

Probabilistic Motion Planning Under Temporal Tasks and Soft Constraints

Meng Guo , Member, IEEE, and Michael M. Zavlanos , Member, IEEE

Abstract—This paper studies motion planning of a mobile robot under uncertainty. The control objective is to synthesize a finite-memory control policy, such that a high-level task specified as a linear temporal logic formula is satisfied with a desired high probability. Uncertainty is considered in the workspace properties, robot actions, and task outcomes, giving rise to a Markov decision process that models the proposed system. Different from most existing methods, we consider cost optimization both in the prefix and suffix of the system trajectory. We also analyze the potential trade-off between reducing the mean total cost and maximizing the probability that the task is satisfied. The proposed solution is based on formulating two coupled linear programs, for the prefix and suffix, respectively, and combining them into a multiobjective optimization problem, which provides provable guarantees on the probabilistic satisfiability and the total cost optimality. We show that our method outperforms relevant approaches that employ Round-Robin policies in the trajectory suffix. Furthermore, we propose a new control synthesis algorithm to minimize the frequency of reaching a bad state when the probability of satisfying the tasks is zero, in which case, most existing methods return no solution. We validate the above-mentioned schemes via both numerical simulations and experimental studies.

Index Terms—Chance-constrained optimization, linear temporal logic (LTL), Markov decision process (MDP), motion planning.

I. INTRODUCTION

IN THIS paper, we study the problem of robot motion planning under uncertainty and temporal task specifications. We consider uncertainty in the workspace properties, robot motion and actions, and outcome of task executions, which gives rise to a Markov decision process (MDP) to model the proposed system. MDPs have been used extensively to model motion and sensing uncertainty in robotics [1], [2], and then solve decision-making problems that optimize a given control objective. The most common objective is to reach a goal state from an initial state while minimizing the cost. The resulting solution is a policy that maps states to actions [2]. On the other hand, a linear

Manuscript received June 21, 2017; revised October 23, 2017; accepted January 11, 2018. Date of publication January 30, 2018; date of current version December 3, 2018. This work was supported by the NSF under Grant IIS #1302283. Recommended by Associate Editor A. Girard. (Corresponding author: Meng Guo.)

The authors are with the Department of Mechanical Engineering and Materials Science, Duke University, Durham, NC 27708 USA (e-mail: meng.guo@duke.edu; michael.zavlanos@duke.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TAC.2018.2799561

temporal logic (LTL) provides a formal language to describe complex high-level tasks beyond the classic start-to-goal navigation. An LTL task formula is usually specified with respect to an abstraction of the robot motion within the allowed workspace [3], modeled by a deterministic finite transition system (FTS). Then, a high-level discrete plan is found using off-the-shelf model-checking algorithms [4], which is then executed through low-level continuous controllers [3], [5]. This framework is extended to allow for both robot motion and actions in the task specification [6] and partially known or dynamic workspaces in [7] and [8].

Recently, there have been many efforts to address the problem of synthesizing a control policy for an MDP that satisfies high-level temporal tasks specified in various formal languages. Different classes of probabilistic computation tree logic (PCTL) formulas have been studied in [9] for abstraction and verification over interval-valued Markov chains (MCs). The work in [10] proposes a control policy for a mobile robot that maximizes the probability of satisfying a bounded LTL formula. Syntactical co-safe LTL formulas are considered in [11] for a deterministic robot that coexists with other robots whose behavior is modeled as an MDP. An FTS with time-varying rewards is controlled to satisfy an LTL formula and maximize the accumulated reward in [12]. A robust control policy for MDPs with uncertain transition probabilities is proposed in [8]. A verification toolbox is provided in [13] for probabilistic discrete-time or continuous-time MC, under a wide variety of quantitative properties expressed in a PCTL, an LTL, a CTL, etc.

In this paper, we study motion planning of a mobile robot under uncertainty in both robot motion and workspace properties. The goal is to synthesize a finite-memory control policy that generates robot trajectories that satisfy a high-level LTL task formula with desired high probability. At the same time, we optimize the total cost *both* in the prefix and suffix parts of the system trajectories. Our proposed approach is based on solving two coupled linear programs, one for the prefix and one for the suffix, over the occupancy measures of the product automaton introduced in [14]. Moreover, we explore cases where the probability of satisfying the LTL tasks is zero, so that an accepting end component (AEC) does not exist in the MDP, where most relevant work returns no solutions. To address such situations, we treat satisfaction of the tasks as soft constraints and propose a relaxed suffix plan that minimizes the frequency with which the system enters bad states that violate the task specifications. We show that our approach outperforms the widely used Round-Robin policy, via both numerical simulations and experimental studies. We also compare our proposed method with the widely used probabilistic model-checking tool PRISM [13].

This paper is related to literature on the following:

- 1) policy synthesis for MDPs under multiple objectives;
- 2) cost optimization within AECs in MDPs; and
- 3) infeasible temporal tasks.

We discuss below this literature and highlight our contributions.

Since we consider both temporal tasks and total-cost criteria over MDPs, this paper is closely related to the policy synthesis of MDPs under multiple objectives. The work in [14] proposes a framework with provable correctness to synthesize a control policy for MDPs under multiple constrained total-cost criteria. A survey on multiobjective decision-making for MDPs can be found in [15]. On the other hand, verification of MDPs under *multiple* high-level tasks is addressed in [16], where the probability of satisfying each subtask is lower bounded by a given value. Moreover, a quantitative multiobjective verification scheme is proposed in [17] and [18] for numerical queries over probabilistic reward predicates. On the other hand, the seminal works [19], [20] consider MDPs with multidimensional weights under multipercentile queries that may be conflicting. However, most of the above-mentioned work does not address cost optimization over the suffix of the system trajectory within the AECs, neither does it address the case where no AECs can be found in the product automaton, which are the main contributions here.

The satisfaction of an LTL formula is associated with reaching the corresponding AECs. In particular, in [4, Ch. 10], a value iteration method is used to solve the maximal reachability problem towards the AECs to obtain a policy for the plan prefix. For planning within the AECs, the Round-Robin policy, which guarantees only correctness but not optimality, is adopted in [4], [17], and [21]. Optimal policies for the plan suffix that keeps the system within the AECs have been proposed in [22]–[25]. Specifically, in [22], the expected cost of satisfying instances of a desired property is minimized, whereas in [23], the minimal bottleneck cost is considered. Both approaches in [22] and [23] require particular types of LTL formulas (such as “always eventually”). The work in [24] and [26] considers MDPs with ω -regular specifications and quantitative resource constraints within the AECs. The work in [25] investigates the Pareto cost of a human-in-the-loop MDP measured by a given discounted cost function. Compared to this literature, the multiobjective optimization problem that we formulate to solve the control synthesis problem allows us to explicitly characterize the trade-off between prefix and suffix optimality. We then extend this methodology to the case where no AECs can be found.

Most aforementioned work [4], [17], [19]–[22], [27] relies on the assumption that the product automaton contains at least one AEC. However, in many situations, this assumption does not hold so that the probability of satisfying the task under any policy is *zero*. In this case, it is still important to identify those policies that minimize the frequency with which the system will reach the bad states that violate the task specifications. Consequently, it is desirable to synthesize a policy with certain risk guarantees even when soft LTL tasks are considered that are only partially feasible. To the best of our knowledge, there is no work on control synthesis for infeasible soft LTL task formulas defined on MDPs, especially when an AEC cannot be found in the resulting product automaton. For deterministic transition systems, a framework for robot motion planning in partially

known workspaces is proposed in [7] that can handle soft LTL task formulas whose satisfiability is improved over time; a least violating control strategy is synthesized in [28] for a set of LTL safety rules. In the case of MDPs, a relevant formulation is considered in [29] where an MDP is controlled to satisfy an ω -regular formula. A policy is proposed to ensure that the MDP enters a failure state relatively late in the prefix. However, a multiobjective criterion of the control policy, especially in the plan suffix, is not considered there. Also, recent work in [30] proposes an approach to increase the satisfaction probability by modifying the task formula which, however, only considers co-safe LTL formulas without cost optimization constraints.

In summary, the main contribution of this paper is threefold, which is given as follows.

- 1) A framework that optimizes the total cost both in the plan prefix and suffix while ensuring that the tasks are satisfied with a desired high probability.
- 2) A new algorithm to synthesize the control policies that have a high probability of satisfying the task over long time intervals, for cases where an AEC does not exist.
- 3) A new method that allows the system to recover from bad states and continue the task.

The rest of this paper is organized as follows. Section II introduces necessary preliminaries. In Section III, we formalize the considered problem. Section IV presents our solution in details, which includes four major parts. Section V demonstrates the feasibility of the results by numerical simulations. Section VI contains the experimental results. We conclude and discuss about future directions in Section VII.

II. PRELIMINARIES

A. Transient MDP

An MDP is defined as a 6-tuple $\mathcal{M} \triangleq (X, U, D, p_D, c_D, x_0)$, where X is the finite state space; U is the finite control action space (with a slight abuse of notation, $U(x)$ also denotes the set of control actions *allowed* at state $x \in X$); $D = \{(x, u) \mid x \in X, u \in U(x)\}$ is the set of possible state-action pairs; $p_D : X \times U \times X \rightarrow [0, 1]$ is the transition probability function so that $p_D(x, u, \tilde{x})$ is the transition probability from state x to state \tilde{x} via control action u and $\sum_{\tilde{x} \in X} p_D(x, u, \tilde{x}) = 1$, $\forall (x, u) \in D$; $c_D : D \rightarrow \mathbb{R}^{>0}$ that $c_D(x, u)$ is the cost of performing action $u \in U(x)$ at state $x \in X$; and $x_0 \in X$ is the initial state. Denote by $\text{Post}(x, u) \triangleq \{\tilde{x} \in X \mid p_D(x, u, \tilde{x}) > 0\}$, $\forall (x, u) \in D$.

The above-mentioned MDP evolves by taking an action $u \in U(x)$ associated with every state $x \in X$. Denote by $R_T = x_0 u_0 x_1 u_1 \dots x_T u_T$ the past run that is a sequence of previous states and actions up to time $T \geq 0$. As defined in [2], a control policy $\mu = \mu_0 \mu_1 \dots$ is a sequence of decision rules μ_t at time $t \geq 0$. A control policy is stationary if $\mu_t = \mu$, $\forall t \geq 0$, where μ can be randomized so that $\mu : X \times U \rightarrow [0, 1]$ or deterministic so that $\mu : X \rightarrow U$, $\forall t \geq 0$. On the other hand, a policy is history dependent or finite memory if $\mu_t : R_t \times U \rightarrow [0, 1]$, where R_t is the past history until time $t \geq 0$.

B. End Components

A *sub-MDP* of \mathcal{M} is a pair (S, A) where $S \subseteq X$ and $A : S \rightarrow 2^U$ such that (i) $S \neq \emptyset$, $\emptyset \neq A(s) \subseteq U(s)$, $\forall s \in S$;

(ii) $\text{Post}(s, u) \subseteq S$, $\forall s \in S$ and $\forall u \in A(s)$. An *end component* (EC) of \mathcal{M} is a sub-MDP (S, A) such that the digraph $G_{(S,A)}$ induced by (S, A) is strongly connected. An EC (S, A) is called *maximal* if there is no other EC (S', A') such that $(S, A) \neq (S', A')$, $S \subseteq S'$, and $A(s) \subseteq A'(s)$, $\forall s \in S$. The set of maximal end components (MECs) of an MDP is finite and can be uniquely determined. The analysis of MECs would include each EC as a special case. We refer the readers to [4, Definitions 10.116, 10.117, and 10.124] for details. Moreover, an *accepting MEC* (AMEC) is an EC that satisfies certain accepting conditions such as the Streett and Robin conditions, which will be defined in the sequel. On the other hand, a *strongly connected component* (SCC) of the digraph $G_{\mathcal{M}}$ induced by \mathcal{M} is a set of states $S \subseteq X$, so that there exists a path in each direction between any pair of states in S . Similarly, an *accepting SCC* (ASCC) is an SCC that satisfies certain accepting conditions. Note that the main difference between an MEC (S, A) and an SCC S is that the SCC does not restrict the set of actions $U(s)$ that can be taken at each state $s \in S$. In other words, there might be paths that start from any state within the SCC and end at states outside the SCC.

C. LTL and Deterministic Rabin Automaton (DRA)

The ingredients of an LTL formula are a set of atomic propositions AP and several Boolean and temporal operators. Atomic propositions are Boolean variables that can be either true or false. An LTL formula is specified according to the syntax [4]: $\varphi \triangleq \top \mid p \mid \varphi_1 \wedge \varphi_2 \mid \neg\varphi \mid \bigcirc\varphi \mid \varphi_1 \mathbf{U}\varphi_2$, where $\top \triangleq \text{True}$, $p \in \text{AP}$, \bigcirc (next), \mathbf{U} (until), and $\perp \triangleq \neg\top$. For brevity, we omit the derivations of other operators like \square (always), \diamond (eventually), and \Rightarrow (implication). The semantics of an LTL is defined over the set of infinite words over 2^{AP} . Intuitively, $p \in \text{AP}$ is satisfied on a word $w = w(1)w(2) \dots$ if it holds at $w(1)$, i.e., if $p \in w(1)$. Formula $\bigcirc\varphi$ holds true if φ is satisfied on the word suffix that begins in the next position $w(2)$, whereas $\varphi_1 \mathbf{U}\varphi_2$ states that φ_1 has to be true until φ_2 becomes true. Finally, $\diamond\varphi$ and $\square\varphi$ are true if φ holds on w eventually and always, respectively. We refer the readers to [4, Ch. 5] for the full definition.

The set of words that satisfy an LTL formula φ over AP can be captured through a DRA \mathcal{A}_φ [4], defined as $\mathcal{A}_\varphi = (Q, 2^{\text{AP}}, \delta, q_0, \text{Acc}_\mathcal{A})$, where Q is a set of states; 2^{AP} is the alphabet; $\delta \subseteq Q \times 2^{\text{AP}} \times Q$ is a transition relation; $q_0 \in Q$ is the initial state; and $\text{Acc}_\mathcal{A} \subseteq 2^Q \times 2^Q$ is a set of accepting pairs, i.e., $\text{Acc}_\mathcal{A} = \{(H_\mathcal{A}^1, I_\mathcal{A}^1), (H_\mathcal{A}^2, I_\mathcal{A}^2), \dots, (H_\mathcal{A}^N, I_\mathcal{A}^N)\}$ where $H_\mathcal{A}^i, I_\mathcal{A}^i \subseteq Q$, $\forall i = 1, 2, \dots, N$. An infinite run $q_0q_1q_2 \dots$ of \mathcal{A} is *accepting* if there exists *at least one* pair $(H_\mathcal{A}^i, I_\mathcal{A}^i) \in \text{Acc}_\mathcal{A}$ such that $\exists n \geq 0$, it holds $\forall m \geq n$, $q_m \notin H_\mathcal{A}^i$, and $\exists \infty n \geq 0$, $q_n \in I_\mathcal{A}^i$, where $\exists \infty$ stands for “existing infinitely many.” Namely, this run should intersect with $H_\mathcal{A}^i$ *finitely* many times, whereas with $I_\mathcal{A}^i$ *infinitely* many times. There are translation tools [31] to obtain \mathcal{A}_φ given φ , which requires the process of translating first the LTL formula to the associated nondeterministic Büchi automaton, and then to the DRA with complexity $2^{2^{\mathcal{O}(n \log n)}}$, where n is the length of φ . Our implementation of the Python interface for [31] can be found in [32]. Note that [31] allows for different levels of automata simplifications to be made regarding the size of \mathcal{A}_φ , and a simplified automation may result in loss of optimality.

