

2012-01-01

Incremental temporal logic synthesis of control policies for robots interacting with dynamic agents

Tichakorn Wongpiromsarn, Alphan Ulusoy, Calin Belta, Emilio Frazzoli, Daniela Rus. 2012. "Incremental Temporal Logic Synthesis of Control Policies for Robots Interacting with Dynamic Agents." 2012 IEEE/RSJ INTERNATIONAL CONFERENCE ON INTELLIGENT ROBOTS AND SYSTEMS (IROS), pp. 229 - 236 (8).

<https://hdl.handle.net/2144/29849>

"Downloaded from OpenBU. Boston University's institutional repository."

Incremental Temporal Logic Synthesis of Control Policies for Robots Interacting with Dynamic Agents

Tichakorn Wongpiromsarn, Alphan Ulusoy, Calin Belta, Emilio Frazzoli and Daniela Rus

Abstract—We consider the synthesis of control policies from temporal logic specifications for robots that interact with multiple dynamic environment agents. Each environment agent is modeled by a Markov chain whereas the robot is modeled by a finite transition system (in the deterministic case) or Markov decision process (in the stochastic case). Existing results in probabilistic verification are adapted to solve the synthesis problem. To partially address the state explosion issue, we propose an incremental approach where only a small subset of environment agents is incorporated in the synthesis procedure initially and more agents are successively added until we hit the constraints on computational resources. Our algorithm runs in an anytime fashion where the probability that the robot satisfies its specification increases as the algorithm progresses.

I. INTRODUCTION

Temporal logics [1], [2], [3] have been recently employed to precisely express complex behaviors of robots. In particular, given a robot specification expressed as a formula in a temporal logic, control policies that ensure or maximize the probability that the robot satisfies the specification can be automatically synthesized based on exhaustive exploration of the state space [4], [5], [6], [7], [8], [9], [10], [11], [12]. Consequently, the main limitation of existing approaches for synthesizing control policies from temporal logic specifications is almost invariably due to a combinatorial blow up of the state space, commonly known as the state explosion problem.

In many applications, robots need to interact with external, potentially dynamic agents, including human and other robots. As a result, the control policy synthesis problem becomes more computationally complex as more external agents are incorporated in the synthesis procedure. Consider, as an example, the problem where an autonomous vehicle needs to go through a pedestrian crossing while there are multiple pedestrians who are already at or approaching the crossing. The state space of the complete system (i.e., the vehicle and all the pedestrians) grows exponentially with the number of the pedestrians. Hence, given a limited budget of computational resources, solving the control policy synthesis problem with respect to temporal logic specifications may not be feasible when there are a large number of pedestrians.

In this paper, we partially address the aforementioned issue and propose an algorithm for computing a robot control policy in an anytime manner. Our algorithm progressively

computes a sequence of control policies, taking into account only a small subset of the environment agents initially and successively adds more agents to the synthesis procedure in each iteration until the computational resource constraints are exceeded. As opposed to existing incremental synthesis approaches that handle temporal logic specifications where representative robot states are incrementally added to the synthesis procedure [8], we consider incrementally adding representative environment agents instead.

The main contribution of this paper is twofold. First, we propose an anytime algorithm for synthesizing a control policy for a robot interacting with multiple environment agents with the objective of maximizing the probability for the robot to satisfy a given temporal logic specification. Second, an incremental construction of various objects needed to be computed during the synthesis procedure is proposed. Such an incremental construction makes our anytime algorithm more efficient by avoiding unnecessary computation and exploiting the objects computed in the previous iteration. Experimental results show that not only we obtain a reasonable solution much faster, but we are also able to obtain an optimal solution faster than existing approaches.

The rest of the paper is organized as follows: We provide useful definitions and descriptions of the formalisms in the following section. Section III is dedicated to the problem formulation. Section IV provides a complete solution to the control policy synthesis problem for robots that interact with environment agents. Incremental computation of control policies is discussed in Section V. Section VI presents experimental results. Finally, Section VII concludes the paper and discusses future work.

II. PRELIMINARIES

We consider systems that comprise multiple (possibly stochastic) components. In this section, we define the formalisms used in this paper to describe such systems and their desired properties. Throughout the paper, we let X^* , X^ω and X^+ denote the set of finite, infinite and nonempty finite strings, respectively, of a set X .

A. Automata

Definition 1: A deterministic finite automaton (DFA) is a tuple $\mathcal{A} = (Q, \Sigma, \delta, q_{init}, F)$ where

- Q is a finite set of states,
- Σ is a finite set called alphabet,
- $\delta : Q \times \Sigma \rightarrow Q$ is a transition function,
- $q_{init} \in Q$ is the initial state, and
- $F \subseteq Q$ is a set of final states.

T. Wongpiromsarn is with the Singapore-MIT Alliance for Research and Technology, Singapore 117543, Singapore. nok@smart.mit.edu

A. Ulusoy and C. Belta are with Boston University, Boston, MA, USA alphan@bu.edu, cbelta@bu.edu

E. Frazzoli and D. Rus are with the Massachusetts Institute of Technology, Cambridge, MA, USA frazzoli@mit.edu, rus@csail.mit.edu

We use the relation notation, $q \xrightarrow{w} q'$ to denote $\delta(q, w) = q'$.

Consider a finite string $\sigma = \sigma_1\sigma_2\dots\sigma_n \in \Sigma^*$. A *run* for σ in a DFA $\mathcal{A} = (Q, \Sigma, \delta, q_{init}, F)$ is a finite sequence of states $q_0q_1\dots q_n$ such that $q_0 = q_{init}$ and $q_0 \xrightarrow{\sigma_1} q_1 \xrightarrow{\sigma_2} q_2 \xrightarrow{\sigma_3} \dots \xrightarrow{\sigma_n} q_n$. A run is *accepting* if $q_n \in F$. A string $\sigma \in \Sigma^*$ is *accepted* by \mathcal{A} if there is an accepting run of σ in \mathcal{A} . The language *accepted* by \mathcal{A} , denoted by $\mathcal{L}(\mathcal{A})$, is the set of all accepted strings of \mathcal{A} .

B. Linear Temporal Logic

Linear temporal logic (LTL) is a branch of logic that can be used to reason about a time line. An LTL formula is built up from a set Π of atomic propositions, the logic connectives \neg , \vee , \wedge and \implies and the temporal modal operators \bigcirc (“next”), \square (“always”), \diamond (“eventually”) and \mathcal{U} (“until”). An LTL formula over a set Π of atomic propositions is inductively defined as

$$\varphi := \text{True} \mid p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \bigcirc\varphi \mid \varphi \mathcal{U} \varphi$$

where $p \in \Pi$. Other operators can be defined as follows: $\varphi \wedge \psi = \neg(\neg\varphi \vee \neg\psi)$, $\varphi \implies \psi = \neg\varphi \vee \psi$, $\diamond\varphi = \text{True} \mathcal{U} \varphi$, and $\square\varphi = \neg\diamond\neg\varphi$.

Semantics of LTL: LTL formulas are interpreted on infinite strings over 2^Π . Let $\sigma = \sigma_0\sigma_1\sigma_2\dots$ where $\sigma_i \in 2^\Pi$ for all $i \geq 0$. The satisfaction relation \models is defined inductively on LTL formulas as follows:

- $\sigma \models \text{True}$,
- for an atomic proposition $p \in \Pi$, $\sigma \models p$ if and only if $p \in \sigma_0$,
- $\sigma \models \neg\varphi$ if and only if $\sigma \not\models \varphi$,
- $\sigma \models \varphi_1 \wedge \varphi_2$ if and only if $\sigma \models \varphi_1$ and $\sigma \models \varphi_2$,
- $\sigma \models \bigcirc\varphi$ if and only if $\sigma_1\sigma_2\dots \models \varphi$, and
- $\sigma \models \varphi_1 \mathcal{U} \varphi_2$ if and only if there exists $j \geq 0$ such that $\sigma_j\sigma_{j+1}\dots \models \varphi_2$ and for all i such all $0 \leq i < j$, $\sigma_i\sigma_{i+1}\dots \models \varphi_1$.

More details on LTL can be found, e.g., in [1], [2], [3].

In this paper, we are particularly interested in a class of LTL known as co-safety formulas. An important property of a co-safety formula is that any word satisfying the formula has a finite good prefix, i.e., a finite prefix that cannot be extended to violate the formula. Specifically, given an alphabet Σ , a language $L \subseteq \Sigma^\omega$ is co-safety if and only if every $w \in L$ has a good prefix $x \in \Sigma^*$ such that for all $y \in \Sigma^\omega$, we have $x \cdot y \in L$. In general, the problem of determining whether an LTL formula is co-safety is PSPACE-complete [13]. However, there is a class of co-safety formulas, known as *syntactically co-safe* LTL formulas, which can be easily characterized. A *syntactically co-safe* LTL formula over Π is an LTL formula over Π whose only temporal operators are \bigcirc , \diamond and \mathcal{U} when written in positive normal form where the negation operator \neg occurs only in front of atomic propositions [3], [13]. It can be shown that for any syntactically co-safe formula φ , there exists a DFA \mathcal{A}_φ that accepts all and only words in $\text{pref}(\varphi)$, i.e., $\mathcal{L}(\mathcal{A}_\varphi) = \text{pref}(\varphi)$, where $\text{pref}(\varphi)$ denote the set of all good prefixes for φ [9].

C. Systems and Control Policies

We consider the case where each component of the system can be modeled by a deterministic finite transition system, Markov chain or Markov decision process, depending on the characteristics of that component. These different models are defined as follows.

Definition 2: A *deterministic finite transition system* (DFTS) is a tuple $\mathcal{T} = (S, \text{Act}, \longrightarrow, s_{init}, \Pi, L)$ where

- S is a finite set of states,
- Act is a finite set of actions,
- $\longrightarrow \subseteq S \times \text{Act} \times S$ is a transition relation such that for all $s \in S$ and $\alpha \in \text{Act}$, $|\text{Post}(s, \alpha)| \leq 1$ where $\text{Post}(s, \alpha) = \{s' \in S \mid (s, \alpha, s') \in \longrightarrow\}$,
- $s_{init} \in S$ is the initial state,
- Π is a set of atomic propositions, and
- $L : S \rightarrow 2^\Pi$ is a labeling function.

$(s, \alpha, s') \in \longrightarrow$ is denoted by $s \xrightarrow{\alpha} s'$. An action α is *enabled* in state s if and only if there exists s' such that $s \xrightarrow{\alpha} s'$.

Definition 3: A (*discrete-time*) *Markov chain* (MC) is a tuple $\mathcal{M} = (S, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$ where S , Π and L are defined as in DFTS and

- $\mathbf{P} : S \times S \rightarrow [0, 1]$ is the transition probability function such that for any state $s \in S$, $\sum_{s' \in S} \mathbf{P}(s, s') = 1$, and
- $\mathbf{l}_{init} : S \rightarrow [0, 1]$ is the initial state distribution satisfying $\sum_{s \in S} \mathbf{l}_{init}(s) = 1$.

Definition 4: A *Markov decision process* (MDP) is a tuple $\mathcal{M} = (S, \text{Act}, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$ where S , Act , \mathbf{l}_{init} , Π and L are defined as in DFTS and MC and $\mathbf{P} : S \times \text{Act} \times S \rightarrow [0, 1]$ is the transition probability function such that for any state $s \in S$ and action $\alpha \in \text{Act}$, $\sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$.

An action α is *enabled* in state s if and only if $\sum_{s' \in S} \mathbf{P}(s, \alpha, s') = 1$. Let $\text{Act}(s)$ denote the set of enabled actions in s .

