

Robust Control of Uncertain Markov Decision Processes with Temporal Logic Specifications

Eric M. Wolff, Ufuk Topcu, and Richard M. Murray

Abstract—We present a method for designing robust controllers for dynamical systems with linear temporal logic specifications. We abstract the original system by a finite Markov Decision Process (MDP) that has transition probabilities in a specified uncertainty set. A robust control policy for the MDP is generated that maximizes the worst-case probability of satisfying the specification over all transition probabilities in the uncertainty set. To do this, we use a procedure from probabilistic model checking to combine the system model with an automaton representing the specification. This new MDP is then transformed into an equivalent form that satisfies assumptions for stochastic shortest path dynamic programming. A robust version of dynamic programming allows us to solve for a ϵ -suboptimal robust control policy with time complexity $O(\log 1/\epsilon)$ times that for the non-robust case. We then implement this control policy on the original dynamical system.

I. INTRODUCTION

As the level of autonomous operation expected of robots, vehicles, and other cyberphysical systems increases, there is a growing need for formal methods that allow desired system properties to be precisely specified and allow performance to be automatically verified. As autonomous systems often operate over long time periods in uncertain environments, it is also important that system performance is robust to both environmental disturbances and modeling errors.

A promising approach for specifying and verifying system properties uses temporal logics such as linear temporal logic (LTL) to specify tasks. LTL provides a natural framework to specify desired properties such as response (if A, then B), liveness (always eventually A), safety (always not B), stability (eventually always A), and priority (first A, then B, then C). This framework requires a finite-state abstraction for continuous dynamical system, which is possible, in the simulation sense, for a wide range of systems [1].

We use a Markov Decision Process (MDP) as a finite-state abstraction of our original dynamical system because it provides a general framework for modeling the non-determinism and probabilistic behavior that is present in many real-world systems. MDPs are also amenable to formal verification techniques for temporal logic properties [4]. These techniques generate a control policy for the MDP that maximizes the probability of satisfying a given LTL specification. However, these techniques assume that the state transition probabilities of the MDP are known exactly, which is often unrealistic. We relax this assumption by allowing the transition probabilities

of our MDP abstraction to lie in uncertainty sets. We then generate a control policy that maximizes the worst-case probability of satisfying a given LTL specification over all transition probabilities in the uncertainty set.

We argue that it is prudent to consider uncertainty in the MDP model. Real systems have unmodeled dynamics and parametric uncertainty that are not fully captured by mathematical models. So, even if exact finite-state abstraction techniques are available for a dynamical system model, the resulting MDP abstraction will only represent the real system to the extent that the dynamical system model does. Furthermore, when approximate abstraction techniques are used for a dynamical system model, the MDP abstraction will be a further approximation of the real system.

There is a large amount of work in addressing the robustness of control policies for MDPs with uncertain transition probabilities. Our approach most closely follows that of Nilim and El Ghaoui [22], which addresses a wide variety of uncertainty models. Other work includes [3], [23], [11].

On the temporal logic side, formal verification is well developed for discrete systems represented as MDPs with exact transition matrices [9], [4] and there exist powerful software tools [16]. While there has been recent work in verification of uncertain MDPs, these works almost exclusively use simple interval models of transition probability uncertainty [27], [7], [6], which our work includes as a special case.

The above references are primarily concerned with discrete models. Controllers that are guaranteed to satisfy a given temporal logic specification can be generated for a variety of non-stochastic hybrid systems [15], [14], [26]. Similar methods for non-stochastic systems have been implemented on stochastic systems [12], [17]. Recent work that explicitly designs controllers for stochastic systems [10], [18] does not consider robustness. Robustness of non-probabilistic discrete transition systems to disturbances is explored in [20].

There are two main contributions in our work. First of all, we create a robust control policy π that maximizes the worst-case probability of satisfying an LTL specification φ for a system represented as a finite labeled MDP \mathcal{M} with transition matrices in an uncertainty set \mathcal{P} . A control policy π is a mapping from each MDP state to an allowable action. The set \mathcal{P} can be non-convex and it includes interval uncertainty sets as a special case. This freedom allows more statistically accurate and less conservative results than interval uncertainty sets. Our second contribution is a method for generating an approximate finite-state MDP abstraction from a possibly nonlinear dynamical system given in state-space form.

Authors are with the Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA, USA. The corresponding author is ewolff@caltech.edu

Due to space constraints, all omitted proofs are included in a technical report [24].

II. PRELIMINARIES

We now present formal definitions for Markov Decision Processes (MDPs), which we use as finite-state abstractions for the original continuous dynamics, and linear temporal logic (LTL), which is the task specification language. Throughout, we will assume that equality and inequality is component-wise for vectors and matrices. Also, $\mathbf{1}$ denotes a vector of ones of appropriate dimension.

A. Markov Decision Processes

Definition 1 (Labeled finite MDP). We specify a *labeled finite MDP*, \mathcal{M} , by the tuple $\mathcal{M} = (S, A, P, s_0, AP, L)$, where S is a finite set of states, A is a finite set of actions, $P : S \times A \times S \rightarrow [0, 1]$ is the transition probability function, s_0 is the initial state, AP is a finite set of atomic propositions, and $L : S \rightarrow 2^{AP}$ is a labeling function. Let $A(s)$ denote the set of available actions at state s . $\sum_{s' \in S} P(s, a, s') = 1$ if $a \in A(s)$ and $P(s, a, s') = 0$ otherwise.

We assume, for notational convenience, that the available actions $A(s)$ are the same for every $s \in S$. We use P_{ij}^a as shorthand for the transition probability from state i to state j when using action a . We call $P^a \in \mathbb{R}^{n \times n}$ a *transition matrix*, where the (i, j) -th entry of P^a is P_{ij}^a . Where it is clear from context, we refer to the row vector P_i^a as p .

To model uncertainty in state transitions, we specify uncertainty sets for the transition matrices of the MDP. We assume that an uncertainty set \mathcal{P}^a corresponding to action $a \in A$ is a set of probability matrices.

Assumption 1. \mathcal{P}^a can be factored as the Cartesian product of its rows, so its rows are uncorrelated. Formally, for every $a \in A$, $\mathcal{P}^a = \mathcal{P}_1^a \times \dots \times \mathcal{P}_n^a$ where each \mathcal{P}_i^a is a subset of the probability simplex in \mathbb{R}^n . We follow Nilim and El Ghaoui and refer to these as *rectangular* uncertainty sets [22].