III. PROBLEM FORMULATION

A. Mathematical Model

In order to model uncertainty in both the robot motion and the workspace properties, we extend the definition of an MDP from Section II-A to include probabilistic labels, as the *probabilistically labeled MDP*

$$\mathcal{M} = (X, U, D, p_D, (x_0, l_0), \text{AP}, L, p_L, c_D) \quad (1)$$

where AP is a set of atomic propositions that capture the properties of interest in the workspace; $L : X \rightarrow 2^{2^{\text{AP}}}$ contains the set of property subsets that can be true at each state; and $p_L : X \times 2^{\text{AP}} \rightarrow [0, 1]$ specifies the associated probability. Particularly, $p_L(x, l)$ denotes the probability that state $x \in X$ satisfies the set of propositions $l \subset \text{AP}$. Note that $\sum_{l \in L(x)} p_L(x, l) = 1$, $\forall x \in X$. Moreover, (x_0, l_0) contains the initial state $x_0 \in X$ and the initial label $l_0 \in L(x_0)$, whereas the rest of the notations in (1) are the same as defined in Section II-A. The probabilistic labeling function provides a way to consider time-varying and dynamic workspace properties. Moreover, there is an LTL task formula φ specified over the same set of atomic propositions AP, as the desired behavior of \mathcal{M} . We assume that the MDP \mathcal{M} in (1) is *fully observable* due to the following assumption.

Assumption 1: At any stage $t \geq 0$, the current robot state $x_t \in X$ and its label $l_t \in L(x_t)$ are fully observable. ■

While the robot is moving within the workspace, it is capable of sensing an actual property and determine the label of the state it is located at. At stage $T \geq 0$, the robot’s past path is given by $X_T = x_0x_1 \dots x_T \in X^{(T+1)}$, the past sequence of observed labels is given by $L_T = l_0l_1 \dots l_T \in (2^{\text{AP}})^{(T+1)}$ and the past sequence of control actions is $U_T = u_0u_1 \dots u_T \in U^{(T+1)}$. It holds that $p_D(x_t, u_t, x_{t+1}) > 0$ and $p_L(x_t, l_t) > 0$, $\forall t \geq 0$. These three sequences can be composed into the complete past run $R_T = x_0l_0u_0x_1l_1u_1 \dots x_Tl_Tu_T$. Denote by \mathbf{X}_T , \mathbf{L}_T , and \mathbf{R}_T the set of all possible past sequences of states, labels, and runs up to stage T . We set $T = \infty$ for infinite sequences.

Definition 1: The mean total cost [2], [33] of an infinite robot run R_∞ of \mathcal{M} is defined as

$$\text{Cost}(R_\infty) = \liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{t=0}^n c_D(x_t, u_t) \quad (2)$$

where $R_\infty = x_0l_0u_0x_1l_1u_1 \dots \in \mathbf{R}_\infty$. ■

As discussed in [2], [20], [24], and [33], the above-mentioned mean total cost is called the *mean-payoff* function (or limit average), where the “lim” operator is needed as the limit-average might not exist for some runs (see [24], [33], and [34]).

Our goal is to find a finite-memory policy for \mathcal{M} , denoted by $\mu = \mu_0\mu_1 \dots$. The control policy at stage $t \geq 0$ is given by $\mu_t : \mathbf{R}_t \times U \rightarrow [0, 1]$, where \mathbf{R}_t is the past run R_t , $\forall t \geq 0$. Denote by $\overline{\mu}$ the set of all such policies. Given a control policy $\mu \in \overline{\mu}$, the probability measure $\text{Pr}_\mu^{\mathcal{M}}(\cdot)$ on the smallest σ -algebra, over all possible infinite sequences \mathbf{R}_∞ that contain R_T , is the unique measure [4] by

$$\text{Pr}_\mu^{\mathcal{M}}(\mathbf{R}_\infty) = \prod_{t=0}^T p_D(x_t, u_t, x_{t+1}) \cdot p_L(x_t, l_t) \cdot \mu_t(\mathbf{R}_t, u_t) \quad (3)$$

where $\mu(\mathbf{R}_t, u_t)$ is defined as the probability of choosing action u_t given the past run \mathbf{R}_t . Then, we define the probability of \mathcal{M} satisfying φ under policy μ by

$$\Pr_{\mathcal{M}}^{\mu}(\varphi) = \Pr_{\mathcal{M}}^{\mu}\{\mathbf{R}_{\infty} \models \varphi\}$$

where the satisfaction relation “ \models ” is introduced in Section II-C, given an infinite word and an LTL formula. Accordingly, the *risk* is defined as the probability that the task formula φ is not satisfied by \mathcal{M} under the policy μ , namely $\text{Risk}_{\mathcal{M}}^{\mu}(\varphi) = 1 - \Pr_{\mathcal{M}}^{\mu}(\varphi)$.

Problem 1: Given the labeled MDP \mathcal{M} defined in (1) and the task specification φ , our goal is to solve

$$\begin{aligned} \min_{\mu \in \bar{\mu}} \mathbb{E}_{\mathcal{M}}^{\mu}\{\text{Cost}(R_{\infty})\} \\ \text{s.t. } \text{Risk}_{\mathcal{M}}^{\mu}(\varphi) \leq \gamma \end{aligned} \quad (4)$$

where $\gamma \geq 0$ is a predefined parameter as the allowed risk; the optimal policy minimizes the mean total cost and ensures that the risk of violating φ remains bounded by γ . ■

Note that the traditional definition of undiscounted expected total cost over an infinite run from [2] and [14] is not used here, as it is infinite except for the special case of transient MDPs defined in Section II-A. However, in this paper, the model \mathcal{M} is not restricted to be transient. Moreover, the discounted total cost in [2] is not used here either due to two reasons: first, it is not obvious how to choose the discount factor for various control tasks φ [25]; and second, we are more interested in optimizing the repetitive long-term behavior of the system, rather than the short-term one [20]. In-depth discussions on the optimization of infinite-horizon undiscounted or discounted total-cost criteria over MDPs with or without constraints can be found in [2].

Remark 1: Different from the maximal reachability problem addressed in [4] and [21], a deterministic policy would not suffice here. Instead, randomization is required due to the mean total-cost criterion and the risk constraint, similar to [14]. ■

IV. SOLUTION

This section contains the three major parts of the proposed solution, which are as follows:

- 1) the construction of the product automaton and its AMECs;
- 2) the algorithms to synthesize the optimal plan prefix and suffix, for both cases where the AMECs exist or not; and
- 3) the complete policy and the online execution algorithm.

A. Product Automaton and AMECs

To begin with, we construct the DRA \mathcal{A}_{φ} associated with the LTL task formula φ via the translation tools [31], [32]. Let it be $\mathcal{A}_{\varphi} = (Q, 2^{AP}, \delta, q_0, \text{Acc}_{\mathcal{A}})$, where the notations are defined in Section II-C. Then, we construct a product automaton between the robot model \mathcal{M} and the DRA \mathcal{A}_{φ} .

Definition 2: Denote by \mathcal{P} the product $\mathcal{M} \times \mathcal{A}_{\varphi}$ as a 7-tuple

$$\mathcal{P} = (S, U, E, p_E, c_E, s_0, \text{Acc}_{\mathcal{P}}) \quad (5)$$

where the state $S \subseteq X \times 2^{AP} \times Q$ is so that $\langle x, l, q \rangle \in S$, $\forall x \in X, \forall l \in L(x)$, and $\forall q \in Q$; the action set U is the same as

in (1) and $U(s) = U(x)$, $\forall s = \langle x, l, q \rangle \in S$; $E = \{(s, u) \mid s \in S, u \in U(s)\}$; the transition probability $p_E : S \times U \times S \rightarrow [0, 1]$ is so that

$$p_E(\langle x, l, q \rangle, u, \langle \check{x}, \check{l}, \check{q} \rangle) = p_D(x, u, \check{x}) \cdot p_L(\check{x}, \check{l}) \quad (6)$$

where (i) $\langle x, l, q \rangle, \langle \check{x}, \check{l}, \check{q} \rangle \in S$; (ii) $(x, u) \in D$; and (iii) $\check{q} \in \delta(q, l)$; the cost function $c_E : E \rightarrow \mathbb{R}^{>0}$ is so that $c_E(\langle x, l, q \rangle, u) = c_D(x, u)$, $\forall (\langle x, l, q \rangle, u) \in E$. Namely, the label l should fulfill the transition condition from q to \check{q} in \mathcal{A}_{φ} ; the single initial state is $s_0 = \langle x_0, l_0, q_0 \rangle \in S$; the accepting pairs are defined as $\text{Acc}_{\mathcal{P}} = \{(H_p^i, I_p^i), i = 1, 2, \dots, N\}$, where $H_p^i = \{\langle x, l, q \rangle \in S \mid q \in H_{\mathcal{A}}^i\}$ and $I_p^i = \{\langle x, l, q \rangle \in S \mid q \in I_{\mathcal{A}}^i\}$, $\forall i = 1, 2, \dots, N$. ■

The product \mathcal{P} computes the intersection between all traces of \mathcal{M} and all words that are accepted to \mathcal{A}_{φ} , to find all admissible robot behaviors that satisfy the task φ . It combines the uncertainty in robot motion and the workspace model by including both x and l in the states. The Rabin accepting condition of \mathcal{P} is defined as follows: An infinite path $R_{\mathcal{P}} = s_0 s_1 \dots$ of \mathcal{P} is accepting if for at least one pair $(H_p^i, I_p^i) \in \text{Acc}_{\mathcal{P}}$, it holds that $R_{\mathcal{P}}$ intersects with H_p^i finitely often, whereas with I_p^i infinitely often. To transform this condition into equivalent graph properties, we need to compute the AMECs of \mathcal{P} associated with its accepting pairs $\text{Acc}_{\mathcal{P}}$. Detailed definition of MECs is given in Section II-B.

In order to find the complete set of AMECs of \mathcal{P} , for each pair $(H_p^i, I_p^i) \in \text{Acc}_{\mathcal{P}}$, perform the following steps:

(i) Build the MDP $\mathcal{Z}_i^{-H} \triangleq (S', U', E', p'_E)$, where $S' = S_i^{-H} \cup \{\nu\}$ is the set of states with $S_i^{-H} = S \setminus H_p^i$ and ν a trap state; $U' = U \cup \{\tau_0\}$ is the set of actions where τ_0 is a pseudoaction; $E' \subset S' \times U$ is the set of transitions with the associated probability p'_E that are defined by three cases: (a) for the transitions within S_i^{-H} , it holds that $(s, u) \in E'$ and $p'_E(s, u, \check{s}) = p_E(s, u, \check{s})$, $\forall (s, u) \in E$, where $s, \check{s} \in S_i^{-H}$; (b) for the transitions from S_i^{-H} to outside S_i^{-H} , it holds that $(s, u) \in E'$ and $p'_E(s, u, \nu) = \sum_{\check{s} \notin S_i^{-H}} p_E(s, u, \check{s})$, $\forall (s, u) \in E$, where $s \in S_i^{-H}$; and (c) the trap state is included in a self-loop such that $(\nu, \tau_0) \in E'$ and $p'_E(\nu, \tau_0, \nu) = 1$. Simply speaking, all transitions from inside S_i^{-H} to outside S_i^{-H} are transformed to transitions to the trap state ν .

(ii) Determine all MECs of \mathcal{Z}_i^{-H} above via [4, Algorithm 47], which is based on splitting the SCCs of \mathcal{Z}_i^{-H} until the conditions of being an EC are fulfilled. Our implementation for this algorithm can be found in [32]. Denote by $\Xi^i = \{(S'_1, U'_1), (S'_2, U'_2), \dots, (S'_{C_i}, U'_{C_i})\}$, the set of MECs, where $S'_c \subset S'$ and $U'_c : S'_c \rightarrow 2^{U'}$, $\forall c = 1, 2, \dots, C_i$. Note that $S'_c \cap S'_{c'} = \emptyset$, $\forall (S'_{c'}, U'_{c'}) \in \Xi^i$.

(iii) Find $(S'_c, U'_c) \in \Xi^i$ that is *accepting*, i.e., it satisfies $\nu \notin S'_c$ and $S'_c \cap I_p^i \neq \emptyset$. Save the AMECs in Ξ_{acc}^i . Since Ξ_{acc}^i is computed for each $(H_p^i, I_p^i) \in \text{Acc}_{\mathcal{P}}$, we denote by $\Xi_{\text{acc}} = \{\Xi_{\text{acc}}^i, i = 1, \dots, N\}$, the complete set of AMECs of \mathcal{P} .

Remark 2: A single state with a self-transition can be an MEC with a proper action set. Therefore, there exist at most $|S'|$ MECs within \mathcal{Z}_i^{-H} , $\forall i = 1, \dots, N$. Thus, Step (ii) above has complexity $\mathcal{O}(|S'|^2)$, as shown in [4, Lemma 10.126], whereas Steps (i) and (iii) have complexity linear with $|S'|$. ■

B. Plan Prefix and Suffix Synthesis

Given the complete set of AMECs Ξ_{acc} of \mathcal{P} , in this section, we show how to synthesize the control policy to drive the system towards Ξ_{acc} and furthermore remain inside Ξ_{acc} while satisfying the accepting condition. As mentioned in Section I, most related work [4], [16], [17], [21] focuses on maximizing the probability of reaching the union of AMECs, i.e., $\cup_{(S'_c, U'_c) \in \Xi_{\text{acc}}} S'_c$, where dynamic programming techniques, such as value or policy iteration, can be applied to obtain the optimal policy. Furthermore, once the system enters any AMEC, e.g., $(S'_c, U'_c) \in \Xi_{\text{acc}}$, it has probability 1 of staying within S'_c by following U'_c (see [4, Lemma 10.119]). The Round-Robin policy is adopted in [4], [17], and [21] that ensures all states in S'_c (including its nonempty intersection with $I_{\mathcal{P}}^i$) are visited infinitely often. As a result, the task φ is satisfied by \mathcal{P} under this policy with the maximal probability.

The above-mentioned solutions may suffice for verification problems that do not optimize cost or for tasks with trivial accepting conditions. However, for the purposes of plan synthesis and for general tasks, it is of practical interest to simultaneously satisfy the probability of reaching *all* the AMECs as well as optimize the mean cost of staying within *any* AMEC and fulfilling the accepting condition. Moreover, when no AECs can be found, instead of simply reporting failure, it is important to obtain a relaxed policy that guarantees high probability of satisfying the task over long time intervals, thus minimizing the frequency of encountering bad events. In what follows, we present a policy synthesis algorithm that consists of four parts.

- 1) The *plan prefix* that drives the system from the initial state to all AMECs while minimizing the expected cost and respecting the risk constraint (see Section IV-B1).
- 2) The *plan suffix* that keeps the system within the AMEC it has reached while satisfying the accepting condition and optimizing the expected suffix cost (see Section IV-B2).
- 3) The *relaxed prefix and suffix plans* for the case where no AECs of \mathcal{P} can be found (see Section IV-B3).
- 4) The complete finite-memory policy for the original MDP \mathcal{M} (see Section IV-C1).

Before stating the solution, we introduce a partition of S , given the initial state s_0 and the set of AMECs Ξ_{acc} . Let $S_r \subseteq S$ be the set of states within S that can be reached from s_0 , which can be derived via a simple graph search in \mathcal{P} .