Given a complete system as the composition of all its components, we are interested in computing a control policy for the system that optimizes certain objectives. We define a control policy for a system modeled by an MDP as follows.

Definition 5: Let $\mathcal{M} = (S, \text{Act}, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$ be a Markov decision process. A *control policy* for \mathcal{M} is a function $\mathcal{C} : S^+ \rightarrow \text{Act}$ such that $\mathcal{C}(s_0s_1\dots s_n) \in \text{Act}(s_n)$ for all $s_0s_1\dots s_n \in S^+$.

Let $\mathcal{M} = (S, \text{Act}, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$ be an MDP and $\mathcal{C} : S^+ \rightarrow \text{Act}$ be a control policy for \mathcal{M} . Given an initial state s_0 of \mathcal{M} such that $\mathbf{l}_{init}(s_0) > 0$, an infinite sequence $r_{\mathcal{M}}^{\mathcal{C}} = s_0s_1\dots$ on \mathcal{M} generated under policy \mathcal{C} is called a *path* on \mathcal{M} if $\mathbf{P}(s_i, \mathcal{C}(s_0s_1\dots s_i), s_{i+1}) > 0$ for all i . The subsequence $s_0s_1\dots s_n$ where $n \geq 0$ is the *prefix* of length n of $r_{\mathcal{M}}^{\mathcal{C}}$. We define $\text{Paths}_{\mathcal{M}}^{\mathcal{C}}$ and $\text{FPaths}_{\mathcal{M}}^{\mathcal{C}}$ as the set of all infinite paths of \mathcal{M} under policy \mathcal{C} and their finite prefixes, respectively, starting from any state s_0 with $\mathbf{l}_{init}(s_0) > 0$. For $s_0s_1\dots s_n \in \text{FPaths}_{\mathcal{M}}^{\mathcal{C}}$, we let $\text{Paths}_{\mathcal{M}}^{\mathcal{C}}(s_0s_1\dots s_n)$ denote the set of all paths in $\text{Paths}_{\mathcal{M}}^{\mathcal{C}}$ with the prefix $s_0s_1\dots s_n$.

The σ -algebra associated with \mathcal{M} under policy \mathcal{C} is defined as the smallest σ -algebra that contains $\text{Paths}_{\mathcal{M}}^{\mathcal{C}}(\hat{r}_{\mathcal{M}}^{\mathcal{C}})$ where

$\hat{r}_{\mathcal{M}}^{\mathcal{C}}$ ranges over all finite paths in $FPaths_{\mathcal{M}}^{\mathcal{C}}$. It follows that there exists a unique probability measure $Pr_{\mathcal{M}}^{\mathcal{C}}$ on the σ -algebra associated with \mathcal{M} under policy \mathcal{C} where for any $s_0 s_1 \dots s_n \in FPaths_{\mathcal{M}}^{\mathcal{C}}$,

$$Pr_{\mathcal{M}}^{\mathcal{C}}\{Paths_{\mathcal{M}}^{\mathcal{C}}(s_0 s_1 \dots s_n)\} = \mathbb{1}_{init}(s_0) \prod_{0 \leq i < n} \mathbf{P}(s_i, \mathcal{C}(s_0 s_1 \dots s_i), s_{i+1}).$$

Given an LTL formula φ , one can show that the set $\{s_0 s_1 \dots \in Paths_{\mathcal{M}}^{\mathcal{C}} \mid L(s_0)L(s_1) \dots \models \varphi\}$ is measurable [3]. The probability for \mathcal{M} to satisfy φ under policy \mathcal{C} is then defined as

$$Pr_{\mathcal{M}}^{\mathcal{C}}(\varphi) = Pr_{\mathcal{M}}^{\mathcal{C}}\{s_0 s_1 \dots \in Paths_{\mathcal{M}}^{\mathcal{C}} \mid L(s_0)L(s_1) \dots \models \varphi\}.$$

For a given (possibly noninitial) state $s \in S$, we let $\mathcal{M}^s = (S, Act, \mathbf{P}, \mathbb{1}_{init}^s, \Pi, L)$ where $\mathbb{1}_{init}^s(t) = 1$ if $s = t$ and $\mathbb{1}_{init}^s(t) = 0$ otherwise. We define $Pr_{\mathcal{M}}^{\mathcal{C}}(s \models \varphi) = Pr_{\mathcal{M}^s}^{\mathcal{C}}(\varphi)$ as the probability for \mathcal{M} to satisfy φ under policy \mathcal{C} , starting from s .

A control policy essentially resolves all nondeterministic choices in an MDP and induces a Markov chain $\mathcal{M}_{\mathcal{C}}$ that formalizes the behavior of \mathcal{M} under control policy \mathcal{C} [3]. In general, $\mathcal{M}_{\mathcal{C}}$ contains all the states in S^+ and hence may not be finite even though \mathcal{M} is finite. However, for a special case where \mathcal{C} is memoryless, it can be shown that $\mathcal{M}_{\mathcal{C}}$ can be identified with a finite MC.

Definition 6: Let $\mathcal{M} = (S, Act, \mathbf{P}, \mathbb{1}_{init}, \Pi, L)$ be a Markov decision process. A control policy \mathcal{C} on \mathcal{M} is *memoryless* if and only if for each sequence $s_0 s_1 \dots s_n$ and $t_0 t_1 \dots t_m \in S^+$ with $s_n = t_m$, $\mathcal{C}(s_0 s_1 \dots s_n) = \mathcal{C}(t_0 t_1 \dots t_m)$. A memoryless control policy \mathcal{C} can be described by a function $\mathcal{C} : S \rightarrow Act$.

III. PROBLEM FORMULATION

Consider a system that comprises the plant (e.g., the robot) and N independent environment agents. We assume that at any time instance, the state of the system, which incorporates the state of the plant and the environment agents, can be precisely observed. The system can regulate the state of the plant but has no control over the state of the environment agents. Hence, we do not distinguish between a control policy *for the system* and a control policy *for the plant* and refer to them as a control policy in general, as there is no confusion that in both cases, only the state of the plant can be regulated and both the system and the plant can precisely observe the current state of the complete system. Hence, even though a control policy may be implemented on the plant, it may be defined over the state of the complete system.

We assume that each environment agent can be modeled by a finite Markov chain. Let $\mathcal{M}_i = (S_i, \mathbf{P}_i, \mathbb{1}_{init,i}, \Pi_i, L_i)$ be the model of the i th environment agent. The plant is modeled either by a deterministic finite transition system or by a finite Markov decision process, depending on whether each control action leads to a deterministic state transition. We use \mathcal{T} to denote the model of the plant and let $\mathcal{T} = (S_0, Act, \longrightarrow, s_{init,0}, \Pi_0, L_0)$ for the case where \mathcal{T} is a DFTS and $\mathcal{T} = (S_0, Act, \mathbf{P}_0, \mathbb{1}_{init,0}, \Pi_0, L_0)$ for the case where \mathcal{T} is an MDP. For the simplicity of the presentation,

we assume that for all $s \in S_0$, $Act(s) \neq \emptyset$. In addition, we assume that all the components $\mathcal{T}, \mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_N$ in the system make a transition simultaneously, i.e., each of them makes a transition at every time step.

Example 1: Consider a problem where an autonomous vehicle (plant) needs to go through a pedestrian crossing while there are N pedestrians (agents) who are already at or approaching the crossing. Suppose the road is discretized into a finite number of cells c_0, c_2, \dots, c_M . The vehicle is modeled by either a DFTS $\mathcal{T} = (S_0, Act, \longrightarrow, s_{init,0}, \Pi_0, L_0)$ or an MDP $\mathcal{T} = (S_0, Act, \mathbf{P}_0, \mathbb{1}_{init,0}, \Pi_0, L_0)$ whose state $s \in S_0$ describes the cell occupied by the vehicle and whose action $\alpha \in Act$ corresponds to a motion primitive of the vehicle (e.g., stop, accelerate, decelerate). If each motion primitive leads to a deterministic change in the vehicle's state, then \mathcal{T} is a DFTS. Otherwise, \mathcal{T} is an MDP. The motion of the i th pedestrian is modeled by an MC $\mathcal{M}_i = (S_i, \mathbf{P}_i, \mathbb{1}_{init,i}, \Pi_i, L_i)$ whose state $s \in S_i$ describes the cell occupied by the i th pedestrian. The labeling function L_i , $i \in \{0, \dots, N\}$ essentially maps each cell to its label, indexed by the agent ID, i.e., $L_i(c_j) = c_j^i$ for all $j \in \{0, \dots, M\}$.

Control Policy Synthesis Problem: Given a system model described by $\mathcal{T}, \mathcal{M}_1, \dots, \mathcal{M}_N$ and a syntactically co-safe LTL formula φ over $\Pi_0 \cup \Pi_1 \cup \dots \cup \Pi_N$, we want to automatically synthesize a control policy that maximizes the probability for the system to satisfy φ .

Example 2: Consider the autonomous vehicle problem described in Example 1 and the desired property stating that the vehicle does not collide with any pedestrian until it reaches cell c_M (e.g., the other side of the pedestrian crossing). In this case, the specification φ is given by $\varphi = (\neg \bigvee_{i \geq 1, j \geq 0} (c_j^0 \wedge c_j^i)) \mathcal{U} c_M^0$. Using simple logic manipulation, it can be checked that φ is a co-safe LTL formula.

IV. CONTROL POLICY SYNTHESIS

We employ existing results in probabilistic verification and consider the following 3 main steps to solve the control policy synthesis problem defined in Section III:

- 1) Compute the composition of all the system components to obtain the complete system.
- 2) Construct the product MDP.
- 3) Extract an optimal control policy for the product MDP.

In this section, we describe these steps in more detail and discuss their connection to our control policy synthesis problem described in Section III.

A. Parallel Composition of System Components

Assuming that all the components of the system make a transition simultaneously, we first construct the synchronous parallel composition of all the components to obtain the complete system. Synchronous parallel composition of different types of components is defined as follows.

Definition 7: Let $\mathcal{M}_1 = (S_1, \mathbf{P}_1, \mathbb{1}_{init,1}, \Pi_1, L_1)$ and $\mathcal{M}_2 = (S_2, \mathbf{P}_2, \mathbb{1}_{init,2}, \Pi_2, L_2)$ be Markov chains. Their synchronous parallel composition, denoted by $\mathcal{M}_1 \parallel \mathcal{M}_2$, is the MC $\mathcal{M} = (S_1 \times S_2, \mathbf{P}, \mathbb{1}_{init}, \Pi_1 \cup \Pi_2, L)$ where:

- For each $s_1, s'_1 \in S_1$ and $s_2, s'_2 \in S_2$, $\mathbf{P}(\langle s_1, s_2 \rangle, \langle s'_1, s'_2 \rangle) = \mathbf{P}_1(s_1, s'_1)\mathbf{P}_2(s_2, s'_2)$.
- For each $s_1 \in S_1$ and $s_2 \in S_2$, $\mathbf{l}_{init}(\langle s_1, s_2 \rangle) = \mathbf{l}_{init,1}(s_1)\mathbf{l}_{init,2}(s_2)$.
- For each $s_1 \in S_1$ and $s_2 \in S_2$, $L(\langle s_1, s_2 \rangle) = L(s_1) \cup L(s_2)$.

