  such that

$$\pi_1^* = \arg\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2^D} E^{(\pi_1, \pi_2)} \left[ \sum_{k=\tau-1}^N r(s_k^{\tau}, \phi) \right], \quad (6)$$

$$\text{where} \quad r(s_k^{\tau},\phi) = \begin{cases} e^{\beta\rho^{\Phi}(s_k^{\tau})} & \text{if } \Phi = F_{[0,T]}\phi \\ -e^{-\beta\rho^{\Phi}(s_k^{\tau})} & \text{if } \Phi = G_{[0,T]}\phi \end{cases},$$

and  $\Pi_2^D$  denotes the set of deterministic policies of the adversarial agent.

The control synthesis problem is in the standard form for multi-agent reinforcement learning where the agents receive an immediate reward during each state transition. The optimal control policy  $\pi_1^*$  maximizes the expected robustness degree with respect to the STL specification  $\Phi$  against the best response of the adversarial agent.

# IV. LEARNING THE OPTIMAL POLICIES

# A. Learning Control Policies Using Minimax Q-Learning

Minimax Q-learning, which was first proposed in [5], is a model-free reinforcement learning method for learning optimal control policies in two-player zero-sum stochastic games. The learning objective in minimax Q-learning is to learn the control policy that ensures that the agents are in adversarial equilibrium provided the adversarial agent plays its optimal policy. In minimax Q-learning, the state-value function of the learning agent (agent 1) is given by

$$V_1(s) = \max_{\pi_1(s) \in \Pi_1} \min_{a_2 \in A_2} \sum_{a_1 \in A_1} Q_1(s, a_1, a_2) \pi_1^{a_1}(s),$$

where  $\pi_1^{a_1}(s)$  is the probability with which action  $a_1$  is chosen at state s under the policy  $\pi_1$ . The learning agent

learns its optimal policy by solving the stage game at each state with the Q-values  $Q_1(s,\cdot,\cdot)$  as the immediate payoffs. The minimax Q-learning algorithm for learning the optimal policy for the approximated control synthesis problem in (6) is presented in Algorithm 1.

In Algorithm 1,  $\alpha_j \in (0,1)$  is the learning rate,  $\gamma < 1$  is the discount factor from the definition of the zero-sum stochastic game,  $N_p$  denotes the number of episodes, and  $T = H(\Phi) = N\Delta t$ . The minimax Q-learning player always plays a safe strategy and therefore does not adapt to the opponent's change in strategies. The following theorem provides conditions for convergence to an adversarial equilibrium for the approximated control synthesis problem defined in (6) under Algorithm 1.

**Algorithm** 1: Minimax Q-learning algorithm

```
1: Initialize Q_1(s^{\tau}, a_1, a_2), V_1(s^{\tau}), \text{ and } \pi_1(s^{\tau})
 2: procedure
 3:
             for j = 1 : N_p do
 4:
                    Initialize the state s_{\tau-1}^{\tau}
                    for k = \tau - 1 : N do
 5:
                           Select an action a_1 based on policy \pi_1 with proba-
      bility \epsilon or a random action from A_1 with probability 1 - \epsilon
 7:
                           At the new state s_{k+1}^{\tau}, observe reward r_1 and
       opponent's action a_2
                          Update Q_1(s_k^{\tau}, a_1, a_2) using the following rule
      Q_1(s_k^{\tau}, a_1, a_2) \leftarrow (1 - \alpha_j) Q_1(s_k^{\tau}, a_1, a_2) + \alpha_j (r_1 + \gamma V_1(s_{k+1}^{\tau}))
                          Update the policy by solving a linear program
           \pi_{1}(s_{k}^{\tau}) \leftarrow \arg\max_{\pi_{1}(s_{k}^{\tau})} \min_{\bar{a}_{2} \in A_{2}} \sum_{\bar{a}_{1} \in A_{1}} Q_{1}(s_{k}^{\tau}, \bar{a}_{1}, \bar{a}_{2}) \pi_{1}^{\bar{a}_{1}}(s_{k}^{\tau})
V_{1}(s_{k}^{\tau}) = \min_{\bar{a}_{2} \in A_{2}} \sum_{\bar{a}_{1} \in A_{1}} Q_{1}(s_{k}^{\tau}, \bar{a}_{1}, \bar{a}_{2}) \pi_{1}^{\bar{a}_{1}}(s_{k}^{\tau})
10:
```