Let F^a denote the nominal transition matrix for action a .

Assumption 2. $F_{ij}^a = 0$ if and only if $P_{ij}^a = 0$ for all $P^a \in \mathcal{P}^a$.

Assumption 2 implies that if a nominal transition is zero (non-zero), then it is zero (non-zero) for all transition matrices in the uncertainty set. This must be enforced as adding or removing a transition to an unsafe set may result in a different satisfaction probability for properties over infinite time.

Definition 2. A *control policy* for an MDP \mathcal{M} is a sequence $\pi = \{\mu_0, \mu_1, \dots\}$, where each $\mu_k : S \rightarrow A$ such that $\mu_k(s) \in A(s)$. We call a policy *stationary* if it is of the form $\pi = \{\mu, \mu, \dots\}$. We let Π be the set of all control policies and Π_s be the set of all stationary control policies.

Definition 3. An *environment policy* for an MDP \mathcal{M} is a sequence of transition matrices $\tau = (P_k^a)_{a \in A, k=1,2,\dots}$. We call an environment policy *stationary* if it is of the form $\tau = (P^a)_{a \in A}$. We let \mathcal{T} be the set of all environment policies and \mathcal{T}_s be the set of all stationary environment policies.

We now associate a cost with each state in \mathcal{M} through a function $c(s, a) : S \times A \rightarrow \mathbb{R}$. This cost is incurred at every stage k and the control policy tries to minimize the total expected cost over the horizon of length N . The total expected cost of an infinite-horizon problem starting from state s when using control policy π under environment policy τ is

$$J^{\pi\tau}(s) := \lim_{N \rightarrow \infty} \mathbb{E}_{\pi\tau} \left[\sum_{k=0}^{N-1} c(s_k, \mu_k(s_k)) \middle| s_0 = s \right], \quad (1)$$

where the expectation $\mathbb{E}_{\pi\tau}$ depends on both the control and environment policies.

The optimal worst-case total expected cost starting from state s is

$$J^*(s) := \min_{\pi \in \Pi} \max_{\tau \in \mathcal{T}} J^{\pi\tau}(s). \quad (2)$$

The optimal worst-case control and environment policies are denoted π^* and τ^* . We call π^* a *robust control policy*.

B. Linear Temporal Logic

We use linear temporal logic (LTL) to concisely and unambiguously specify the desired system behavior. We will only touch on the key aspects of LTL for our problem; we defer the reader to [4] for details.

We describe an *execution* of a discrete-time system by an infinite sequence of its states, $\sigma = s_0 s_1 s_2 \dots$ where $s_t \in S$ is the state of the system at time t .

An *atomic proposition* is a statement that has a unique truth value (*True* or *False*).

LTL is built up from (a) a set of atomic propositions, (b) the logic connectives: negation (\neg), disjunction (\vee), conjunction (\wedge) and material implication (\implies), and (c) the temporal modal operators: next (\circ), always (\square), eventually (\diamond) and until (\mathcal{U}). An LTL formula is defined inductively as follows: (1) any atomic proposition p is an LTL formula; and (2) given LTL formulas φ and ψ , $\neg\varphi$, $\varphi \vee \psi$, $\circ\varphi$ and $\varphi \mathcal{U} \psi$ are also LTL formulas.

The semantics of LTL is defined inductively as follows: (a) For an atomic proposition p , $s_i \models p$ if and only if (iff) $s_i \models p$; (b) $s_i \models \neg\varphi$ iff $s_i \not\models \varphi$; (c) $s_i \models \varphi \vee \psi$ iff $s_i \models \varphi$ or $s_i \models \psi$; (d) $s_i \models \circ\varphi$ iff $s_{i+1} \models \varphi$; and (e) $s_i \models \varphi \mathcal{U} \psi$ iff there exists $j \geq i$ such that $s_j \models \psi$ and $\forall k \in [i, j), s_k \models \varphi$. Based on this definition, $\circ\varphi$ holds at position s_i iff φ holds at the next state s_{i+1} , $\square\varphi$ holds at position i iff φ holds at every position in σ starting at position i , and $\diamond\varphi$ holds at position i iff φ holds at some position $j \geq i$ in σ .

Definition 4. An execution $\sigma = s_0 s_1 s_2 \dots$ *satisfies* φ , denoted by $\sigma \models \varphi$, if $s_0 \models \varphi$.

Remark 1. Properties typically studied in the control and hybrid systems domains are safety and stability. LTL also allows properties such as guarantee, liveness, and response.

An automaton representation of any LTL formula can be given by a deterministic Rabin automaton.

Definition 5. A deterministic Rabin automaton is a tuple $\mathcal{A} = (Q, \Sigma, \delta, q_0, Acc)$ where Q is a finite set of states, Σ is an

alphabet, $\delta : Q \times \Sigma \rightarrow Q$ is the transition function, $q_0 \in Q$ is the initial state, and accepting state pairs $Acc \subseteq 2^Q \times 2^Q$.

Let Σ^ω be the set of infinite words over Σ . A run for $\sigma = \mathcal{A}_0 \mathcal{A}_1 \mathcal{A}_2 \dots \in \Sigma^\omega$ denotes an infinite sequence $q_0 q_1 q_2 \dots$ of states in \mathcal{A} such that $q_{i+1} \in \delta(q_i, \mathcal{A}_i)$ for $i \geq 0$. The run $q_0 q_1 q_2 \dots$ is accepting if there exists a pair $(L, K) \in Acc$ and an $n \geq 0$, such that for all $m \geq n$ we have $q_k \notin L$ and there exist infinitely many k such that $q_k \in K$.

Intuitively, a run is accepted by a deterministic Rabin automaton if the set of states L is visited finitely often and the set K is visited infinitely often.

III. PROBLEM STATEMENT

We now provide a formal statement of the problem and an overview of our approach.

The system evolves with the discrete-time dynamics

$$x_{t+1} = f(x_t, u_t, w_t) \quad \text{for } t = 0, 1, \dots, \quad (3)$$

where $x_t \in \mathbb{R}^l$ is the system state, $u_t \in \mathbb{R}^m$ is the control input, and $w_t \in \mathbb{R}^p$ is a stochastic disturbance input.