Definition 8: Let $\mathcal{T}_1 = (S_1, Act, \longrightarrow, s_{init}, \Pi_1, L_1)$ be a deterministic finite transition system and $\mathcal{M}_2 = (S_2, \mathbf{P}_2, \mathbf{l}_{init,2}, \Pi_2, L_2)$ be a Markov chain. Their synchronous parallel composition, denoted by $\mathcal{T}_1 || \mathcal{M}_2$, is the MDP $\mathcal{M} = (S_1 \times S_2, Act, \mathbf{P}, \mathbf{l}_{init}, \Pi_1 \cup \Pi_2, L)$ where:

- For each $s_1, s'_1 \in S_1$, $s_2, s'_2 \in S_2$ and $\alpha \in Act$, $\mathbf{P}(\langle s_1, s_2 \rangle, \alpha, \langle s'_1, s'_2 \rangle) = \mathbf{P}_2(s_2, s'_2)$ if $s_1 \xrightarrow{\alpha} s'_1$ and $\mathbf{P}(\langle s_1, s_2 \rangle, \alpha, \langle s'_1, s'_2 \rangle) = 0$ otherwise.
- For each $s_2 \in S_2$, $\mathbf{l}_{init}(\langle s_{init}, s_2 \rangle) = \mathbf{l}_{init,2}(s_2)$ and $\mathbf{l}_{init}(\langle s_1, s_2 \rangle) = 0$ for all $s_1 \in S \setminus \{s_{init}\}$.
- For each $s_1 \in S_1$ and $s_2 \in S_2$, $L(\langle s_1, s_2 \rangle) = L(s_1) \cup L(s_2)$.

Definition 9: Let $\mathcal{M}_1 = (S_1, Act, \mathbf{P}_1, \mathbf{l}_{init,1}, \Pi_1, L_1)$ be a Markov decision process and $\mathcal{M}_2 = (S_2, \mathbf{P}_2, \mathbf{l}_{init,2}, \Pi_2, L_2)$ be a Markov chain. Their synchronous parallel composition, denoted by $\mathcal{M}_1 || \mathcal{M}_2$, is the MDP $\mathcal{M} = (S_1 \times S_2, Act, \mathbf{P}, \mathbf{l}_{init}, \Pi_1 \cup \Pi_2, L)$ where:

- For each $s_1, s'_1 \in S_1$, $s_2, s'_2 \in S_2$ and $\alpha \in Act$, $\mathbf{P}(\langle s_1, s_2 \rangle, \alpha, \langle s'_1, s'_2 \rangle) = \mathbf{P}_1(s_1, \alpha, s'_1)\mathbf{P}_2(s_2, s'_2)$.
- For each $s_1 \in S_1$ and $s_2 \in S_2$, $\mathbf{l}_{init}(\langle s_1, s_2 \rangle) = \mathbf{l}_{init,1}(s_1)\mathbf{l}_{init,2}(s_2)$.
- For each $s_1 \in S_1$ and $s_2 \in S_2$, $L(\langle s_1, s_2 \rangle) = L(s_1) \cup L(s_2)$.

From the above definitions, our complete system can be modeled by the MDP $\mathcal{T} || \mathcal{M}_1 || \dots || \mathcal{M}_N$, regardless of whether \mathcal{T} is a DFTS or an MDP. We denote this MDP by $\mathcal{M} = (S, Act, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$.

B. Construction of Product MDP

Let $\mathcal{A}_\varphi = (Q, 2^\Pi, \delta, q_{init}, F)$ be a DFA that recognizes the good prefixes of φ . Such \mathcal{A}_φ can be automatically constructed using existing tools [14]. Our next step is to obtain a finite MDP $\mathcal{M}_p = (S_p, Act_p, \mathbf{P}_p, \mathbf{l}_{p,init}, Q, L_p)$ as the product of \mathcal{M} and \mathcal{A}_φ , defined as follows.

Definition 10: Let $\mathcal{M} = (S, Act, \mathbf{P}, \mathbf{l}_{init}, \Pi, L)$ be an MDP and let $\mathcal{A} = (Q, 2^\Pi, \delta, q_{init}, F)$ be a DFA. The product of \mathcal{M} and \mathcal{A} is the MDP $\mathcal{M}_p = \mathcal{M} \otimes \mathcal{A}$ defined by¹ $\mathcal{M}_p = (S_p, Act, \mathbf{P}_p, \mathbf{l}_{p,init}, \Pi, L_p)$ where $S_p = S \times Q$ and $L_p(\langle s, q \rangle) = L(s)$. \mathbf{P}_p is defined as

$$\mathbf{P}_p(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \begin{cases} \tilde{\mathbf{P}}_p(\langle s, q \rangle, \alpha, \langle s', q' \rangle) & \text{if } q' = \delta(q, L(s')) \\ 0 & \text{otherwise} \end{cases}, \quad (1)$$

where $\tilde{\mathbf{P}}_p(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \mathbf{P}(s, \alpha, s')$. For the rest of the paper, we refer to $\tilde{\mathbf{P}}_p : S_p \times Act \times S_p \rightarrow [0, 1]$ as the

intermediate transition probability function for \mathcal{M}_p . Finally,

$$\mathbf{l}_{p,init}(\langle s, q \rangle) = \begin{cases} \tilde{\mathbf{l}}_{p,init}(\langle s, q \rangle) & \text{if } q = \delta(q_{init}, L(s)) \\ 0 & \text{otherwise} \end{cases}, \quad (2)$$

where $\tilde{\mathbf{l}}_{p,init}(\langle s, q \rangle) = \mathbf{l}_{init}(s)$. For the rest of the paper, we refer to $\tilde{\mathbf{l}}_{p,init} : S_p \rightarrow [0, 1]$ as the intermediate initial state distribution for \mathcal{M}_p .

Stepping through the above definition shows that given a path $r_{\mathcal{M}_p}^{\mathcal{C}_p} = \langle s_0, q_0 \rangle \langle s_1, q_1 \rangle \dots$ on \mathcal{M}_p generated under some control policy \mathcal{C}_p , the corresponding path $s_0 s_1 \dots$ on \mathcal{M} generates a word $L(s_0)L(s_1) \dots$ that satisfies φ if and only if there exists $n \geq 0$ such that $q_n \in F$ (and hence $q_0 q_1 \dots q_n$ is an accepting run on \mathcal{A}_φ), in which case we say that $r_{\mathcal{M}_p}^{\mathcal{C}_p}$ is accepting. Therefore, each accepting path of \mathcal{M}_p uniquely corresponds to a path of \mathcal{M} whose word satisfies φ . In addition, a control policy \mathcal{C}_p on \mathcal{M}_p induces the corresponding control policy \mathcal{C} on \mathcal{M} . The details for generating \mathcal{C} from \mathcal{C}_p can be found, e.g. in [3], [10].

Based on this argument, our control policy synthesis problem defined in Section III can be reduced to computing a control policy for \mathcal{M}_p that maximizes the probability of reaching a state in $B_p = \{\langle s, q \rangle \in S_p \mid q \in F\}$.

C. Control Policy Synthesis for Product MDP

For each $s \in S_p$, let x_s denote the maximum probability of reaching a state in B_p , starting from s . Formally, $x_s = \sup_{\mathcal{C}} \Pr_{\mathcal{M}_p}^{\mathcal{C}}(s \models \diamond B_p)$, where, with an abuse of notation, $B_p \in \diamond B_p$ is a proposition that is satisfied by all states in B_p . There are two main techniques for computing the probability x_s for each $s \in S_p$: linear programming (LP) and value iteration. LP-based techniques yield an exact solution but it typically does not scale as well as value iteration. On the other hand, value iteration is an iterative numerical technique. This method works by successively computing the probability vector $(x_s^{(k)})_{s \in S_p}$ for increasing $k \geq 0$ such that $\lim_{k \rightarrow \infty} x_s^{(k)} = x_s$ for all $s \in S_p$. Initially, we set $x_s^{(0)} = 1$ if $s \in B_p$ and $x_s^{(0)} = 0$ otherwise. In the $(k+1)$ th iteration where $k \geq 0$, we set

$$x_s^{(k+1)} = \begin{cases} 1 & \text{if } s \in B_p \\ \max_{\alpha \in Act_p(s)} \sum_{t \in S_p} \mathbf{P}_p(s, \alpha, t) x_t^{(k)} & \text{otherwise.} \end{cases} \quad (3)$$

In practice, we terminate the computation and say that $x_s^{(k)}$ converges when a termination criterion such as $\max_{s \in S_p} |x_s^{(k+1)} - x_s^{(k)}| < \epsilon$ is satisfied for some fixed (typically very small) threshold ϵ .

As discussed in [15], [16], decomposition of \mathcal{M}_p into strongly connected components (SCC) can help speed up value iteration. $C \subseteq S_p$ is an SCC of \mathcal{M}_p if there is a path in \mathcal{M}_p between any two states in C and C is maximal (i.e., there does not exist any $\tilde{C} \subseteq S_p$ such that $C \subset \tilde{C}$ and \tilde{C} is an SCC). The algorithm proposed in [17] allows us to identify all the SCCs of \mathcal{M}_p with time and space complexity that is linear in the size of \mathcal{M}_p .

¹We slightly modify the definition of atomic propositions and labeling function of the product MDP from the definition often used in literature to facilitate incremental construction of product MDP, which is explained in Section V-B.

The SCC-based value iteration works as follows. First, we set $x_s^{(0)} = 1$ if $s \in B_p$ and $x_s^{(0)} = 0$ otherwise.² Next, we identify all the SCCs $C_1^{\mathcal{M}_p}, \dots, C_m^{\mathcal{M}_p}$ of \mathcal{M}_p . From the definition of SCC, we get that $C_i^{\mathcal{M}_p} \cap C_j^{\mathcal{M}_p} = \emptyset, \forall i \neq j$ and $\bigcup_i C_i^{\mathcal{M}_p} = S_p$. For each SCC $C_i^{\mathcal{M}_p}$, we define $Succ(C_i^{\mathcal{M}_p}) \subseteq S_p \setminus C_i^{\mathcal{M}_p}$ to be the set of all the immediate successors of states in $C_i^{\mathcal{M}_p}$ that are not in $C_i^{\mathcal{M}_p}$. A (strict) partial order, $\prec_{\mathcal{M}_p}$, among $C_1^{\mathcal{M}_p}, \dots, C_m^{\mathcal{M}_p}$ can be defined such that $C_j^{\mathcal{M}_p} \prec_{\mathcal{M}_p} C_i^{\mathcal{M}_p}$ if $Succ(C_i^{\mathcal{M}_p}) \cap C_j^{\mathcal{M}_p} \neq \emptyset$. (Note that from the definition of SCC and $Succ$, there cannot be cyclic dependency among SCCs; hence, such a partial order can always be defined.)

An important property of SCCs and their partial order that we will exploit in the computation of the probability vector $(x_s)_{s \in S_p}$ is that the probability values of states in $C_i^{\mathcal{M}_p}$ can be affected only by the probability values of states in $C_i^{\mathcal{M}_p}$ and all $C_j^{\mathcal{M}_p} \prec_{\mathcal{M}_p} C_i^{\mathcal{M}_p}$. Thus, our next step is to generate an order $\mathbb{O}^{\mathcal{M}_p}$ among $C_1^{\mathcal{M}_p}, \dots, C_m^{\mathcal{M}_p}$ such that $C_i^{\mathcal{M}_p}$ appears before $C_j^{\mathcal{M}_p}$ in $\mathbb{O}^{\mathcal{M}_p}$ if $C_i^{\mathcal{M}_p} \prec_{\mathcal{M}_p} C_j^{\mathcal{M}_p}$. We can then process each SCC separately, according to the order in $\mathbb{O}^{\mathcal{M}_p}$, since the probability values of states in $C_j^{\mathcal{M}_p}$ that appears after $C_i^{\mathcal{M}_p}$ in $\mathbb{O}^{\mathcal{M}_p}$ cannot affect the probability values of states in $C_i^{\mathcal{M}_p}$. Processing of SCC $C_i^{\mathcal{M}_p}$ terminates at the k th iteration where all $x_s^{(k)}, s \in C_i^{\mathcal{M}_p}$ converges. Let x_s be the value to which $x_s^{(k)}$ converges. When processing $C_i^{\mathcal{M}_p}$, we exploit the order in $\mathbb{O}^{\mathcal{M}_p}$ and existing values of x_t for all $t \in Succ(C_i^{\mathcal{M}_p})$ to determine the set of $s \in C_i^{\mathcal{M}_p}$ where $x_s^{(k+1)}$ needs to be updated from $x_s^{(k)}$. The formula in (3) with $x_t^{(k)}$ replaced by x_t for all $t \in Succ(C_i^{\mathcal{M}_p})$ can be used to update those $x_s^{(k+1)}$. We refer the reader to [15], [16] for more details.

