

Temporal Logic Motion Control using Actor-Critic Methods

– Technical Report – *

Xu Chu Ding[†], Jing Wang[‡], Morteza Lahijanian[‡], Ioannis Ch. Paschalidis[‡], and Calin A. Belta[‡]

Abstract—In this paper, we consider the problem of deploying a robot from a specification given as a temporal logic statement about some properties satisfied by the regions of a large, partitioned environment. We assume that the robot has noisy sensors and actuators and model its motion through the regions of the environment as a Markov Decision Process (MDP). The robot control problem becomes finding the control policy maximizing the probability of satisfying the temporal logic task on the MDP. For a large environment, obtaining transition probabilities for each state-action pair, as well as solving the necessary optimization problem for the optimal policy are usually not computationally feasible. To address these issues, we propose an approximate dynamic programming framework based on a least-square temporal difference learning method of the actor-critic type. This framework operates on sample paths of the robot and optimizes a randomized control policy with respect to a small set of parameters. The transition probabilities are obtained only when needed. Hardware-in-the-loop simulations confirm that convergence of the parameters translates to an approximately optimal policy.

Index Terms—Motion planning, Markov Decision Processes, dynamic programming, actor-critic methods.

I. INTRODUCTION

One major goal in robot motion planning and control is to specify a mission task in an expressive and high-level language and convert the task automatically to a control strategy for the robot. The robot is subject to mechanical constraints, actuation and measurement noise, and limited communication and sensing capabilities. The challenge in this area is the development of a computationally efficient framework accommodating both the robot constraints and the uncertainty of the environment, while allowing for a large spectrum of task specifications.

In recent years, temporal logics such as Linear Temporal Logic (LTL) and Computation Tree Logic (CTL) have been promoted as formal task specification languages for robotic applications [1]–[6]. They are appealing due to their high expressivity and closeness to human language. Moreover, several existing formal verification [7], [8] and synthesis [8] tools can be adapted to generate motion plans and provably correct control strategies for the robots.

* Research partially supported by the NSF under grant EFRI-0735974, by the DOE under grant DE-FG52-06NA27490, by the ODDR&E MURI10 program under grant N00014-10-1-0952, and by ONR MURI under grant N00014-09-1051.

[†] Xu Chu Ding is with Embedded Systems and Networks group, United Technologies Research Center, East Hartford, CT 06108 (dingx@utrc.utc.com).

Jing Wang, Morteza Lahijanian, Ioannis Ch. Paschalidis, and Calin A. Belta are with the Division of System Eng., Dept. of Mechanical Eng., Dept. of Electrical & Computer Eng., and Dept. of Mechanical Eng., Boston University, Boston, MA 02215 ({wangjing,morteza,yannisp,cbelta}@bu.edu), respectively.

In this paper, we assume that the robot model in the environment is described by a (finite) Markov Decision Process (MDP). In this model, the robot can precisely determine its current state, and by applying an action (corresponding to a motion primitive) enabled at each state, it triggers a transition to an adjacent state with a fixed probability. We are interested in controlling the MDP robot model such that it maximizes the probability of satisfying a temporal logic formula over a set of properties satisfied at the states of the MDP. By adapting existing probabilistic model checking [8]–[10] and synthesis [11], [12] algorithms, we recently developed such computational frameworks for formulas of LTL [13] and a fragment of probabilistic CTL [14].

With the above approaches, an optimal control policy can be generated to maximize the satisfying probability, given that the transition probabilities are known for each state-action pair of the MDP, which can be computed by using a Monte-Carlo method and repeated forward simulations. However, it is often not feasible for realistic robotic applications to obtain the transition probabilities for each state-action pair, even if an accurate model or a simulator of the robot in the environment is available. Moreover, the problem size is even larger when considering temporal logic specifications. For example, in order to find an optimal policy for an MDP satisfying an LTL formula, one need to solve a dynamical programming problem on the product between the original MDP and a Rabin automaton representing the formula. As such, exact solution can be computationally prohibitive for realistic settings.

In this paper, we show that approximate dynamic programming [15] can be effectively used to address the above limitations. For large dynamic programming problems, an approximately optimal solution can be provided using actor-critic algorithms [16]. In particular, actor-critic algorithms with Least Squares Temporal Difference (LSTD) learning have been shown recently to be a powerful tool to solve large-sized problems [17], [18]. This paper extends from [19], in which we proposed an actor-critic method for maximal reachability (MRP) problems, *i.e.*, maximizing the probability of reaching a set of states, to a computational framework that finds a control policy such that the probability of its paths satisfying an arbitrary LTL formula is locally optimal over a set of parameters. This set of parameters is designed to tailor to this class of approximate dynamical programming problems.

Our proposed algorithm produces a *randomized policy*, which gives a probability distribution over enabled actions at a state. Our method requires transition probabilities to be generated only along sample paths, and is therefore

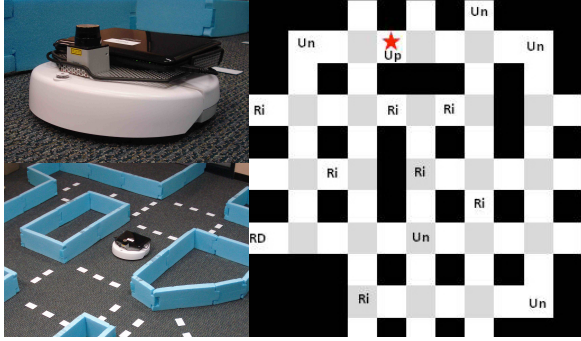


Fig. 1. Robotic InDoor Environment (RIDE) platform. Left: An iCreate mobile platform moving autonomously through the corridors and intersections of an indoor-like environment. Right: The partial schematics of the environment. The black blocks represent walls, and the grey and white regions are intersection and corridors, respectively. The labels inside a region represents observations associated with regions, such as **Un** (unsafe regions) and **Ri** (risky regions).

particularly suitable for robotic applications. To the best of our knowledge, this is the first of combining temporal logic formal synthesis with actor-critic type methods. We illustrate the algorithms with hardware-in-the-loop simulations using an accurate simulator of our Robotic InDoor Environment (RIDE) platform [20].

Notation: We use bold letters to denote sequences and vectors. Vectors are assumed to be column vectors. Transpose of a vector \mathbf{x} is denoted by \mathbf{x}^T . $\|\cdot\|$ stands for the Euclidean norm. $|S|$ denotes the cardinality of a set S .

II. PROBLEM FORMULATION AND APPROACH

We consider a robot moving in an environment partitioned into regions such as the Robotic Indoor Environment (RIDE) (see Fig. 1). Each region in the environment is associated with a set of observations. Observations can be **Un** for unsafe regions, or **Up** for a region where the robot can upload data. We assume that the robot can detect its current region. Moreover, the robot is programmed with a set of motion primitives allowing it to move from a region to an adjacent region. To capture noise in actuation and sensing, we make the natural assumption that, at a given region, a motion primitive designed to take the robot to a specific adjacent region may take the robot to a different adjacent region.