Theorem 1: Given an STL specification  $\Phi$  with  $H(\Phi) = T$  and a zero-sum stochastic game denoted by the tuple  $G = \langle S, A_1, A_2, P, R, \gamma \rangle$ , let  $G^{\tau} = \langle S^{\tau}, A_1, A_2, P^{\tau}, R^{\tau}, \gamma \rangle$  be the extended zero-sum stochastic game. Suppose that the optimal Q-function and the optimal policy are denoted by  $Q_1^*(s^{\tau}, a_1, a_2)$  and  $\pi_1^*(s^{\tau})$ , respectively, and that the following assumptions hold:

- i) Every state s and action  $a_i \in A_i$ , for i = 1, 2, are visited infinitely often and
- ii) The learning rate  $\alpha$  satisfies the conditions  $\sum_{j=0}^{\infty} \alpha_j = \infty$  and  $\sum_{j=0}^{\infty} (\alpha_j)^2 < \infty$ .

Then, Algorithm 1 with the update rules

$$\begin{split} Q_1^{j+1}(s_k^\tau, a_1, a_2) &= (1 - \alpha_j) Q_1^j(s_k^\tau, a_1, a_2) + \\ \alpha_j \left( r_1 + \gamma \max_{\pi_{1_j}(s_{k+1}^\tau)} \min_{\bar{a}_2 \in A_2} \sum_{\bar{a}_1 \in A_1} Q_1^j(s_{k+1}^\tau, \bar{a}_1, \bar{a}_2) \pi_{1_j}^{\bar{a}_1}(s_{k+1}^\tau) \right), \\ \pi_{1_{j+1}}(s_k^\tau) &= \arg\max_{\pi_{1_j}(s_k^\tau)} \min_{\bar{a}_2 \in A_2} \sum_{\bar{a}_1 \in A_1} Q_i^j(s_k^\tau, \bar{a}_1, \bar{a}_2) \pi_{1_j}^{\bar{a}_1}(s_k^\tau), \\ \\ \text{where } r_1 &= \begin{cases} e^{\beta \rho^\Phi(s_{k+1}^\tau)} & \text{if } \Phi = F_{[0,T]}\phi \\ -e^{-\beta \rho^\Phi(s_{k+1}^\tau)} & \text{if } \Phi = G_{[0,T]}\phi \end{cases}, \end{split}$$

for some  $\beta>0$ , ensures convergence of  $Q_1^j(\cdot,\cdot,\cdot)$  and  $\pi_{1_j}(\cdot)$  to  $Q_1^*(\cdot,\cdot,\cdot)$  and  $\pi_1^*(\cdot)$ , respectively, with probability 1 as  $j\to\infty$ .

*Proof.* The proof is a generalization of Theorem 3.1 in [22] and is omitted here due to paucity of space.