Definition 3: Given s_0 and Ξ_{acc} , S is partitioned as $S = S_o \cup S_c \cup S_d \cup S_n$, where $S_o \triangleq S \setminus S_r$ is the set of states that can *not* be reached from s_0 ; S_c is the union of all goal states in Ξ_{acc} , i.e., $S_c \triangleq \cup_{(S'_c, U'_c) \in \Xi_{\text{acc}}} S'_c$; $S_d \subseteq S_r$ can be reached from s_0 but cannot reach any state in S_c ; and $S_n \triangleq S_r \setminus (S_c \cup S_d)$. ■

The set S_d can be derived through a simple graph search, e.g., by reversing the directed graph associated with \mathcal{P} , finding all reachable nodes of any state within each $(S'_c, U'_c) \in \Xi_{\text{acc}}$ (as any AMEC is strongly connected), and finally computing its cross intersection with S_r (see [32] for implementation details). Roughly speaking, S_n is the set of states related to the plan prefix, S_c is the set of goal states related to the plan suffix, and S_d is set of bad states to be avoided during the prefix. Since S_o contains the states that cannot be reached from s_0 , it is neglected hereafter for the purpose of plan synthesis.

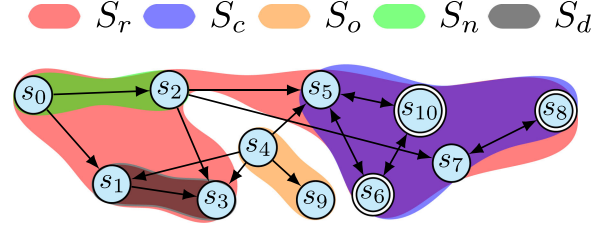


Fig. 1. Illustration of the partition of S in Definition 3, where S_r , S_c , S_o , S_n , and S_d are highlighted by red, blue, orange, green, and black areas, respectively. Details can be found in Example 1.

Example 1: This example illustrates the partition in Definition 3. Consider the toy product automaton \mathcal{P} in Fig. 1. For state s_0 , the set of reachable states is $S_r = \{s_0, s_1, s_2, s_3, s_4, s_5, s_6, s_7, s_8, s_9, s_{10}\}$, the set of unreachable states is $S_o = \{s_4, s_9\}$, the states within an AMEC are $S'_{c_1} = \{s_5, s_6, s_{10}\}$, and another AMEC $S'_{c_2} = \{s_7, s_8\}$, thus $S_c = S'_{c_1} \cup S'_{c_2} = \{s_5, s_6, s_7, s_8, s_{10}\}$, the states that can be reached from s_0 but cannot reach S_c are $S_d = \{s_1, s_3\}$, and the states that s_0 can reach outside $S_c \cup S_d$ are $S_n = \{s_0, s_2\}$. ■

1) Plan Prefix: Similar to [17] and [18], we first construct a modified sub-MDP \mathcal{Z}_{pre} of \mathcal{P} as $\mathcal{Z}_{\text{pre}} \triangleq (S_p, U_p, E_p, s_0, p_p, c_p)$, where the set of states is given by $S_p = S_n \cup S_c$ with S_n, S_c being defined in Definition 3. The set of actions is given by $U_p = U \cup \{\tau_0\}$ where τ_0 is a self-loop action. The set of transitions E_p is the subset of E associated with S_p . Moreover, the transition probability p_p is defined by (i) $p_p(s, u, \tilde{s}) = p_E(s, u, \tilde{s})$, $\forall s, \tilde{s} \in S_p$ where $s \notin S_c$ and $\forall u \in U(s)$; and (ii) $p_p(s, \tau_0, s) = 1$, $\forall s \in S_c$. Finally, the cost function c_p is defined by (i) $c_p(s, u) = c_E(s, u)$, $\forall s \in S_n$ and $\forall u \in U(s)$; and (ii) $c_p(s, \tau_0) = 0$, $\forall s \in S_c$.

Then, we find a policy for \mathcal{Z}_{pre} such that, starting from s_0 , it can reach the set of goal states S_c with a probability larger than $1 - \gamma$, while at the same time, minimizing the expected total cost. Formally, consider the problem as follows.

Problem 2: Given the sub-MDP \mathcal{Z}_{pre} , compute an optimal stationary prefix policy $\pi_{\text{pre}}^* \in \bar{\pi}$ that solves the problem

$$\begin{aligned} \min_{\pi \in \bar{\pi}} \left[C_{\text{pre}}(S_c) \triangleq \mathbb{E}_{\mathcal{Z}_{\text{pre}}}^{\pi} \left\{ \sum_{t=0}^{\infty} c_p(s_t, u_t) \right\} \right] \\ \text{s.t. } \Pr_{s_0}^{\pi}(\diamond S_c) \geq 1 - \gamma \end{aligned} \quad (7)$$

where $s_0 u_0 s_1 u_1 \dots$ is a run of \mathcal{Z}_{pre} , $\bar{\pi}$ is the set of all stationary policies, the objective function is the expected total cost, $\Pr_{s_0}^{\pi}(\diamond S_c)$ is the probability of reaching S_c from the initial state s_0 , under the policy π ; and $\gamma > 0$ is from (4). ■

Note that the objective function in (7) is well defined and finite due to the fact that \mathcal{Z}_{pre} is transient with respect to S_n , and is equal to the expected total cost of reaching S_c since the cost of staying within S_c is zero. We omit the proof that \mathcal{Z}_{pre} is transient here and refer the interested readers to [2] and [14]. Our proposed solution to Problem 2 is based on transforming it into a constrained optimization problems for MDPs, which can be then solved using linear programming (LP). The approach is inspired by [14], [16], and [17]. Particularly, denote by $y_{s,u}$, the *expected number of times* over the infinite horizon that the sys-

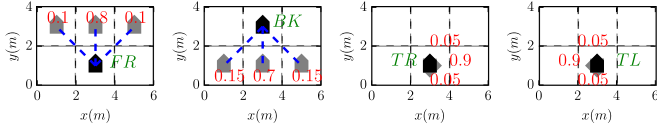


Fig. 2. Uncertainty of each action primitive, see Section V for details. Possible post states are in gray from the starting state in black, where the associated possibilities are marked in red.

tem is at state s and action u is taken, $\forall s \in S_n$ and $\forall u \in U(s)$, which are often referred to as occupancy measures [14] as it holds $y_{s,u} = \sum_{t=0}^{\infty} Pr_{s_0}^{\pi} [s_t = s, u_t = u]$, where the probability is conditioned on a policy π and the initial state s_0 . Note that an occupancy measure is a sum of probabilities, but not a probability itself. Consider the linear program as follows:

$$\min_{\{y_{s,u}\}} \left[C_{\text{pre}}(S_c) \triangleq \sum_{(s,u)} \sum_{\tilde{s} \in S_p} y_{s,u} p_p(s, u, \tilde{s}) c_p(s, u) \right] \quad (8a)$$

$$\text{s.t.} \quad \sum_{(s,u)} \sum_{\tilde{s} \in S_c} y_{s,u} p_p(s, u, \tilde{s}) \geq 1 - \gamma \quad (8b)$$

$$\sum_{u \in U(\tilde{s})} y_{\tilde{s},u} = \sum_{(s,u)} y_{s,u} p_p(s, u, \tilde{s}) + \mathbb{1}(\tilde{s} = s_0) \quad \forall \tilde{s} \in S_n \quad (8c)$$

$$y_{s,u} \geq 0 \quad \forall s \in S_n \quad \forall u \in U(s) \quad (8d)$$

where $\sum_{(s,u)} \triangleq \sum_{s \in S_n} \sum_{u \in U(s)}$, the indicator function $\mathbb{1}(\tilde{s} = s_0) = 1$ if $\tilde{s} = s_0$, and $\mathbb{1}(\tilde{s} = s_0) = 0$, otherwise. Denote by $C_{\text{pre}}(S_c)$, the objective function associated with S_c . Let the solution of (8) be $y_{\text{pre}}^* = \{y_{s,u}^*, s \in S_n, u \in U(s)\}$. Then, the optimal *stationary* policy for the plan prefix, denoted by π_{pre}^* , can be derived as follows: the probability of choosing action u at state s equals to $\pi_{\text{pre}}^*(s, u) = y_{s,u}^* / (\sum_{u \in U(s)} y_{s,u}^*)$ if $\sum_{u \in U(s)} y_{s,u}^* \neq 0$; otherwise, the action at s can be chosen randomly, $\forall s \in S_c$.

Lemma 1: Given an optimal solution y_{pre}^* of (8), the associated policy π_{pre}^* ensures that $Pr_{s_0}^{\pi_{\text{pre}}^*}(\diamond S_c) \geq 1 - \gamma$.

Proof: First, $y_{s,u}$ is finite and well defined since Z_{pre} is transient with respect to S_n . The second part of the proof is similar to [16, Lemma 3.3]. The summation $\sum_{(s,u)} \sum_{\tilde{s} \in S_c} y_{s,u} p_p(s, u, \tilde{s})$ is the expected number of times that Z_{pre} transitions from any state in S_n into S_c for the *first time*, under policy π_{pre}^* from the initial state s_0 . Since the system remains within S_c once it enters S_c , the summation equals the probability of eventually reaching the set S_c , which is lower bounded by $1 - \gamma$. ■

Example 2: This example illustrates the important role of γ in the tradeoff between reducing the expected total cost and minimizing the risk in Problem 2. Consider the unicycle robot with action primitives illustrated in Fig. 2 and defined in Section V. The robot moves within partitioned cells, as shown in Fig. 3, where the red cell has probability 0.9 to be occupied by an obstacle. Consider the task: $\varphi = (\diamond \square b) \wedge (\square \neg \text{obs})$, i.e., to reach the yellow base without crossing any obstacle. In what follows, we solve (8) under risk factors $\gamma = 0$ and $\gamma = 0.4$ to derive two different optimal policies. Fig. 3 shows a shorter trajectory with lower expected total cost of about 12.6 when a

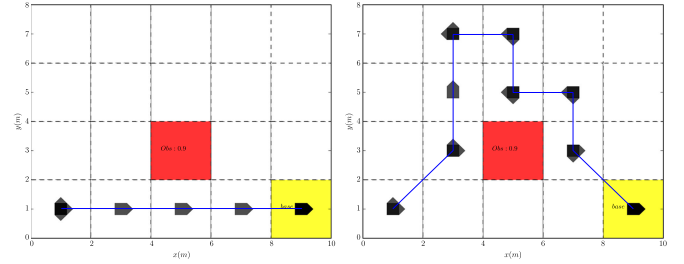


Fig. 3. Trajectories when setting $\gamma = 0.4$ (left) and $\gamma = 0$ (right). The task is to reach the yellow base while avoiding the red cell.

larger risk is allowed, compared with the right trajectory that avoids completely colliding with the obstacle, but with a much higher total cost of about 33.7. ■

2) Plan Suffix With AMECs: In this section, we present an algorithm to synthesize the *plan suffix* that minimizes the mean total cost within the AMECs while ensuring that the system trajectory satisfies the accepting condition of \mathcal{P} . Note that the plan prefix π_{pre}^* from the previous section guarantees that the system enters S_c from s_0 with probability higher than $1 - \gamma$. Recall also that $S_c = \cup_{(S'_c, U'_c) \in \Xi_{\text{acc}}} S'_c$. Thus, it is *possible* that the system enters *any* set S'_c within Ξ_{acc} . For this reason, we propose to treat each AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$ *separately*, as each S'_c is associated with different U'_c , and thus a different accepting condition for $S'_c \cap I_p^i$. Specifically, consider any AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$ and let $I'_c \triangleq S'_c \cap I_p^i$, which is nonempty by the definition of an AMEC.

Once the system enters any AMEC, most related work [4], [17], [21] adopts the Round-Robin policy defined as follows.

Definition 4: For each state $s_t \in S'_c$, create *any* ordered sequence of actions from $U'_c(s_t)$, denoted by $\bar{U}(s_t)$ and its infinite repetition by $\bar{U}^\omega(s_t)$. Then, at any stage $t > 0$, whenever the system reaches $s_t \in S'_c$, the *Round-Robin policy* instructs the system to take the *next* action in $\bar{U}^\omega(s_t)$, starting from the first action in $\bar{U}^\omega(s_t)$. ■

Namely, once the system enters S'_c , the Round-Robin policy iterates over the allowed actions for each state, which in turn ensures that all states in S'_c (which include I'_c) are visited infinitely often. Detailed can be found in [4, Lemma 10.119].

Definition 5: An *accepting cyclic path* of \mathcal{P} , associated with S'_c and I'_c , is a finite path that starts from any state $s_f \in I'_c$ and ends in any state $s_g \in I'_c$ while remaining within S'_c . ■

Note that an accepting cyclic path does not necessarily start and end at the same state in I'_c . Furthermore, we can define the mean cyclic cost of \mathcal{P} under a stationary policy.

Definition 6: The total cost of a cyclic path $P_a = s_0 u_0 s_1 u_1 \dots s_{N_a} u_{N_a}$ is defined as

$$\bar{C}_{\text{suf}}(P_a) \triangleq \sum_{t=0}^{N_a} c_D(s_t, u_t) \quad (9)$$

where $N_a \geq 1$ is the length of the path and $s_0, s_{N_a} \in I'_c$. Then, its mean total cost is defined as $C_{\text{suf}}(P_a) \triangleq \frac{1}{N_a} \bar{C}_{\text{suf}}(P_a)$. ■

Problem 3: Find a *stationary suffix policy* π_{suf}^* for \mathcal{P} within S'_c that minimizes the *mean cyclic cost*

$$C_{\text{suf}}(S'_c, U'_c) = \mathbb{E}_{P_a \in \mathcal{P}_a}^{\pi_{\text{suf}}^*} \{C_{\text{suf}}(P_a)\} \quad (10)$$

where \mathcal{P}_a is the set of all accepting cyclic paths associated with the AMEC (S'_c, U'_c) . ■

Inspired by [20], [24], and [33], we formulate a linear program to solve the mean-payoff optimization problem. First, we construct a modified sub-MDP \mathcal{Z}_{suf} of \mathcal{P} over S'_c by splitting I'_c into two virtual copies: I_{in} that only has incoming transitions into I'_c and I_{out} that has only outgoing transitions from I'_c . Formally, we define $\mathcal{Z}_{\text{suf}} \triangleq (S_e, U_e, E_e, y_0, p_e, c_e)$, where the set of states is $S_e = (S'_c \setminus I'_c) \cup I_{\text{in}} \cup I_{\text{out}}$ with $I_{\text{in}} = \{s_f^{\text{in}}, \forall s_f \in I'_c\}$ and $I_{\text{out}} = \{s_f^{\text{out}}, \forall s_f \in I'_c\}$, the virtual copies of I'_c . The set of control actions is $U_e = U \cup \{\tau_0\}$, where τ_0 is a self-loop action. The set of state-action pairs $E_e \subset S_e \times U_e$ is defined by (i) $(s, u) \in E_e, \forall s \in S'_c \setminus I'_c$ and $u \in U'_c(s)$; (ii) $(s, \tau_0) \in E_e, \forall s \in I_{\text{in}}$; and (iii) $(s_f^{\text{out}}, u) \in E_e, \forall s_f \in I'_c$ and $u \in U'_c(s_f)$. Moreover, y_0 is the initial distribution of all states in S'_c that can be reached by taking a transition from states in S'_n defined by

$$y_0(s) = \sum_{\check{s} \in S'_n} \sum_{u \in U_p(\check{s})} p_p(\check{s}, u, s) y_{\text{pre}}(\check{s}, u) \quad \forall s \in (S'_c \setminus I'_c) \cup I_{\text{out}}$$

where $\{y_{\text{pre}}(s, u)\}$ are the variables of (8). Furthermore, the transition probability p_e is defined in five cases as follows: (a) for transitions within $S'_c \setminus I'_c$, it holds that $p_e(s, u, \check{s}) = p_E(s, u, \check{s}), \forall s, \check{s} \in S'_c \setminus I'_c, \forall u \in U_e(s)$; (b) for transitions originated from I_{out} , it holds that $p_e(s_f^{\text{out}}, u, \check{s}) = p_E(s_f, u, \check{s}), \forall s_f^{\text{out}} \in I_{\text{out}}, \forall u \in U_e(s_f^{\text{out}})$, and $\forall \check{s} \in S'_c \setminus I'_c$; (c) for transitions into I_{in} , it holds that $p_e(s, u, s_f^{\text{in}}) = p_E(s, u, s_f), \forall s \in S'_c \setminus I'_c, \forall u \in U_e(s)$, and $\forall s_f^{\text{in}} \in I_{\text{in}}$; (d) for transitions from I_{out} to I_{in} , it holds that $p_e(s_f^{\text{out}}, u, s_f^{\text{in}}) = p_E(s_f, u, s_f), \forall s_f^{\text{out}} \in I_{\text{out}}$ and $\forall u \in U_e(s_f^{\text{out}})$; and (e) for transitions within I_{in} , $p_e(s_f^{\text{in}}, \tau_0, s_f^{\text{in}}) = 1, \forall s_f^{\text{in}} \in I_{\text{in}}$. Finally, the cost function satisfies $c_e(s, u) = c_E(s, u), \forall s \in (S_e \setminus I_{\text{in}}), \forall u \in U_e(s)$, and $c_e(s_f^{\text{in}}, \tau_0) = 0, \forall s_f^{\text{in}} \in I_{\text{in}}$.