We assume that relevant system properties are specified by a given LTL formula φ over a finite set AP of atomic propositions. Let $\mathcal{X} \subset \mathbb{R}^l$ be partitioned into a finite set \mathcal{R} of regions. Each region $r \in \mathcal{R}$ has an associated set of *True* atomic propositions associated with it, which we will informally call a *label*. Let $R : \mathcal{X} \rightarrow \mathcal{R}$ be a surjective mapping that associates each system state $x \in \mathcal{X}$ with a region $r \in \mathcal{R}$. Furthermore, we associate each region with a labeling function $L : \mathcal{R} \rightarrow 2^{AP}$. We assume a bijection between regions $r \in \mathcal{R}$ and MDP states $s \in S$, hence we use the same notation for this labeling function as in Definition 1. Since every system state $x \in \mathcal{X}$ in a given region has the same label, this is a *proposition-preserving* partition. Finally, let $\mathcal{S} : 2^{AP} \rightarrow \mathcal{X}$ which maps labels to system states.

Given a control policy $\pi = \{\mu_0, \mu_1, \dots\}$ the MDP abstraction of the dynamical system (3), we can induce a control policy π_{proj}^* on (3) by using action $\mu_k(R(x))$ for the system state $x \in \mathcal{X}$ at time $k = 0, 1, \dots$. This simply maps the system state to the corresponding region with R , and thus a unique MDP state due to the assumed bijection.

Problem 1. Given a labeled finite MDP \mathcal{M} with transition matrices in an uncertainty set \mathcal{P}^a for each action $a \in A$ and an LTL specification φ over AP , create a robust control policy π^* for \mathcal{M} . Furthermore, if \mathcal{M} is an abstraction of a dynamical system (3), create the induced control policy π_{proj}^* for (3).

We present a method for abstracting the original system as a finite MDP \mathcal{M} with transition matrix uncertainty sets in Section VI. We then solve Problem 1 by first creating the product MDP \mathcal{M}_p which contains system trajectories that satisfy both \mathcal{M} and the deterministic Rabin automaton representing the given LTL specification. We modify \mathcal{M}_p so that all policies for it are proper, and thus it satisfies the stochastic shortest path assumptions. We then show that maximizing the worst-case probability of satisfying the specification is equivalent

to a creating a control policy that maximizes the worst-case probability of reaching a certain set of states in \mathcal{M}_p . We solve for this policy using a robust variant of dynamic programming. Finally, we map the robust control policy back to \mathcal{M} and the original dynamical system (3).

IV. ROBUST DYNAMIC PROGRAMMING

We now prove when a robust form of dynamic programming holds for MDPs with uncertain transition matrices.

A. Dynamic Programming

We consider a MDP \mathcal{M} with a finite set of states $S = \{1, 2, \dots, n, t\}$ and actions $a \in A(s)$ for all $s \in S$. We assume that t is a special terminal state, which is absorbing ($P_{tt}^a = 1$) and incurs zero cost ($c(t) = 0$) for all $a \in A(t)$ and all $P \in \mathcal{P}_t^a$ [5]. It follows that the total expected cost for t , $J(t) = 0$ for all control and environment policies.

We require policies to be proper, i.e. they almost surely reach the terminal state t for all transition matrices in the uncertainty set [5].

Definition 6. A stationary control policy μ is *proper* if, under that policy, there is positive probability that the terminal state will be reached after at most n stages, regardless of the initial state and transition matrices, that is, if

$$\rho_{\mu\tau} := \max_{s=1, \dots, n} \max_{P \in \mathcal{P}^{\mu(s)}} P(s_n \neq t | s_0 = s, \mu) < 1. \quad (4)$$

Assumption 3. All stationary control policies are proper.

Remark 2. This assumption implies that the terminal state will eventually be reached under any stationary policy. This will allow us to make statements regarding convergence rates. While this is usually a rather strong condition, these assumptions are not restrictive for Problem 1, as in Section V.

In preparation for the main result of this section, we give the following classical theorem [21].

Theorem 1 (Contraction Mapping Theorem). *Let (M, d) be a complete metric space and let $f : M \rightarrow M$ be a contraction, i.e., there is a real number β , $0 \leq \beta < 1$, such that $d(f(x), f(y)) \leq \beta d(x, y)$ for all x and y in M . Then there exists a unique point x^* in M such that $f(x^*) = x^*$. Additionally, if x is any point in M , then $\lim_{k \rightarrow \infty} f^k(x) = x^*$, where f^k is the composition of f with itself k times.*

We now define mappings that play an important role in the rest of this section. We loosely follow the notation of [5]. The scalar $J(s)$ corresponds to the total expected cost to go when starting at state $s \in S$. The shorthand $J \in \mathbb{R}^n$ represents the cost vector for all $s \in S \setminus t$. Since the cost is always zero at the terminal state t , we do not include it. The T and $T_{\mu\tau}$ operators are mappings from \mathbb{R}^n to \mathbb{R}^n . For each state $s \in S \setminus t$, define the s -th component of TJ and $T_{\mu\tau}J$ respectively as

$$(TJ)(s) := \min_{a \in A(s)} [c(s, a) + \max_{P \in \mathcal{P}_s^a} p^T J], \quad (5)$$

$$(T_{\mu\tau}J)(s) := c(s, \mu(s)) + p^T J. \quad (6)$$

In the following two lemmas, we show that these mappings are monotonic and contractive. We prove these for (5); the proofs for (6) follow by limiting the actions and transition probabilities at each state s to $\mu(s)$ and P_s^a respectively. T^k is the composition of T with itself k times.

Lemma 1 (Monotonicity). *For any vectors $u, v \in \mathbb{R}^n$, such that $u \leq v$, we have that $T^k u \leq T^k v$ for $k = 1, 2, \dots$*

Proof: Immediate from (5) and the definition of \mathcal{P}_s^a . ■

Definition 7. The *weighted maximum norm* $\|\cdot\|_w$ of a vector $u \in \mathbb{R}^n$ is defined by $\|u\|_w = \max_{i=1, \dots, n} |u(i)|/w(i)$ where vector $w \in \mathbb{R}^n$ and $w > 0$.