Such a robot model naturally leads to a labeled Markov Decision Process (MDP), which is defined below.

Definition II.1 (Labeled Markov Decision Process). A labeled Markov decision process (MDP) is a tuple $\mathcal{M} = (Q, q_0, U, A, P, \Pi, h)$, where

- (i) $Q = \{1, \dots, n\}$ is a finite set of states;
- (ii) $q_0 \in Q$ is the initial state;
- (iii) U is a finite set of actions;
- (iv) $A : Q \rightarrow U$ maps a state $q \in Q$ to actions enabled at q ;
- (v) $P : Q \times U \times Q \rightarrow [0, 1]$ is the transition probability function such that for all $q \in Q$, $\sum_{q' \in Q} P(q, u, q') = 1$ if $u \in A(q)$, and $P(q, u, q') = 0$ for all $q' \in Q$ if $u \notin A(q)$;

- (vi) Π is a set of observations;
- (vii) $h : Q \rightarrow 2^\Pi$ is the observation map.

Each state of the MDP \mathcal{M} modeling the robot in the environment corresponds to an ordered set of regions in the environment, while the actions label the motion primitives that can be applied at a region. For example, a state of \mathcal{M} may be labelled as I_1-C_1 , which means that the robot is currently at region C_1 , coming from region I_1 . Each ordered set of regions corresponds to a recent history of the robot trajectory, and is needed to ensure the Markov property (more details on such MDP abstraction of the robot in the environment can be found in e.g., [14]). The transition probability function P can be obtained through extensive simulations of the robot in the environment. We assume that there exists an accurate simulator that is capable of generating (computing) the transition probability $P(q, u, \cdot)$ for each state-action pair $q \in Q$ and $u \in A(q)$. More details of the construction of the MDP model for a robot in the RIDE platform are included in Sec. IV.

If the exact transition probabilities are not known, \mathcal{M} can be seen as a labeled non-deterministic transition system (NTS) $\mathcal{M}^\mathcal{N} = (Q, q_0, U, A, P^\mathcal{N}, \Pi, h)$, where P in \mathcal{M} is replaced by $P^\mathcal{N} : Q \times U \times Q \rightarrow \{0, 1\}$, and $P^\mathcal{N}(q, u, q') = 1$ indicates a possible transition from q to q' applying an enabled action $u \in A(q)$; if $P^\mathcal{N}(q, u, q') = 0$, then the transition from q to q' is not possible under u .

A path on \mathcal{M} is a sequence of states $\mathbf{q} = q_0q_1\dots$ such that for all $k \geq 0$, there exists $u_k \in A(q_k)$ such that $P(q_k, u_k, q_{k+1}) > 0$. Along a path $\mathbf{q} = q_0q_1\dots, q_k$ is said to be the state at time k . The trajectory of the robot in the environment is represented by a path \mathbf{q} on \mathcal{M} (which corresponds to a sequence of regions in the environment). A path $\mathbf{q} = q_1q_2\dots$ generates a sequence of observations $\mathbf{h}(\mathbf{q}) := o_1o_2\dots$, where $o_k = h(q_k)$ for all $k \geq 0$. We call $\mathbf{o} = \mathbf{h}(\mathbf{q})$ the word generated by \mathbf{q} .

Definition II.2 (Policy). A control policy for an MDP \mathcal{M} is an infinite sequence $M = \mu_0\mu_1\dots$, where $\mu_k : Q \times U \rightarrow [0, 1]$ is such that $\sum_{u \in A(q)} \mu_k(q, u) = 1$, for all $k \geq 0$.

Namely, at time k , $\mu_k(q, \cdot)$ is a discrete probability distribution over $A(q)$. If $\mu = \mu_k$ for all $k \geq 0$, then $M = \mu\mu\dots$ is called a *stationary* policy. If for all $k \geq 0$, $\mu_k(q, u) = 1$ for some u , then M is *deterministic*; otherwise, M is *randomized*. Given a policy M , we can then generate a set of paths on \mathcal{M} , by applying u_k with probability $\mu_k(q_k, u_k)$ at state q_k for all time k .

We require the trajectory of the robot in the environment to satisfy a rich task specification given as a Linear Temporal Logic (LTL) (see, e.g., [7], [8]) formula over a set of observations Π . An LTL formula over Π is evaluated over an (infinite) sequence $\mathbf{o} = o_0o_1\dots$ (e.g., a word generated by a path on \mathcal{M}), where $o_k \subseteq \Pi$ for all $k \geq 0$. We denote $\mathbf{o} \models \phi$ if word \mathbf{o} satisfies the LTL formula ϕ , and we say \mathbf{q} satisfies ϕ if $\mathbf{h}(\mathbf{q}) \models \phi$. Roughly, ϕ can be constructed from a set of observations Π , Boolean operators \neg (negation), \vee (disjunction), \wedge (conjunction), \longrightarrow (implication), and temporal operators X (next), U (until), F (eventually), G

(always). A variety of robotic tasks can be easily translated to LTL formulas. For example, the following complex task command in natural language: “*Gather data at locations **Da** infinitely often. Only reach a risky region **Ri** if valuable data at **VD** can be gathered, and always avoid unsafe regions (**Un**)*” can be translated to the LTL formula:

$$\phi := \text{G F Da} \wedge \text{G (Ri} \longrightarrow \text{VD)} \wedge \text{G } \neg \text{Un.}$$

In this paper, we consider the following problem.

Problem II.3. *Given a labeled MDP $\mathcal{M} = (Q, q_0, U, A, P, \Pi, h)$ modeling the motion of a robot in a partitioned environment and a mission task specified as an LTL formula ϕ over Π , find a control policy that maximizes the probability of its path satisfying ϕ .*

The probability that paths generated under a policy M satisfy an LTL formula ϕ is well defined with a suitable measure over the set of all paths generated by M [8].

In [13], we proposed a computational framework to solve Prob. II.3, by adapting methods from the area of probabilistic model checking [8]–[10]. However, this framework relies upon the fact that the transition probabilities are known for all state-action pairs. These transition probabilities are typically not available for robotic applications and computationally expensive to compute. Moreover, even if the transition probabilities are obtained for each state-action pair, this method still requires solving a linear program on the product of the MDP and the automata representing the formula, which can be very large (thousands or even millions of states). In this case an approximate method might be more desirable. For these reasons, we instead focus on the following problem.

Problem II.4. *Given a labeled NTS $\mathcal{M}^N = (Q, q_0, U, A, P^N, \Pi, h)$ modeling a robot in a partitioned environment, a mission task specified as an LTL formula ϕ over Π , and an accurate simulator to compute transition probabilities $P(q, u, \cdot)$ given a state-action pair (q, u) , find a control policy that approximately maximizes the probability of its path satisfying ϕ .*

In many robotic applications, the NTS model $\mathcal{M}^N = (Q, q_0, U, A, P^N, \Pi, h)$ can be quickly constructed for the robot in the environment. Our approach to Prob. II.4 can be summarized as follows: First, we proceed to translate the problem to a maximal reachability probability (MRP) problem using \mathcal{M}^N and ϕ (Sec. III-A). We then use an actor critic framework to find a randomized policy giving an approximate solution to the MRP problem (Sec. III-B). The randomized policy is constructed to be a function of a small set of parameters and we find a policy that is locally optimal with respect to these parameters. The construction of a class of policies suitable for MRP problems without using the transition probabilities is explained in Sec. III-C. The algorithmic framework presented in this paper is summarized in Sec. III-D.