Remark 3: The initial distribution y_0 of \mathcal{Z}_{suf} indicates how likely it is that the system controlled by the plan prefix π_{pre}^* will enter the AMEC (S'_c, U'_c) via each state inside S'_c . ■

Let also $S'_e \triangleq S_e \setminus I_{\text{in}}$ and denote by $z_{s,u}$ the *long-run frequency* with which the system is at state s and the action u is applied, $\forall s \in S'_e$ and $\forall u \in U_e(s)$. Then, we can formulate the following linear program to solve Problem 3:

$$\min_{\{z_{s,u}\}} \left[\mathbf{C}_{\text{suf}}(S'_c, U'_c) \triangleq \sum_{(s,u)} \sum_{\check{s} \in S_e} z_{s,u} p_e(s, u, \check{s}) c_e(s, u) \right] \quad (11a)$$

$$\text{s.t.} \quad \sum_{(s,u)} \sum_{\check{s} \in I_{\text{in}}} z_{s,u} p_e(s, u, \check{s}) = \sum_{s \in S'_e} y_0(s) \quad (11b)$$

$$\sum_{u \in U_e(s)} z_{s,u} = \sum_{(\check{s},u)} z_{\check{s},u} p_e(\check{s}, u, s) + y_0(s) \quad \forall s \in S'_e \quad (11c)$$

$$z_{s,u} \geq 0, \quad \forall s \in S'_e \quad \forall u \in U_e(s) \quad (11d)$$

where $\sum_{(s,u)} \triangleq \sum_{s \in S'_e} \sum_{u \in U_e(s)}$, the first constraint ensures that I_{in} is eventually reached, whereas the second constraint balances the incoming and outgoing flow at each state. Let its solution be $z_{\text{suf}}^* = \{z_{s,u}^*, \forall s \in S'_e, \forall u \in U_e(s)\}$. Then, the optimal *stationary* policy for the plan suffix, denoted by π_{suf}^* , can be derived as follows: the probability of choosing ac-

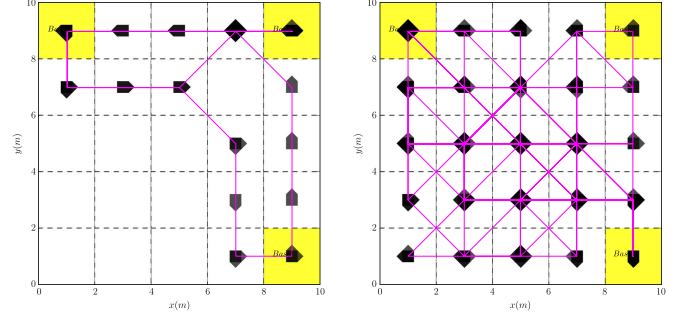


Fig. 4. Simulated trajectory under π_{suf}^* (left) and under the Round-Robin policy (right), see Example 3.

tion u at state s equals to $\pi_{\text{suf}}^*(s, u) = z_{s,u}^* / (\sum_{u \in U_e(s)} z_{s,u}^*)$ if $\sum_{u \in U_e(s)} z_{s,u}^* \neq 0$; otherwise the action at s is chosen randomly, $\forall s \in S'_e$. Note that $\pi_{\text{suf}}^*(s_f, u) = \pi_{\text{suf}}^*(s_f^{\text{out}}, u), \forall s_f \in I'_c$ and $\forall u \in U'_c(s_f)$. Namely, once the system reaches any state $s_g \in I'_c$, the control policy at s_g will be the control policy for $s_g^{\text{out}} \in I_{\text{out}}$, according to the solution of (11).

Remark 4: The initial distribution is derived from (8), instead of being arbitrarily set as in [25]; Moreover, (11b) ensures that only I'_c is intersected infinitely often, instead of enforcing that all states in the set S'_c are visited infinitely often as in [25]. ■

Lemma 2: If (11) has a solution, then the plan suffix π_{suf}^* solves Problem 3 for the chosen AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$.

Proof: First, by Definition 5, the objective in (11) equals the mean cyclic cost of all accepting cyclic paths for I'_c . Moreover, by the definition of an AMEC, any path remains within S'_e by choosing only actions within $U'_c(s)$ at each state $s \in S'_e$.

Lemma 3: Let $\tau_{\mathcal{P}}$ be the set of all accepting runs of \mathcal{P} that enter S'_c after a finite number of steps. If $\tau_{\mathcal{P}} \in \tau_{\mathcal{P}}$ is generated under π_{suf}^* , then $\tau_{\mathcal{P}}$ satisfies the accepting condition of \mathcal{P} . Moreover, the mean total cost in (2) equals the mean cyclic cost in (10), i.e., $\mathbb{E}_{\tau_{\mathcal{P}} \in \tau_{\mathcal{P}}} \{\text{Cost}(\tau_{\mathcal{P}})\} = \mathbf{C}_{\text{suf}}(S'_c, U'_c)$.

Proof: By (11), any system trajectory of \mathcal{P} under π_{suf}^* contains infinite occurrences of accepting cyclic paths. Since any accepting cyclic path starts from and ends in I'_c (which is finite), $\tau_{\mathcal{P}}$ intersects with I'_c infinitely often. Moreover, since any accepting cyclic path remains within S'_c , $\tau_{\mathcal{P}}$ remains within S'_c for all time after entering S'_c . In other words, $\tau_{\mathcal{P}}$ intersects with $H_{\mathcal{P}}^i$ a finite number of times before entering S'_c and then intersects $I_{\mathcal{P}}^i$ infinitely often after entering S'_c , which satisfies the Rabin accepting condition of \mathcal{P} . To show the second part, notice that the product \mathcal{P} under π_{suf}^* evolves as an MC and the set of all accepting cyclic paths within S'_c has a stationary distribution. By viewing any accepting run $\tau_{\mathcal{P}}$ as the *concatenation* of an infinite number of cyclic paths, the mean total cost of $\tau_{\mathcal{P}}$ defined in (4) over an infinite time horizon equals the mean cyclic cost in (10) of all cyclic paths contained in $\tau_{\mathcal{P}}$. This result is important in showing the equivalence between Problems 1 and 3 later in Theorem 6.

Example 3: This example illustrates the difference between the plan suffix obtained by (11) and the Round-Robin policy. Consider the same robot model from Example 2 and the partitioned workspace in Fig. 4. The task is to surveil three base stations in the corners, i.e., $\varphi = (\square \diamond b1) \wedge (\square \diamond b2) \wedge (\square \diamond b3)$. The plan prefix is derived by solving (8) but two different plan

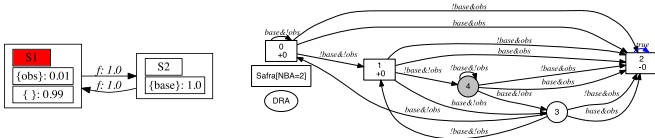


Fig. 5. MDP \mathcal{M} (left) and DRA \mathcal{A}_φ (right, derived via [31] and [32]) described in Example 4, with one accepting pair $(\{2\}, \{0, 1\})$.

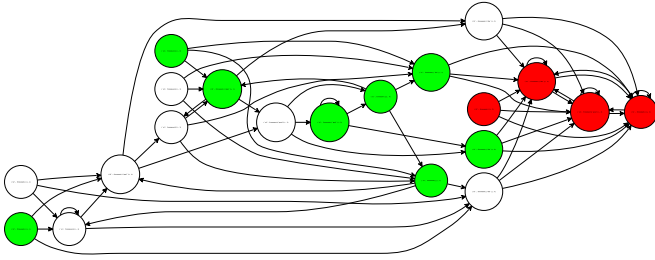


Fig. 6. Product \mathcal{P} of \mathcal{M} and \mathcal{A}_φ in Fig. 5. The state and edge names are omitted as the structure is of importance here. At least one green state should be visited infinitely often while avoiding all red states. Note all transitions are driven by the action f .

suffixes are used: one using (11) and the Round-Robin policy. Fig. 4 shows the simulated trajectory under these two policies. It can be seen that the trajectory under the optimal plan suffix approximates the shortest route to cross all base stations, whereas the trajectory under the Round-Robin policy exhibits a rather random behavior. ■

3) Plan Synthesis When AECs Do Not Exist: The synthesis algorithms proposed in Sections IV-B1 and IV-B2 rely on the assumption that the set of AMECs Ξ_{acc} of \mathcal{P} is nonempty which, however, might not hold in many scenarios. In this case, most existing techniques proposed in [4], [17], [21], and [22] cannot be applied. In this section, we first provide a simple example where no AECs exist, and then propose an approach to synthesize a *relaxed* plan prefix and suffix.

Example 4: This example provides a robot model \mathcal{M} and its task φ for which no AECs exist in the product automaton \mathcal{P} . Consider the MDP \mathcal{M} in Fig. 5 that transitions between two states (S_1, S_2) with probability 1 using the action f . Note that S_1 has only probability 0.01 of being occupied by an obstacle and S_2 is the base station. The task is to surveil the base station while avoiding obstacles, i.e., $\varphi = (\square \diamond b) \wedge (\square \neg \text{obs})$. The associated DRA is shown in Fig. 5. The resulting \mathcal{P} is shown in Fig. 6, where the set of states $H_i^{\mathcal{P}}$ to avoid in the suffix is in red and the set of states $I_i^{\mathcal{P}}$ to intersect infinitely often is in green. The reason that no AECs exist in \mathcal{P} is because by definition, an AEC $(S', \{f\})$ should include *all* successor states that are reachable by the single action f . Then, starting from any green state in $I_i^{\mathcal{P}}$, the set of reachable states eventually intersect with the red states in $H_i^{\mathcal{P}}$. ■

When no AECs exist in \mathcal{P} , the probability of satisfying the task under *any* policy is *zero*. However, it is still important to identify those policies that ensure high probability of avoiding bad states over long time intervals. Consequently, we propose to use an ASCC of \mathcal{P} as the *relaxed* AMEC, due to the following lemma.

Lemma 4: Assume that there exists one infinite path of \mathcal{P} that is accepting. Then, there exists at least one SCC of \mathcal{P} that intersects with $I_{\mathcal{P}}^i$ but not with $H_{\mathcal{P}}^i$, for at least one pair $(H_{\mathcal{P}}^i, I_{\mathcal{P}}^i) \in \text{Acc}_{\mathcal{P}}$.

Proof: As mentioned before, an infinite path of \mathcal{P} , denoted by $R_{\mathcal{P}}$, is accepting if for at least one pair $(H_{\mathcal{P}}^i, I_{\mathcal{P}}^i) \in \text{Acc}_{\mathcal{P}}$, it holds that $R_{\mathcal{P}}$ intersects with all states in $H_{\mathcal{P}}^i$ finitely often, whereas with $I_{\mathcal{P}}^i$ infinitely often. Since both $H_{\mathcal{P}}^i$ and $I_{\mathcal{P}}^i$ are finite, there exists a cyclic path $s_k \dots s_f \dots s_k$ of \mathcal{P} that contains at least one $s_f \in I_{\mathcal{P}}^i$ and does not contain any state within $H_{\mathcal{P}}^i$. By definition, this cyclic path is an SCC of \mathcal{P} that intersects with $I_{\mathcal{P}}^i$ but not with $H_{\mathcal{P}}^i$.

Denote the set of SCCs in \mathcal{P} as $\Omega \triangleq \{S'_1, S'_2, \dots, S'_C\}$, where $S'_c \subseteq S$. This set can be derived using Tarjan's algorithm [4], [32]. Moreover, denote by $\Omega_{\text{acc}}^i = \{S'_c \in \Omega \mid S'_c \cap I_{\mathcal{P}}^i \neq \emptyset, S'_c \cap H_{\mathcal{P}}^i = \emptyset\}$, the set of SCCs that satisfy the accepting conditions associated with $(H_{\mathcal{P}}^i, I_{\mathcal{P}}^i) \in \text{Acc}_{\mathcal{P}}$. Lemma 4 ensures that $\Omega_{\text{acc}}^i \neq \emptyset$ for at least one pair $(H_{\mathcal{P}}^i, I_{\mathcal{P}}^i) \in \text{Acc}_{\mathcal{P}}$. Therefore, the union $\Omega_{\text{acc}} \triangleq \cup_{i=1, \dots, N} \Omega_{\text{acc}}^i$ is not empty.

Now, the union $S_c \triangleq \cup_{S'_c \in \Omega_{\text{acc}}} S'_c$ serves as the set of states the system should enter, starting from the initial state, and then remain inside any of the ASCC to satisfy the accepting condition. Again, the first step is to formulate a linear program that minimizes the expected total cost of reaching S_c from s_0 while ensuring the risk is upper bounded by the chosen $\gamma_{\text{prex}} > 0$. It can be done analogously as in (8) but over $S_n \triangleq S \setminus S_c$ (which is omitted here). Denote the objective function by $\mathbf{C}_{\text{prex}}(S_c)$ and its set of variables by $\{y_{\text{prex}}(s, u)\}$ and the associated relaxed plan prefix as π_{prex} . Similar to Section IV-B2, it is possible that the system under the policy π_{prex} can enter any ASCC in Ω_{acc} . Assume that the system enters $S'_c \in \Omega_{\text{acc}}$. Different from an AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$, the action set at each state of $S'_c \in \Omega_{\text{acc}}$ is *not* constrained. Thus, there is no guarantee that the system will stay within S'_c after entering it.

