

Incentive Design for Temporal Logic Objectives

Yagiz Savas, Vijay Gupta, Melkior Ornik, Lillian J. Ratliff, Ufuk Topcu

Abstract—We study the problem of designing an optimal sequence of incentives that a principal should offer to an agent so that the agent’s optimal behavior under the incentives realizes the principal’s objective expressed as a temporal logic formula. We consider an agent with a finite decision horizon and model its decision-making process as a Markov decision process (MDP). Under certain assumptions, we present a polynomial-time algorithm to synthesize an incentive sequence that minimizes the cost to the principal. We show that if the underlying MDP has only deterministic transitions, the principal can hide its objective from the agent and still realize the desired behavior through incentives. On the other hand, an MDP with stochastic transitions may require the principal to share its objective with the agent. Finally, we demonstrate the proposed method in motion planning examples where a principal changes the optimal trajectory of an agent by providing incentives.

I. INTRODUCTION

Consider a scenario where a principal provides incentives to an agent so that the optimal behavior of the agent under the provided incentives satisfies the principal’s objective. If the principal had enough resources to provide arbitrarily large incentives, it would be straightforward to obtain the desired agent behaviour. However, since the resources are limited in practice, it is important to establish the minimum amount of incentives that leads to the desired behavior. In this paper, we are interested in designing a sequence of incentives that minimizes the cost to the principal while guaranteeing the realization of its objective by the agent with maximum probability.

We model the sequential decision-making process of the agent as a Markov decision process (MDP) [1], and assume that the agent’s objective is to maximize its expected total reward at the end of a finite planning horizon. Although each planning horizon is finite, the agent plans its future decisions infinitely many times. Examples of such an agent can be a person who plans her schedule on a weekly basis or an autonomous system with a limited computational power which plans its route by considering only a small subset of all possible environment states.

The principal’s objective is described by a syntactically co-safe linear temporal logic (LTL) formula. LTL specifications are widely used to describe complex tasks for autonomous

robots [2], design security protocols [3] and check the reliability of software [4]. For example, in a navigation scenario, syntactically co-safe LTL formulae allow one to specify tasks such as liveness (eventually visit the region A) or priority (first visit the region A and then B).

We assume that the principal is aware of the agent’s reward function and the length of its planning horizon. In many real-world applications, the decision horizon and the reward structure of an agent can be known or at least inferred through observations. For example, a manufacturing company is generally interested in maximizing its profit at the end of a fiscal year, and an autonomous car aims to reach its destination within certain time interval.

From a practical point of view, an interesting question is whether an adversarial principal can convince an agent to satisfy its objective through incentives. In such a scenario, if the agent knows the principal’s objective explicitly, it will reject the provided incentives because the resulting behavior under the incentives will serve to the benefit of the enemy. However, if the principal can design an incentive sequence without sharing its objective with the agent, then the incentives may lead to the desired agent behavior. Therefore, it is important to establish the conditions under which the principal can actually hide its objective from the agent.

The contributions of this paper can be summarized as follows. First, we present an algorithm, based on a series of linear optimization problems, to synthesize a sequence of incentives that minimizes the cost to the principal while ensuring that the optimal agent behavior under the provided incentives satisfies a syntactically co-safe LTL formula with maximum probability. Second, we present an example scenario where the principal has to share its objective with the agent to induce the desired behavior. Third, we provide sufficient conditions on the structure of the MDP and the length of the agent’s decision horizon under which there exists an optimal incentive design that allows the principal to hide its objective from the agent.