III. CONTROL SYNTHESIS

A. Formulation of the MRP Problem

The formulation of the MRP problem is based on [8]–[10], [13] with modification if needed when using the NTS \mathcal{M}^N instead of \mathcal{M} . We start by converting the LTL formula ϕ over Π to a so-called deterministic *Rabin automaton*, which is defined as follows.

Definition III.1 (Deterministic Rabin Automaton). *A deterministic Rabin automaton (DRA) is a tuple $\mathcal{R} = (S, s_0, \Sigma, \delta, F)$, where*

- (i) S is a finite set of states;
- (ii) $s_0 \in S$ is the initial state;
- (iii) Σ is a set of inputs (alphabet);
- (iv) $\delta : S \times \Sigma \rightarrow S$ is the transition function;
- (v) $F = \{(L(1), K(1)), \dots, (L(M), K(M))\}$ is a set of pairs of sets of states such that $L(i), K(i) \subseteq S$ for all $i = 1, \dots, M$.

A run of a Rabin automaton \mathcal{R} , denoted by $\mathbf{r} = s_0 s_1 \dots$, is an infinite sequence of states in \mathcal{R} such that for each $k \geq 0$, $s_{k+1} \in \delta(s_k, \alpha)$ for some $\alpha \in \Sigma$. A run \mathbf{r} is *accepting* if there exists a pair $(L, K) \in F$ such that \mathbf{r} intersects with L finitely many times and K infinitely many times. For any LTL formula ϕ over Π , one can construct a DRA (for which we denote by \mathcal{R}_ϕ) with input alphabet $\Sigma = 2^\Pi$ accepting all and only words over Π that satisfy ϕ (see [21]).

We then obtain an MDP as the product of a labeled MDP \mathcal{M} and a DRA \mathcal{R}_ϕ , which captures all paths of \mathcal{M} satisfying ϕ . Note that this product MDP can only be constructed from an MDP and a deterministic automaton, this is why we require a DRA instead of, *e.g.*, a (generally non-deterministic) Büchi automaton (see [8]).

Definition III.2 (Product MDP). *The product MDP $\mathcal{M} \times \mathcal{R}_\phi$ between a labeled MDP $\mathcal{M} = (Q, q_0, U, A, P, \Pi, h)$ and a DRA $\mathcal{R}_\phi = (S, s_0, 2^\Pi, \delta, F)$ is an MDP $\mathcal{P} = (S_{\mathcal{P}}, s_{\mathcal{P}0}, U_{\mathcal{P}}, A_{\mathcal{P}}, P_{\mathcal{P}}, \Pi, h_{\mathcal{P}})$, where*

- (i) $S_{\mathcal{P}} = Q \times S$ is a set of states;
- (ii) $s_{\mathcal{P}0} = (q_0, s_0)$ is the initial state;
- (iii) $U_{\mathcal{P}} = U$ is a set of actions inherited from \mathcal{M} ;
- (iv) $A_{\mathcal{P}}$ is also inherited from \mathcal{M} and $A_{\mathcal{P}}((q, s)) := A(q)$;
- (v) $P_{\mathcal{P}}$ gives the transition probabilities:

$$P_{\mathcal{P}}((q, s), u, (q', s')) = \begin{cases} P(q, u, q') & \text{if } q' = \delta(s, h(q)) \\ 0 & \text{otherwise;} \end{cases}$$

Note that $h_{\mathcal{P}}$ is not used in the product MDP. Moreover, \mathcal{P} is associated with pairs of accepting states (similar to a DRA) $F_{\mathcal{P}} := \{(L_{\mathcal{P}}(1), K_{\mathcal{P}}(1)), \dots, (L_{\mathcal{P}}(M), K_{\mathcal{P}}(M))\}$ where $L_{\mathcal{P}}(i) = Q \times L(i)$, $K_{\mathcal{P}}(i) = Q \times K(i)$, for $i = 1, \dots, M$;

The product MDP is constructed in a ways such that, given a path $(s_0, q_0)(s_1, q_1) \dots$, the corresponding path $s_0 s_1 \dots$ on \mathcal{M} satisfies ϕ if and only if there exists a pair $(L_{\mathcal{P}}, K_{\mathcal{P}}) \in F_{\mathcal{P}}$ satisfying the Rabin acceptance condition, *i.e.*, the set $K_{\mathcal{P}}$ is visited infinitely often and the set $L_{\mathcal{P}}$ is visited finitely often.

We can make a very similar product between a labeled NTS $\mathcal{M}^N = (Q, q_0, U, A, P^N, \Pi, h)$ and \mathcal{R}_ϕ . This

product is also an NTS, which we denote by $\mathcal{P}^{\mathcal{N}} = (S_{\mathcal{P}}, s_{\mathcal{P}0}, U_{\mathcal{P}}, A_{\mathcal{P}}, P_{\mathcal{P}}^{\mathcal{N}}, \Pi, h_{\mathcal{P}}) := \mathcal{M}^{\mathcal{N}} \times \mathcal{R}_{\phi}$, associated with accepting sets $F_{\mathcal{P}}$. The definition (and the accepting condition) of $\mathcal{P}^{\mathcal{N}}$ is exactly the same as for the product MDP. The only difference between $\mathcal{P}^{\mathcal{N}}$ and \mathcal{P} is in $P_{\mathcal{P}}^{\mathcal{N}}$, which is either 0 or 1 for every state-action-state tuple.

From the product \mathcal{P} or equivalently $\mathcal{P}^{\mathcal{N}}$, we can proceed to construct the MRP problem. To do so, it is necessary to produce the so-called *accepting maximum end components* (AMECs). An end component is a subset of an MDP (consisting of a subset of states and a subset of enabled actions at each state) such that for each pair of states (i, j) in \mathcal{P} , there is a sequence of actions such that i can be reached from j with positive probability, and states outside the component cannot be reached. An AMEC of \mathcal{P} is the largest end component containing at least one state in $K_{\mathcal{P}}$ and no state in $L_{\mathcal{P}}$, for a pair $(K_{\mathcal{P}}, L_{\mathcal{P}}) \in F_{\mathcal{P}}$.

A procedure to obtain all AMECs of an MDP is outlined in [8]. This procedure is intended to be used for the product MDP \mathcal{P} , but it can be used without modification to find all AMECs associated with \mathcal{P} when $\mathcal{P}^{\mathcal{N}}$ is used instead of \mathcal{P} . This is because the information needed to construct the AMECs is the set of all possible state transitions at each state, and this information is already contained in $\mathcal{P}^{\mathcal{N}}$.