Lemma 2 (Contraction). *If all stationary control policies are proper, then there exists a vector $w > 0$ and a scalar $\gamma \in [0, 1)$ such that $\|Tu - Tv\|_w \leq \gamma \|u - v\|_w$ for all $u, v \in \mathbb{R}^n$.*

Sketch of proof: The proof of Lemma 2 closely follows that in [5] (Vol. II, Section 2.4) where the environment policy is fixed. More specifically, the proof in [5] is modified to allow maximization over environment policies. This modification holds due to Assumption 3 and Lemma 1. Details are available in a technical report [24]. ■

We now prove the main result of this section. We remind the reader that the vector $J^* \in \mathbb{R}^n$, defined in (2), is the optimal worst-case total expected cost starting from state $s \in S$.

Theorem 2 (Robust Dynamic Programming). *Under the assumption that all stationary control policies μ are proper for a finite MDP \mathcal{M} with transition matrices in the uncertainty set \mathcal{P}^a for $a \in A$, the following statements hold.*

(a) *The optimal worst-case cost vector J^* is the unique fixed-point of T ,*

$$J^* = TJ^*. \quad (7)$$

(b) *The optimal worst-case cost vector J^* is given by,*

$$J^* = \lim_{k \rightarrow \infty} T^k J, \quad (8)$$

for all $J \in \mathbb{R}^n$. This limit is unique.

(c) *A stationary control policy μ and a stationary environment policy τ are optimal if and only if*

$$T_{\mu\tau} J^* = TJ^*. \quad (9)$$

Proof: Parts (a) and (b) follow immediately from Theorem 1 and Lemmas 1 and 2.

Part (c): First assume that $T_{\mu\tau} J^* = TJ^*$. Thus $T_{\mu\tau} J^* = TJ^* = J^*$ from (7) and $J^{\mu\tau} = J^*$ from the uniqueness of the fixed-point. Thus, μ and τ are optimal policies. Now assume that μ and τ are optimal policies so that $J^{\mu\tau} = J^*$. We then have that $T_{\mu\tau} J^* = T_{\mu\tau} J^{\mu\tau} = J^{\mu\tau} = J^*$. ■

Corollary 1. *Given the optimal worst-case cost vector J^* , the optimal control actions a^* satisfy*

$$a^*(s) \in \arg \min_{a \in A(s)} [c(s, a) + \max_{p \in \mathcal{P}_s^a} p^T J^*], \quad s \in S. \quad (10)$$

and, with some abuse of notation, the optimal transition vectors (for the environment) are

$$P_s^{*a} \in \arg \max_{p \in \mathcal{P}_s^a} p^T J^*, \quad s \in S, a \in A(s). \quad (11)$$

Proof: Follows from Part (c) in Theorem 2 and (5). ■

To recap, we showed that T is monotone and a contraction with respect to a weighted max norm. This allowed us to prove in Theorem 2 that T has a unique fixed-point that can be found by an iterative procedure (often called *value iteration*). We further gave conditions on the optimality of stationary policies and showed how to determine optimal actions for the system and the environment.

B. Uncertainty Set Representations

Referring back to the T operator (5), we see that it is composed of two nested optimization problems—the outer problem for the system and the inner problem for the environment. To be clear, the environment optimization problem for a given state $s \in S$ and control action $a \in A(s)$ refers to $\max_{p \in \mathcal{P}_s^a} p^T J$.

The tractability of this optimization problem depends on the structure of the uncertainty set \mathcal{P}_s^a . In the remainder of this section, we investigate interval and likelihood uncertainty sets, as these are both statistically meaningful and computationally efficient. Due to lack of space, we do not discuss maximum a priori, entropy, scenario, or ellipsoidal uncertainty models, even though these are included in this framework. The reader should refer to Nilim and El Ghaoui for details [22].

1) *Interval Models:* A common description of uncertainty for transition matrices corresponding to action $a \in A$ is by intervals $\mathcal{P}^a = \{P^a \mid \underline{P}^a \leq P^a \leq \overline{P}^a, P^a \mathbf{1} = \mathbf{1}\}$, where \underline{P}^a and \overline{P}^a are nonnegative matrices $\underline{P}^a \leq \overline{P}^a$. This representation is motivated by statistical estimates of confidence intervals on the individual components of the transition matrix [19]. The environmental optimization problem can be solved in $O(n \log(n))$ time using a bisection method [22].

2) *Likelihood Models:* The likelihood uncertainty model is motivated by determining the transition probabilities between states through experiments. We denote the experimentally measured transition probability matrix corresponding to action a by F^a and the optimal log-likelihood by β_{\max} .

Uncertainty in the transition matrix for each action $a \in A$ is described by the *likelihood region* [19]

$$\mathcal{P}^a = \{P^a \in \mathbb{R}^{n \times n} \mid P^a \geq 0, P^a \mathbf{1} = \mathbf{1}, \sum_{i,j} F_{ij}^a \log P_{ij}^a \geq \beta^a\}, \quad (12)$$

where $\beta^a < \beta_{\max}^a$ and can be estimated for a desired confidence level by using a large sample Gaussian approximation [22]. As described in Assumption 2, we enforce that $F_{ij}^a = 0$ if and only if $P_{ij}^a = 0$ for all $i, j \in S$ and all $a \in A$.

Since the likelihood region above is not rectangular (contrary to Assumption 1), it must be approximated by projection onto each row of the transition matrix. Even with the approximation, likelihood regions are less conservative uncertainty representations than intervals, which arise from further projection onto the row's components. A bisection algorithm can

approximate the environment optimization problem to within an accuracy δ in $O(\log(J_{\max}/\delta))$ time, where J_{\max} is the maximum value of the cost vector, J [22].

V. CREATING THE ROBUST CONTROL POLICY

We now solve Problem 1 by creating a robust control policy. First, we form the product MDP \mathcal{M}_p which allows behaviors that satisfy both the system MDP \mathcal{M} and the LTL specification φ . We show that \mathcal{M}_p can be transformed into an equivalent form \mathcal{M}_{ssp} where all stationary control policies are proper. Next, we use the robust dynamic programming developed in Section IV to find a control policy that maximizes the probability of satisfying φ . Finally, we project this control policy to a policy for \mathcal{M} .

A. Forming the product MDP

The product MDP \mathcal{M}_p restricts behaviors to those that satisfy both the system transitions and the deterministic Rabin automaton \mathcal{A} representing the LTL specification.