One of the major shortcomings of Algorithm 1, or Qlearning in general, is scalability as the size of the state space increases exponentially with  $\tau$  and the number of partitions of the continuous state space. The deep Q-learning algorithm proposed in [9] overcomes the issue of scalability by representing the Q-function as a deep neural network, where the network weights are learned during the training process. Inspired by deep Q-learning, we propose an algorithm called the minimax deep Q-learning algorithm, which is an extension of the minimax Q-learning algorithm presented in Algorithm 1, to learn the optimal policies. Similar to the input of the Q-function in [9], the input of the O-function in Algorithm 1 can be thought of as a series of image frames of length  $\tau+2$ , where the first  $\tau$  image frames correspond to  $s_{(\cdot)}^{\tau}$ , and the last two frames encode the actions of the agents. This will enable us to use convolutional neural networks to approximate the Q-function; such neural networks have been found to provide promising results for applications that involve recognizing patterns in data [8], [9].

```
Algorithm 2: Minimax deep Q-learning algorithm
1: Initialize the replay buffer D of length M; initialize the Q-
    network Q_1 with random weights \theta; initialize the target network
    Q_1 with weights \theta = \theta; initialize \pi_1(s^{\tau})
2: procedure
3:
         for k = 1 : N_p do
               Initialize the state s_{\tau-1}^{\tau}
4:
               for t = \tau - 1 : N do
5:
                     Select an action a_1^t based on policy \pi_1 with proba-
    bility \epsilon or a random action from A_1 with probability 1 - \epsilon
                     At new state s_{t+1}^{\tau}, observe reward r_t and opponent's
7:
    action a_2^t; store transition (s_t^{\tau}, a_1^t, a_2^t, s_{t+1}^{\tau}, r_t) in D
    Sample random minibatch of transitions (s_j^{\tau}, a_1^j, a_2^j, s_{j+1}^{\tau}, r_j) from D; set the target y_j as r_j if
    episode terminates at j+1, or else set y_j as
     y_j \leftarrow r_j + \max_{\pi_1(s_{j+1}^{\tau})} \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \hat{Q}_1(s_{j+1}^{\tau}, a_1, a_2; \hat{\theta}) \pi_1^{a_1}(s_{j+1}^{\tau})
                     Perform a gradient descent
     (y_j - Q_1(s_j^\tau, a_j^1, a_2^j; \theta))^2 \text{ with respect to } \theta 
 \pi_1(s_t^\tau) \leftarrow \arg\max_{\pi_1(s_t^\tau)} \min_{a_2 \in A_2} Q_1(s_t^\tau, a_1, a_2; \theta) \pi_1^{a_1}(s_t^\tau) 
                     s_t^{\tau} \leftarrow s_{t+1}^{\tau}; every C steps, reset \hat{Q}_1 = Q_1
```

It is known that using a nonlinear function approximator such as a neural network to represent the Q-function might result in instability during the learning process [23]. The deep Q-learning algorithm proposed in [9] makes use of two heuristics to overcome the issue of instability. In the proposed minimax deep O-learning algorithm presented in Algorithm 2, we take advantage of these heuristics in order to stabilize the learning process. The first heuristic involves storing the experiences of the learning agent at each time step in a replay buffer  $D = \{e_1, \dots, e_t\}$ , where  $e_t = (s_i^{\tau}, a_1, a_2, s_i^{\tau}, r_1)$ . During the Q-learning update, the algorithm randomly samples the experience from the replay buffer and uses it to update the parameters of the network. The second heuristic makes use of a separate target Q-network for generating the targets  $y_i$ during the Q-learning update. The target network stabilizes the learning process by ensuring that the Q-function to be learned is not continuously changing.

In this section, Algorithms 1 and 2 are applied to a simulation case study. The case study involves an agent moving in a  $4\times4$  spatial grid and tasked to reach a goal and stay within it for a stipulated number of time steps in the presence of an adversary. The agents are allowed to take a step in either the north (N), the east (E), the west (W), the south (S), the north-east (NE), the north-west (NW), the south-east (SE), or the south-west (SW) direction. The uncertainty in the dynamics is modeled as follows: if an agent selects an action a, then the agent executes the selected action with probability 0.95 or executes a randomly selected action from the set  $\{a_-, a_+\}$  with probability 0.05. If one visualizes the agents' possible actions to be rays emanating from the center of a circle,  $a_{-}$  and  $a_{+}$  are actions that are immediately to the left and the right of a, respectively. For instance, if the selected action is N, then  $a_{-}$  and  $a_{+}$  are the actions NW and NE, respectively. When an agent in state s takes an action that drives it outside the boundary of the grid, the agent remains in the same state s after executing the action. If the actions of the agents in state s are such that they move into the same grid, then both agents remain in safter the transition.