If we denote $S_{\mathcal{P}}^*$ as the union of all states in all AMECs associated with \mathcal{P} , it has been shown in probabilistic model checking (see *e.g.*, [8]) that the probability of satisfying the LTL formula is given by the maximal probability of reaching the set $S_{\mathcal{P}}^*$ from the initial state $s_{\mathcal{P}0}$. The desired optimal policy can then be obtained as the policy maximizing this probability. If transition probabilities are available for each state-action pair, then the solution to this MRP problem can be solved as by a linear program (see [8], [22]). The resultant optimal policy is deterministic (i.e., $M = \mu\mu\dots$) on the product MDP \mathcal{P} . To implement this policy on \mathcal{M} , it is necessary to use the DRA as a feedback automaton to keep track of the current state $s_{\mathcal{P}}$ on \mathcal{P} , and apply the action u where $\mu(s_{\mathcal{P}}, u) = 1$ (since μ is deterministic).

Remark III.3. *It is only necessary to find the optimal policy for states not in the set $S_{\mathcal{P}}^*$. This is because by construction, there exists a policy inside any AMEC that almost surely satisfies the LTL formula ϕ by reaching a state in $K_{\mathcal{P}}$ infinitely often. This policy can be obtained by simply choosing an action (among the subset of actions retained by the AMEC) at each state randomly, i.e., a trivial randomized stationary policy exists that almost surely satisfies ϕ .*

B. LSTD Actor-Critic Method

We now describe how relevant results in [19] can be applied to solve Prob. II.4. An approximate dynamic programming algorithm of the actor-critic type was presented in [19], which obtains a stationary randomized policy (RSP) (see Def. II.2) $M = \mu_{\theta}\mu_{\theta}\dots$, where $\mu_{\theta}(q, u)$ is a function of the state-action pair (q, u) and $\theta \in \mathbb{R}^n$, which is a vector of parameters. For the convenience of notations, we denote an RSP $\mu_{\theta}\mu_{\theta}\dots$ simply by μ_{θ} . In this sub-section we assume

that the RSP $\mu_{\theta}(q, u)$ to be given, and we will describe in Sec. III-C on how to design a suitable RSP.

Given an RSP μ_{θ} , actor-critic algorithms can be applied to optimize the parameter vector θ by policy gradient estimations. The basic idea is to use stochastic learning techniques to find θ that locally optimizes a cost function. In particular, the algorithm presented in [19] is targeted at Stochastic Shortest Path (SSP) problems commonly studied in literature (see *e.g.*, [22]). Given an MDP $\mathcal{M} = (Q, q_0, U, A, P, \Pi, h)$, a termination state $q^* \in Q$ and a function $g(q, u)$ defining the one-step cost of applying action u at state q , the *expected total cost* is defined as:

$$\bar{\alpha}(\theta) = \lim_{N \rightarrow \infty} E \left\{ \sum_{k=0}^{N-1} g(q_k, u_k) \right\}, \quad (1)$$

where (q_k, u_k) is the state-action pair at time k along a path under RSP μ_{θ} .

The SSP problem is formulated as the problem of finding θ^* minimizing (1). Note that, in general, we assume q^* to be cost-free and absorbing, i.e., $g(q^*, u) = 0$ and $P(q^*, u, q^*) = 1$ for all $u \in A(q^*)$. Under these conditions, the expected total cost (1) is finite.

We note that an MRP problem as described in Sec. III-A can be immediately converted to an SSP problem.

Definition III.4 (Conversion from MRP to SSP). *Given the product MDP $\mathcal{P} = (S_{\mathcal{P}}, s_{\mathcal{P}0}, U_{\mathcal{P}}, A_{\mathcal{P}}, P_{\mathcal{P}}, F_{\mathcal{P}})$ and a set of states $S_{\mathcal{P}}^* \subseteq S_{\mathcal{P}}$, the problem of maximizing the probability of reaching $S_{\mathcal{P}}^*$ can be converted to an SSP problem by defining a new MDP $\tilde{\mathcal{P}} = (\tilde{S}_{\mathcal{P}}, \tilde{s}_{\mathcal{P}0}, \tilde{U}_{\mathcal{P}}, \tilde{A}_{\mathcal{P}}, \tilde{P}_{\mathcal{P}}, g_{\mathcal{P}})$, where*

- (i) $\tilde{S}_{\mathcal{P}} = (S_{\mathcal{P}} \setminus S_{\mathcal{P}}^*) \cup \{s_{\mathcal{P}}^*\}$, where $s_{\mathcal{P}}^*$ is a “dummy” terminal state;
- (ii) $\tilde{s}_{\mathcal{P}0} = s_{\mathcal{P}0}$ (without the loss of generality, we exclude the trivial case where $s_{\mathcal{P}0} \in S_{\mathcal{P}}^*$);
- (iii) $\tilde{U}_{\mathcal{P}} = U_{\mathcal{P}}$;
- (iv) $\tilde{A}_{\mathcal{P}}(s_{\mathcal{P}}) = A_{\mathcal{P}}(s_{\mathcal{P}})$ for all $s_{\mathcal{P}} \in S_{\mathcal{P}}$, and for the dummy state we set $\tilde{A}_{\mathcal{P}}(s_{\mathcal{P}}^*) = \tilde{U}_{\mathcal{P}}$;
- (v) The transition probability is redefined as follows. We first define $\tilde{S}_{\mathcal{P}}^*$ as the set of states on $\tilde{\mathcal{P}}$ that cannot reach $S_{\mathcal{P}}^*$ under any policy. We then define:

$$\tilde{P}_{\mathcal{P}}(s_{\mathcal{P}}, u, s'_{\mathcal{P}}) = \begin{cases} \sum_{s''_{\mathcal{P}} \in S_{\mathcal{P}}^*} P_{\mathcal{P}}(s_{\mathcal{P}}, u, s''_{\mathcal{P}}), & \text{if } s'_{\mathcal{P}} = s_{\mathcal{P}}^* \\ P_{\mathcal{P}}(s_{\mathcal{P}}, u, s'_{\mathcal{P}}), & \text{if } s'_{\mathcal{P}} \in S_{\mathcal{P}} \setminus S_{\mathcal{P}}^* \end{cases}$$

for all $s_{\mathcal{P}} \in S_{\mathcal{P}} \setminus (S_{\mathcal{P}}^* \cup \tilde{S}_{\mathcal{P}}^*)$ and $u \in \tilde{U}_{\mathcal{P}}$. Moreover, for all $s_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}}^*$ and $u \in \tilde{U}_{\mathcal{P}}$, we set $\tilde{P}_{\mathcal{P}}(s_{\mathcal{P}}^*, u, s_{\mathcal{P}}^*) = 1$ and $\tilde{P}_{\mathcal{P}}(s_{\mathcal{P}}, u, s_{\mathcal{P}0}) = 1$;

- (vi) For all $s_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}}$ and $u \in \tilde{U}_{\mathcal{P}}$, we define the one-step cost $g_{\mathcal{P}}(s_{\mathcal{P}}, u) = 1$ if $s_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}}^*$, and $g_{\mathcal{P}}(s_{\mathcal{P}}, u) = 0$ otherwise.