Definition 8. For labeled finite MDP $\mathcal{M} = (S, A, P, s_0, AP, L)$ and deterministic Rabin automaton $\mathcal{A} = (Q, 2^{AP}, \delta, q_0, Acc)$, the product MDP $\mathcal{M}_p = (S_p, A, P_p, s_{0p}, Q, L_p)$ with $S_p = S \times Q$,

- $P_p((s, q), \alpha, (s', q')) = \begin{cases} P(s, \alpha, s') & \text{if } q' = \delta(q, L(s')) \\ 0 & \text{otherwise,} \end{cases}$
- $s_{0p} = (s_0, q)$ such that $q = \delta(q_0, L(s_0))$,
- $L_p((s, q)) = \{q\}$.

The accepting product state pairs $Acc_p = \{(L_1^p, K_1^p), \dots, (L_k^p, K_k^p)\}$ are lifted directly from Acc . Formally, for every $(L_i, K_i) \in Acc$, state $(s, q) \in S_p$ is in L_i^p if $q \in L_i$, and $(s, q) \in K_i^p$ if $q \in K_i$.

There is a one-to-one correspondence between the paths on \mathcal{M}_p and \mathcal{M} , which induces a one-to-one correspondence for policies on \mathcal{M}_p and \mathcal{M} . So, given a policy $\pi^p = \{\mu_0^p, \mu_1^p, \dots\}$ on \mathcal{M}_p , we can induce a policy $\pi = \{\mu_0, \mu_1, \dots\}$ on \mathcal{M} by setting $\mu_i(s_i) = \mu_i^p((s_i, q_i))$ for every stage $i = 0, 1, \dots$. This is always valid since \mathcal{M}_p and \mathcal{M} have the same action set A . If π^p is stationary, π is guaranteed to be finite-memory, but not necessarily stationary [4].

B. Reachability in product MDP

We now show how to use the product MDP \mathcal{M}_p to determine a robust control policy that maximizes the worst-case probability that a given LTL specification is satisfied. Given a control and environment policy, the probability of satisfying an LTL formula is equivalent to the probability of reaching a certain set of states in \mathcal{M}_p , which are *accepting maximal end components* [4]. We call this probability the *reachability probability*. Informally, accepting maximal end components are sets of states that the system can remain in forever and where the acceptance condition of the deterministic Rabin automaton is satisfied.

Definition 9. A *sub-MDP* of a MDP is a pair of states and action sets (C, D) where: (1) $\emptyset \neq C \subseteq S$ and $D : C \rightarrow 2^A$ is a

function such that $\emptyset \neq D(s) \subseteq A(s)$ for all states $s \in C$ and (2) $s \in C$ and $a \in D(s)$ implies $Post(s, a) = \{t \in S | P_{st}^a > 0\} \subseteq C$.

Definition 10. An *end component* is a sub-MDP (C, D) such that the digraph $G_{(C, D)}$ induced by (C, D) is strongly connected.

An end component (C, D) is *maximal* if there is no end component (C', D') such that $(C, D) \neq (C', D')$ and $C \subseteq C'$ and $D(s) \subseteq D'(s)$ for all $s \in C$. Furthermore, (C, D) is *accepting* for the deterministic Rabin automaton \mathcal{A} if for some $(L, K) \in Acc$, $L \subseteq C$ and $K \in C$.

After computing the accepting maximal end components of \mathcal{M}_p , we need to determine a control policy that maximizes the worst-case probability of reaching an accepting maximal end component from our initial state. Without considering transition probability uncertainty, this policy can now be computed using either linear or dynamic programming methods [4]. We need to use the robust dynamic programming approach from Section IV. Note that this approach does not directly apply to \mathcal{M}_p , as not all stationary control policies on \mathcal{M}_p are proper. We can transform \mathcal{M}_p into an equivalent MDP where all stationary policies are proper, as shown in the next subsection.

C. Transformation to Stochastic Shortest Path

We now show how to transform the product MDP \mathcal{M}_p into an equivalent form \mathcal{M}_{ssp} where all stationary control policies μ are proper (cf. Assumption 3). It is important to note that \mathcal{M}_p and \mathcal{M}_{ssp} are equivalent only in terms of the probability of reaching an accepting maximal end component—both the states and the transition probabilities will likely change.

In the remainder of this section, we use the simplified notation $\mathcal{M}_p = (S, A, P)$ to describe the states, actions, and transition matrices of the product MDP. We refer to the state of \mathcal{M}_p as s instead of (s, q) when clear from context.

We first partition the states S of \mathcal{M}_p into three disjoint sets. We let B be the union of all accepting maximal end components in \mathcal{M}_p . By definition, every state $s \in B$ has a maximum reachability probability of 1. We define the set of states that have zero probability of reaching B by S_0 . These can be found efficiently by graph algorithms [4]. Finally, we define the set $S_r = S - (B \cup S_0)$ as all of states not in an accepting maximal end component with non-zero maximum reachability probability. It is easy to see that these three sets form a partition of S .

Algorithm 1 Appending the terminal state

Require: $\mathcal{M}_p = (S, A, P)$ and S_r, S_0, B
 $S := S \cup \{t\}$ and $A(t) := \{u\}$ and $c(t, u) := 0$
 $A(s) := \{u\}$ and $P_{st}^u := 1$ for all $s \in B \cup S_0$

In Algorithm 1, we augment S with a terminal state t which is absorbing and incurs zero cost (cf. Section IV).

Remark 3. Algorithm 1 does not change the probability of reaching an accepting maximal end component for any state $s \in S$ under any control and environment policies.

Algorithm 2 End component elimination (de Alfaro [9])

Require: MDP $\mathcal{M}_p = (S, A, P)$ and S_r, S_0, B

Ensure: MDP $\mathcal{M}_{ssp} = (\hat{S}, \hat{A}, \hat{P})$

$\{(C_1, D_1), \dots, (C_k, D_k)\}$ max end components in S_r

$\hat{S}_0 := S_0$ and $\hat{B} := B$

$\hat{S} := S \cup \{\hat{s}_1, \dots, \hat{s}_k\} - \cup_{i=1}^k C_i$

$\hat{S}_r := S_r \cup \{\hat{s}_1, \dots, \hat{s}_k\} - \cup_{i=1}^k C_i$

$\hat{A}(s) := \{(s, a) \mid a \in A(s)\}$ for $s \in S - \cup_{i=1}^k C_i$

$\hat{A}(\hat{s}_i) := \{(s, a) \mid s \in C_i \wedge a \in A(s) - D(s)\}$ for $1 \leq i \leq k$

For $s \in \hat{S}, t \in S - \cup_{i=1}^k C_i$ and $(u, a) \in \hat{A}(s), \hat{P}_{st}^{(u,a)} := P_{ut}^a$
and $\hat{P}_{s\hat{s}_i}^{(u,a)} := \sum_{t \in C_i} P_{ut}^a$

Algorithm 2 eliminates the maximal end components in S_r and replaces them with new states. This procedure is from Algorithm 3.3 of [9], where it is also proven (Theorem 3.8 in [9]) that the reachability probabilities are unchanged by this procedure. The intuition behind this result is that one can move between any two states r and s in a maximal end component in S_r with probability one and zero cost.