Therefore, the second step is to synthesize the relaxed plan suffix that keeps the system inside S'_c to satisfy the accepting condition with the maximal probability. Define the set $I'_c = S'_c \cap I_{\mathcal{P}}^i$, which is not empty for an ASCC S'_c . Then, an accepting cyclic path of \mathcal{P} associated with I'_c , and the cyclic cost associated with S'_c and I'_c can be defined similarly as in Definition 5. Formally, we consider the following problem.

Problem 4: Find a control policy for \mathcal{P} that minimizes the mean cyclic cost associated with the ASCC S'_c : $\mathbb{E}_{P_a}^{\pi} \{\mathbf{C}_{\text{suffix}}(P_a)\}$, where \mathbf{P}_a is the set of all accepting cyclic paths associated with S'_c and $\mathbf{C}_{\text{suffix}}$ is defined as in Definition 6; while *at the same time*, maximizing the probability that the cyclic paths stay within S'_c . ■

In Problem 4, the first objective of minimizing the mean cyclic cost corresponds to minimizing the mean total cost in (4) in Problem 1. The objective of maximizing the probability of the system staying within the ASCC S'_c corresponds to minimizing the frequency with which the system will reach the bad states that violate the task specifications. It constitutes a relaxation of the risk constraint (4) in Problem 1. To solve Problem 4, first we construct a modified MDP $\mathcal{Z}_{\text{suffix}}$ over S'_c , which is similar to $\mathcal{Z}_{\text{suffix}}$ in Section IV-B2. The set I'_c is split into two virtual copies: I_{in} that only has incoming transitions and I_{out} that has only outgoing transitions. Formally, we define $\mathcal{Z}_{\text{suffix}} = (S_{\mathcal{R}}, U_{\mathcal{R}}, E_{\mathcal{R}}, y_0, p_{\mathcal{R}}, c_{\mathcal{R}})$, where

the set of states is $S_r = (S'_c \setminus I'_c) \cup I_{\text{in}} \cup I_{\text{out}} \cup \{s_{\text{bad}}\}$, with $I_{\text{in}} = \{s_f^{\text{in}}, \forall s_f \in I'_c\}$ and $I_{\text{out}} = \{s_f^{\text{out}}, \forall s_f \in I'_c\}$, the two virtual copies of I'_c , and s_{bad} is a virtual bad state. The set of control actions is given by $U_r = U \cup \{\tau_0\}$, where τ_0 is a self-loop action. The set of transition is $E_r \subset S_r \times U_r$, which satisfies that (i) $(s, u) \in E_r, \forall s \in S'_c$ and $u \in U(s)$; (ii) $(s, \tau_0) \in E_r, \forall s \in I_{\text{in}}$; and (iii) $(s_{\text{bad}}, \tau_0) \in E_r$. Moreover, y_0 is the initial distribution of states in S'_c based on the incoming transitions from states in S'_n

$$y_0(s) = \sum_{(\check{s}, u)} p_p(\check{s}, u, s) y_{\text{prex}}(\check{s}, u), \quad \forall s \in (S'_c \setminus I'_c) \cup I_{\text{out}}$$

where $\sum_{(\check{s}, u)} \triangleq \sum_{\check{s} \in S'_n} \sum_{u \in U_p(\check{s})}$ and $\{y_{\text{prex}}(s, u)\}$ are the variable solutions from the synthesis of the relaxed plan prefix, and $y_0(s_{\text{bad}}) = 0$. Furthermore, the transition probability p_r is defined in seven cases as follows: (a) for transitions within $S'_c \setminus I'_c$, it holds that $p_r(s, u, \check{s}) = p_E(s, u, \check{s})$, $\forall s, \check{s} \in S'_c \setminus I'_c, \forall u \in U_r(s)$; (b) for transitions originated from I_{out} , it holds that $p_r(s_f^{\text{out}}, u, \check{s}) = p_E(s_f, u, \check{s})$, $\forall s_f^{\text{out}} \in I_{\text{out}}, \forall u \in U_r(s_f^{\text{out}})$, and $\forall \check{s} \in S'_c \setminus I'_c$; (c) for transitions into I_{in} , it holds that $p_r(s, u, s_f^{\text{in}}) = p_E(s, u, s_f)$, $\forall s \in S'_c \setminus I'_c, \forall u \in U_r(s)$, and $\forall s_f^{\text{in}} \in I_{\text{in}}$; (d) for transitions from I_{out} to I_{in} , it holds that $p_r(s_f^{\text{out}}, u, s_f^{\text{in}}) = p_E(s_f, u, s_f)$, $\forall s_f^{\text{out}} \in I_{\text{out}}$ and $\forall u \in U_r(s_f^{\text{out}})$; (e) for transitions into the bad state s_{bad} , it holds that $p_r(s, u, s_{\text{bad}}) = p_E(s, u, \check{s})$, $\forall s \in S'_c \setminus I_{\text{in}}, \forall \check{s} \in S' \setminus S'_c$, and $u \in U_r(s)$; (f) each state within I_{in} is included in a self-loop such that $p_r(s_f^{\text{in}}, \tau_0, s_f^{\text{in}}) = 1, \forall s_f^{\text{in}} \in I_{\text{in}}$; (g) the bad state is included in a self-loop such that $p_r(s_{\text{bad}}, \tau_0, s_{\text{bad}}) = 1$. Finally, the cost function c_r is defined in two cases: (i) $c_r(s, u) = c_E(s, u)$, $\forall s \in S_r \setminus I_{\text{in}}, \forall u \in U_r(s)$; and (ii) $c_r(s_f^{\text{in}}, \tau_0) = 0, \forall s_f^{\text{in}} \in I_{\text{in}}$ and $c_r(s_{\text{bad}}, \tau_0) = 0$.

Remark 5: Note that E_r contains all actions for each state in S'_c , compared with E_e as allowed by the AMEC. ■

Let $S'_r \triangleq S_r \setminus (I_{\text{in}} \cup \{s_{\text{bad}}\})$ and $S''_r \triangleq S_r \setminus \{s_{\text{bad}}\}$. We can also show that $\mathcal{Z}_{\text{suffix}}$ above is S'_r -transient. Then, to solve Problem 4, we rely on a technique proposed in [35] to deal with dead ends in stochastic shortest path problems. First, we introduce a large positive penalty for reaching the dead state, denoted by $d > 0$. Then, we modify (11) as follows: denote by $z_{s,u}$ the long-run frequency with which the system is at state s and the action u is taken, $\forall s \in S'_r$ and $\forall u \in U_r(s)$. We want to minimize the mean total cost of reaching I_{in} from I_{out} while minimizing the probability of leaving S''_e . In particular, we consider the following optimization:

$$\min_{\{z_{s,u}\}} \left[\mathbf{C}_{\text{suffix}}(S'_c, d) \triangleq \sum_{(\check{s}, u)} \left(\sum_{s \in S''_e} \eta(\check{s}, u, s) c_r(\check{s}, u) + \eta(\check{s}, u, s_{\text{bad}}) d \right) \right] \quad (12a)$$

$$\text{s.t.} \quad \sum_{u \in U_r(s)} z_{s,u} = \sum_{(\check{s}, u)} \eta(\check{s}, u, s) + y_0(s) \quad \forall s \in S'_r \quad (12b)$$

$$\sum_{(\check{s}, u)} \left(\sum_{s \in I_{\text{in}}} \eta(\check{s}, u, s) + \eta(\check{s}, u, s_{\text{bad}}) \right) = \sum_{s \in S'_c} y_0(s) \quad (12c)$$

$$z_{s,u} \geq 0 \quad \forall s \in S'_r \quad \forall u \in U_r(s) \quad (12d)$$

where the notation $\sum_{(\check{s}, u)} \triangleq \sum_{\check{s} \in S'_e} \sum_{u \in U_r(\check{s})}$, the variables satisfy that $\eta(\check{s}, u, s) \triangleq z_{\check{s},u} p_r(\check{s}, u, s)$, $\eta(\check{s}, u, s_{\text{bad}}) \triangleq z_{\check{s},u} p_r(\check{s}, u, s_{\text{bad}})$, and $\mathbf{C}_{\text{suffix}}(S'_c, d)$ denotes the objective function as the summation of the mean cost of reaching I_{in} and the expected penalty of reaching s_{bad} . The first constraint balances the incoming and outgoing flow at each state, whereas the second constraint ensures that $I_{\text{in}} \cup \{s_{\text{bad}}\}$ are eventually reached. Let the optimal solution of (12) be $z^*_{\text{suffix}} = \{z^*_{s,u}, s \in S'_r, u \in U_r(s)\}$. Then, the optimal stationary policy for the relaxed plan suffix, denoted by π^*_{suffix} , can be derived as follows: for states in S'_r , the optimal policy is given by $\pi^*_{\text{suffix}}(s, u) = z^*_{s,u} / (\sum_{u \in U_r(s)} z^*_{s,u})$ if $\sum_{u \in U_r(s)} z^*_{s,u} \neq 0$; otherwise the action at s is chosen randomly, $\forall s \in S'_r$. Note that $\pi^*_{\text{suffix}}(s_f, u) = \pi^*_{\text{suffix}}(s_f^{\text{out}}, u)$, $\forall s_f \in I'_c$ and $\forall u \in U(s_f)$.

Lemma 5: Under the relaxed plan suffix π^*_{suffix} , the probability of $\mathcal{Z}_{\text{suffix}}$ reaching I_{in} from I_{out} while staying within S''_r over an infinite horizon is lower bounded by $1 - \gamma_{\text{suffix}}(d)$, where $\gamma_{\text{suffix}}(d) \triangleq \sum_{\check{s} \in S'_e} \sum_{u \in U_r(\check{s})} z^*_{\text{suffix}}(\check{s}, u) p_r(\check{s}, u, s_{\text{bad}})$.

Proof: The proof is a simple inference from (12c). ■

Remark 6: A lower bound can be enforced on γ_{suffix} as in (8). However, this bound is hard to estimate and a large bound can yield the problem infeasible. In contrast, (12) always has a solution and $\gamma_{\text{suffix}}(d)$ is tunable by varying d .

C. Complete Policy

In this section, we present how to combine the stationary plan prefix and plan suffix of \mathcal{P} into the complete finite-memory policy of the original MDP \mathcal{M} . Furthermore, we show how to execute this finite-memory policy online.

1) Combining the Plan Prefix and Suffix: When AMECs of \mathcal{P} exist, we can combine the plan prefix synthesis and the plan suffix synthesis for each AMEC into one linear program as follows:

$$\min_{\{y_{s,u}, z_{s,u}\}} \beta \cdot \mathbf{C}_{\text{pre}}(S_c) + (1 - \beta) \sum_{(S'_c, U'_c) \in \Xi_{\text{acc}}} \mathbf{C}_{\text{suf}}(S'_c, U'_c) \quad (13)$$

s.t. Constraints (8b)–(8d) and (11c)–(11d)

where $\mathbf{C}_{\text{pre}}(S_c)$ and $\mathbf{C}_{\text{suf}}(S'_c, U'_c)$ are defined in (8a) and (11a), respectively, the variables $\{y_{s,u}\}$ satisfy the constraints (8b)–(8d) and (11c), and the variables $z_{s,u} \triangleq \{z_{s,u}(S'_c), \forall (S'_c, U'_c) \in \Xi_{\text{acc}}\}$, where $z_{s,u}(S'_c)$, satisfy the constraints (11c)–(11d) for the AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$. The parameter $0 \leq \beta \leq 1$ captures the importance of minimizing the expected total cost to reach S_c versus stay in S_c . Note that the initial conditions y_0 in (11c) for each state in the suffix are expressed over the variables $\{y_{s,u}\}$. In other words, the initial conditions of each AMEC are now optimized to solve the combined objective function (13). It can be solved via any LP solver, e.g., ‘‘Gurobi’’ [36] and ‘‘CPLEX.’’ Once the optimal solution $\{y^*_{s,u}\}$ and $z^*_{s,u}$ is obtained, the optimal plan prefix π^*_{pre} can be constructed as described in Section IV-B1 and the plan suffix π^*_{suf} as in Section IV-B2.

On the other hand, when no AECs of \mathcal{P} exist, as discussed in Section IV-B3, we can combine the relaxed plan prefix and suffix

synthesis for each ASCC into one linear program as follows:

$$\begin{aligned} \min_{\{y_{s,u}, z_{s,u}\}} & \beta \cdot \mathbf{C}_{\text{prex}}(S_c) + (1 - \beta) \sum_{S'_c \in \Omega_{\text{acc}}} \mathbf{C}_{\text{sufx}}(S'_c, d) \\ \text{s.t.} & \text{Constraints (8b)–(8d) and (12b)–(12d)} \quad (14) \end{aligned}$$

where $\mathbf{C}_{\text{prex}}(S_c)$ and $\mathbf{C}_{\text{sufx}}(S'_c, d)$ are defined in (8a) and (12a), respectively, the variables $\{y_{s,u}\}$ satisfy the constraints (8b)–(8d), and the variables $z_{s,u} \triangleq \{z_{s,u}(S'_c), \forall S'_c \in \Omega_{\text{acc}}\}$, where $z_{s,u}(S'_c)$, satisfy the constraints (12b)–(12d) for the ASCC $S'_c \in \Omega_{\text{acc}}$. The parameter $0 \leq \beta \leq 1$ captures the importance of minimizing the expected total cost to reach S_c versus stay in S_c . Similar to the previous case, the initial conditions y_0 in (12b) for each state in the ASCCs are expressed over the variables $\{y_{s,u}\}$. Thus, the initial conditions are now optimized to solve the combined objective function (14). Again, it can be solved via any LP solver. Once the optimal $\{y_{s,u}^*\}$ and $z_{s,u}^*$ is obtained, the optimal relaxed plan prefix π_{prex}^* and relaxed plan suffix π_{sufx}^* can be constructed as described in Section IV-B3.

Note that the size of both linear programs in (13) and (14) is linear with respect to the number of transitions in \mathcal{P} and can be solved in polynomial time [37]. Note also that the multiobjective costs introduced in (13) and (14) provide a balance between optimizing the plan prefix and suffix. Compared to only optimizing the plan suffix, i.e., for $\beta = 0$ as required to solve Problems 3 and 4, increasing slightly the value of β can lead to a significant decrease in the total cost of the plan prefix, without sacrificing much the optimality in the plan suffix.

Observe that the optimal policy derived above only includes the states within $S_n \cup S_c$. Thus, no policy is specified for the *bad states* in S_d . Once the system reaches any bad state, it has violated the formula φ and cannot satisfy it anymore. Thus, it is common practice to stop the system once that happens [4], [21]. We propose here a new method that allows the system to *recover* from the bad state in S_d and continue performing the task, which could be useful for partially feasible tasks with soft constraints, as discussed in [7].