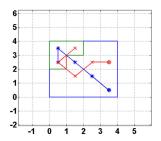
The objective of agent 1 is eventually to reach the region shown in Figure 1 with a green outline within 5 time steps and stay within that region for at least 2 consecutive time steps. The task objective is formally specified using the following STL specification:

 $\Phi = F_{[0,5]}G_{[0,1]}(\text{region } A \land \neg \text{region } B),$ (7) where region A represents  $(x_1 > 0 \land x_1 < 2 \land y_1 >$  $2 \wedge y_1 < 4$ ) and region B represents  $(x_1 > 1 \wedge x_1 <$  $2 \wedge y_1 > 2 \wedge y_1 < 3$ ). The specification results in  $\tau = 2$ with a state space of dimension  $|S^2| = 7056$ . We use Algorithms 1 and 2 presented in Section IV to learn the optimal policies that enforce the STL formula  $\Phi$  given in (7). During the implementation of Algorithm 2, a convolutional neural network (CNN) is used to represent the O-function. The input to the CNN consists of  $\tau + 2$  frames, where the first  $\tau$  frames represent the state of the system and the last two frames encode the actions taken by each of the agents. The frames that correspond to the states take values in the set  $\{-1,0,1\}$ , where a value of 1 at a cell corresponds to the presence of the agent and a value of -1 corresponds to the presence of the adversary. A value of 0 indicates that the cell is empty. The frames that encode the agents' actions take values in the set  $\{1, \dots, 8\}$ . The output of the network is the value of the Q-function corresponding to the state of the system and the actions executed by the two agents. The CNN used to represent the Q-function is composed of two convolution layers having 15 filters each and a stride of 1, with the second convolution layer followed by a hyperbolic tangent nonlinearity and two fully connected layers. The parameters used in the implementation of Algorithms 1 and 2 are M=1000, C=70,  $\gamma=0.9999$ ,  $\alpha=0.99\times10^{(\log(0.01)/N_p)}$ ,  $\epsilon$ =0.9,  $\beta$ =7, and  $N_p$ =50000. The size of the minibatch used in Algorithm 2 is 25.

In this case study, two types of adversaries were considered, namely the random adversary (RA) that plays a

	Algorithm 1		Algorithm 2	
	$\pi_A^*$	$\pi_B^*$	$\pi_A^*$	$\pi_B^*$
RA	0.812	0.846	0.938	0.919
OA	0.271	0.394	0.353	0.481

TABLE I:  $Pr[s_{[0:6]} \models \Phi]$  computed from 1000 simulations.



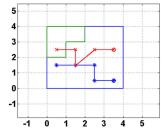


Fig. 1: Sample trajectories using the optimal strategy  $\pi_B^*$  learned using Algorithm 2; the trajectories of the defensive agent and the adversarial agent are shown in blue and red, respectively; the adversary is playing its optimal strategy; the figure on the left shows a simulation where the defensive agent successfully satisfies the specification, whereas in the figure on the right, the agent fails to satisfy the specification.

stationary random strategy, and the optimal adversary (OA) that plays the optimal strategy. Let  $\pi_A^*$  and  $\pi_B^*$  denote the optimal policies of the agent learned against RA and OA, respectively. The optimal strategy of the adversary is first learned against a deterministic hand-built strategy of the agent. The performance measure used to evaluate the learned optimal policies is the ratio of the number of trajectories that satisfy the specification  $\Phi$  to the total number of simulated trajectories, and is denoted by  $Pr[s_{[0:7]} \models \Phi]$ . The results are summarized in Table I. It is observed that the optimal policies learned using Algorithm 2 outperform the optimal policies learned using Algorithm 1. The reason for the improvement in performance when applying Algorithm 2 is the use of a convolution network to represent the Qfunction. In Algorithm 1, the Q-function is updated only for the state-action combinations that the agents visit. However, in Algorithm 2, the Q-network that represents the Q-function is updated at all state-action combinations and therefore provides a better representation of the O-values for the stateaction combinations that are not often visited. It is noted that in an ideal case when both the agent and the adversary are playing their optimal policies, the performance measure should be close to 0.5 since the agent and the adversary are in an adversarial equilibrium. However, in this case study, the optimal policy used by the adversary is learned against a hand-built policy of the agent and therefore might not have converged to the adversarial equilibrium, thereby explaining the performance values shown in the second row of Table I. Results from two representative simulations where the agent plays the optimal policy  $\pi_B^*$ , which is learned using Algorithm 2, against the optimal adversary are shown in Figure 1.