Note that computation of an order $\mathbb{O}^{\mathcal{M}_p}$ requires $O(|S_p|^2)$ time. Thus, the pre-computation required by the SCC-based value iteration can be computationally expensive, unless all the SCCs of \mathcal{M}_p and an order $\mathbb{O}^{\mathcal{M}_p}$ are provided a-priori. As a result, the SCC-based value iteration may require more computation time than the normal value iteration, if the pre-computation time is also taken into account.

Once the vector $(x_s)_{s \in S_p}$ is computed, a memoryless control policy \mathcal{C} such that for any $s \in S_p$, $\Pr_{\mathcal{M}}^{\mathcal{C}}(s \models \diamond B_p) = x_s$ can be constructed as follows. For each state $s \in S_p$, let $Act_p^{max}(s) \subseteq Act_p(s)$ be the set of actions such that for all $\alpha \in Act_p^{max}(s)$, $x_s = \sum_{t \in S_p} \mathbf{P}(s, \alpha, t) x_t$. For each $s \in S_p$ with $x_s > 0$, let $\|s\|$ be the length of a shortest path from s to a state in B_p , using only actions in Act_p^{max} . $\mathcal{C}(s) \in Act_p^{max}(s)$ for a state $s \in S_p \setminus B_p$ with $x_s > 0$ is then chosen such that $\mathbf{P}_p(s, \mathcal{C}(s), t) > 0$ for some $t \in S_p$ with $\|t\| = \|s\| - 1$. For a state $s \in S_p$ with $x_s = 0$ or a state $s \in B_p$, $\mathcal{C}(s) \in Act_p(s)$ can be chosen arbitrarily.

²In the original algorithm, all the states $s \in S_p$ with $x_s = 1$ and all the states that cannot reach B_p under any control policy need to be identified but it has been shown in [16] that this step is not necessary for the correctness of the algorithm.

V. INCREMENTAL COMPUTATION OF CONTROL POLICIES

Automatic synthesis described in the previous section suffers from the state explosion problem as the composition of \mathcal{T} and all $\mathcal{M}_1, \dots, \mathcal{M}_N$ needs to be constructed, leading to an exponential blow up of the state space. In this section, we propose an incremental synthesis approach where we progressively compute a sequence of control policies, taking into account only a small subset of the environment agents initially and successively add more agents to the synthesis procedure in each iteration until we hit the computational resource constraints. Hence, even though the complete synthesis problem cannot be solved due to the computational resource limitation, we can still obtain a reasonably good control policy.

A. Overview of Incremental Computation of Control Policies

Initially, we consider a small subset $\mathbf{M}_0 \subset \{\mathcal{M}_1, \dots, \mathcal{M}_N\}$ of the environment agents. For each $\mathcal{M}_i = (S_i, \mathbf{P}_i, i_{init,i}, \Pi_i, L_i) \notin \mathbf{M}_0$, we consider a simplified model $\tilde{\mathcal{M}}_i$ that essentially assumes that the i th environment agent is stationary (i.e., we take into account their presence but do not consider their full model). Formally, $\tilde{\mathcal{M}}_i = (\{s_i\}, \tilde{\mathbf{P}}_i, \tilde{i}_{init,i}, \Pi_i, \tilde{L}_i)$ where $s_i \in S_i$ can be chosen arbitrarily, $\tilde{\mathbf{P}}_i(s_i, s_i) = 1$, $\tilde{i}_{init,i}(s_i) = 1$ and $\tilde{L}_i(s_i) = L_i(s_i)$. Note that the choice of $s_i \in S_i$ may affect the performance of our incremental synthesis algorithm; hence, it should be chosen such that it is the most likely state of \mathcal{M}_i . We let $\tilde{\mathbf{M}}_0 = \{\tilde{\mathcal{M}}_i \mid \mathcal{M}_i \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_0\}$.

The composition of \mathcal{T} , all $\mathcal{M}_i \in \mathbf{M}_0$ and all $\tilde{\mathcal{M}}_j \in \tilde{\mathbf{M}}_0$ is then constructed. We let $\mathcal{M}^{\mathbf{M}_0}$ be the MDP that represents such composition. Note that since $\tilde{\mathcal{M}}_i$ is typically smaller \mathcal{M}_i , $\mathcal{M}^{\mathbf{M}_0}$ is typically much smaller than the composition of $\mathcal{T}, \mathcal{M}_1, \dots, \mathcal{M}_N$. We identify all the SCCs of $\mathcal{M}^{\mathbf{M}_0}$ and their partial order. Following the steps for synthesizing a control policy described in Section IV, we construct $\mathcal{M}_p^{\mathbf{M}_0} = \mathcal{M}^{\mathbf{M}_0} \otimes \mathcal{A}_\varphi$ where $\mathcal{A}_\varphi = (Q, 2^\Pi, \delta, q_{init}, F)$ is a DFA that recognizes the good prefixes of φ . We also store the intermediate transition probability function and the intermediate initial state distribution for $\mathcal{M}_p^{\mathbf{M}_0}$ and denote these functions by $\tilde{\mathbf{P}}_p^{\mathbf{M}_0}$ and $\tilde{i}_{p,init}^{\mathbf{M}_0}$, respectively.

At the end of the initialization period (i.e., the 0th iteration), we obtain a control policy $\mathcal{C}^{\mathbf{M}_0}$ that maximizes the probability for $\mathcal{M}^{\mathbf{M}_0}$ to satisfy φ . $\mathcal{C}^{\mathbf{M}_0}$ resolves all nondeterministic choices in $\mathcal{M}^{\mathbf{M}_0}$ and induces a Markov chain, which we denote by $\mathcal{M}_{\mathbf{C}^{\mathbf{M}_0}}^{\mathbf{M}_0}$.

Our algorithm then successively adds more full models of the rest of the environment agents to the synthesis procedure at each iteration. In the $(k+1)$ th iteration where $k \geq 0$, we consider $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$ for some $\mathcal{M}_l \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_k$. Such \mathcal{M}_l may be picked such that the probability for $\mathcal{M}_{\mathbf{C}^{\mathbf{M}_0}}^{\mathbf{M}_0} \parallel \mathcal{M}_l$ to satisfy φ is the minimum among all $\mathcal{M}_i \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_k$. This probability can be efficiently computed using probabilistic verification [3]. (As an MC can be considered a special case of MDP with exactly one action enabled in each state, we can easily adapt the techniques for computing the probability vector of a product MDP described in Section IV-C to compute the

probability that $\mathcal{M}_{\mathcal{C}_{M_0}^{M_0}} \parallel \mathcal{M}_l$ satisfies φ .) We let $\tilde{\mathbf{M}}_{k+1} = \tilde{\mathbf{M}}_k \setminus \{\tilde{\mathcal{M}}_l\}$ and let $\mathcal{M}^{\mathbf{M}_{k+1}}$ be the MDP that represents the composition of \mathcal{T} , all $\mathcal{M}_i \in \mathbf{M}_{k+1}$ and all $\tilde{\mathcal{M}}_j \in \tilde{\mathbf{M}}_{k+1}$. Next, we construct $\mathcal{M}_p^{\mathbf{M}_{k+1}} = \mathcal{M}^{\mathbf{M}_{k+1}} \otimes \mathcal{A}_\varphi$ and obtain a control policy $\mathcal{C}^{\mathbf{M}_{k+1}}$ that maximizes the probability for $\mathcal{M}^{\mathbf{M}_{k+1}}$ to satisfy φ . Similar to the initialization step, during the construction of $\mathcal{M}_p^{\mathbf{M}_{k+1}}$, we store the intermediate transition probability function and the intermediate initial state distribution for $\mathcal{M}_p^{\mathbf{M}_{k+1}}$ and denote these functions by $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}$ and $\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_{k+1}}$, respectively.

The process outlined in the previous paragraph terminates at the K th iteration where $\mathbf{M}_K = \{\mathcal{M}_1, \dots, \mathcal{M}_N\}$ or when the computational resource constraints are exceeded. To make this process more efficient, we avoid unnecessary computation and exploit the objects computed in the previous iteration. Consider an arbitrary iteration $k \geq 0$. In Section V-B, we show how $\mathcal{M}_p^{\mathbf{M}_{k+1}}$, $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}$, and $\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_{k+1}}$ can be incrementally constructed from $\mathcal{M}_p^{\mathbf{M}_k}$, $\tilde{\mathbf{P}}_p^{\mathbf{M}_k}$ and $\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_k}$. Hence, we can avoid computing $\mathcal{M}^{\mathbf{M}_{k+1}}$. In addition, as previously discussed in Section IV-C, generating an order of SCCs can be computationally expensive. Hence, we only compute the SCCs and their order for $\mathcal{M}^{\mathbf{M}_0}$ and all $\mathcal{M}_j \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_0$, which are typically small. Incremental construction of SCCs of $\mathcal{M}^{\mathbf{M}_{k+1}}$ and their order from those of $\mathcal{M}^{\mathbf{M}_k}$ is considered in Section V-C. (Note that we do not compute $\mathcal{M}^{\mathbf{M}_k}$ but only maintain its SCCs and their order, which are incrementally constructed using the results from the previous iteration.) Finally, Section V-D describes computation of $\mathcal{C}^{\mathbf{M}_k}$, using a method adapted from SCC-based value iteration where we avoid having to identify the SCCs of $\mathcal{M}_p^{\mathbf{M}_k}$ and their order. Instead, we exploit the SCCs of $\mathcal{M}^{\mathbf{M}_k}$ and their order, which can be incrementally constructed using the approach described in Section V-C.

B. Incremental Construction of Product MDP

For an iteration $k \geq 0$, let $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$ for some $\mathcal{M}_l \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_k$. In general, one can construct $\mathcal{M}_p^{\mathbf{M}_{k+1}}$ by first computing $\mathcal{M}^{\mathbf{M}_{k+1}}$, which requires taking the composition of a DFTS or an MDP with N MCs, and then constructing $\mathcal{M}_p^{\mathbf{M}_{k+1}} \otimes \mathcal{A}_\varphi$. To accelerate the process of computing $\mathcal{M}_p^{\mathbf{M}_{k+1}}$, we exploit the presence of $\mathcal{M}_p^{\mathbf{M}_k}$, its intermediate transition probability function $\tilde{\mathbf{P}}_p^{\mathbf{M}_k}$ and intermediate initial state distribution $\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_k}$, which are computed in the previous iteration.

First, note that a state s_p of $\mathcal{M}_p^{\mathbf{M}_k}$ is of the form $s_p = \langle s, q \rangle$ where $s = \langle s_0, s_1, \dots, s_N \rangle \in S_0 \times S_1 \times \dots \times S_N$ and $q \in Q$. For $s = \langle s_0, s_1, \dots, s_N \rangle \in S_0 \times S_1 \times \dots \times S_N$, $i \in \{0, \dots, N\}$ and $r \in S_i$, we define $s|_{i \leftarrow r} = \langle s_0, \dots, s_{i-1}, r, s_{i+1}, \dots, s_N \rangle$, i.e., $s|_{i \leftarrow r}$ is obtained by replacing the i th element of s by r .