Definition 7: The *projected distance* of a bad state $s_d = \langle x, l, q \rangle \in S_d$ onto $S_c \cup S_n$ via $u \in U(s_d)$ is defined as follows:

$$\kappa(s_d, u) \triangleq \sum_{\check{s} \in S_c \cup S_n} \frac{D(l, \chi(q, \check{q}))}{|\chi(q, \check{q})|} \cdot p_E(x, u, \check{x}) \cdot p_L(\check{x}, \check{l}) \quad (15)$$

where $\check{s} \triangleq \langle \check{x}, \check{l}, \check{q} \rangle$ and function $D: 2^{\text{AP}} \times 2^{\text{AP}} \rightarrow \mathbb{N}$ returns the distance between an element $l \in 2^{\text{AP}}$ and a set $\chi \subseteq 2^{\text{AP}}$, was first introduced in [7] and restated as follows. ■

Simply speaking, $\kappa(s_d, u)$ evaluates how much the product automaton \mathcal{P} is violated on the average if the bad state $s_d \in S_d$ is projected into the set of good states $S_c \cup S_n$ using action $u \in U(s_d)$. Function $D(l, \chi) = 0$ if $l \in \chi$ and $D(l, \chi) = \min_{\ell' \in \chi} |\{a \in \text{AP} \mid a \in \ell, a \notin \ell'\}|$ otherwise. Namely, it returns the minimal difference between ℓ and any element in χ . Given $\kappa(\cdot)$, the policy at $s_d \in S_d$ is given by

$$\pi^*(s_d, u) = \begin{cases} 1 & \text{for } u = \arg \min_{u \in U(s_d)} \kappa(s_d, u) \\ 0 & \text{other } u \in U(s_d) \end{cases} \quad (16)$$

which chooses the single action that minimizes (15). Combing (13), (14), and (16) provides the complete policy for \mathcal{P} . The above-mentioned discussions are summarized in Algorithm 1.

Algorithm 1: Complete Policy Synthesis.

Input: \mathcal{P} by Definition 2, γ, β

Output: the complete policy π^*, μ^*

if $\Xi_{\text{acc}} \neq \emptyset$ **then**

1. Construct $\mathcal{Z}_{\text{prex}}$, and \mathcal{Z}_{suf} for

each $(S'_c, U'_c) \in \Xi_{\text{acc}}$.

2. Derive π^* via solving (13), and (16).

else

1. Construct $\mathcal{Z}_{\text{prex}}$, and $\mathcal{Z}_{\text{sufx}}$ for each $S'_c \in \Omega_{\text{acc}}$.

2. Derive π^* via solving (14), and (16).

3. Construct μ^* from π^* by (17)

2) Mapping π^* to μ^* : Finally, we need to map the optimal stationary policy π^* of \mathcal{P} above to the optimal finite-memory policy μ^* of \mathcal{M} . Starting from stage $t = 0$, the initial state $s_0 = \langle x_0, l_0, q_0 \rangle \in S_n$ and the optimal action to take is given by the distribution $\pi^*(s_0)$. Assume that $u \in U(s_0)$ is taken. Then, at stage $t = 1$, the robot observes its resulting state x_1 and the label l_1 . Thus, the subsequent state in \mathcal{P} is $s_1 = \langle x_1, l_1, q_1 \rangle$, where $q_1 = \delta(q_0, l_0)$ is unique as \mathcal{A}_φ is deterministic. The optimal action to take now is given by the distribution $\pi^*(s_1)$. This process repeats itself indefinitely. Denote by $s_t \in S$, the *reachable* state at stage $t \geq 0$, which is always unique given the robot's past sequence of states $X_t = x_0 x_1 \dots x_t$ and labels $L_t = l_0 l_1 \dots l_t$. Thus, the optimal policy μ^* at stage $t \geq 0$, given X_t and L_t , is

$$\mu^*(X_t, L_t) = \pi^*(s_t) \quad (17)$$

i.e., the control policy at the reachable state s_t in \mathcal{P} is the best control policy in \mathcal{M} at stage $t, \forall t \geq 0$. Last but not least, if the system reaches a bad state at stage $t - 1$, i.e., $s_{t-1} \in S_d$, according to policy (16), the robot will take action u^* and more importantly, the next reachable state is *set to be* $s_t \triangleq \langle x_t, l_t, q'_t \rangle \in (S_c \cup S_n)$, where x_t and l_t are the observed robot location and label at stage t and $q'_t \triangleq \arg \min_{\check{q} \in \text{Post}(q_{t-1})} D(l_{t-1}, \chi(q_{t-1}, \check{q}))$.

Theorem 6: Algorithm 1 solves Problem 1, if AECs of \mathcal{P} exist and $\beta = 0$. Otherwise, if no AECs of \mathcal{P} exist, then Problem 1 has no solution. In this case, Algorithm 1 provides a relaxed policy that minimizes the *relaxed* suffix cost $\mathbf{C}_{\text{sufx}}(S'_c, d)$ defined in (12). Moreover, given any finite run $S_T = s_0 s_1 \dots s_T$ of \mathcal{P} under the optimal policy π^* , the probability that S_T does not intersect with the set of bad states S_d for all time $t \in [0, T]$ is bounded as

$$Pr(s_t \notin S_d, \forall t \in [0, T]) \geq (1 - \gamma_{\text{prex}}) \cdot (1 - \gamma_{\text{sufx}}(d))^{N_s}$$

where $N_s \geq 0$ is the number of accepting cyclic paths contained in S_T that depends on T .

Proof: To show the *first* part of this theorem, similar to Lemma 1, the constraints of (8b)–(8d) ensure that the total probability of reaching the union of all AMECs is lower bounded by $1 - \gamma$. Moreover, the first part of Lemma 3 shows that any infinite run $\tau_{\mathcal{P}}$ of \mathcal{P} would satisfy φ once it enters any AMEC $(S'_c, U'_c) \in \Xi_{\text{acc}}$, by following the plan suffix. The fact that π^* also minimizes the mean total cost in (4) when $\beta = 0$ in (13) can be shown as follows: as discussed in [24], [33], and [34], the mean payoff objective depends on how the system suffix behaves within the AMECs. The second part of Lemma 3 guarantees that the derived plan suffix π_{suf}^* minimizes the mean total

Algorithm 2: Policy Execution.

Input: \mathcal{M} , φ , observed state x_t and label l_t at stage $t \geq 0$

Output: μ^* and u_t at stage $t \geq 0$

1. **Offline:** Construct \mathcal{P} and synthesize π^* by Alg. 1.
2. At $t = 0$: set $s_0 = \langle x_0, l_0, q_0 \rangle$ and apply $u_0 \sim \pi^*(s_0)$.
3. **while** $t = 1, 2, \dots$ **do**
 - observe x_t and l_t .
 - if** $s_{t-1} \notin S_d$ **then**
 - Set $s_t = \langle x_t, l_t, q_t \rangle$, where $q_t = \delta(q_{t-1}, l_{t-1})$.
 - else**
 - Set $s_t = \langle x_t, l_t, q'_t \rangle \in (S_n \cup S_c)$.
 - Apply action $u_t \sim \pi^*(s_t)$.

cost of staying within any of the AMECs while satisfying the accepting condition.

To show the *second* part of the theorem, no solution to Problem 1 exists regardless of the choice of γ , as the probability of satisfying the task is zero. Instead, when $\beta = 0$, the optimal policy π^* obtained by Algorithm 1 minimizes the *relaxed* suffix cost $C_{\text{suffx}}(S'_c, d)$. At the same time, due to the constraints in (8) that are also present in (13), the plan prefix π_{prex}^* ensures that all runs stay within S_n with at least probability $(1 - \gamma_{\text{prex}})$ before entering any ASCC $S'_c \in \Omega_{\text{acc}}$, whereas the relaxed plan suffix π_{suffx}^* ensures that the runs stay within S'_c with at least probability $(1 - \gamma_{\text{suffx}}(d))$ for one execution of any accepting cyclic path. Consequently, if the finite run contains N_s accepting cyclic paths, the probability of avoiding S_d is lower bounded by $(1 - \gamma_{\text{prex}}) \cdot (1 - \gamma_{\text{suffx}}(d))^{N_s}$. Even though this probability approaches zero as N_s approaches infinity, this result still ensures that the frequency of visiting bad states over finite intervals is minimized.

3) Policy Execution: Clearly, the optimal policy μ^* from (17) requires only a finite memory to save the current reachable state s_t and the optimal policy π^* . It is synthesized offline once via Algorithm 1 and its online execution involves observing the current state x_t and label l_t , updating the reachable state s_t , and applying the action according to $\pi^*(s_t)$. Details are given in Algorithm 2.

V. SIMULATION RESULTS

In this section, we present simulation results to validate the scheme. All algorithms are implemented in Python 2.7 and available online [32]. All simulations are carried out on a laptop (3.06-GHz Duo CPU and 8 GB of RAM).

A. Model Description

We consider a partitioned 10 m \times 10 m workspace as shown in Fig. 8, where each cell is a 2 m \times 2 m area. The properties of interest are $\{0\text{bs}, \text{b1}, \text{b2}, \text{b3}, \text{Sp1}\}$. The properties satisfied at each cell are probabilistic: three cells at the corners satisfy b1 , b2 , and b3 , respectively, with probability 1. Four cells at (1 m, 5 m), (5 m, 3 m), (9 m, 5 m), and (5 m, 9 m) satisfy Sp1 with probabilities ranging from 0.2 to 0.8, modeling the likelihood that a supply appears at that particular cell. One cell

at (5 m, 1 m) satisfies 0bs with probability 0.7. Other obstacles will be described later upon different task scenarios.

The robot motion follows the unicycle model, i.e., $\dot{x} = v \cos(\theta)$, $\dot{y} = v \sin(\theta)$, and $\dot{\theta} = \omega$, where $p(t) = (x(t), y(t)) \in \mathbb{R}^2$, $\theta(t) \in (-\mathbf{pi}, \mathbf{pi}]$ are the robot's position and orientation at time $t \geq 0$. The control input is $u(t) = (v(t), \omega(t))$ and contains the linear and angular velocities. Due to actuation noise and drifting, the robot's motion is subject to uncertainty. The action primitives and the associated uncertainties are shown in Fig. 2 and described as follows: action "FR" means driving forward for 2 m by setting $v(t) = v_0$ and $\omega(t) = 0$, $\forall t = [0, 2/v_0]$. This action has probability 0.8 of reaching 2 m forward and probability 0.1 of drifting to the left or right by 2 m, respectively; action "BK" can be defined analogously to "FR"; action "TR" means turning right by an angle of $\mathbf{pi}/2$ by setting $v(t) = 0$ and $\omega(t) = -\omega_0$, $\forall t = [0, \mathbf{pi}/(2\omega_0)]$. This action has probability 0.9 of turning to the right by $\mathbf{pi}/2$, probability 0.05 of turning less than $\mathbf{pi}/4$ due to undershoot, and probability 0.05 of turning more than $3\mathbf{pi}/4$ due to overshoot; action "TL" can be defined analogously to "TR"; finally, action "ST" means staying still by setting $v(t) = \omega(t) = 0$, $\forall t = [0, T_0]$, where T_0 is the chosen waiting time. It has probability 1.0 of staying where it is. The cost of each action is given by $[2, 4, 3, 3, 1]$, respectively, where the cost of "ST" is set to 1 as it consumes a unit time to wait at one cell.

With the above-mentioned model, we can abstract the robot state by the cell coordinate in which it belongs, namely, $(x_c, y_c) \in \{1, 3, \dots, 9\}^2$ and its four possible orientations (N, E, S, W). The transition relation and probability can be built following the aforementioned description. The resulting probabilistically labeled MDP has 100 states and 816 edges.

In the sequel, we consider *three* different task formulas in the order of increasing complexity. We used "Gurobi" [36] to solve the linear programs in (13) and (14). When comparing the performance in the plan suffix, we also use the total cost in (9) as an indicator, especially when the difference in the mean total cost in (10) is too small to measure.

B. Ordered Reachability

In this case, we show the tradeoff between reducing the expected total cost and decreasing the risk factor in the plan prefix synthesis using (8). In particular, the robot needs to reach b1 , b2 , and b3 (in this order) from the initial cell while avoiding obstacles for all time. Afterward, it should stay at b3 . The LTL formula for this task is

$$\varphi_1 = (\diamond(\text{b1} \wedge \diamond(\text{b2} \wedge \diamond\text{b3}))) \wedge (\square\neg 0\text{bs}) \wedge (\diamond\square\text{b3}). \quad (18)$$

The associated DRA derived using [31] has 7 states, 24 transitions, and 1 accepting pair. An additional obstacle is added that has probability 0.7 of appearing in the cell (5 m, 9 m).

It took 10.9 s to construct the product automaton that has 840 states, and 7280 transitions. Since one AMEC exists, we synthesize the optimal policy using Algorithm 1 via solving (13) under $\beta = 0.5$ and different risk factors γ chosen from $\{0, 0.1, \dots, 0.4\}$, which took on average 0.1 s. Then, we perform 1000 Monte Carlo simulations of 500 time steps each, where we evaluate the total cost in (7) and whether the task is satisfied. As shown in Table I, the total cost increases when

TABLE I
STATISTICS OF 1000 MONTE CARLO SIMULATIONS OF 500 TIME STEPS,
UNDER DIFFERENT γ FOR TASK (18)

γ	Total cost	Failure	Success	Unfinished
0	132.2	0	910	90
0.1	118.1	99	872	29
0.2	110.5	219	770	11
0.3	104.6	308	692	0
0.4	98.3	417	583	0

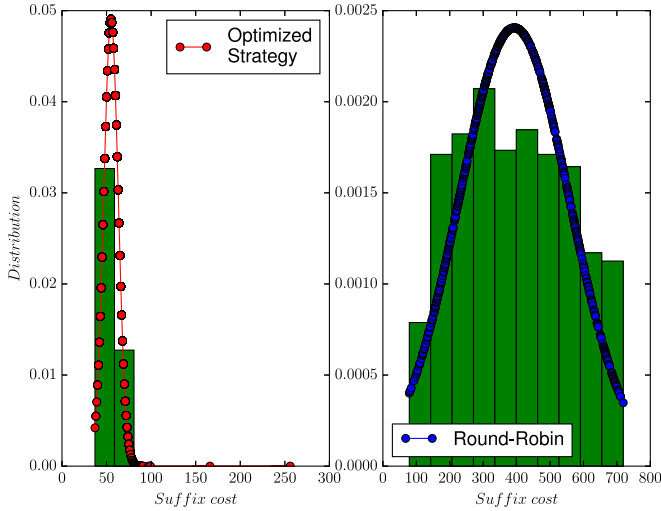


Fig. 7. Normalized distribution of the total cost of accepting cyclic paths from 1000 Monte Carlo simulations under the optimal plan suffix (left) and the Round-Robin policy (right), for task (19).

the allowed risk factor γ is decreased. The percentage of simulated runs that collide with an obstacle is approximately $(1 - \gamma)$, which verifies the risk constraint in Lemma 1.

C. Surveillance

In this case, we compare the efficiency of the optimal plan suffix from Algorithm 1 and the Round-Robin policy. Particularly, the robot should visit b1, b2, and b3 infinitely often for surveillance and avoid all obstacles

$$\varphi_2 = (\square \diamond b1) \wedge (\square \diamond b2) \wedge (\square \diamond b3) \wedge (\square \neg obs). \quad (19)$$

The associated DRA has 8 states, 30 transitions, and 1 accepting pair. It took 5.8 s to construct the product \mathcal{P} that has 700 states, 5712 transitions, and 1 accepting pair. Since one AMEC exists in the product, we synthesize the optimal policy using Algorithm 1 via solving (13) under $\gamma = 0$ and $\beta = 0.1$, which took 0.2 s. We conducted 1000 Monte Carlo simulations and Fig. 7 shows that the total cost from (9) of accepting cyclic paths in the plan suffix under the optimal policy is much lower than the Round-Robin policy (50 versus 400). Moreover, Fig. 9 shows that the average number of times each base station is visited by the robot under the optimal policy is much higher than under the Round-Robin policy.