Related work. The problem of obtaining desired agent behavior through a sequence of incentives has been extensively studied in the literature. In [5] and [6], the authors present methods to design incentive sequences with *limited* resources that maximizes the value of the principal’s objective function. They employ techniques from inverse reinforcement learning literature and prove NP-hardness of the considered design problem [5]. The work [7] provides a polynomial-time algorithm to synthesize minimum incentives for inducing a *specific* agent policy. Reference [8] considers a bandit model and presents methods to induce desired agent actions under different constraints on the incentives. Although it is quite

Y. Savas and U. Topcu are with the Department of Aerospace Engineering, University of Texas at Austin, TX, USA. E-mail: {yagiz.savas, utopcu}@utexas.edu

V. Gupta is with the Department of Electrical Engineering, University of Notre Dame, IN, USA. E-mail: vgupta@nd.edu

M. Ornik is with the Department of Aerospace Engineering and the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, IL, USA. E-mail: mornik@illinois.edu

L. J. Ratliff is with the Department of Electrical Engineering, University of Washington, WA, USA. E-mail: ratliff@uw.edu

different from the problem considered here, the design of feasible incentives that aligns the objectives of an agent and a principal is discussed in [9] from a control theoretic perspective. Unlike the references mentioned above, in this paper, we consider the problem of designing minimum incentives that maximizes the value of the principal's objective function expressed as a temporal logic formula. We also note that establishing the complexity of the design problem considered in this paper is mentioned as an open problem in [5].

II. PRELIMINARIES

For a set S , we denote its power set and cardinality by 2^S and $|S|$, respectively. Additionally, $\mathbb{N}=\{1, 2, \dots\}$, $\mathbb{N}_0=\{0, 1, 2, \dots\}$ and $\mathbb{R}_{\geq 0}=[0, \infty)$.

A. Markov Decision Processes

Definition 1: A *Markov decision process* (MDP) is a tuple $\mathcal{M}=(S, s_0, \mathcal{A}, \mathcal{P}, \mathcal{AP}, \mathcal{L}, \mathcal{R})$ where S is a finite set of states, $s_0 \in S$ is an initial state, \mathcal{A} is a finite set of actions, $\mathcal{P}: S \times \mathcal{A} \times S \rightarrow [0, 1]$ is a transition function such that $\sum_{s' \in S} \mathcal{P}(s, a, s')=1$ for all $s \in S$ and $a \in \mathcal{A}(s)$ where $\mathcal{A}(s)$ denote the available actions in s , \mathcal{AP} is a set of atomic propositions, $\mathcal{L}: S \rightarrow 2^{\mathcal{AP}}$ is a function that labels each state with a subset of atomic propositions, and $\mathcal{R}: S \times \mathcal{A} \rightarrow \mathbb{R}$ is a reward function.

We denote the transition probability $\mathcal{P}(s, a, s')$ by $\mathcal{P}_{s,a,s'}$.

Definition 2: For an MDP \mathcal{M} , a *decision rule* $d: S \times \mathcal{A} \rightarrow [0, 1]$ is a function such that $\sum_{a \in \mathcal{A}(s)} d(s, a)=1$ for all $s \in S$. A decision rule d is said to be *deterministic* if for all $s \in S$ there exists $a \in \mathcal{A}(s)$ such that $d(s, a)=1$, and *randomized* otherwise. For an MDP \mathcal{M} , we denote the set of all (deterministic) decision rules by $(\mathcal{D}^D(\mathcal{M})) \mathcal{D}(\mathcal{M})$.

For an MDP \mathcal{M} , a decision-maker, i.e., an agent, chooses a decision rule $d \in \mathcal{D}(\mathcal{M})$ at each *stage*.

Definition 3: An N -stage *policy* for an MDP \mathcal{M} is a sequence $\pi=(d_1, d_2, \dots, d_N)$ where $N \leq \infty$ and $d_t \in \mathcal{D}(\mathcal{M})$ for all $t \leq N$. A *stationary* policy is a policy such that $d_t=d_1$ for all $t \leq N$. A policy is said to be *deterministic* if $d_t \in \mathcal{D}^D(\mathcal{M})$ for all t , and *randomized* otherwise. For an MDP \mathcal{M} , we denote the set of all N -stage policies by $\Pi_N(\mathcal{M})$. For notational simplicity, we denote the set of ∞ -stage policies by $\Pi(\mathcal{M})$.