We have shown in [19] that the policy minimizing (1) for the SSP problem with MDP $\tilde{\mathcal{P}}$ and the termination state $s_{\mathcal{P}}^*$ is a policy maximizing the probability of reaching the set $S_{\mathcal{P}}^*$ on \mathcal{P} , i.e., a solution to the MRP problem formulated in Sec. III-A.

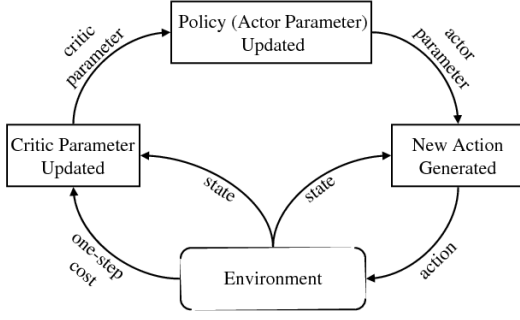


Fig. 2. Diagram illustrating an actor-critic algorithm.

The SSP problem can also be constructed from the NTS \mathcal{P}^N . In this case we obtain an NTS $\tilde{\mathcal{P}}^N(\tilde{S}_{\mathcal{P}}, \tilde{s}_{\mathcal{P}0}, \tilde{U}_{\mathcal{P}}, \tilde{A}_{\mathcal{P}}, \tilde{P}_{\mathcal{P}}^N, g_{\mathcal{P}})$, using the exact same construction as Def. III.4, except for the definition of $\tilde{P}_{\mathcal{P}}^N$. The transition function $\tilde{P}_{\mathcal{P}}^N(s_{\mathcal{P}}, u, s'_{\mathcal{P}})$ is instead defined as:

$$\tilde{P}_{\mathcal{P}}^N(s_{\mathcal{P}}, u, s'_{\mathcal{P}}) = \begin{cases} \max_{s''_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}}^*} P_{\mathcal{P}}^N(s_{\mathcal{P}}, u, s''_{\mathcal{P}}), & \text{if } s'_{\mathcal{P}} = s_{\mathcal{P}}^* \\ P_{\mathcal{P}}^N(s_{\mathcal{P}}, u, s'_{\mathcal{P}}), & \text{if } s'_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}} \setminus \tilde{S}_{\mathcal{P}}^* \end{cases}$$

for all $s_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}} \setminus (\tilde{S}_{\mathcal{P}}^* \cup \tilde{S}_{\mathcal{P}}^*)$ and $u \in \tilde{U}_{\mathcal{P}}$. Moreover, for all $s_{\mathcal{P}} \in \tilde{S}_{\mathcal{P}}^*$ and $u \in \tilde{U}_{\mathcal{P}}$, we set $\tilde{P}_{\mathcal{P}}^N(s_{\mathcal{P}}^*, u, s_{\mathcal{P}}^*) = 1$ and $\tilde{P}_{\mathcal{P}}^N(s_{\mathcal{P}}, u, s_{\mathcal{P}0}) = 1$.

Once the SSP problem is constructed, the algorithm presented in [19] is an iterative procedure that obtains a policy that locally minimizes the cost function (1) by simulating sample paths on $\tilde{\mathcal{P}}$. Each sample paths on $\tilde{\mathcal{P}}$ starts at $s_{\mathcal{P}0}$ and ends when the termination state $s_{\mathcal{P}}^*$ is reached. Since the probabilities is needed only along the sample path, we do not require the MDP $\tilde{\mathcal{P}}$, but only $\tilde{\mathcal{P}}^N$.

An actor-critic algorithm operates in the following way: the critic observes state and one-step cost from MDP and uses observed information to update the critic parameters, then the critic parameters are used to update the policy; the actor generates the action based on the policy and applies the action to the MDP. The algorithm stops when the gradient of $\bar{\alpha}(\theta)$ is small enough (*i.e.*, θ is locally optimal). The actor-critic update mechanism is shown in Fig. 2.

We summarize the actor-critic update algorithm in Alg. 1, and we note that it does not depend on the form of RSP μ_{θ} . The vectors $\mathbf{z}_k \in \mathbb{R}^n, \mathbf{b}_k \in \mathbb{R}^n, \mathbf{r}_k \in \mathbb{R}^n$ and the matrix $\mathbf{A}_k \in \mathbb{R}^{n \times n}$ are updated during each critic update, while simultaneously, the vector $\theta_k \in \mathbb{R}^n$ is updated during each actor update. Both the critic and actor update depend on

$$\psi_{\theta}(x, u) := \nabla_{\theta} \ln(\mu_{\theta}(x, u)), \quad (2)$$

which is the gradient of the logarithm of $\mu_{\theta}(x, u)$, to estimate the gradient $\nabla \bar{\alpha}(\theta)$. Lastly, sequence $\{\gamma_k\}$ controls the critic step-size, while $\{\beta_k\}$ and $\Gamma(\mathbf{r}_k)$ control the actor step-size. We note that all step-size parameters are positive, and their effect on the convergence rate is discussed in [19].

The critic update algorithm in Alg. 1 is of the LSTD type, which has shown to be superior to other approximate

dynamic programming methods in terms of the convergence rate [18]. More detail of this algorithm can be found in [19].

Algorithm 1 LSTD Actor-critic algorithm for SSP problems

Input: The NTS $\tilde{\mathcal{P}}^N(\tilde{S}_{\mathcal{P}}, \tilde{s}_{\mathcal{P}0}, \tilde{U}_{\mathcal{P}}, \tilde{A}_{\mathcal{P}}, \tilde{P}_{\mathcal{P}}^N, g_{\mathcal{P}})$ with the terminal state $s_{\mathcal{P}}^*$, the RSP μ_{θ} , and a computation tool to obtain $\tilde{P}_{\mathcal{P}}(s_{\mathcal{P}}, u, \cdot)$ for a given $(s_{\mathcal{P}}, u)$ state-action pair.

- 1: **Initialization:** Set all entries in $\mathbf{z}_0, \mathbf{A}_0, \mathbf{b}_0$ and \mathbf{r}_0 to zeros. Let θ_0 take some initial value. Set initial state $x_0 := \tilde{s}_{\mathcal{P}0}$. Obtain action u_0 using the RSP μ_{θ_0} .
- 2: **repeat**
- 3: Compute the transition probabilities $\tilde{P}(x_k, u_k, \cdot)$.
- 4: Obtain the simulated subsequent state x_{k+1} using the transition probabilities $\tilde{P}(x_k, u_k, \cdot)$. If $x_k = s_{\mathcal{P}}^*$, set $x_{k+1} := x_0$.
- 5: Obtain action u_{k+1} using the RSP μ_{θ_k} .
- 6: **Critic Update:**

$$\begin{aligned} \mathbf{z}_{k+1} &= \lambda \mathbf{z}_k + \psi_{\theta_k}(x_k, u_k) \\ \mathbf{b}_{k+1} &= \mathbf{b}_k + \gamma_k (g(x_k, u_k) \mathbf{z}_k - \mathbf{b}_k) \\ \mathbf{A}_{k+1} &= \mathbf{A}_k + \gamma_k (\mathbf{z}_k (\psi_{\theta_k}^T(x_{k+1}, u_{k+1}) - \psi_{\theta_k}^T(x_k, u_k)) \\ &\quad - \mathbf{A}_k), \\ \mathbf{r}_{k+1} &= -\mathbf{A}_k^{-1} \mathbf{b}_k. \end{aligned}$$

- 7: **Actor Update:**
 - 8: **until** $\|\nabla \bar{\alpha}(\theta_k)\| \leq \epsilon$ for some given ϵ
-

C. Designing an RSP

In this section we describe a randomized policy suitable to be used in Alg. 1 for MRP problems, and do not require the transition probabilities. We propose a family of RSPs that perform a “ t steps look-ahead”. This class of policies consider all possible sequences of actions in t steps and obtain a probability for each action sequence.