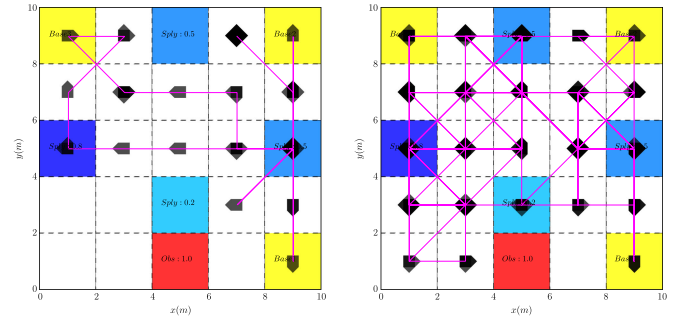


Fig. 8. Simulated trajectory suffix under the optimized plan suffix (left) and the Round-Robin policy (right).

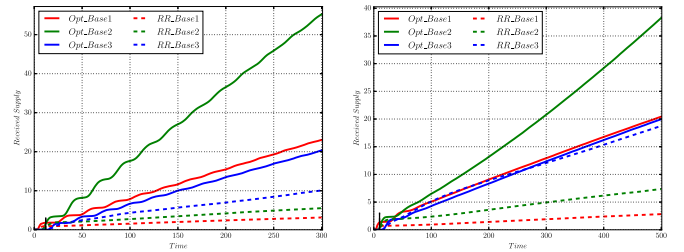


Fig. 9. Left: The average number of times each base is visited, for task (19). Right: The average number of supplies received at each base for task (20). The optimal policy is shown in solid lines while the Round-Robin policy is shown in dashed lines.

TABLE II
OPTIMAL PREFIX COST, SUFFIX COST, AND THE BALANCED COST AS
DEFINED IN (13) OF TASK (20) UNDER DIFFERENT β WITH $\gamma = 0$

β	Prefix cost	Suffix cost		Balanced cost by (13)
		Total	Mean	
0	180.7	66.1	2.524	66.1
0.2	62.4	67.1	2.533	65.2
0.4	50.5	72.9	2.551	64.1
0.6	49.8	73.5	2.552	59.3
0.8	49.5	74.3	2.554	54.4
1.0	49.5	246.7	2.817	49.5

D. Ordered Supply Delivery

In this case, we demonstrate the reactivity of the derived optimal policy. The robot needs to collect supplies from the cells that are marked by Sp1, where supplies appear probabilistically. Then, it needs to transport these supplies to each base station. Furthermore, the robot should *not* visit two base stations consecutively without collecting a supply first. It should always avoid obstacles. The LTL task formula is

$$\varphi = \varphi_{\text{all.base}} \wedge \varphi_{\text{order}} \wedge (\square \neg obs) \quad (20)$$

where $\varphi_{\text{all.base}} = (\square \diamond b1) \wedge (\square \diamond b2) \wedge (\square \diamond b3)$ means that all base stations should be visited infinitely often and $\varphi_{\text{order}} = \square(\varphi_{\text{one}} \rightarrow \bigcirc((\neg \varphi_{\text{one}}) \cup Sp1))$, with $\varphi_{\text{one}} = (b1 \vee b2 \vee b3)$ means that when one base station is visited, then no base can be visited until a supply has been collected. The associated DRA

TABLE III
SIZE AND COMPUTATION TIME OF VARIOUS MODELS \mathcal{M} AS DESCRIBED IN SECTION V-E UNDER TASK (20)

\mathcal{M}		\mathcal{P}		AMECs Ξ_{acc}		π^* via (13)		
Size	Time [s]	Size	Time [s]	Size	Time [s]	Size of (8)	Size of (11)	Time to solve (13) [s]
(100, 816)	0.13	(4.2e3, 4.1e4)	16.3	1.2e3	4.15	(443, 2.0e3, 8.3e3)	(1.2e3, 4.9e3, 2.1e4)	0.21
(324, 2.8e3)	1.69	(1.1e4, 1.0e5)	41.2	3.6e3	29.4	(1.3e3, 6.3e3, 2.2e4)	(3.6e3, 1.7e4, 5.9e4)	0.72
(900, 8.4e3)	24.2	(2.9e4, 2.8e5)	106.8	1.0e4	337.1	(3.6e3, 1.7e4, 6.0e4)	(9.9e3, 4.8e4, 1.6e5)	16.74
(1.4e3, 1.3e4)	88.7	(4.7e4, 4.5e5)	391.7	1.6e4	1.1e3	(5.8e3, 2.8e4, 9.7e4)	(1.5e4, 7.7e4, 2.6e5)	20.81
(2.5e3, 2.4e4)	326.9	(8.1e4, 7.8e5)	290.1	2.7e4	4.8e3	(1.0e4, 4.9e4, 1.6e5)	(2.7e4, 1.3e5, 4.5e5)	15.74
(3.3e3, 3.2e4)	558.3	(1.0e5, 1.1e6)	380.1	3.7e4	9.4e3	(1.3e4, 6.6e4, 2.2e5)	(3.7e4, 1.8e5, 6.1e5)	32.04

The notation $\text{aeb} \triangleq a \times 10^b$ for $a, b > 0$. The size of \mathcal{M} , \mathcal{A}_φ , and \mathcal{P} includes the number of states and transitions. The size of LP problems (13), which contains (8) and (11), includes the number of rows, columns, and variables in the linear equations, as indicated by the ‘‘Gurobi’’ solver [36].

is derived using [31] and [32] in 0.05 s, which has 32 states, 298 transitions, and 1 accepting pair.

It took around 16 s to construct the product automaton that has 4224 states, 41344 transitions, and 1 accepting pair. Since two AMECs exist in the product, we synthesize the optimal policy using Algorithm 1 via solving (13) under $\gamma = 0$ and $\beta = 0.1$, which took around 0.2 s, given the complexity of task (20). Notice that the optimal plan sometimes requires the robot to wait at a cell marked by Sp1 by taking action ‘‘ST’’ since the expected cost of traveling to another cell with supply might be higher than waiting there for the supply to appear. Fig. 8 compares the simulated trajectories under the optimal policy and the Round-Robin policy. Based on 1000 Monte Carlo simulations, the total cost of accepting cyclic paths is much lower under the optimal policy than the Round-Robin policy (70 versus 550). Furthermore, Fig. 9 shows the average number of supplies received at each base under these two policies. It can be seen that much more supplies are received at each base station under the optimal policy. Simulation videos of both cases can be found in [38]. Finally, to show how the choice of β in (13) affects the optimal prefix and suffix cost, we repeat the above-mentioned procedure for different β and the results are summarized in Table II. In the table, the prefix cost equals to $\mathbf{C}_{\text{pre}}(S_c)$, and the mean suffix cost equals to $\sum_{(S'_c, U'_c) \in \Xi_{\text{acc}}} \mathbf{C}_{\text{suf}}(S'_c, U'_c)$ from (13). The total suffix cost is computed based on (9) in order to magnify the changes in the suffix cost. It can be noticed that for small nonzero values of β , less than 0.2, the optimal prefix cost is reduced *dramatically* (from 180.7 to 62.4), without increasing much the optimal suffix cost (from 66.1 to 67.1).

In order to demonstrate scalability and computational complexity of the proposed algorithm, we repeat the policy synthesis under the same task (20), but for workspaces of various sizes. Particularly, we increase the number of cells from 5^2 to 9^2 , 15^2 , 19^2 , 25^2 , and 29^2 . The size of resulting \mathcal{M} , \mathcal{P} , and Ξ_{acc} , and the time taken to compute them are shown in Table III, where we also list the complexity of the LP (13), which consists of (8) and (11), and the time taken to solve (13). It can be seen from Table III that solving (13) requires a small fraction of total time, compared to the construction of \mathcal{M} , \mathcal{P} , and Ξ_{acc} .

E. Surveillance With Clustered Obstacles

In this case, we demonstrate how the relaxed plan prefix and suffix can be synthesized under scenarios where no AECs can be found. In particular, we consider the surveillance task in (19)

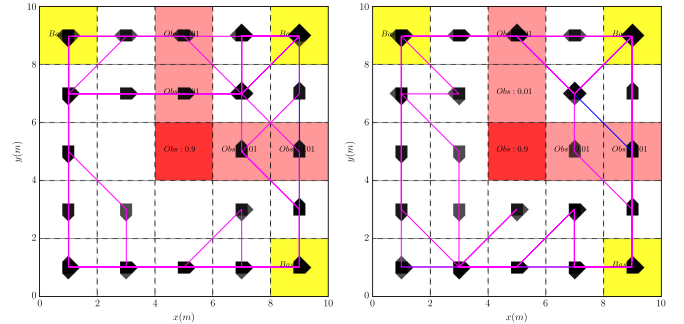


Fig. 10. Two simulated trajectories of 200 time steps for the surveillance task (19), under the relaxed optimal policy.

TABLE IV
STATISTICS OF 1000 MONTE CARLO SIMULATIONS UNDER DIFFERENT γ_{prex} AND d , FOR TASK (19) IN SECTION V-E

γ_{prex}	d	γ_{sufx}	Failure	Pre. Success	Suf. Success
0.1	300	0.05	106	894	852
0.2	300	0.05	169	831	785
0.3	300	0.05	318	682	650
0.4	300	0.05	409	591	549
0.1	280	0.85	888	901	117
0.1	270	0.98	997	903	4

but more obstacles are placed in the workspace as shown in Fig. 10. The center cell (5 m, 5 m) has probability 0.9 of being occupied by an obstacle and the four cells above and on the left have probability 0.01 of being occupied by an obstacle. Thus, b1 is surrounded by possible obstacles around it, even though the probability is very low.

The resulting product automaton has 1184 states, 13888 transitions, and 1 accepting pair. It can be verified that *no AECs exist* in \mathcal{P} , and thus the second case of Algorithm 1 is activated, where the optimal solution is derived by solving (14). We synthesize the relaxed optimal policy under different γ_{prex} and d , as shown in Table IV. It took in average 37 s to synthesize the complete policy for $\beta = 0.1$ and any chosen γ_{prex} and d in this case. Recall that d is a large positive penalty for entering the set of bad states in (12). In particular, we first choose $\gamma_{\text{prex}} = 0.1$ and $d = 300$. Two simulated trajectories under the derived policy are shown in Fig. 10. Furthermore, we perform 1000 Monte Carlo

TABLE V
SIZE AND COMPUTATION TIME OF VARIOUS MODELS \mathcal{M} UNDER TASK (19) WHERE NO AECs EXIST IN \mathcal{P}

\mathcal{M}		\mathcal{P}		ASCCs Ω_{acc}		π^* via (14)		
Size	Time [s]	Size	Time [s]	Size	Time [s]	Size of (8)	Size of (12)	Time to solve (14) [s]
(100, 816)	0.13	(1.0e3, 1.1e4)	0.9	3.1e2	0.66	(202, 920, 3.4e3)	(301, 1.4e3, 4.9e3)	0.45
(324, 2.8e3)	1.57	(2.9e3, 3.1e4)	3.39	9.8e2	1.84	(6.5e2, 3.1e3, 1.1e4)	(9.7e2, 4.7e3, 1.6e4)	2.41
(900, 8.4e3)	23.9	(7.7e3, 7.9e4)	7.04	2.7e3	5.09	(1.8e3, 8.7e3, 3.0e4)	(2.7e3, 1.3e4, 4.5e4)	9.89
(1.4e3, 1.3e4)	92.2	(1.2e4, 1.2e5)	9.78	4.3e3	8.41	(2.9e3, 1.4e4, 4.9e4)	(4.3e3, 2.1e4, 7.2e4)	22.94
(2.5e3, 2.4e4)	322.1	(2.1e4, 2.1e5)	20.1	7.5e3	17.1	(5.1e3, 2.5e4, 8.5e4)	(7.5e3, 3.7e4, 1.3e5)	83.33
(3.3e3, 3.2e4)	625.2	(2.8e4, 2.9e5)	23.1	1.0e4	19.6	(6.7e3, 3.3e4, 1.1e5)	(1.0e4, 4.9e4, 1.7e5)	145.8

The notations are defined similarly as in Table III. In this case, the combined LP in (14) contains (8) and (12) instead.

simulation under γ_{prex} and d listed in Table IV, where we compare the number of times that the robot fails the task by colliding with obstacles (the failure), the number of times that the robot successfully reaches the set of ASCC S_c (the prefix success), and the number of times that the robot successfully executes one accepting cyclic path associated with S'_c and I'_c of one ASCC (the suffix success). It can be seen that $(1 - (1 - \gamma_{\text{prex}})(1 - \gamma_{\text{sufx}}))$, $(1 - \gamma_{\text{prex}})$, and $(1 - \gamma_{\text{sufx}})$ match very well the probability of failure, the prefix success, and the suffix success, respectively, as discussed in Theorem 6. Also, it can be seen that the system can recover from the bad states and continue executing the task if the recovery policy proposed in (16) is activated. It can also be seen that increasing γ_{prex} leads to a lower prefix success rate and decreasing d leads to a lower suffix success rate.

To demonstrate scalability and computational complexity of the proposed algorithm when AMECs do not exist, we repeat the policy synthesis under the same task (19) but for different workspaces of various sizes, as in Section V-D. We set $\gamma = 0.3$, $d = 300$, and $\beta = 0.1$. The size of resulting \mathcal{M} , \mathcal{P} , and Ω_{acc} , and the time taken to compute them are shown in Table V, where we also list the complexity of (14), which consists of (8) and (12), and the time taken to solve (14). It can be seen above that solving (14) now requires a larger fraction of total time, compared to the construction of \mathcal{M} , \mathcal{P} , and Ω_{acc} . However, it requires much less time to compute the set of ASCCs Ω_{acc} than the set of AMECs Ξ_{acc} . For instance, in the case of 29^2 cells in the workspace, it took around 23.1 s to construct \mathcal{P} (which has approximately 2.8×10^4 states and 2.9×10^5 transitions) and 19.6 s to construct its ASCCs (compared with 160 min in Table III). Once (14) is constructed, it took around 2.5 min to solve it.

F. Comparison With PRISM

In this section, we compare the proposed algorithm to the widely used model-checking tool PRISM [13]. The following results were obtained using PRISM 4.3.1, where LP is chosen as the solution method. First, since PRISM does not take the probabilistically labeled MDP in (1) as inputs, we translate the product automaton in (5) into PRISM language and verify its Rabin accepting condition directly. Implementation details can be found in [32]. For tasks (18), (19), and (20), PRISM verifies that the probability of satisfying each of them is 1.0, within time 0.46 s, 0.38 s, and 6.4 s, respectively. The difference in computation time is likely due to the difference in the LP solvers. Second, in order to test different values of γ , we use the “multi-

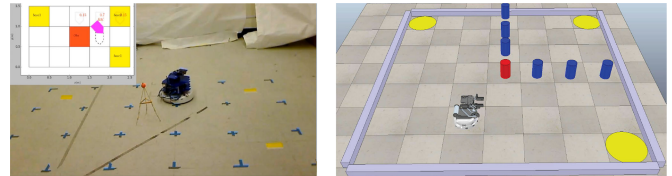


Fig. 11. Experiment workspace (left) with the monitoring panel. Three bases are marked by yellow tapes, whereas the tripod represents the obstacle. The monitoring panel displays the real-time position, the control policy, the motion uncertainty, and the robot status [being in prefix (green) or in suffix (magenta)]. Customizable virtual experiment platform (right) in V-REP for task (19) where no AMECs can be found in the product (see [32] and [40]).

objective property” to find the minimal cumulative reward while ensuring that the risk of violating the task is bounded by γ . Note that the associated model has to be the modified product model \mathcal{Z}_{pre} defined in Section IV-B1 as PRISM does not currently support multiobjective property with the “F target” operator (i.e., $\diamond S_c$). The computation time is approximately the same as in the previous cases. Last, the current PRISM version does not support the mean-payoff optimization in the AMECs, nor does it generate the relaxed control policy for the case where no AMECs exist in the product automaton. In fact, PRISM will simply return that the maximal probability of satisfying the task is 0. The MultiGain tool recently proposed in [34] can handle multiple mean-payoff constraints but does not allow the tuning of the satisfaction probability $(1 - \gamma)$.