Lemma 1: Consider an arbitrary iteration $k \geq 0$. Let $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$ where $\mathcal{M}_l \in \{\mathcal{M}_1, \dots, \mathcal{M}_N\} \setminus \mathbf{M}_k$. Suppose $\mathcal{M}_p^{\mathbf{M}_k} = (S_p^{\mathbf{M}_k}, Act_p^{\mathbf{M}_k}, \mathbf{P}_p^{\mathbf{M}_k}, \mathbf{1}_{p,init}^{\mathbf{M}_k}, \Pi_p^{\mathbf{M}_k}, L_p^{\mathbf{M}_k})$ and $\mathcal{M}_l = (S_l, \mathbf{P}_l, \mathbf{1}_{init,l}, \Pi_l, L_l)$. Assuming that for any $i, j \in \{0, \dots, N\}$, $\Pi_i \cap \Pi_j = \emptyset$, then $\mathcal{M}_p^{\mathbf{M}_{k+1}} = (S_p^{\mathbf{M}_{k+1}}, Act_p^{\mathbf{M}_{k+1}}, \mathbf{P}_p^{\mathbf{M}_{k+1}}, \mathbf{1}_{p,init}^{\mathbf{M}_{k+1}}, \Pi_p^{\mathbf{M}_{k+1}}, L_p^{\mathbf{M}_{k+1}})$ where

$$S_p^{\mathbf{M}_{k+1}} = \{\langle s|_{l \leftarrow r}, q \rangle \mid \langle s, q \rangle \in S_p^{\mathbf{M}_k} \text{ and } r \in S_l\},$$

$$Act_p^{\mathbf{M}_{k+1}} = Act_p^{\mathbf{M}_k}, \Pi_p^{\mathbf{M}_{k+1}} = \Pi_p^{\mathbf{M}_k}, \text{ and for any } s = \langle s_0, \dots, s_N \rangle, s' = \langle s'_0, \dots, s'_N \rangle \in S_0 \times \dots \times S_N \text{ and } q, q' \in Q,$$

$$\bullet \mathbf{P}_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \begin{cases} \tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) & \text{if } q' = \delta(q, L_p^{\mathbf{M}_{k+1}}(\langle s', q' \rangle)) \\ 0 & \text{otherwise} \end{cases},$$

where the intermediate transition probability function is given by

$$\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \mathbf{P}_l(s_l, s'_l) \tilde{\mathbf{P}}_p^{\mathbf{M}_k}(\langle \tilde{s}, q \rangle, \alpha, \langle \tilde{s}', q' \rangle) \quad (4)$$

for any $\langle \tilde{s}, q \rangle, \langle \tilde{s}', q' \rangle \in S_p^{\mathbf{M}_k}$ such that $\tilde{s}|_{l \leftarrow s_l} = s$ and $\tilde{s}'|_{l \leftarrow s'_l} = s'$,

$$\bullet \mathbf{1}_{p,init}^{\mathbf{M}_{k+1}}(\langle s, q \rangle) = \begin{cases} \tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_{k+1}}(\langle s, q \rangle) & \text{if } q = \delta(q_{init}, L_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle)) \\ 0 & \text{otherwise} \end{cases}$$

where the intermediate initial state distribution is given by

$$\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_{k+1}}(\langle s, q \rangle) = \mathbf{1}_{init,l}(s_l) \tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_k}(\langle \tilde{s}, q \rangle) \quad (5)$$

for any $\langle \tilde{s}, q \rangle \in S_p^{\mathbf{M}_k}$ such that $\tilde{s}|_{l \leftarrow s_l} = s$, and

$$\bullet L_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle) = (L_p^{\mathbf{M}_k}(\langle \tilde{s}, q \rangle) \setminus L_l(\tilde{s}_l)) \cup L_l(s_l) \text{ for any } \langle \tilde{s}, q \rangle \in S_p^{\mathbf{M}_k} \text{ such that } \tilde{s}|_{l \leftarrow s_l} = s.$$

Proof: The correctness of $S_p^{\mathbf{M}_{k+1}}$, $Act_p^{\mathbf{M}_{k+1}}$, $\Pi_p^{\mathbf{M}_{k+1}}$ and $L_p^{\mathbf{M}_{k+1}}$ is straightforward to verify. Hence, we will only provide the proof for the correctness of $\mathbf{P}_p^{\mathbf{M}_{k+1}}$ and $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}$. The correctness of $\mathbf{1}_{p,init}^{\mathbf{M}_{k+1}}$ and $\tilde{\mathbf{1}}_{p,init}^{\mathbf{M}_{k+1}}$ can be proved in a similar way.

Consider an arbitrary iteration $k \geq 0$ and let $\mathcal{M}^{\mathbf{M}_k} = (S^{\mathbf{M}_k}, Act^{\mathbf{M}_k}, \mathbf{P}^{\mathbf{M}_k}, \mathbf{1}_{init}^{\mathbf{M}_k}, \Pi^{\mathbf{M}_k}, L^{\mathbf{M}_k})$ and $\mathcal{M}^{\mathbf{M}_{k+1}} = (S^{\mathbf{M}_{k+1}}, Act^{\mathbf{M}_{k+1}}, \mathbf{P}^{\mathbf{M}_{k+1}}, \mathbf{1}_{init}^{\mathbf{M}_{k+1}}, \Pi^{\mathbf{M}_{k+1}}, L^{\mathbf{M}_{k+1}})$. It is obvious from the definition of product MDP that $\mathbf{P}_p^{\mathbf{M}_{k+1}}$ is correct as long as $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}$ is correct, i.e., $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \mathbf{P}^{\mathbf{M}_{k+1}}(s, \alpha, s')$ for all $\langle s, q \rangle, \langle s', q' \rangle \in S_p^{\mathbf{M}_{k+1}}$ and $\alpha \in Act_p^{\mathbf{M}_{k+1}}$. Hence, we only need to prove the correctness of $\tilde{\mathbf{P}}_p^{\mathbf{M}_{k+1}}$.

Assume that $\tilde{\mathbf{P}}_p^{\mathbf{M}_k}$ is correct, i.e., $\tilde{\mathbf{P}}_p^{\mathbf{M}_k}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) = \mathbf{P}^{\mathbf{M}_k}(s, \alpha, s')$ for all $\langle s, q \rangle, \langle s', q' \rangle \in S_p^{\mathbf{M}_k}$ and $\alpha \in Act_p^{\mathbf{M}_k}$. Let l be the index such that $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$. Consider arbitrary $\langle s, q \rangle, \langle s', q' \rangle \in S_p^{\mathbf{M}_{k+1}}$ and $\alpha \in Act_p^{\mathbf{M}_{k+1}}$. Suppose $s = \langle s_0, \dots, s_N \rangle$ and $s' = \langle s'_0, \dots, s'_N \rangle$. Note that since $\tilde{\mathcal{M}}_l$ only contains one state, there exists exactly one $\langle \tilde{s}, q \rangle \in S_p^{\mathbf{M}_k}$ and exactly one $\langle \tilde{s}', q' \rangle \in S_p^{\mathbf{M}_k}$ such that $\tilde{s}|_{l \leftarrow s_l} = s$ and $\tilde{s}'|_{l \leftarrow s'_l} = s'$. Since \mathbf{M}_k is the composition of \mathcal{T} , all $\mathcal{M}_i \in \mathbf{M}_k$ and all $\tilde{\mathcal{M}}_j \in \tilde{\mathbf{M}}_k$ and since $\mathcal{M}_l \notin \mathbf{M}_k$ and $\tilde{\mathbf{P}}_l(\cdot, \cdot) = 1$, it follows that if \mathcal{T} is a DFTS, then

$$\mathbf{P}^{\mathbf{M}_k}(\tilde{s}, \alpha, \tilde{s}') = \begin{cases} \prod_{i \in \{1, \dots, N\} \setminus \{l\}} \mathbf{P}_i(s_i, s'_i) & \text{if } s_0 \xrightarrow{\alpha} s'_0 \\ 0 & \text{otherwise} \end{cases},$$

and if \mathcal{T} is an MDP, then

$$\mathbf{P}^{\mathbf{M}_k}(\tilde{s}, \alpha, \tilde{s}') = \mathbf{P}_0(s_0, \alpha, s'_0) \prod_{i \in \{1, \dots, N\} \setminus \{l\}} \mathbf{P}_i(s_i, s'_i).$$

Thus, $\mathbf{P}^{\mathbf{M}^{k+1}}(s, \alpha, s') = \mathbf{P}_l(s_l, s'_l) \mathbf{P}^{\mathbf{M}^k}(\tilde{s}, \alpha, \tilde{s}')$. Combining this with (4), we get

$$\begin{aligned} & \tilde{\mathbf{P}}_p^{\mathbf{M}^{k+1}}(\langle s, q \rangle, \alpha, \langle s', q' \rangle) \\ &= \mathbf{P}_l(s_l, s'_l) \tilde{\mathbf{P}}_p^{\mathbf{M}^k}(\langle \tilde{s}, q \rangle, \alpha, \langle \tilde{s}', q' \rangle) \\ &= \mathbf{P}_l(s_l, s'_l) \mathbf{P}^{\mathbf{M}^k}(\tilde{s}, \alpha, \tilde{s}') \\ &= \mathbf{P}^{\mathbf{M}^{k+1}}(s, \alpha, s'). \end{aligned}$$

By definition, we can conclude that $\tilde{\mathbf{P}}_p^{\mathbf{M}^{k+1}}$ is correct. ■

C. Incremental Construction of SCCs

Consider an arbitrary iteration $k \geq 0$. Let l be the index of the environment agent such that $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$. In this section, we first provide a way to incrementally identify all the SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$ from all the SCCs of $\mathcal{M}^{\mathbf{M}^k}$ and \mathcal{M}_l . We conclude the section with incremental construction of the partial order over the SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$ from the partial order defined over the SCCs of $\mathcal{M}^{\mathbf{M}^k}$ and \mathcal{M}_l .

Lemma 2: Let $C^{\mathbf{M}^k}$ be an SCC of $\mathcal{M}^{\mathbf{M}^k}$ and C^l be an SCC of \mathcal{M}_l where $\mathbf{M}_{k+1} = \mathbf{M}_k \cup \{\mathcal{M}_l\}$. Suppose either of the following conditions holds:

Cond 1: $|C^{\mathbf{M}^k}| = 1$ and the state in $C^{\mathbf{M}^k}$ does not have a self-loop in $\mathcal{M}^{\mathbf{M}^k}$.

Cond 2: $|C^l| = 1$ and the state in C^l does not have a self-loop in \mathcal{M}_l .

Then, for any $s \in C^{\mathbf{M}^k}$ and $r \in C^l$, $\{s|_{l \leftarrow r}\}$ is an SCC of $\mathcal{M}^{\mathbf{M}^{k+1}}$. Otherwise, $\{s|_{l \leftarrow r} \mid s \in C^{\mathbf{M}^k}, r \in C^l\}$ is an SCC of $\mathcal{M}^{\mathbf{M}^{k+1}}$.