#### VI. CONCLUSIONS

This work addressed the problem of learning optimal control policies for a multi-agent system in an adversarial environment where the mission objective is expressed as an STL specification. The control synthesis problem was formulated as a zero-sum discounted stochastic game. It was shown that the control synthesis problem is very similar to the deep Q-learning problem where the inputs are in the form of a sequence of image frames. A multi-agent deep Q-learning algorithm was proposed to learn the optimal policies. The effectiveness of the proposed approach was illustrated using a simulation case study.

#### REFERENCES

- [1] J. van der Wal, "Stochastic dynamic programming," Ph.D. dissertation, Methematisch Centrum, 1980.
- [2] C. Belta, B. Yordanov, and E. A. Gol, "Formal methods for discretetime dynamical systems," 2017.
- [3] O. Maler and D. Nickovic, "Monitoring temporal properties of continuous signals," in FORMATS/FTRTFT, vol. 3253, 2004, pp. 152–166.
- [4] D. Aksaray, A. Jones, Z. Kong, M. Schwager, and C. Belta, "Q-Learning for robust satisfaction of Signal Temporal Logic specifications," 2016, pp. 6565–6570.
- [5] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Proceedings of the eleventh international* conference on machine learning, vol. 157, 1994, pp. 157–163.
- [6] —, "Value-function reinforcement learning in Markov games," Cognitive Systems Research, vol. 2, no. 1, pp. 55–66, 2001.
- [7] J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," *Journal of Machine Learning Research*, vol. 4, no. Nov, pp. 1039–1069, 2003.
- [8] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [9] 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.
- [10] A. Nikou, J. Tumova, and D. V. Dimarogonas, "Cooperative task planning of multi-agent systems under timed temporal specifications," in *Proceedings of the American Control Conference (ACC)*. IEEE, 2016, pp. 7104–7109.
- [11] Z. Liu, J. Dai, B. Wu, and H. Lin, "Communication-aware motion planning for multi-agent systems from signal temporal logic specifications," arXiv preprint arXiv:1705.11085, 2017.
- [12] S. Andersson, A. Nikou, and D. V. Dimarogonas, "Control synthesis for multi-agent systems under metric interval temporal logic specifications," arXiv preprint arXiv:1703.02780, 2017.
- [13] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," *PloS one*, vol. 12, no. 4, p. e0172395, 2017.
- [14] J. K. Gupta, M. Egorov, and M. Kochenderfer, "Cooperative multiagent control using deep reinforcement learning," in *Proceedings of the Adaptive and Learning Agents workshop (at AAMAS 2017)*, 2017.
- [15] G. Palmer, K. Tuyls, D. Bloembergen, and R. Savani, "Lenient multiagent deep reinforcement learning," arXiv preprint arXiv:1707.04402, 2017.
- [16] A. Pnueli, "The temporal logic of programs," in Foundations of Computer Science, 1977., 18th Annual Symposium on, 1977, pp. 46– 57.
- [17] A. Dokhanchi, B. Hoxha, and G. Fainekos, "On-line monitoring for temporal logic robustness," in *International Conference on Runtime* Verification. Springer, 2014, pp. 231–246.
- [18] J. Filar and K. Vrieze, Competitive Markov decision processes. Springer Science & Business Media, 2012.
- [19] M. L. Littman, "Friend-or-foe q-learning in general-sum games," in ICML, vol. 1, 2001, pp. 322–328.
- [20] A. Neyman and S. Sorin, *Stochastic games and applications*. Springer Science & Business Media, 2003, vol. 570.
- [21] M. G. Lagoudakis and R. Parr, "Learning in zero-sum team markov games using factored value functions," in *Advances in Neural Infor*mation Processing Systems, 2003, pp. 1659–1666.
- [22] C. Szepesvári and M. L. Littman, "Generalized markov decision processes: Dynamic-programming and reinforcement-learning algorithms," in *Proceedings of International Conference of Machine Learn*ing, vol. 96, 1996.
- [23] J. N. Tsitsiklis and B. Van Roy, "Analysis of temporal-diffference learning with function approximation," in Advances in neural information processing systems, 1997, pp. 1075–1081.