For an MDP \mathcal{M} and a policy $\pi \in \Pi(\mathcal{M})$, let $\mu_t^\pi(s, a)$ be the joint probability of being in state $s \in S$ and taking the action $a \in \mathcal{A}(s)$ at stage t , which is uniquely determined through the recursive formula

$$\mu_{t+1}^\pi(s', a') = \sum_{s \in S} \sum_{a \in \mathcal{A}(s)} \mathcal{P}_{s,a,s'} \mu_t^\pi(s, a) d_{t+1}(s', a') \quad (1)$$

where $\mu_1^\pi(s, a)=d_1(s, a)\mu_0(s)$ and $\mu_0: S \rightarrow \{0, 1\}$ is a function such that $\mu_0(s_0)=1$ and $\mu_0(s)=0$ for all $s \in S \setminus \{s_0\}$.

Definition 4: For an MDP \mathcal{M} and a policy $\pi \in \Pi(\mathcal{M})$, the *expected residence time* in a state-action pair (s, a) is

$$\xi^\pi(s, a) := \sum_{t=1}^{\infty} \mu_t^\pi(s, a). \quad (2)$$

An infinite sequence $\varrho^\pi=s_0 s_1 s_2 \dots$ of states generated in \mathcal{M} under a policy $\pi \in \Pi(\mathcal{M})$, which starts from the initial

state s_0 and satisfies $\sum_{a_t \in \mathcal{A}(s_t)} d_k(s_t, a_t) \mathcal{P}_{s_t, a_t, s_{t+1}} > 0$ for all $t \geq 0$, is called a *path*. Any finite prefix of ϱ^π a finite path fragment. We define the set of all paths and finite path fragments in \mathcal{M} under the policy π by $Paths^\pi(\mathcal{M})$ and $Paths_{fin}^\pi(\mathcal{M})$, respectively. We use the standard probability measure over the outcome set $Paths^\pi(\mathcal{M})$ [10].

Definition 5: An *incentive design* for an MDP \mathcal{M} is a sequence $\Gamma=(\gamma_1, \gamma_2, \dots)$ where $\gamma_t: S \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$. A *stationary* incentive design is a design such that $\gamma_t=\gamma_1$ for all $t \in \mathbb{N}$. For an MDP \mathcal{M} , we denote the set of all incentive designs by $\Theta(\mathcal{M})$.

B. Linear temporal logic

We consider syntactically co-safe linear temporal logic (scLTL) formulae to specify tasks and refer the reader to [10], [11] for the syntax and semantics of scLTL.

An scLTL formula is built up from a set \mathcal{AP} of atomic propositions, logical connectives such as conjunction (\wedge) and negation (\neg), and temporal modal operators such as until (\mathcal{U}) and eventually (\diamond). An infinite sequence of subsets of \mathcal{AP} defines an infinite *word*, and an scLTL formula is interpreted over infinite words on $2^{\mathcal{AP}}$. We denote by $w \models \varphi$ that a word $w=w_0 w_1 w_2 \dots$ satisfies an scLTL formula φ .

For an MDP \mathcal{M} under a policy π , a path $\varrho^\pi=s_0 s_1 \dots$ generates a word $w=w_0 w_1 \dots$ where $w_k=\mathcal{L}(s_k)$ for all $k \geq 0$. With a slight abuse of notation, we use $\mathcal{L}(\varrho^\pi)$ to denote the word generated by ϱ^π . For an scLTL formula φ , the set $\{\varrho^\pi \in Paths^\pi(\mathcal{M}) : \mathcal{L}(\varrho^\pi) \models \varphi\}$ is measurable [10]. Hence, we define

$$\Pr_{\mathcal{M}}^\pi(s_0 \models \varphi) := \Pr_{\mathcal{M}}^\pi\{\varrho^\pi \in Paths^\pi(\mathcal{M}) : \mathcal{L}(\varrho^\pi) \models \varphi\}$$

as the probability of satisfying the scLTL formula φ for an MDP \mathcal{M} under the policy $\pi \in \Pi(\mathcal{M})$.

III. PROBLEM STATEMENT

We consider an *agent* whose sequential decision-making process is modeled as an MDP \mathcal{M} , and a *principal* that provides the agent a sequence of incentives $\Gamma \in \Theta(\mathcal{M})$.