VI. EXPERIMENTAL STUDY

In this section, we present an experimental study. We use a differential-driven “iRobot” whose position we track in real time via an Optitrack motion capture system. The communication among the planning module, the robot actuation module, and the Optitrack is handled by the Robot Operating System. The software implementation for this experiment is available in [39]. The experiment videos are online [40].

A. Model Description

Consider the $2.5 \text{ m} \times 1.5 \text{ m}$ experiment workspace as shown in Fig. 11, with three base stations located at the corners and one obstacle region. It consists of 5×3 square cells of dimension $0.5 \text{ m} \times 0.5 \text{ m}$ each. The robot’s motion within the workspace

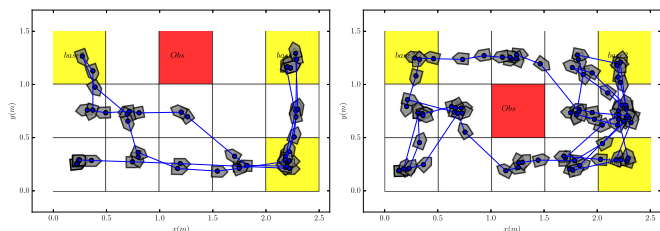


Fig. 12. Robot trajectory to satisfy task (18) (left) and (19) (right) when $\gamma = 0$, sampled at every 15 s.

is abstracted similarly as in Section V-A. The resulting MDP has 60 states and 456 edges.

B. Experimental Results

We consider two different tasks: first, the sequential visiting task (18), and then, the surveillance task (19).

1) Sequential Visiting Task: The LTL task formula is given in (18) and the associated DRA is constructed in Section V-B. The obstacle has probability 0.1 of appearing in the cell (1.25 m, 1.25 m). The resulting product automaton in this case has 532 states, 4228 edges, and 1 accepting pair. For $\gamma = 0$ and $\beta = 0.1$, it took 3.16 s to synthesize the complete policy using Algorithm 1, resulting in an average prefix cost 47.72 and suffix cost 1.0. Then, the robot was controlled in real time using Algorithm 2. The robot state was retrieved using the motion capture system and the observed label was generated randomly. The complete video is online [40] and the resulting trajectory is shown in Fig. 12. Notice that the robot avoids complete collision with the obstacle.

2) Surveillance Task: The LTL task formula is given in (19) and the associated DRA is constructed in Section V-B. The obstacle has probability 0.1 of appearing in the cell (1.25 m, 0.75 cm). The resulting product automaton in this case has 608 states, 4992 edges, and 1 accepting pair.

In the *first* experiment, we choose $\gamma = 0$ and $\beta = 0.1$, so that there is no risk allowed in the plan prefix. It took 5.2 s to synthesize the complete plan offline using Algorithm 1. The real-time execution of the system followed Algorithm 2. The resulting trajectory is shown in Fig. 12. In the *second* experiment, we selected $\gamma = 0.1$ and $\beta = 0.1$ to allow risk in the plan prefix. It took 4.9 s to synthesize the complete policy. Compared to the case where $\gamma = 0$, the optimal policy instructs the robot to move forward, straight to the base station at (2.25 m, 0.25 m), even though there is a risk of colliding with the obstacle at (1.25 m, 0.75 m) due to the uncertainty in its forward action. Both experiment videos are online [40].

Finally, to demonstrate the proposed scheme for much larger workspaces and more complex tasks, particularly when no AMECs can be found in the product automaton, we create a virtual experiment platform based on V-REP [41], which is available in [32]. A snapshot is shown in Fig. 12. The user can easily change the configuration of the workspace and the robot task specification. Once the control policy is synthesized via Algorithm 1 and saved, the user can perform any number of test runs in this environment. Demonstration videos are online [40] where we replicate the surveillance task with clustered

obstacles from Section V-E. It can be seen that the relaxed control policy can ensure high probability of avoiding bad states over long time intervals.

VII. CONCLUSION AND FUTURE WORK

In this paper, we propose a plan synthesis algorithm for probabilistic motion planning, subject to high-level LTL task formulas and risk constraints. Uncertainties in both the robot motion and the workspace properties are considered. We obtain optimal policies that optimize the total cost both in the prefix and suffix of the system trajectory. We also address the case where no AECs exist in the product automaton, in which case, the probability of satisfying the task is zero. The proposed solution provides provable guarantees on the probabilistic satisfiability and the mean total-cost optimality, and is verified via both numerical simulations and experimental studies. Future work involves extensions to multirobot systems.

REFERENCES

- [1] S. Thrun, W. Burgard, and D. Fox, *Probabilistic Robotics*. Cambridge, MA, USA: MIT Press, 2005.
- [2] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. New York, NY, USA: Wiley, 2014.
- [3] G. E. Fainekos, A. Girard, H. Kress-Gazit, and G. J. Pappas, "Temporal logic motion planning for dynamic robots," *Automatica*, vol. 45, no. 2, pp. 343–352, 2009.
- [4] C. Baier and J.-P. Katoen, *Principles of Model Checking*. Cambridge, MA, USA: MIT Press, 2008.
- [5] C. Belta, A. Bicchi, M. Egerstedt, E. Frazzoli, E. Klavins, and G. J. Pappas, "Symbolic planning and control of robot motion," *IEEE Robot. Autom. Mag.*, vol. 14, no. 1, pp. 61–70, Mar. 2007.
- [6] M. Guo, M. Egerstedt, and D. V. Dimarogonas, "Hybrid control of multi-robot systems using embedded graph grammars," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2016, pp. 5242–5247.
- [7] M. Guo and D. V. Dimarogonas, "Multi-agent plan reconfiguration under local LTL specifications," *Int. J. Robot. Res.*, vol. 34, no. 2, pp. 218–235, 2015.
- [8] E. M. Wolff, U. Topcu, and R. M. Murray, "Robust control of uncertain Markov decision processes with temporal logic specifications," in *Proc. IEEE 51st Annu. Conf. Decis. Control*, 2012, pp. 3372–3379.
- [9] M. Lahijanian, S. B. Andersson, and C. Belta, "Formal verification and synthesis for discrete-time stochastic systems," *IEEE Trans. Autom. Control*, vol. 60, no. 8, pp. 2031–2045, Aug. 2015.
- [10] I. Cizelj and C. Belta, "Control of noisy differential-drive vehicles from time-bounded temporal logic specifications," *Int. J. Robot. Res.*, vol. 33, no. 8, pp. 1112–1129, 2014.
- [11] A. Ulusoy, T. Wongpiromsarn, and C. Belta, "Incremental controller synthesis in probabilistic environments with temporal logic constraints," *Int. J. Robot. Res.*, vol. 33, no. 8, pp. 1130–1144, 2014.
- [12] X. Ding, M. Lazar, and C. Belta, "LTL receding horizon control for finite deterministic systems," *Automatica*, vol. 50, no. 2, pp. 399–408, 2014.
- [13] M. Kwiatkowska, G. Norman, and D. Parker, "Prism 4.0: Verification of probabilistic real-time systems," in *Computer Aided Verification*. New York, NY, USA: Springer, 2011, pp. 585–591.
- [14] E. Altman, "Constrained Markov decision processes with total cost criteria: Occupation measures and primal LP," *Math. Methods Oper. Res.*, vol. 43, no. 1, pp. 45–72, 1996.
- [15] D. M. Roijers, P. Vamplew, S. Whiteson, and R. Dazeley, "A survey of multi-objective sequential decision-making," *J. Artif. Intell. Res.*, vol. 48, pp. 67–113, 2013.
- [16] K. Etessami, M. Kwiatkowska, M. Y. Vardi, and M. Yannakakis, "Multi-objective model checking of Markov decision processes," in *Tools and Algorithms for the Construction and Analysis of Systems*. Springer, pp. 50–65, 2007.
- [17] V. Forejt, M. Kwiatkowska, G. Norman, D. Parker, and H. Qu, "Quantitative multi-objective verification for probabilistic systems," in *Tools and Algorithms for the Construction and Analysis of Systems*. New York, NY, USA: Springer, 2011, pp. 112–127.

- [18] V. Forejt, M. Kwiatkowska, and D. Parker, "Pareto curves for probabilistic model checking," in *International Symposium on Automated Technology for Verification and Analysis*. New York, NY, USA: Springer, 2012, pp. 317–332.
- [19] M. Randour, J.-F. Raskin, and O. Sankur, "Variations on the stochastic shortest path problem," in *Proc. Int. Workshop Verification, Model Check. Abstract Interpretation*. New York, NY, USA: Springer, 2015, pp. 1–18.
- [20] M. Randour, J.-F. Raskin, and O. Sankur, "Percentile queries in multi-dimensional Markov decision processes," in *International Conference on Computer Aided Verification*. New York, NY, USA: Springer, 2015, pp. 123–139.
- [21] X. C. Ding, S. L. Smith, C. Belta, and D. Rus, "MDP optimal control under temporal logic constraints," in *Proc. 50th IEEE Conf. Decis. Control*, 2011, pp. 532–538.
- [22] X. Ding, S. L. Smith, C. Belta, and D. Rus, "Optimal control of Markov decision processes with linear temporal logic constraints," *IEEE Trans. Automat. Control*, vol. 59, no. 5, pp. 1244–1257, May 2014.
- [23] S. L. Smith, J. Tumova, C. Belta, and D. Rus, "Optimal path planning for surveillance with temporal-logic constraints," *Int. J. Robot. Res.*, vol. 30, no. 14, pp. 1695–1708, 2011.
- [24] K. Chatterjee and L. Doyen, "Energy and mean-payoff parity Markov decision processes," in *International Symposium on Mathematical Foundations of Computer Science*. New York, NY, USA: Springer, 2011, pp. 206–218.
- [25] J. Fu and U. Topcu, "Pareto efficiency in synthesizing shared autonomy policies with temporal logic constraints," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2015, pp. 361–368.
- [26] V. Bruyere, E. Filiot, M. Randour, and J.-F. Raskin, "Meet your expectations with guarantees: Beyond worst-case synthesis in quantitative games," *Inf. Comput.*, vol. 254, no. 2, pp. 259–295, Jun. 2017.
- [27] R. Dimitrova, J. Fu, and U. Topcu, "Robust optimal policies for Markov decision processes with safety-threshold constraints," in *Proc. IEEE 55th Conf. Decis. Control*, 2016, pp. 7081–7086.
- [28] J. Tumova, G. C. Hall, S. Karaman, E. Frazzoli, and D. Rus, "Least-violating control strategy synthesis with safety rules," in *Proc. Int. Conf. Hybrid Syst., Comput. Control*, 2013.
- [29] R. Ehlers, S. Moaref, and U. Topcu, "Risk-averse ω -regular Markov decision process control," in *Proc. IEEE Conf. Decis. Control*, 2016, pp. 426–433.
- [30] M. Lahijanian and M. Kwiatkowska, "Specification revision for Markov decision processes with optimal trade-off," in *Proc. IEEE Conf. Decis. Control*, 2016, pp. 7411–7418.
- [31] J. Klein, "ltl2dstar-LTL to deterministic Streett and Rabin automata," 2007. [Online]. Available: <http://www.ltl2dstar.de>
- [32] GitHub, Inc., MDP_TG. 2017. [Online]. Available: https://github.com/MengGuo/P_MDP_TG
- [33] T. Brázdil, V. Brozek, K. Chatterjee, V. Forejt, and A. Kucera, "Two views on multiple mean-payoff objectives in Markov decision processes," in *Proc. 26th Annu. IEEE Symp. Logic Comput. Sci.*, 2011, pp. 33–42.
- [34] T. Brázdil, K. Chatterjee, V. Forejt, and A. Kučera, "Multigain: A controller synthesis tool for MDPs with multiple mean-payoff objectives," in *Proc. Int. Conf. Tools Algorithms Construction Anal. Syst.* New York, NY, USA: Springer, 2015, pp. 181–187.
- [35] F. Trevizan, S. Thiébaux, P. Santana, and B. Williams, "Heuristic search in dual space for constrained stochastic shortest path problems," *Proc. Assoc. Adv. Artif. Intell.*, 2016, pp. 326–334.
- [36] Gurobi. 2017. [Online]. Available: <https://www.gurobi.com/>
- [37] G. Dantzig, *Linear Programming and Extensions*. Princeton, NJ, USA: Princeton Univ. Press, 2016.
- [38] Simulation_Videos. [Online]. Available: <https://vimeo.com/169438447>; <https://vimeo.com/169438832>; <https://vimeo.com/174351505>; <https://vimeo.com/175143095>
- [39] GitHub, Inc., Py_iRobot_OptiTrack. [Online]. Available: https://github.com/MengGuo/Py_iRobot_OptiTrack
- [40] Experiment_Videos. [Online]. Available: <https://vimeo.com/180983006>; <https://vimeo.com/180985419>; <https://vimeo.com/180987471>; <https://vimeo.com/222038744>
- [41] E. Rohmer, S. P. Singh, and M. Freese, "V-REP: A versatile and scalable robot simulation framework," in *Proc. IEEE Int. Conf. Intell. Robot. Syst.*, 2013, pp. 1321–1326.



Meng Guo (S'14–M'16) received the M.Sc. degree in system, control, and robotics in 2011, and the Ph.D. degree in electrical engineering in 2016, both from KTH Royal Institute of Technology, Stockholm, Sweden.

He is currently a Postdoctoral Associate with the Department of Mechanical Engineering and Materials Science, Duke University, Durham, NC, USA. His research interests include distributed motion and task planning of multiagent systems and formal control synthesis.



Michael M. Zavlanos (S'05–M'09) received the Diploma in mechanical engineering from the National Technical University of Athens, Athens, Greece, in 2002, and the M.S.E. and Ph.D. degrees in electrical and systems engineering from the University of Pennsylvania, Philadelphia, PA, USA, in 2005 and 2008, respectively.

He is currently an Assistant Professor with the Department of Mechanical Engineering and Materials Science, Duke University, Durham, NC, USA. He also holds a secondary appointment with the Department of Electrical and Computer Engineering and the Department of Computer Science. Prior to joining Duke University, he was an Assistant Professor with the Department of Mechanical Engineering, Stevens Institute of Technology, Hoboken, NJ, USA, and a Postdoctoral Researcher in the GRASP Lab, University of Pennsylvania, Philadelphia, PA, USA. His research interests include a wide range of topics in the emerging discipline of networked systems, with applications in robotic, sensor, and communication networks. He is particularly interested in hybrid solution techniques, on the interface of control theory, distributed optimization, estimation, and networking.

Dr. Zavlanos is the recipient of various awards including the 2014 Naval Research Young Investigator Program Award and the 2011 National Science Foundation Faculty Early Career Development (CAREER) Award.