Proof: First, we consider the case where Cond 1 or Cond 2 holds and consider arbitrary $s \in C^{\mathbf{M}^k}$ and $r \in C^l$. To show that $\{s|_{l \leftarrow r}\}$ is an SCC of $\mathcal{M}^{\mathbf{M}^{k+1}}$, we will show that there is no path from $s|_{l \leftarrow r}$ to itself in $\mathcal{M}^{\mathbf{M}^{k+1}}$. Since condition (1) or condition (2) holds, either there is no path from s to itself in $\mathcal{M}^{\mathbf{M}^k}$ or there is no path from r to itself in C^l . Assume, by contradiction, that there is a path from $s|_{l \leftarrow r}$ to itself in $\mathcal{M}^{\mathbf{M}^{k+1}}$. Let this path be $s|_{l \leftarrow r}, s^1, s^2, \dots, s^n, s|_{l \leftarrow r}$ where for each $i \in \{1, \dots, n\}$, $s^i = \langle s_0^i, \dots, s_N^i \rangle$. From the proof of Lemma 1, we get that $\mathbf{P}^{\mathbf{M}^{k+1}}(s|_{l \leftarrow r}, \alpha, s^1) = \mathbf{P}_l(r, s_l^1) \mathbf{P}^{\mathbf{M}^k}(s, \alpha, \tilde{s}^1)$, $\mathbf{P}^{\mathbf{M}^{k+1}}(s^n, \alpha, s|_{l \leftarrow r}) = \mathbf{P}_l(s_l^n, r) \mathbf{P}^{\mathbf{M}^k}(\tilde{s}^n, \alpha, s)$ and $\mathbf{P}^{\mathbf{M}^{k+1}}(s^i, \alpha, s^{i+1}) = \mathbf{P}_l(s_l^i, s_l^{i+1}) \mathbf{P}^{\mathbf{M}^k}(\tilde{s}^i, \alpha, \tilde{s}^{i+1})$ for all $\alpha \in Act^{\mathbf{M}^{k+1}}$ where for each $i \in \{1, \dots, n\}$, $\tilde{s}^i \in S^{\mathbf{M}^k}$ such that $\tilde{s}^i|_{l \leftarrow s_l^i} = s^i$.

Since $s|_{l \leftarrow r}, s^1, s^2, \dots, s^n, s|_{l \leftarrow r}$ is a path in $\mathcal{M}^{\mathbf{M}^{k+1}}$, there exist $\alpha_0, \dots, \alpha_n \in Act^{\mathbf{M}^{k+1}}$ such that $\mathbf{P}^{\mathbf{M}^{k+1}}(s|_{l \leftarrow r}, \alpha_0, s^1)$, $\mathbf{P}^{\mathbf{M}^{k+1}}(s^n, \alpha_n, s|_{l \leftarrow r})$, $\mathbf{P}^{\mathbf{M}^{k+1}}(s^i, \alpha_i, s^{i+1}) > 0$ for all $i \in \{1, \dots, n\}$. Thus, it must be the case that $\mathbf{P}_l(r, s_l^1)$, $\mathbf{P}_l(s_l^n, r)$, $\mathbf{P}_l(s_l^i, s_l^{i+1}) > 0$ and $\mathbf{P}^{\mathbf{M}^k}(s, \alpha, \tilde{s}^1)$, $\mathbf{P}^{\mathbf{M}^k}(\tilde{s}^n, \alpha, s)$, $\mathbf{P}^{\mathbf{M}^k}(\tilde{s}^i, \alpha, \tilde{s}^{i+1}) > 0$ for all $i \in \{1, \dots, n\}$. But then, $r, s_l^1, \dots, s_l^n, r$ is a path in C^l and $s, \tilde{s}^1, \dots, \tilde{s}^n, s$ is a path in $\mathcal{M}^{\mathbf{M}^k}$, leading to a contradiction.

Next, consider the case where both Cond 1 and Cond 2 do not hold. To show that $C^{\mathbf{M}^{k+1}} = \{s|_{l \leftarrow r} \mid s \in C^{\mathbf{M}^k}, r \in C^l\}$ is an SCC of $\mathcal{M}^{\mathbf{M}^{k+1}}$, we need to show that for any $s, \tilde{s} \in C^{\mathbf{M}^{k+1}}$ and any $s' \notin C^{\mathbf{M}^{k+1}}$, (1) there is a path in $\mathcal{M}^{\mathbf{M}^{k+1}}$ from s to \tilde{s} , and (2) there is no path in $\mathcal{M}^{\mathbf{M}^{k+1}}$

either from s to s' or from s' to s . Both of these statements can be proved by contradiction, using the same reasoning as in the proof above for the case where either Cond 1 or Cond 2 holds. ■

We say that an SCC $C^{\mathbf{M}^{k+1}}$ of $\mathcal{M}^{\mathbf{M}^{k+1}}$ is *derived* from $\langle C^{\mathbf{M}^k}, C^l \rangle$, where $C^{\mathbf{M}^k}$ is an SCC of $\mathcal{M}^{\mathbf{M}^k}$ and C^l is an SCC of \mathcal{M}_l , if $C^{\mathbf{M}^{k+1}}$ is constructed from $C^{\mathbf{M}^k}$ and C^l according to Lemma 2, i.e., $C^{\mathbf{M}^{k+1}} = \{s|_{l \leftarrow r}\}$ for some $s \in C^{\mathbf{M}^k}$ and $r \in C^l$ if Cond 1 or Cond 2 in Lemma 2 holds; otherwise, $C^{\mathbf{M}^{k+1}} = \{s|_{l \leftarrow r} \mid s \in C^{\mathbf{M}^k}, r \in C^l\}$.

Lemma 3: For each SCC $C^{\mathbf{M}^{k+1}}$ of $\mathcal{M}^{\mathbf{M}^{k+1}}$, there exists a unique $\langle C^{\mathbf{M}^k}, C^l \rangle$ from which $C^{\mathbf{M}^{k+1}}$ is derived.

Proof: Similar to Lemma 1, it can be checked that $S^{\mathbf{M}^{k+1}} = \{s|_{l \leftarrow r} \mid s \in S^{\mathbf{M}^k} \text{ and } r \in S_l\}$ is the set of states of $\mathcal{M}^{\mathbf{M}^{k+1}}$. Consider an arbitrary SCC $C^{\mathbf{M}^{k+1}}$ of $\mathcal{M}^{\mathbf{M}^{k+1}}$ and an arbitrary $s = \langle s_0, \dots, s_N \rangle \in C^{\mathbf{M}^{k+1}}$.

By definition, for any arbitrary SCC $C^{\mathbf{M}^k}$ of $\mathcal{M}^{\mathbf{M}^k}$ and arbitrary SCC C^l of \mathcal{M}_l , $C^{\mathbf{M}^{k+1}}$ is derived from $\langle C^{\mathbf{M}^k}, C^l \rangle$ only if $s_l \in C^l$ and there exist $s' \in C^{\mathbf{M}^k}$ such that $s'|_{l \leftarrow s_l} = s$. But since $\tilde{\mathcal{M}}_l$ contains exactly one state, there exists a unique $s' \in S^{\mathbf{M}^k}$ such that $s'|_{l \leftarrow s_l} = s$. Also, from the definition of SCC, there exist a unique SCC $C^{\mathbf{M}^k}$ of $\mathcal{M}^{\mathbf{M}^k}$ and a unique SCC C^l of \mathcal{M}_l such that $s' \in C^{\mathbf{M}^k}$ and $s_l \in C^l$. Thus, it cannot be the case that $C^{\mathbf{M}^{k+1}}$ is derived from $\langle \tilde{C}^{\mathbf{M}^k}, \tilde{C}^l \rangle$ where $\tilde{C}^{\mathbf{M}^k} \neq C^{\mathbf{M}^k}$ or $\tilde{C}^l \neq C^l$. Applying Lemma 2, we get that there exists an SCC $\tilde{C}^{\mathbf{M}^{k+1}}$ of $\mathcal{M}^{\mathbf{M}^{k+1}}$ that is derived from $\langle C^{\mathbf{M}^k}, C^l \rangle$ and contains s . Since $s \in C^{\mathbf{M}^{k+1}}$ and $s \in \tilde{C}^{\mathbf{M}^{k+1}}$, from the definition of SCC, it must be the case that $C^{\mathbf{M}^{k+1}} = \tilde{C}^{\mathbf{M}^{k+1}}$; thus, $C^{\mathbf{M}^{k+1}}$ must be derived from $\langle C^{\mathbf{M}^k}, C^l \rangle$. ■

Lemma 2 and Lemma 3 provide a way to generate all the SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$ from all the SCCs of $\mathcal{M}^{\mathbf{M}^k}$ and \mathcal{M}_l as formally stated below.

Corollary 1: The set of all the SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$ is given by

$$\{C^{\mathbf{M}^{k+1}} \text{ derived from } \langle C^{\mathbf{M}^k}, C^l \rangle \mid C^{\mathbf{M}^k} \text{ is an SCC of } \mathcal{M}^{\mathbf{M}^k} \text{ and } C^l \text{ is an SCC of } \mathcal{M}_l\}.$$

Finally, in the following lemma, we provide a necessary condition, based on the partial order over the SCCs of $\mathcal{M}^{\mathbf{M}^k}$ and \mathcal{M}_l , for the existence of the partial order between two SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$.

Lemma 4: Let $C_1^{\mathbf{M}^{k+1}}$ and $C_2^{\mathbf{M}^{k+1}}$ be SCCs of $\mathcal{M}^{\mathbf{M}^{k+1}}$. Suppose $C_1^{\mathbf{M}^{k+1}}$ is derived from $\langle C_1^{\mathbf{M}^k}, C_1^l \rangle$ and $C_2^{\mathbf{M}^{k+1}}$ is derived from $\langle C_2^{\mathbf{M}^k}, C_2^l \rangle$ where $C_1^{\mathbf{M}^k}$ and $C_2^{\mathbf{M}^k}$ are SCCs of $\mathcal{M}^{\mathbf{M}^k}$ and C_1^l and C_2^l are SCCs of \mathcal{M}_l . Then, $C_1^{\mathbf{M}^{k+1}} \prec_{\mathcal{M}^{\mathbf{M}^{k+1}}} C_2^{\mathbf{M}^{k+1}}$ only if $C_1^{\mathbf{M}^k} \prec_{\mathcal{M}^{\mathbf{M}^k}} C_2^{\mathbf{M}^k}$ and $C_1^l \prec_{\mathcal{M}_l} C_2^l$.

Proof: Consider the case where $C_1^{\mathbf{M}^{k+1}} \prec_{\mathcal{M}^{\mathbf{M}^{k+1}}} C_2^{\mathbf{M}^{k+1}}$. By definition, $Succ(C_2^{\mathbf{M}^{k+1}}) \cap C_1^{\mathbf{M}^{k+1}} \neq \emptyset$. Consider a state $s' = \langle s'_0, \dots, s'_N \rangle \in Succ(C_2^{\mathbf{M}^{k+1}}) \cap C_1^{\mathbf{M}^{k+1}}$. Since $s' \in Succ(C_2^{\mathbf{M}^{k+1}})$, there exists $s = \langle s_0, \dots, s_N \rangle \in C_2^{\mathbf{M}^{k+1}}$ and $\alpha \in Act^{\mathbf{M}^{k+1}}$ such that $\mathbf{P}^{\mathbf{M}^{k+1}}(s, \alpha, s') > 0$. But from the proof of Lemma 1, $\mathbf{P}^{\mathbf{M}^{k+1}}(s, \alpha, s') = \mathbf{P}_l(s_l, s'_l) \mathbf{P}^{\mathbf{M}^k}(\tilde{s}, \alpha, \tilde{s}')$ where \tilde{s} and \tilde{s}' are unique states in $S^{\mathbf{M}^k}$ such that $\tilde{s}|_{l \leftarrow s_l} = s$ and $\tilde{s}'|_{l \leftarrow s'_l} = s'$. Thus, it must be the case that $\mathbf{P}_l(s_l, s'_l) > 0$ and $\mathbf{P}^{\mathbf{M}^k}(\tilde{s}, \alpha, \tilde{s}') > 0$.