The agent's objective is to maximize its expected total reward after N stages. However, since the incentive sequence offered by the principal might be non-stationary, the agent computes an N -stage policy every N stages. A graphical illustration of the agent's planning method is shown in Fig. 1. Formally, let $N \in \mathbb{N}$ be a constant, and $\mathcal{R}(S_t, A_t)$ and $\gamma_t(S_t, A_t)$ be the random reward and incentive received in stage $t \leq N$. Additionally, let $J:= (J_0, J_1, \dots)$ be a sequence of objective functions where $J_k: \Pi_N(\mathcal{M}) \times \Theta(\mathcal{M}) \rightarrow \mathbb{R}^{|S|}$ is such that

$$J_k(\pi, \Gamma)(s) := \mathbb{E}_s^\pi \left[\sum_{t=1}^N (\mathcal{R}(S_t, A_t) + \gamma_{kN+t}(S_t, A_t)) \right]$$

for all $s \in S$ where the expectation is taken over the finite path fragments that are generated by the policy $\pi \in \Pi_N(\mathcal{M})$ and start from the state s . Then, for a given incentive design

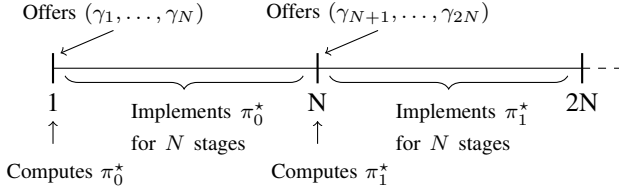


Fig. 1: An illustration of the incentive implementation and the agent's decision-making process. The principal offers incentives for the next N stages. After receiving the incentive offers, the agent computes and implements its optimal decisions for the next N stages.

$\Gamma \in \Theta(\mathcal{M})$, the agent's optimal ∞ -stage policy is given by $\pi^* := (\pi_0^*, \pi_1^*, \dots)$ where π_k^* is such that

$$\pi_k^* \in \arg \max_{\pi \in \Pi_N(\mathcal{M})} J_k(\pi, \Gamma)(s) \quad (3)$$

for all $s \in S$ and $k \in \mathbb{N}_0$. Note that the agent's policy π_k^* maximizes the total reward starting from any $s \in S$.

The principal's objective is to design an incentive sequence such that the agent's optimal policy under the provided incentives satisfies an scLTL formula φ with maximum probability.

The problem that we consider is the synthesis of an incentive design that minimizes the cost to the principal while realizing its objective. We make the following assumptions:

- (i) Agent's reward function \mathcal{R} is known by the principal.
- (ii) Agent's decision horizon N is known by the principal.
- (iii) The principal pays the offered incentives if and only if the agent takes the incentivized action.

Then, the optimization problem that we are interested in to solve is the following:

$$\min_{\Gamma \in \Theta(\mathcal{M})} \mathbb{E}_{s_0}^{\pi^*} \left[\sum_{t=1}^{\infty} \gamma_t(s, a) \right] \quad (4a)$$

$$\text{subject to: } \pi^* = (\pi_0^*, \pi_1^*, \dots) \quad (4b)$$

$$\pi_k^* \in \arg \max_{\pi \in \Pi_N(\mathcal{M})} J_k(\pi, \Gamma)(s) \quad \forall s \in S, \forall k \in \mathbb{N}_0 \quad (4c)$$

$$\Pr_{\mathcal{M}}^{\pi^*}(s_0 \models \varphi) = \max_{\pi \in \Pi(\mathcal{M})} \Pr_{\mathcal{M}}^{\pi}(s_0 \models \varphi) \quad (4d)$$

where $\Gamma = (\gamma_1, \gamma_2, \dots)$.

IV. THE DESIGN OF INCENTIVE SEQUENCES

In this section, we provide a method to synthesize an incentive design that solves the problem (4a)-(4d). For simplicity, we restrict our attention to reachability specifications, i.e., $\varphi = \diamond p$ where $p \in \mathcal{AP}$. The incentive design for general scLTL specifications is discussed in Section VI.