To simplify notation, for a pair of states $i, j \in \tilde{S}_{\mathcal{P}}$, we denote $i \xrightarrow{t} j$ if there is a positive probability of reaching j from i in t step. This can be quickly verified given $\tilde{P}_{\mathcal{P}}^N$ without transition probabilities. At state $i \in \tilde{S}_{\mathcal{P}}$, we denote an action sequence from i with t steps look-ahead as $e = u_1 u_2 \dots u_t$, where $u_k \in \tilde{A}_{\mathcal{P}}(j)$ for some j such that $i \xrightarrow{k} j$, for all $k = 1, \dots, t$. We denote the set of all action sequences from state i as $E(i)$. Given $e \in E(i)$, we denote $\tilde{P}_{\mathcal{P}}^N(i, e, j) = 1$ if there is a positive probability of reaching j from i with the action sequence e . This can also be recursively obtained given $\tilde{P}_{\mathcal{P}}^N(i, u, \cdot)$.

For each pair of states $i, j \in \tilde{S}_{\mathcal{P}}$, we define $d(i, j)$ as the minimum number of steps from i to reach j (this again can be obtained quickly from $\tilde{P}_{\mathcal{P}}^N$ without transition probabilities). We denote $j \in N(i)$ if and only if $d(i, j) \leq r_N$, where r_N is a fixed integer given apriori. If $j \in N(i)$, then we say i is in the neighborhood of j , and r_N represents the radius of the neighborhood around each state.

For each state $i \in \tilde{S}_{\mathcal{P}}$, We define the safety score $\text{safe}(i)$ as the ratio of the neighboring states not in $\tilde{S}_{\mathcal{P}}^*$ over all neighboring states of i . Recall that $\tilde{S}_{\mathcal{P}}^*$ is the set of states with 0 probability of reaching the goal states $S_{\mathcal{P}}^*$. To be

more specific, we define:

$$\text{safe}(i) := \frac{\sum_{j \in N(i)} I(j)}{|N(i)|}, \quad (3)$$

where $I(i)$ is an indicator function such that $I(i) = 1$ if and only if $i \in \tilde{S}_{\mathcal{P}} \setminus \bar{S}_{\mathcal{P}}^*$ and $I(i) = 0$ if otherwise. A higher safety score for the current state implies that it is less likely to reach $\bar{S}_{\mathcal{P}}^*$ in the near future. Furthermore, we define the progress score of a state $i \in \tilde{S}_{\mathcal{P}}$ as $\text{progress}(i) := \min_{j \in S_{\mathcal{P}}^*} d(i, j)$, which is the minimum number of transitions from i to any goal state.

We can now present the definition of our RSP. Let $\theta := [\theta_1, \theta_2]^T$. We define:

$$\begin{aligned} & a(\theta, i, e) \\ = & \exp\left(\theta_1 \sum_{j \in N(i)} \text{safe}(j) \tilde{P}_{\mathcal{P}}^N(i, e, j) \right. \\ & \left. + \theta_2 \sum_{j \in N(i)} (\text{progress}(j) - \text{progress}(i)) \right. \\ & \left. \tilde{P}_{\mathcal{P}}^N(i, e, j) \right), \end{aligned} \quad (4)$$

where \exp is the exponential function. Note that $a(\theta, i, e)$ is the combination of the expected safety score of the next state applying the action sequence e , and the expected improved progress score from the current state applying e , weighted by θ_1 and θ_2 . We assign the probability of pick the action sequence e at i proportional to the combined score $a(\theta, i, e)$. Hence, the probability to pick action sequence e at state i is defined as:

$$\tilde{\mu}_{\theta}(i, e) = \frac{a(\theta, i, e)}{\sum_{e \in E(i)} a(\theta, i, e)}. \quad (5)$$

Note that, if the action sequence $e = u_1 u_2 \dots u_t$ is picked, only the first action u_1 is applied. Hence, at stat i , the probability that an action $u \in \tilde{A}_{\mathcal{P}}(i)$ can be derived from Eq. (5):

$$\mu_{\theta}(i, u) = \sum_{\{e \in E(i) \mid e = uu_2 \dots u_t\}} \tilde{\mu}_{\theta}(i, e), \quad (6)$$

which completes the definition of the RSP.

D. Overall Algorithm

We now connect all the pieces together and present the overall algorithm giving a solution to Prob. II.4.

Proposition III.5. *Alg. 2 returns in finite time with θ^* locally maximizing the probability of the RSP μ_{θ} satisfying the LTL formula ϕ .*

Proof. In [19], we have shown that the actor-critic algorithm used in this paper returns in finite time with a locally optimal θ^* such that $\|\nabla \bar{\alpha}(\theta^*)\| \leq \epsilon$ for a given ϵ . We have shown throughout the paper that the optimal policy maximizing the probability of reaching $S_{\mathcal{P}}^*$ on \mathcal{P} is a policy maximizing the probability of satisfying ϕ . We also showed throughout the paper that the SSP problem, as well as the RSP μ_{θ} can be constructed without the transition probabilities, and only with \mathcal{M}^N . Therefore, Alg. 2 produces an RSP maximizing

Algorithm 2 Overall algorithm providing a solution to Prob. II.4

Input: A labeled NTS $\mathcal{M}^N = (Q, q_0, U, A, P^N, \Pi, h)$ modeling a robot in a partitioned environment, LTL formula ϕ over Π , and a simulator to compute $P(q, u, \cdot)$ given a state-action pair (q, u)

- 1: Translate the LTL formula ϕ to a DRA \mathcal{R}_{ϕ}
- 2: Generate the product NTS $\mathcal{P}^N = \mathcal{M}^N \times \mathcal{R}_{\phi}$
- 3: Find the union of all AMECs $S_{\mathcal{P}}^*$ associated with \mathcal{P}^N
- 4: Convert from an MRP to an SSP and generate $\tilde{\mathcal{P}}^N$
- 5: Obtained the RSP μ_{θ} with \mathcal{P}^N
- 6: Execute Alg. 1 with $\tilde{\mathcal{P}}^N$ and μ_{θ} as inputs until $\|\nabla \bar{\alpha}(\theta^*)\| \leq \epsilon$ for a θ^* and a given ϵ

Output: RSP μ_{θ} and θ^* locally maximizing the probability of satisfying ϕ with respect to θ up to a threshold ϵ

the probability of satisfying ϕ with respect to θ up to a threshold ϵ . ■

IV. HARDWARE-IN-THE-LOOP SIMULATION

We test the algorithms proposed in this paper through hardware-in-the-loop simulation for the RIDE environment (as shown in Fig. 1). The transition probabilities are computed by an accurate simulator of RIDE as needed. We apply both LTL control synthesis methods of linear programming (exact solution) and actor-critic (approximate solution) and compare the results.