In addition, since $C_1^{M_{k+1}}$ is derived from $\langle C_1^{M_k}, C_1^l \rangle$ and $C_2^{M_{k+1}}$ is derived from $\langle C_2^{M_k}, C_2^l \rangle$, from Lemma 2 and Lemma 3, it must be the case that $\tilde{s} \in C_2^{M_k}$, $\tilde{s}' \in C_1^{M_k}$, $s_l \in C_2^l$ and $s'_l \in C_1^l$. Since $\tilde{s} \in C_2^{M_k}$, $\tilde{s}' \in C_1^{M_k}$ and $\mathbf{P}^{M_k}(\tilde{s}, \alpha, \tilde{s}') > 0$, we can conclude that $\tilde{s}' \in \text{Succ}(C_2^{M_k}) \cap C_1^{M_k}$, and therefore, by definition, $C_1^{M_k} \prec_{\mathcal{M}^{M_k}} C_2^{M_k}$. Similarly, since $s_l \in C_2^l$, $s'_l \in C_1^l$ and $\mathbf{P}_l(s_l, s'_l) > 0$, we can conclude that $s'_l \in \text{Succ}(C_2^l) \cap C_1^l$, and therefore, by definition, $C_1^l \prec_{\mathcal{M}_l} C_2^l$. ■

D. Computation of Probability Vector and Control Policy for $\mathcal{M}_p^{M_k}$ from SCCs of \mathcal{M}^{M_k}

Consider an arbitrary iteration $k \geq 0$ and the associated product MDP $\mathcal{M}_p^{M_k} = (S_p^{M_k}, Act_p^{M_k}, \mathbf{P}_p^{M_k}, \mathbf{I}_{p,init}^{M_k}, \Pi_p^{M_k}, L_p^{M_k})$. Similar to the SCC-based value iteration, we want to generate a partition $\{D_{p,1}^{M_k}, \dots, D_{p,m_k}^{M_k}\}$ of $S_p^{M_k}$ with a partial order $\prec_{\mathcal{M}_p^{M_k}}$ such that $D_{p,j}^{M_k} \prec_{\mathcal{M}_p^{M_k}} D_{p,i}^{M_k}$ if $\text{Succ}(D_{p,i}^{M_k}) \cap D_{p,j}^{M_k} \neq \emptyset$. However, we relax the condition that each $D_{p,i}^{M_k}, i \in \{1, \dots, m_k\}$ is an SCC of $\mathcal{M}_p^{M_k}$ and only require that if $D_{p,i}^{M_k}$ contains a state in an SCC $C_p^{M_k}$ of $\mathcal{M}_p^{M_k}$, then it has to contain all the states in $C_p^{M_k}$. Hence, $D_{p,i}^{M_k}$ may include all the states in multiple SCCs of $\mathcal{M}_p^{M_k}$. The following lemmas provide a method for constructing $\{D_{p,1}^{M_k}, \dots, D_{p,m_k}^{M_k}\}$ and their partial order from SCCs of \mathcal{M}^{M_k} and their partial order, which can be incrementally constructed as described in Section V-C.

Lemma 5: Let $C_p^{M_k}$ be an SCC of $\mathcal{M}_p^{M_k}$. Then, there exists a unique SCC C^{M_k} of \mathcal{M}^{M_k} such that $C_p^{M_k} \subseteq C^{M_k} \times Q$.

Proof: This follows from the definition of product MDP that for any $s, s' \in S^{M_k}$ and $q, q' \in Q$, there is a path from $\langle s, q \rangle$ to $\langle s', q' \rangle$ in $\mathcal{M}_p^{M_k}$ only if there is a path from s to s' in \mathcal{M}^{M_k} . ■

Lemma 6: Let $C_p^{M_k}$ and $\tilde{C}_p^{M_k}$ be SCCs of $\mathcal{M}_p^{M_k}$. Suppose C^{M_k} and \tilde{C}^{M_k} are unique SCCs of \mathcal{M}^{M_k} such that $C_p^{M_k} \subseteq C^{M_k} \times Q$ and $\tilde{C}_p^{M_k} \subseteq \tilde{C}^{M_k} \times Q$. Then, $C_p^{M_k} \prec_{\mathcal{M}_p^{M_k}} \tilde{C}_p^{M_k}$ only if $C^{M_k} \prec_{\mathcal{M}^{M_k}} \tilde{C}^{M_k}$.

Proof: This follows from the definition of product MDP since for any $\langle s, q \rangle \in S_p^{M_k}$, $\langle \tilde{s}, \tilde{q} \rangle \in S_p^{M_k}$ is a successor of $\langle s, q \rangle$ in $\mathcal{M}_p^{M_k}$ only if \tilde{s} is a successor of s in \mathcal{M}^{M_k} . ■

Lemma 7: Let $C_1^{M_k}, \dots, C_{m_k}^{M_k}$ be all the SCCs of \mathcal{M}^{M_k} and for each $i \in \{1, \dots, m_k\}$, let $D_{p,1}^{M_k} = C_i^{M_k} \times Q$. Then, $\{D_{p,1}^{M_k}, \dots, D_{p,m_k}^{M_k}\}$ is a partition of $S_p^{M_k}$. In addition, the following statements hold for all $i, j \in \{1, \dots, m_k\}$.

- If $D_{p,i}^{M_k}$ contains a state in an SCC $C_p^{M_k}$ of $\mathcal{M}_p^{M_k}$, then it contains all the states in $C_p^{M_k}$.
- $\text{Succ}(D_{p,i}^{M_k}) \cap D_{p,j}^{M_k} \neq \emptyset$ only if $C_j^{M_k} \prec_{\mathcal{M}^{M_k}} C_i^{M_k}$.

Proof: Consider arbitrary $i, j \in \{1, \dots, m_k\}$. It follows directly from Lemma 5 that if $D_{p,i}^{M_k}$ contains a state in an SCC $C_p^{M_k}$ of $\mathcal{M}_p^{M_k}$, then it contains all the states in $C_p^{M_k}$. Next, consider the case where $\text{Succ}(D_{p,i}^{M_k}) \cap D_{p,j}^{M_k} \neq \emptyset$. Then, from Lemma 5, there exist SCCs $C_{p,i}^{M_k} \subseteq D_{p,i}^{M_k}$ and $C_{p,j}^{M_k} \subseteq D_{p,j}^{M_k}$ of $\mathcal{M}_p^{M_k}$ such that $\text{Succ}(C_{p,i}^{M_k}) \cap C_{p,j}^{M_k} \neq$

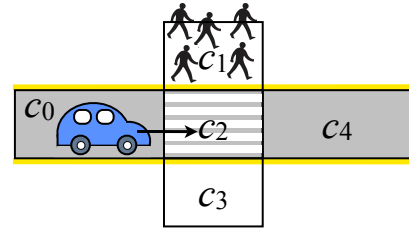


Fig. 1. The road and its partition used in the autonomous vehicle example.

\emptyset . Thus, $C_{p,j}^{M_k} \prec_{\mathcal{M}_p^{M_k}} C_{p,i}^{M_k}$. Applying Lemma 6, we get $C_j^{M_k} \prec_{\mathcal{M}^{M_k}} C_i^{M_k}$. ■

Applying Lemma 7, we generate a partition $\{D_{p,1}^{M_k}, \dots, D_{p,m_k}^{M_k}\}$ of $S_p^{M_k}$ where for each $i \in \{1, \dots, m_k\}$, $D_{p,1}^{M_k} = C_i^{M_k} \times Q$ and $C_1^{M_k}, \dots, C_{m_k}^{M_k}$ are all the SCCs of \mathcal{M}^{M_k} . A partial order $\prec_{\mathcal{M}_p^{M_k}}$ over this partition is defined such that $D_{p,j}^{M_k} \prec_{\mathcal{M}_p^{M_k}} D_{p,i}^{M_k}$ if $C_j^{M_k} \prec_{\mathcal{M}^{M_k}} C_i^{M_k}$. Hence, an order $\mathbb{O}_p^{M_k}$ among $D_{p,1}^{M_k}, \dots, D_{p,m_k}^{M_k}$ can be simply derived from the order of $C_1^{M_k}, \dots, C_{m_k}^{M_k}$, which can be incrementally constructed based on Lemma 4. This order $\mathbb{O}_p^{M_k}$ has the property that the probability values of states in $D_{p,j}^{M_k}$ that appears after $D_{p,i}^{M_k}$ in $\mathbb{O}_p^{M_k}$ cannot affect the probability values of states in $D_{p,i}^{M_k}$. Hence, we can follow the SCC-based value iteration and process each $D_{p,i}^{M_k}$ separately, according to the order in $\mathbb{O}_p^{M_k}$ to compute the probability x_s for all $s \in D_{p,i}^{M_k}$. Finally, we generate a memoryless control policy \mathcal{C}^{M_k} from the probability vector $(x_s)_{s \in S_p^{M_k}}$ as described at the end of Section IV.

VI. EXPERIMENTAL RESULTS

Consider, once again, the autonomous vehicle problem described in Example 1 and Example 2. Suppose the road is discretized into 5 cells c_0, \dots, c_4 where c_2 is the pedestrian crossing area as shown in Figure 1. The vehicle starts in cell c_0 and has to reach cell c_4 . There are 5 pedestrians, modeled by MCs $\mathcal{M}_1, \dots, \mathcal{M}_5$, initially at cell c_1 . The models of the vehicle and the pedestrians are shown in Figure 2. A DFA \mathcal{A}_φ that accepts all and only words in $\text{pref}(\varphi)$ where $\varphi = (\neg \bigvee_{i \geq 1, j \geq 0} (c_j^0 \wedge c_j^i)) \mathcal{U} c_4^0$ is shown in Figure 3.

First, we apply the LP-based, value iteration and SCC-based value iteration techniques described in Section IV to synthesize a control policy that maximizes the probability that the complete system $\mathcal{M} = \mathcal{T} \parallel \mathcal{M}_1 \parallel \mathcal{M}_2 \parallel \dots \parallel \mathcal{M}_5$ satisfies φ . The time required for each step of computation is summarized in Table I. All the approaches yield the probability of 0.8 that \mathcal{M} satisfies φ under the synthesized control policy. The comparison of the total computation time required for these different approaches is shown in Figure 4. As discussed in Section IV-C, although the SCC-based value iteration itself takes significantly less computation time than the LP-based technique or value iteration, the time spent in identifying SCCs and their order renders the total

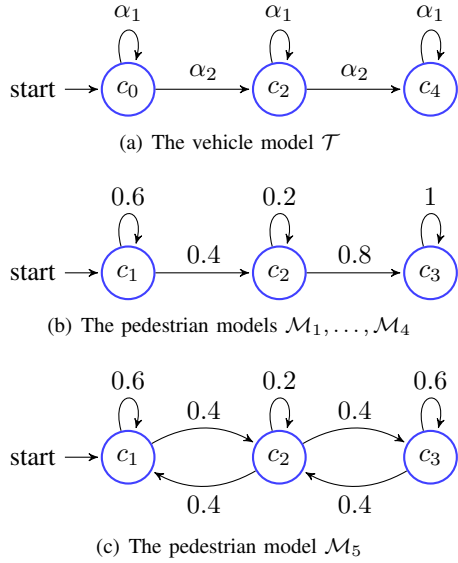


Fig. 2. The models of vehicle and pedestrians.