We first partition the states into three disjoint sets as follows. Let $B \subseteq S$ be the set of all states such that $\{p\} \subseteq \mathcal{L}(s)$, i.e., the set of states that the principal wants the agent to reach, and $S_0 \subseteq S$ be the set of states that have zero probability of reaching the states in B under any policy. More precisely, $s \in S_0$ if $\Pr_{\mathcal{M}}^{\pi}(s \models \diamond p) = 0$ for all $\pi \in \Pi(\mathcal{M})$. Finally, we let $S_r = S \setminus B \cup S_0$ be the set of all states that

are not in B and have nonzero probability of reaching a state in B under some policy. These sets can be found in time polynomial in the size of the MDP using graph search algorithms [10].

The agent's initial state $s_0 \in S$ can belong to either B , S_0 or S_r . However, we only consider the case $s_0 \in S_r$ since otherwise the optimal incentive design is trivially $\gamma_t(s, a) = 0$ for all $t \in \mathbb{N}$.

A. The cost of control

Recall that the agent's first objective function $J_0: \Pi_N(\mathcal{M}) \times \Theta(\mathcal{M}) \rightarrow \mathbb{R}^{|S|}$ is

$$J_0(\pi, \Gamma)(s) = \mathbb{E}_s^{\pi} \left[\sum_{t=1}^N (\mathcal{R}(S_t, \mathcal{A}_t) + \gamma_t(S_t, \mathcal{A}_t)) \right]$$

for all $s \in S$. Let $V_n: S \rightarrow \mathbb{R}$ be the agent's value function at stage n such that

$$V_n(s) := \max_{\pi \in \Pi_N(\mathcal{M})} \mathbb{E}_s^{\pi} \left[\sum_{t=n}^N (\mathcal{R}(S_t, \mathcal{A}_t) + \gamma_t(S_t, \mathcal{A}_t)) \right]$$

for all $s \in S$, where the expectation is taken over the paths that occupy s at stage n . Then, we have the recursive formula

$$V_n(s) = \max_{a \in \mathcal{A}(s)} \mathcal{R}(s, a) + \gamma_n(s, a) + \sum_{s' \in S} \mathcal{P}_{s, a, s'} V_{n+1}(s')$$

for all $1 \leq n \leq N$, where $V_{N+1}(s) = 0$ for all $s \in S$. Let $Q_n: S \times \mathcal{A} \rightarrow \mathbb{R}$ be the agent's Q -function at stage n such that

$$Q_n(s, a) := \mathcal{R}(s, a) + \gamma_n(s, a) + \sum_{s' \in S} \mathcal{P}_{s, a, s'} V_{n+1}(s').$$

By the principle of optimality [1], [12], the agent's optimal policy $\pi_0^* = (d_1^*, d_2^*, \dots, d_N^*)$ is such that, for all $1 \leq n \leq N$, $d_n(s, a') > 0$ only if

$$a' \in \arg \max_{a \in \mathcal{A}(s)} Q_n(s, a).$$

We recursively define

$$\bar{Q}_n(s, a) := \mathcal{R}(s, a) + \sum_{s' \in S} \mathcal{P}_{s, a, s'} \bar{V}_{n+1}(s'), \quad (5)$$

$$\bar{V}_n(s) := \max_{a \in \mathcal{A}(s)} \bar{Q}_n(s, a), \quad (6)$$

for all $s \in S$ and $a \in \mathcal{A}(s)$. For a given $\epsilon \geq 0$, we finally define a real-valued function $\phi_n^\epsilon: S \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\phi_n^\epsilon(s, a) := \begin{cases} \bar{V}_n(s) - \bar{Q}_n(s, a) + \epsilon & \text{if } s \in S_r, a \in \mathcal{A}(s) \\ 0 & \text{otherwise.} \end{cases}$$

For an arbitrarily small $\epsilon > 0$, the value of $\phi_n^\epsilon(s, a)$, referred as the cost of control for the state-action pair (s, a) , is the minimum incentive that should be offered to the agent in order to make the action $a \in \mathcal{A}(s)$ uniquely optimal at stage t . It is worth noting that although the cost of control $\phi_n^\epsilon(s, a)$ depends on the stage number n , it is independent of the objective number, i.e., it is the same for all J_k . This is because the agent's reward function \mathcal{R} is stationary, and therefore, $\bar{V}_n(s)$ and $\bar{Q}_n(s, a)$ do not change with the objective number k as can be seen from (5)-(6).