After applying Algorithms 1 and 2, we call the resulting MDP \mathcal{M}_{ssp} . We note that \hat{S}_r, \hat{B} , and \hat{S}_0 form a disjoint partition of \hat{S} . All stationary control policies for \mathcal{M}_{ssp} are proper, i.e., they will almost surely reach the terminal state t .

Theorem 3. *All stationary control policies for \mathcal{M}_{ssp} are proper.*

Thus, \mathcal{M}_{ssp} is equivalent in terms of reachability probabilities to the original product MDP \mathcal{M}_p . Furthermore, all stationary control policies are proper.

D. Computing the optimal control policy

We now perform robust value iteration as described in Section IV on the transformed product MDP \mathcal{M}_{ssp} . This method applies because \mathcal{M}_{ssp} has been constructed so that all stationary policies are proper.

As we formulated the dynamic programming approach in terms of total expected cost minimization, we define the total expected cost as the negative of the reachability probability, which is equivalent to the probability of satisfying the LTL formula. We refer to cost and probability interchangeably in the remainder. Thus, for all $a \in \hat{A}$, the appropriate costs are $c(s, a) = -1$ for all $s \in \hat{B}$ and $c(s, a) = 0$ for all $s \in \hat{S}_0$. For the remaining states, $s \in \hat{S}_r$, we initialize the costs arbitrarily in $[-1, 0]$ and compute the optimal worst-case cost vector using the iteration presented in Theorem 2. The resulting cost vector is $J_{ssp}^* \in \mathbb{R}^m$, where m is the number of states in \mathcal{M}_{ssp} .

We first determine the cost vector $J_p^* \in \mathbb{R}^n$ for \mathcal{M}_p from $J_{ssp}^* \in \mathbb{R}^m$. For $s_p \in S_p$, we determine the corresponding state $s_{ssp} \in \hat{S}$ and let $J_p^*(s_p) = J_{ssp}^*(s_{ssp})$. This mapping is surjective, as there is at least one s_p for each s_{ssp} .

Now that we have the optimal worst-case cost vector J_p^* for the original product MDP \mathcal{M}_p , we need to determine the optimal actions $a^*(s) \in A(s)$ for each $s \in S_r$. We do not consider actions for states in $S_0 \cup B$ at this time. However,

we cannot simply use the approach for selecting actions given by (10), because not all stationary control policies on \mathcal{M}_p are proper. For states in a maximal end component in S_r , there may be multiple actions that satisfy (10). Arbitrarily selecting actions can lead to situations where the stationary control policy stays in the maximal end component forever and thus never satisfies the specification. We avoid this situation by only selecting an action if it is both optimal (i.e., satisfies (10)) and it has a non-zero probability of transitioning to a state that is not in a maximal end component in S_r . Algorithm 3 selects the action with the highest probability of transitioning to a state not in a maximal end component in S_r .

Algorithm 3 Product MDP Control Policy

Require: $J_p^* \in \mathbb{R}^n, \mathcal{M}_p = (S, A, P), S_r, S_0, B$

$visited := S_0 \cup B$

$possAct(s) := \{a \in A(s) \mid (T_a J_p^*)(s) = J_p^*(s)\}$

for $s \in S_r$ **do**

if $|possAct(s)| = 1$ **then**

$\mu(s) := possAct(s)$ and $visited := visited \cup \{s\}$

end if

end for

while $visited \neq S$ **do**

for $s \in S_r \setminus visited$ **do**

$maxLeaveProb := 0$

$leaveProb := \max_{a \in possAct(s)} \sum_{t \in visited} P_{st}^a$

if $leaveProb > maxLeaveProb$ **then**

$optAct := a$ and $optState := s$

end if

end for

$\mu(s) := optAct$ and $visited := visited \cup \{optState\}$

end while

Theorem 4. *Algorithm 3 returns a robust control policy μ that satisfies $J_p^{\mu\tau} = J_p^*$ for the worst-case environmental policy τ .*

Sketch of Proof: By construction, every state has at least one action that satisfies (10). For states in a maximal end component in S_r , some actions may cause the system to stay in S_r . This is avoided by always selecting optimal actions leading outside of maximal end components formed by states where an action has not already been selected. At least one such action always exists from the definition of end components. ■

The optimal control policy for satisfying the LTL specification φ consists of two parts: a stationary deterministic policy for reaching an accepting maximal end component, and a finite-memory deterministic policy for staying there. The former policy is given by Algorithm 3 and denoted μ_{reach} . The latter policy is a finite-memory controller π_B that selects actions in a round-robin fashion to ensure that the system stays inside the accepting maximal end component for all time and thus satisfies φ [4]. The overall optimal policy is $\pi_p^* = \mu_{reach}$ if $s \notin B$ and $\pi_p^* = \pi_B$ if $s \in B$. We induce an optimal policy π^* on \mathcal{M} as described in Section V-A.

E. Complexity

We now detail the complexity of our approach for solving Problem 1. The number of states in the system MDP \mathcal{M} is the number of labeled regions in \mathcal{R} due to the bijection between MDP states and regions. The number of transitions for each state is given by the reachability between regions $r \in \mathcal{R}$, and is usually sparse for robot motion planning problems.

The (worst-case) size of the deterministic Rabin automaton is doubly-exponential in the length of the LTL formula [8]. Experimental work in [13] has shown that deterministic Rabin automaton sizes are often exponential or lower for many common types of LTL formulae. Also, there are fragments of LTL, which include all safety and guarantee properties, that generate a deterministic Rabin automaton whose size is singly-exponential in the length of the formula [2].