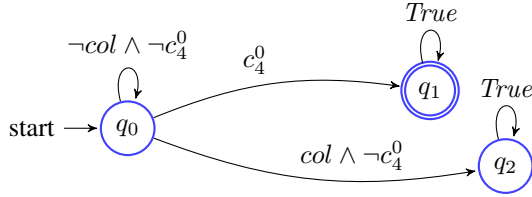


Fig. 3. A DFA \mathcal{A}_φ that recognizes the prefixes of $\varphi = \neg col \cup c_4^0$ where col is defined as $col = \bigvee_{i \geq 1, j \geq 0} (c_j^0 \wedge c_j^i)$. q_1 is the accepting state.

computation time of the SCC-based value iteration more than the other two approaches.

Technique	\mathcal{M}_p	SCCs & order of \mathcal{M}_p	Prob vector	Control policy	Total
LP	156.3	-	8.8	6.8	171.9
Value iteration	156.3	-	31.3	6.8	194.4
SCC-based value iteration	156.3	71.1	1.9	6.8	236.1

TABLE I

TIME REQUIRED (IN SECONDS) FOR COMPUTING VARIOUS OBJECTS USING DIFFERENT TECHNIQUES WHEN THE FULL MODELS OF ALL THE ENVIRONMENT AGENTS ARE CONSIDERED.

Next, we apply the incremental technique where we progressively compute a sequence of control policies as more agents are added to the synthesis procedure in each iteration as described in Section V. We let $\mathbf{M}_0 = \emptyset$, $\mathbf{M}_1 = \{\mathcal{M}_1\}$, $\mathbf{M}_2 = \{\mathcal{M}_1, \mathcal{M}_2\}, \dots, \mathbf{M}_6 = \{\mathcal{M}_1, \dots, \mathcal{M}_5\}$, i.e., we successively add each pedestrian $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_5$, respectively, in each iteration. We consider 2 cases: (1) no incremental construction of various objects is employed (i.e., when $\mathcal{M}^{\mathbf{M}_{k+1}}$ and $\mathcal{M}_p^{\mathbf{M}_{k+1}}$, $k \geq 0$ are computed from scratch in every iteration), and (2) incremental construction of various objects as described in Section V-B-V-D is

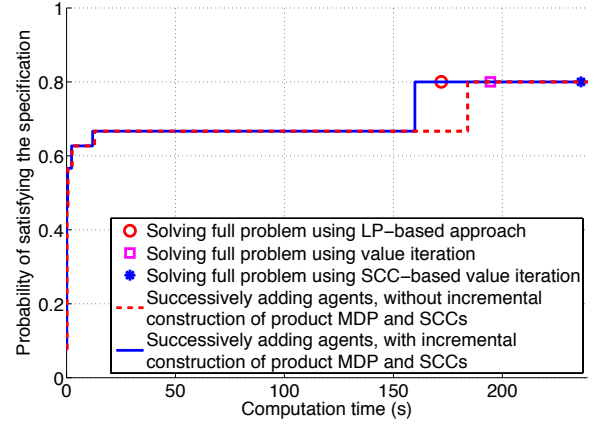


Fig. 4. Comparison of the computation time and the probability for the system to satisfy the specification computed using different techniques.

applied. For the first case, we apply the LP-based technique to compute the probability vector as it has been shown to be the fastest technique when applied to this problem, taking into account the required pre-computation, which needs to be done in every iteration. For both cases, 6 control policies $\mathcal{C}^{\mathbf{M}_0}, \dots, \mathcal{C}^{\mathbf{M}_5}$ are generated for $\mathcal{M}^{\mathbf{M}_0}, \dots, \mathcal{M}^{\mathbf{M}_5}$, respectively. For each policy $\mathcal{C}^{\mathbf{M}_k}$, we compute the probability $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_k}}(\varphi)$ that the complete system \mathcal{M} satisfies φ under policy $\mathcal{C}^{\mathbf{M}_k}$. (Note that $\mathcal{C}^{\mathbf{M}_k}$, when applied to \mathcal{M} , is only a function of states of $\mathcal{M}_i \in \mathbf{M}_k$ since it assumes that the other agents $\mathcal{M}_j \notin \mathbf{M}_k$ are stationary.) These probabilities are given by $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_0}}(\varphi) = 0.08$, $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_1}}(\varphi) = 0.46$, $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_2}}(\varphi) = 0.57$, $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_3}}(\varphi) = 0.63$, $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_4}}(\varphi) = 0.67$ and $\Pr_{\mathcal{M}}^{\mathcal{C}^{\mathbf{M}_5}}(\varphi) = 0.8$.

The comparison of the cases where the incremental construction of various objects is not and is employed is shown in Figure 4. A jump in the probability occurs each time a new control policy is computed. The time spent during each step of computation is summarized in Table II and Table III for the first and the second case, respectively. Notice that the time required for identifying the SCCs and their order when the incremental approach is applied is significantly less than when the full model of all the pedestrians is considered in one shot since $\mathcal{M}^{\mathbf{M}_0}, \mathcal{M}_1, \dots, \mathcal{M}_5$, each of which contains 3 states, are much smaller than \mathcal{M}_p , which contains 2187 states.

From Figure 4, our incremental approach is able to obtain an optimal control policy faster than any other techniques. This is mainly due to the efficiency of our incremental construction of SCCs and their order. In addition, we are able to obtain a reasonable solution, with 0.67 probability of satisfying φ , within 12 seconds while the maximum probability of satisfying φ is 0.8, which requires 160 seconds of computation (or 171.9 seconds without employing the incremental approach).

Iteration	\mathcal{M}^{M_k}	$\mathcal{M}_p^{M_k}$	Prob vector	Control policy	Total
0	0.0064	0.0185	0.0464	0.0084	0.08
1	0.0123	0.0762	0.0203	0.0104	0.12
2	0.0154	0.3383	0.0231	0.0296	0.41
3	0.0357	1.7055	0.0542	0.1503	1.95
4	0.1393	9.1950	0.2155	0.7975	10.35
5	3.1836	152.86	8.2302	6.8938	171.17

TABLE II

TIME REQUIRED (IN SECONDS) FOR COMPUTING VARIOUS OBJECTS IN EACH ITERATION WHEN INCREMENTAL CONSTRUCTION IS NOT APPLIED.

Iteration	\mathcal{M}^{M_0}	SCCs & order of \mathcal{M}^{M_0} , $\mathcal{M}_1, \dots, \mathcal{M}_5$	$\mathcal{M}_p^{M_k}$, partition & order	Prob vector	Control policy	Total
0	0.0055	0.0043	0.0203	0.0112	0.0036	0.04
1	-	-	0.0726	0.0102	0.0087	0.09
2	-	-	0.3239	0.0193	0.0282	0.37
3	-	-	1.6036	0.0567	0.1424	1.80
4	-	-	8.6955	0.1876	0.7755	9.66
5	-	-	139.27	1.6122	7.0125	147.89

TABLE III

TIME REQUIRED (IN SECONDS) FOR COMPUTING VARIOUS OBJECTS IN EACH ITERATION WHEN INCREMENTAL CONSTRUCTION IS APPLIED.

VII. CONCLUSIONS AND FUTURE WORK

An anytime algorithm for synthesizing a control policy for a robot interacting with multiple environment agents with the objective of maximizing the probability for the robot to satisfy a given temporal logic specification was proposed. Each environment agent is modeled by a Markov chain whereas the robot is modeled by a finite transition system (in the deterministic case) or Markov decision process (in the stochastic case). The proposed algorithm progressively computes a sequence of control policies, taking into account only a small subset of the environment agents initially and successively adding more agents to the synthesis procedure in each iteration until we hit the constraints on computational resources. Incremental construction of various objects needed to be computed during the synthesis procedure was proposed. Experimental results showed that not only we obtain a reasonable solution much faster than existing approaches, but we are also able to obtain an optimal solution faster than existing approaches.

Future work includes extending the algorithm to handle full LTL specifications. This direction appears to be promising because the remaining step is only to incrementally construct *accepting maximal end components* of an MDP. We are also examining an effective approach to determine an agent to be added in each iteration. As mentioned in Section V-A, such an agent may be picked based on the result from probabilistic verification but this comes at the extra cost of adding the verification phase.

REFERENCES

[1] E. A. Emerson, "Temporal and modal logic," *Handbook of Theoretical Computer Science (Vol. B): Formal Models and Semantics*, pp. 995–1072, 1990.

[2] Z. Manna and A. Pnueli, *The temporal logic of reactive and concurrent systems*. Springer-Verlag, 1992.

[3] C. Baier and J.-P. Katoen, *Principles of Model Checking (Representation and Mind Series)*. The MIT Press, 2008.

[4] G. Fainekos, H. Kress-Gazit, and G. Pappas, "Temporal logic motion planning for mobile robots," in *IEEE International Conference on Robotics and Automation*, pp. 2020–2025, 2005.

[5] H. Kress-Gazit, G. Fainekos, and G. Pappas, "Where's Waldo? Sensor-based temporal logic motion planning," in *IEEE International Conference on Robotics and Automation*, pp. 3116–3121, 2007.

[6] C. Belta, A. Bicchi, M. Egerstedt, E. Frazzoli, E. Klavins, and G. Pappas, "Symbolic planning and control of robot motion [grand challenges of robotics]," *IEEE Robotics & Automation Magazine*, vol. 14, no. 1, pp. 61–70, 2007.

[7] D. Conner, H. Kress-Gazit, H. Choset, A. Rizzi, and G. Pappas, "Valet parking without a valet," in *IEEE/RSJ International Conference on Intelligent Robots and Systems, 2007*, pp. 572–577, 2007.

[8] S. Karaman and E. Frazzoli, "Sampling-based motion planning with deterministic μ -calculus specifications," in *Proc. of IEEE Conference on Decision and Control*, 2009.

[9] A. Bhatia, L. E. Kavrakli, and M. Y. Vardi, "Sampling-based motion planning with temporal goals," in *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2689–2696, 2010.

[10] X. C. Ding, S. L. Smith, C. Belta, and D. Rus, "LTL control in uncertain environments with probabilistic satisfaction guarantees," in *IFAC World Congress*, 2011.

[11] A. I. Medina Ayala, S. B. Andersson, and C. Belta, "Temporal logic control in dynamic environments with probabilistic satisfaction guarantees," in *IEEE/RSJ International Conference on Intelligent Robots and Systems, 2007*, pp. 3108–3113, 2011.

[12] H. Kress-Gazit, T. Wongpiromsarn, and U. Topcu, "Correct, reactive robot control from abstraction and temporal logic specifications," *Special Issue of the IEEE Robotics & Automation Magazine on Formal Methods for Robotics and Automation*, vol. 18, pp. 65–74, 2011.

[13] O. Kupferman and M. Y. Vardi, "Model checking of safety properties," *Formal Methods in System Design*, vol. 19, pp. 291–314, 2001.

[14] T. Latvala, "Efficient model checking of safety properties," in *Model Checking Software. 10th International SPIN Workshop*, pp. 74–88, Springer, 2003.

[15] F. Ciesinski, C. Baier, M. Größer, and J. Klein, "Reduction techniques for model checking markov decision processes," in *Proceedings of the 2008 Fifth International Conference on Quantitative Evaluation of Systems*, pp. 45–54, 2008.

[16] M. Kwiatkowska, D. Parker, and H. Qu, "Incremental quantitative verification for markov decision processes," in *IEEE/IFIP International Conference on Dependable Systems & Networks*, pp. 359–370, 2011.

[17] R. Tarjan, "Depth-first search and linear graph algorithms," *SIAM Journal on Computing*, vol. 1, pp. 146–160, 1972.