B. An $\bar{\epsilon}$ -optimal incentive design

To synthesize the minimum incentive sequence, we should specify the actions to be incentivized by the principal at each state for each stage. To this aim, we modify the MDP \mathcal{M} by considering the agent's decision horizon N as another dimension in the state-space.

Definition 6: For an MDP \mathcal{M} and $T=\{1, 2, \dots, N\}$, the *expanded MDP* is a tuple $\bar{\mathcal{M}}=(\bar{S}, \bar{s}_0, \mathcal{A}, \bar{\mathcal{P}}, \mathcal{A}\mathcal{P}, \bar{\mathcal{L}}, \mathcal{R})$ where

- $\bar{S}=S \times T$,
- $\bar{s}_0=(s_0, 1)$ is the initial state,
- $\bar{\mathcal{P}}:\bar{S} \times \mathcal{A} \times \bar{S} \rightarrow [0, 1]$ is such that

$$\bar{\mathcal{P}}_{(s,n),a,(s',n')} = \begin{cases} \mathcal{P}_{s,a,s'} & \text{if } 1 \leq n \leq N-1 \text{ and } n' = n+1 \\ \mathcal{P}_{s,a,s'} & \text{if } n = N \text{ and } n' = 1 \\ 0 & \text{otherwise,} \end{cases}$$

- $\bar{\mathcal{L}}:\bar{S} \rightarrow 2^{\mathcal{A}\mathcal{P}}$ is such that $\bar{\mathcal{L}}((s,t))=\mathcal{L}(s)$ for all $s \in S$ and for all $t \in T$,

and \mathcal{A} , $\mathcal{A}\mathcal{P}$ and \mathcal{R} are as defined for \mathcal{M} .

We note that the transition function $\bar{\mathcal{P}}$ is defined such that the agent's initial state while computing the k -th N stage policy is the state occupied by the agent at $kN+1$ -st stage on the expanded MDP.

Let $\bar{B} \cup \bar{S}_0 \cup \bar{S}_r$ be the partition of the states of $\bar{\mathcal{M}}$ such that if $s \in \bar{B}$, then $(s,n) \in \bar{B}$ for all $n \in T$, and the sets \bar{S}_0 and \bar{S}_r are defined similarly. Then, the principal's objective on $\bar{\mathcal{M}}$ is to induce an agent policy that reaches the set \bar{B} with maximum probability. To synthesize an incentive design under which the optimal agent policy satisfies the desired property, we modify the expanded MDP $\bar{\mathcal{M}}$ by making its states $s \in \bar{B} \cup \bar{S}_0$ absorbing, and denote the resulting MDP by $\bar{\mathcal{M}}'$. Then, for a given $\epsilon \geq 0$, we define the cost of control for a state-action pair on $\bar{\mathcal{M}}$ through the function $\underline{\phi}^\epsilon: \bar{S} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\underline{\phi}^\epsilon((s,n), a) := \begin{cases} \bar{V}_n(s) - \bar{Q}_n(s,a) + \epsilon & \text{if } s \in \bar{S}_r, a \in \mathcal{A}(s) \\ 0 & \text{otherwise.} \end{cases}$$

Let $\Xi(\bar{\mathcal{M}}') \subseteq \Pi(\bar{\mathcal{M}}')$ be a subset of the set of ∞ -stage policies such that $\pi' \in \Xi(\bar{\mathcal{M}}')$ if and only if

$$\pi' \in \arg \max_{\pi \in \Pi(\bar{\mathcal{M}}')} \Pr^\pi(\bar{s}_0 \models \varphi), \quad (7)$$

and for $\epsilon \geq 0$, $f_\epsilon: \Xi(\bar{\mathcal{M}}') \rightarrow \mathbb{R}$ be a function such that

$$f_\epsilon(\pi) := \mathbb{E}_{\bar{s}_0}^\pi \left[\sum_{t=1}^{\infty} \underline{\phi}^\epsilon(S_t, \mathcal{A}_t) \right]. \quad (8)$$

Then, for an arbitrarily small $\bar{\epsilon} > 0$, an $\bar{\epsilon}$ -optimal incentive sequence can be designed in two steps as follows.