The size of the product MDP \mathcal{M}_p is equal to the size of \mathcal{M} times the size of \mathcal{A} . \mathcal{M}_p has n states and m transitions. Maximal end components can be found in $O(n^2)$ time. Since T is a contraction, an ϵ -suboptimal control policy can be computed in $O(n^2 m \log(1/\epsilon))$ time without uncertainty sets [5] and $O(n^2 m \log(1/\epsilon)^2)$ when using likelihood transition uncertainty sets. Thus, the computational cost for incorporating robustness is $O(\log(1/\epsilon))$ times the non-robust case.

VI. CREATION OF MDP AND UNCERTAINTY SETS

We now discuss a simulation-based approach for creating a finite MDP abstraction \mathcal{M} of a dynamical system (3). It should be noted that the main contribution of our work, the solution of Problem 1, does not depend on the method for creating \mathcal{M} and its corresponding transition matrix uncertainty sets \mathcal{P} . Thus, if specific systems are more amenable to other abstraction methods, these can be used to create \mathcal{M} and \mathcal{P} , after which our methods can compute a robust control policy. An approach using control primitives is given in [17].

Given a dynamical system (3) and a set \mathcal{R} of labeled regions, where the labels correspond to sets of *True* atomic propositions that an LTL formula is defined over, we want to compute reachability for the system between all regions in \mathcal{R} . We say that for two regions $a, b \in \mathcal{R}$, region b is *reachable* from region a , if there exists a control sequence $u_{1:t}$ for every system state in a that evolves the system to some system state in region b in finite time.

Computing reachability is hard in general, so it is often necessary to consider a simplified model of (3) where reachability can be efficiently computed. Given candidate control sequences based on reachability analysis of the simplified system, one can use simulations to estimate the effect of these control actions on the full system dynamics and use these estimates of reachability to create an MDP abstraction. For every region, randomly sample N initial system states inside it and apply the candidate control sequences. Based on the average of all N simulated trajectories of the system (3) between regions i and j , create the estimated transition probability matrix F_{ij}^a entries. It is important to consider a transition between regions i and j successful only if the trajectory remains in the union of the two regions. This is because one cannot

guarantee that a LTL specification is satisfied if the system visits other regions in between. Trajectories that visit multiple regions can transition to a special 'catch-all' state that has zero probability of satisfying the specification. The designer can tradeoff between the number of samples N and the size of the transition matrix uncertainty sets.

Remark 4. Even if a method can create an exact MDP from (3), transition matrix uncertainty sets should still be considered, as (3) may not fully model the physical system.

An uncertainty set for the transition matrices of \mathcal{M} can be generated using the simulation data. Methods for determining these sets in a principled manner can be found in [22].

VII. EXAMPLE

We demonstrate our robust control framework on a mobile robot, where there is uncertainty in the control inputs from stochastic disturbances caused by terrain variation.

We use the discrete-time form of the standard unicycle model for our system

$$\begin{aligned} x_{k+1} &= x_k + (v + d_v) \cos(\theta) \Delta t \\ y_{k+1} &= y_k + (v + d_v) \sin(\theta) \Delta t \\ \theta_{k+1} &= \theta_k + (\omega + d_\omega) \Delta t \end{aligned}$$

where $\Delta t = t_{k+1} - t_k$, v is the velocity input, ω is the angular velocity input, and $d = (d_v, d_\omega)^T$ is a stochastic disturbance.

In order to determine candidate control sequences, we ignore the stochastic disturbances and use feedback linearization to map between the original control inputs $u = (v_x, v_y)^T$ and new inputs $u_{lin} = (v, \omega)^T$. This map is given by $v = v_x \cos(\theta) + v_y \sin(\theta)$ and $\omega = 1/\epsilon (-v_x \sin(\theta) + v_y \cos(\theta))$, where ϵ is a tuning parameter related to the wheel base length.

We then check reachability for the simplified model $x_{k+1} = x_k + v_x$ and $y_{k+1} = y_k + v_y$, where v_x and v_y are the velocities in the x and y directions. This system is linear, so we can use TuLiP to check reachability [25]. We ignore θ here because our labels only depend on x and y .

The task for the robot is to sequentially visit two regions of interest while always remaining safe. Once the robot has visited the regions in order, it should return to the start and remain there. The atomic propositions are $\{home, unsafe, R1, R2\}$. The LTL specification formalizing this task is $\varphi = home \wedge \diamond \square home \wedge \square \neg unsafe \wedge \diamond (R1 \wedge \diamond R2)$.

We used the procedure outlined in Section VI to create a finite MDP abstraction \mathcal{M} of (13), where each state of \mathcal{M} is a square region in \mathcal{R} . The actions at each state include transitioning to a neighbor, rotating in place, or not moving. As the system does not satisfy the Markov property (due to the θ term), we add four MDP states per region as a coarse discretization for θ . Due to symmetry of the regions and robot, we only need to calculate transitions for one region. An additional 16-fold speedup is possible by also exploiting rotational symmetries, but this is not implemented.

All computations were run on an 2.4 GHz dual core desktop with 2 GB of RAM. It took 20 minutes to construct the MDP abstraction and uncertainty sets using 400

samples for the Monte Carlo simulation. We used large sample approximations to estimate β^a (12) for an uncertainty level of 0.9 [22]. The deterministic Rabin automaton representing φ has seven states. The product MDP has 597 states and 18455 transitions and took 1.1 seconds to compute. It took 25.4 seconds to compute \mathcal{M}_{ssp} . It took 44.3 seconds to generate an ϵ -suboptimal robust control policy with likelihood uncertainty sets and 3.4 seconds to generate a non-robust control policy. In both cases, $\epsilon = 10^{-5}$.

The calculated maximum probability of satisfying the specification was 0.763 for the non-robust policy and 0.526 for the robust policy. However, as the non-robust policy does not account for transition matrix uncertainties, it is likely to be overly optimistic when implemented on a real system. In this case, the non-robust policy only had a 0.491 probability of satisfying the specification given the worst-case transition matrices in the uncertainty set. Figure 1 shows sample trajectories of the robot using the robust and non-robust control policies.