A. Environment

In this case study, we consider an environment whose topology is shown in Fig. 3. This environment is made of square blocks forming 164 corridors and 84 intersections. The corridors (C_1, C_2, \dots, C_{164}) shown as white regions in Fig. 3 are of three different lengths, one-, two-, and three-unit lengths. The three-unit corridors are used to build corners in the environment. The intersections (I_1, I_2, \dots, I_{84}) are of two types, three-way and four-way, and are shown as grey blocks in Fig. 3. The black regions in this figure represent the walls of the environment. Note that there is always a corridor between two intersections.

There are five properties of interest (observations) associated with the regions of the environment. These properties are: **VD** = *ValuableData* (regions containing valuable data to be collected), **RD** = *RegularData* (regions containing regular data to be collected), **Up** = *Upload* (regions where data can be uploaded), **Ri** = *Risky* (regions that could pose a threat to the robot), and **Un** = *Unsafe* (regions that are unsafe for the robot).

B. Construction of the MDP model

The robot is equipped with a set of feedback control primitives (actions) - FollowRoad, GoRight, GoLeft, and GoStraight. The controller FollowRoad is only available (enabled) at the corridors. At four-way intersections, controllers are GoRight, GoLeft, and GoStraight. At three-way intersections, depending on the shape of the intersection, two of the four controllers are available. Due to the presence of noise in the actuators and sensors, however, the resulting

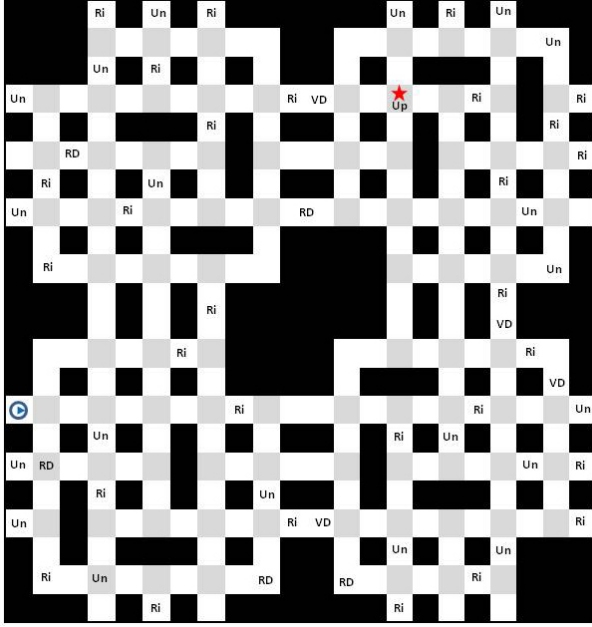


Fig. 3. Schematic representation of the environment with 84 intersections and 164 corridors. The black blocks represent walls, and the grey and white regions are intersection and corridors, respectively. There are five properties of interest in the regions indicated with **VD** = *ValuableData*, **RD** = *RegularData*, **Up** = *Upload*, **Ri** = *Risky*, and **Un** = *Unsafe*. The initial position of the robot is shown with a blue disk and the upload region is indicated with a red star.

motion may be different than intended. Thus, the outcome of each control primitive is characterized probabilistically.

To create an MDP model of the robot in RIDE, we define each state of the MDP as a collection of two adjacent regions (a corridor and an intersection). For instance the pairs C_1-I_2 and I_3-C_4 are two states of the MDP. Through this pairing of regions, it was shown that the Markov property (*i.e.*, the result of an action at a state depends only on the current state) can be achieved [14]. The resulting MDP has 608 states.

The set of actions available at a state is the set of controllers available at the last region corresponding to the state. For example, when in state C_1-I_2 only those actions from region I_2 are allowed. Each state of the MDP whose second region satisfies an observation in Π is mapped to that observation.

To obtain transition probabilities, we use an accurate simulator (see Fig. 4) incorporating the motion and sensing of an iRobot Create platform with a Hokuyo URG-04LX laser range finder, APSX RW-210 RFID reader, and an MSI Wind U100-420US netbook (the robot is shown in Fig. 1) in RIDE. Specifically, it emulates experimentally measured response times, sensing and control errors, and noise levels and distributions in the laser scanner readings. More detail for the software implementation of the simulator can be found in [14]. We perform a total of 1,000 simulations for each action available in each MDP state.

C. Task specification and results

We consider the following mission task:

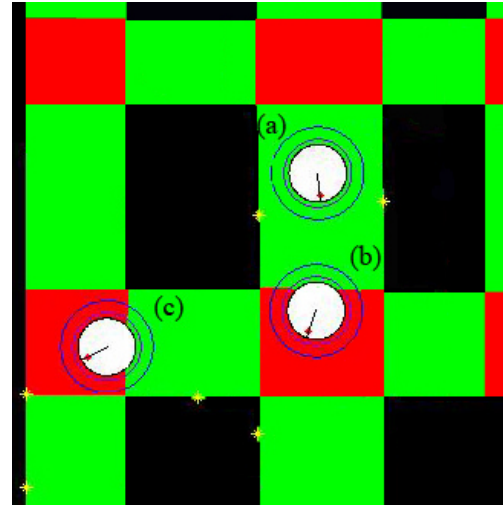


Fig. 4. Simulation snapshots. The white disk represents the robot and the different circles around it indicate different "zones" in which different controllers are activated. The yellow dots represent the laser readings used to define the target angle. (a) The robot centers itself on a stretch of corridor by using *FollowRoad*; (b) The robot applies *GoRight* in an intersection; (c) The robot applies *GoLeft*.

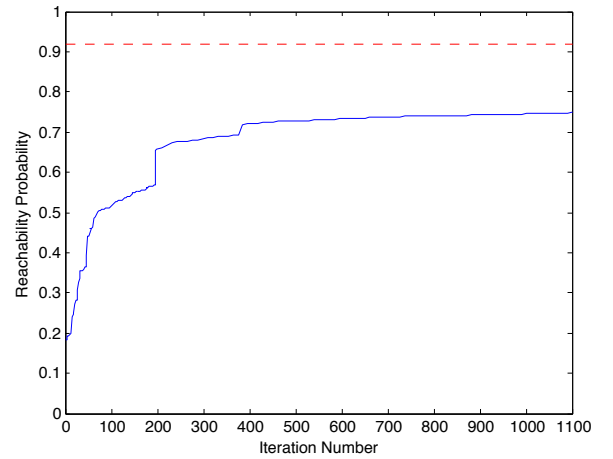


Fig. 5. The optimal solution (the maximal probability of satisfying the specification) is shown with the dashed line, and the solid line represents the exact reachability probability for the RSP as a function of the number of iterations applying the proposed algorithm.