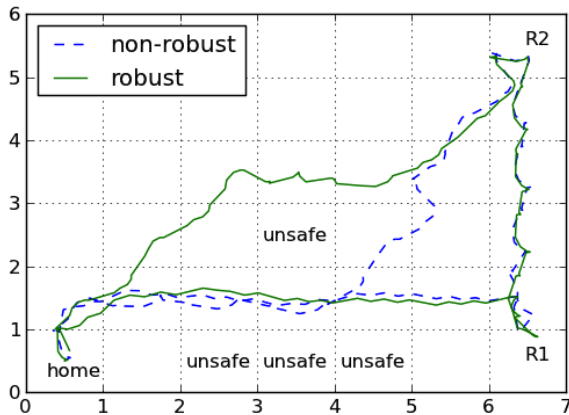


Fig. 1. Sample trajectories of the robot using the robust and non-robust control policies. The robust control policy takes a more conservative path.

VIII. CONCLUSIONS AND FUTURE DIRECTIONS

We have presented a method for creating robust control policies for finite MDP abstractions of dynamical systems with temporal logic specifications. This robustness is useful both when the abstraction is not exact or when the dynamical system model does not adequately represent the physical system. Designing an ϵ -suboptimal robust control policy increases the time complexity only by a factor of $O(\log(1/\epsilon))$ compared to the non-robust policy for statistically relevant uncertainty sets.

In the future, we plan to extend these results to other temporal logics, such as probabilistic computational tree logic (PCTL). We are also looking into weakening the rectangularity assumption on the transition matrix uncertainty sets to allow for correlation. Finally, further work needs to be done on principled methods for abstraction of selected classes of dynamical systems as finite MDPs.

ACKNOWLEDGEMENTS

The authors would like to thank Scott Livingston for helpful comments. This work was supported by an NSF Graduate Research Fellowship and the Boeing Corporation.

REFERENCES

- [1] R. Alur, T. A. Henzinger, G. Lafferriere, and G. J. Pappas. Discrete abstractions of hybrid systems. *Proc. IEEE*, 88(7):971–984, 2000.
- [2] R. Alur and S. La Torre. Deterministic generators and games for LTL fragments. *ACM Trans. Comput. Logic*, 5(1):1–25, 2004.
- [3] J. A. Bagnell, A. Y. Ng, and J. G. Schneider. Solving uncertain Markov decision processes. Technical report, Carnegie Mellon University, 2001.
- [4] C. Baier and J.-P. Katoen. *Principles of Model Checking*. MIT Press, 2008.
- [5] D. P. Bertsekas. *Dynamic Programming and Optimal Control (Vol. I and II)*. Athena Scientific, 2001.
- [6] O. Buffet. Reachability analysis for uncertain ssp. In *Proceedings of the IEEE International Conf. on Tools with Artificial Intelligence*, 2005.
- [7] K. Chatterjee, K. Sen, and T. Henzinger. Model-checking omega-regular properties of interval Markov chains. In R. M. Amadio, editor, *Foundations of Software Science and Computation Structure (FoSSaCS)*, pages 302–317, March 2008.
- [8] C. Courcoubetis and M. Yannakakis. The complexity of probabilistic verification. *Journal of the Association for Computing Machinery*, 42:857–907, 1995.
- [9] L. de Alfaro. *Formal Verification of Probabilistic Systems*. PhD thesis, Stanford University, 1997.
- [10] X. C. Ding, S. L. Smith, C. Belta, and D. Rus. LTL control in uncertain environments with probabilistic satisfaction guarantees. In *Proceedings of 18th IFAC World Congress*, 2011.
- [11] R. Givan, S. Leach, and T. Dean. Bounded-parameter Markov decision processes. *Artificial Intelligence*, 122:71–109, 2000.
- [12] B. Johnson and H. Kress-Gazit. Probabilistic analysis of correctness of high-level robot behavior with sensor error. In *Proceedings of Robotics: Science and Systems*, 2011.
- [13] J. Klein and C. Baier. Experiments with deterministic omega-automata for formulas of linear temporal logic. *Theoretical Computer Science*, 363:182–195, 2006.
- [14] M. Kloetzer and C. Belta. A fully automated framework for control of linear systems from temporal logic specifications. *IEEE Transaction on Automatic Control*, 53(1):287–297, 2008.
- [15] H. Kress-Gazit, G. Fainekos, and G. Pappas. Where’s Waldo? Sensor-based temporal logic motion planning. In *Proc. of IEEE International Conference on Robotics and Automation*, pages 3116–3121, 2007.
- [16] M. Kwiatkowska, G. Norman, and D. Parker. PRISM 4.0: verification of probabilistic real-time systems. In *Proceedings of 23rd International Conference on Computer Aided Verification*, 2011.
- [17] M. Lahijanian, S. B. Andersson, and C. Belta. A probabilistic approach for control of a stochastic system from LTL specifications. In *IEEE Conference on Decision and Control*, 2009.
- [18] M. Lahijanian, S. B. Andersson, and C. Belta. Control of Markov decision processes from PCTL specifications. In *Proceedings of the American Control Conference*, 2011.
- [19] E. L. Lehmann and J. P. Romano. *Testing Statistical Hypotheses*. Springer, 2005.
- [20] R. Majumdar, E. Render, and P. Tabuada. Robust discrete synthesis against unspecified disturbances. In *Proc. Hybrid Systems: Computation and Control*, pages 211–220, 2011.
- [21] A. W. Naylor and G. R. Sell. *Linear Operator Theory in Engineering and Science*. Springer-Verlag, 1982.
- [22] A. Nilim and L. El Ghaoui. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53:780–798, 2005.
- [23] J. K. Satia and R. E. L. Jr. Markovian decision processes with uncertain transition probabilities. *Operations Research*, 21(3):pp. 728–740, 1973.
- [24] E. M. Wolff, U. Topcu, and R. Murray. Robust Control of Uncertain Markov Decision Processes with Temporal Logic Specifications. Technical report, 2011. Available at <http://users.cms.caltech.edu/~utopcu/accl2mdp.html>.
- [25] T. Wongpiromsarn, U. Topcu, and R. M. Murray. Receding horizon control for temporal logic specifications. In *Proc. of the 13th International Conference on Hybrid Systems: Computation and Control*, 2010.
- [26] T. Wongpiromsarn, U. Topcu, and R. M. Murray. Receding horizon temporal logic planning. *IEEE Transactions on Automatic Control*, 2010, submitted.
- [27] D. Wu and X. Koutsoukos. Reachability analysis of uncertain systems using bounded-parameter Markov decision processes. *Artif. Intell.*, 172:945–954, 2008.