Specification: Reach a location with *ValuableData* (**VD**) or *RegularData* (**RD**), and then reach *Upload* (**Up**). Do not reach *Risky* (**Ri**) regions unless eventually reach a location with *ValuableData* (**VD**). Always avoid *Unsafe* (**Un**) regions until *Upload* (**Up**) is reached (and mission completed).

The above task specification can be translated to the LTL formula:

$$\phi := \mathbf{FUp} \wedge (\neg \mathbf{Un} \mathbf{U Up}) \wedge \mathbf{G}(\mathbf{Ri} \longrightarrow \mathbf{FVD}) \wedge \mathbf{G}(\mathbf{VD} \vee \mathbf{RD} \longrightarrow \mathbf{X FUp}) \quad (7)$$

The initial position of the robot is shown as a blue circle in Fig. 3 with the orientation towards the neighboring intersection. We used the computational frameworks described

in this paper to find the control strategy maximizing the probabilities of satisfying the specification. The size of the DRA is 17 which results in the product MDP with 10336 states. By applying both methods of linear programming (exact solution) and actor-critic (approximate solution), we found the maximum probabilities of satisfying the specification were 92% and 75%, respectively. The graph of the convergence of the actor-critic solution is shown in Fig. 5. The parameters for this examples are: $\lambda = 0.9$, and the initial $\theta = [5, -0.5]^T$. The look-ahead window t for the RSP is 2.

It should be emphasized that, we only compute the transition probabilities along the sample path. Thus, when Alg. 2 is completed (at iteration 1100), at most 1100 transition probabilities of state-action pairs were computed. In comparison, in order to solve the probability exactly, around 30000 transition probabilities of state-action pairs must be computed.

V. CONCLUSIONS

We presented a framework that brings together an approximate dynamic programming computational method of the actor critic type, with formal control synthesis for Markov Decision Processes (MDPs) from temporal logic specifications. We show that this approach is particular suitable for problems where the transition probabilities of the MDP are difficult or computationally expensive to compute, such as for many robotic applications. We show that this approach effectively finds an approximate optimal policy within a class of randomized stationary policies maximizing the probability of satisfying the temporal logic formula. Future direction includes extending this result to multi-robot teams, examining exactly how to choose an appropriate look-ahead window when designing the RSP, and applying the result to more realistic problem settings with the MDP containing possibility millions of states.

REFERENCES

- [1] H. Kress-Gazit, G. Fainekos, and G. J. Pappas, "Where's Waldo? Sensor-based temporal logic motion planning," in *IEEE Int. Conf. on Robotics and Automation*, Rome, Italy, 2007, pp. 3116–3121.
- [2] S. Karaman and E. Frazzoli, "Sampling-based motion planning with deterministic μ -calculus specifications," in *IEEE Conf. on Decision and Control*, Shanghai, China, 2009, pp. 2222 – 2229.
- [3] S. G. Loizou and K. J. Kyriakopoulos, "Automatic synthesis of multiagent motion tasks based on LTL specifications," in *43rd IEEE Conference on Decision and Control*, December 2004.
- [4] M. M. Quottrup, T. Bak, and R. Izadi-Zamanabadi, "Multi-robot motion planning: A timed automata approach," in *IEEE Int. Conf. on Robotics and Automation*, New Orleans, LA, Apr. 2004, pp. 4417–4422.
- [5] T. Wongpiromsarn, U. Topcu, and R. M. Murray, "Receding horizon temporal logic planning for dynamical systems," in *IEEE Conference on Decision and Control*, Shanghai, China, 2009.
- [6] A. Bhatia, L. Kavraki, and M. Vardi, "Sampling-based motion planning with temporal goals," in *Robotics and Automation (ICRA), 2010 IEEE International Conference on*. IEEE, 2010, pp. 2689–2696.
- [7] E. M. M. Clarke, D. Peled, and O. Grumberg, *Model checking*. MIT Press, 1999.
- [8] C. Baier, J.-P. Katoen, and K. G. Larsen, *Principles of Model Checking*. MIT Press, 2008.
- [9] L. De Alfaro, "Formal verification of probabilistic systems," Ph.D. dissertation, Stanford University, 1997.
- [10] M. Vardi, "Probabilistic linear-time model checking: An overview of the automata-theoretic approach," *Formal Methods for Real-Time and Probabilistic Systems*, pp. 265–276, 1999.
- [11] C. Courcoubetis and M. Yannakakis, "Markov decision processes and regular events," *IEEE Transactions on Automatic Control*, vol. 43, no. 10, pp. 1399–1418, 1998.
- [12] C. Baier, M. Größer, M. Leucker, B. Bollig, and F. Ciesinski, "Controller synthesis for probabilistic systems," in *In Proceedings of IFIP TCS'2004*. Citeseer, 2004.
- [13] X. C. Ding, S. L. Smith, C. Belta, and D. Rus, "LTL planning for groups of robots control in uncertain environments with probabilistic satisfaction guarantees," in *18th IFAC World Congress*, 2011.
- [14] M. Lahijanian, J. Wasniewski, S. B. Andersson, and C. Belta, "Motion planning and control from temporal logic specifications with probabilistic satisfaction guarantees," in *IEEE Int. Conf. on Robotics and Automation*, Anchorage, AK, 2010, pp. 3227 – 3232.
- [15] J. Si, *Handbook of learning and approximate dynamic programming*. Wiley-IEEE Press, 2004, vol. 2.
- [16] A. Barto, R. Sutton, and C. Anderson, "Neuronlike adaptive elements that can solve difficult learning control problems," *IEEE Transactions on Systems, Man, & Cybernetics*, 1983.
- [17] I. Paschalidis, K. Li, and R. Estanjini, "An actor-critic method using least squares temporal difference learning," in *Decision and Control, 2009 held jointly with the 2009 28th Chinese Control Conference. CDC/CCC 2009. Proceedings of the 48th IEEE Conference on*. IEEE, pp. 2564–2569.
- [18] V. Konda and J. Tsitsiklis, "On actor-critic algorithms," *SIAM Journal on Control and Optimization*, vol. 42, no. 4, pp. 1143–1166, 2004.
- [19] R. Estanjini, X. Ding, M. Lahijanian, J. Wang, C. Belta, and I. Paschalidis, "Least squares temporal difference actor-critic methods with applications to robot motion control," in *IEEE Conference on Decision and Control (CDC), Orlando, FL, December 2011*.
- [20] "Robotic indoor environment." [Online]. Available: www.hyness.bu.edu/ride
- [21] E. Gradel, W. Thomas, and T. Wilke, *Automata, logics, and infinite games: A guide to current research*, ser. Lecture Notes in Computer Science. Springer, 2002, vol. 2500.
- [22] M. Puterman, *Markov decision processes: Discrete stochastic dynamic programming*. John Wiley & Sons, Inc. New York, NY, USA, 1994.