Step 1: Compute $\bar{V}_n(s)$ and $\bar{Q}_n(s,a)$ given in (5)-(6), and construct the cost of control function $\underline{\phi}^\epsilon$. Then for the modified expanded MDP $\bar{\mathcal{M}}'$, compute a *stationary deterministic* policy $\tilde{\pi}=(\tilde{d}, \tilde{d}, \dots)$ such that

$$\tilde{\pi} \in \arg \min_{\pi \in \Xi(\bar{\mathcal{M}}')} f_\epsilon(\pi). \quad (9)$$

Step 2: Let $\varrho^\pi \in \text{Paths}^\pi(\mathcal{M})$ be the path followed by the agent. At stage kN where N is the agent's decision horizon and $k \in \mathbb{N}_0$, provide the agent with the incentive sequence $\{\tilde{\gamma}_1, \tilde{\gamma}_2, \dots, \tilde{\gamma}_N\}$ such that

- if $\varrho^\pi[n] \notin B \cup S_0$ for all $n \leq kN$

$$\tilde{\gamma}_n(s,a) := \begin{cases} \underline{\phi}^\epsilon((s,n), a) & \text{if } s \in S_r \text{ and } \tilde{d}((s,n))(a) > 0, \\ \epsilon & \text{if } s \notin S_r \text{ and } \tilde{d}((s,n))(a) > 0, \\ 0 & \text{otherwise,} \end{cases} \quad (10)$$

- $\tilde{\gamma}_n(s,a) := 0$ otherwise.

Under the proposed incentive design (10), the agent's value function V_n satisfies $V_n(s) = \bar{V}_n(s) + (N+1-n)\epsilon$ for all $s \in S, n \leq N$. Additionally, if $\varrho^\pi[n] \notin B \cup S_0$ for all $n \leq kN$, then for all $s \in S, \tilde{d}((s,n))(a) > 0$ implies that the agent's Q -function satisfies

$$\begin{aligned} Q_n(s,a) &= \mathcal{R}(s,a) + \tilde{\gamma}_n(s,a) + \sum_{s' \in S} \mathcal{P}_{s,a,s'} V_{n+1}(s') \\ &= \tilde{\gamma}_n(s,a) + \bar{Q}_n(s,a) + (N-n)\epsilon \\ &= (N+1-n)\epsilon + \bar{V}_n(s) \\ &> (N-n)\epsilon + \bar{V}_n(s) = \max_{a' \in \mathcal{A}(s) \setminus \{a\}} Q_n(s,a'). \end{aligned}$$

Consequently, the agent is guaranteed to take the incentivized actions at each stage until reaching the set $B \cup S_0$.

We now show $\bar{\epsilon}$ -optimality of the proposed incentive design. Note that an optimal incentive design, i.e., $\bar{\epsilon}=0$, does not exist since choosing $\epsilon=0$ in the cost of control function $\underline{\phi}_n^\epsilon$ may not make the incentivized action uniquely optimal for the agent. As a result, the principal may not be able to control the agent's actions by offering such incentives.

We need the following technical lemma to state the main result.

Lemma 1: There exists a policy $\tilde{\pi} \in \arg \min_{\pi \in \Xi(\bar{\mathcal{M}}')} f_0(\pi)$ such that $\xi^{\tilde{\pi}}(s,a) < \infty$ for all $s \in \bar{S}_r$ and $a \in \mathcal{A}(s)$.

Proof (Sketch): The problem of synthesizing a policy $\tilde{\pi}$ such that $\tilde{\pi} \in \arg \min_{\pi \in \Xi(\bar{\mathcal{M}}')} f_0(\pi)$ can be recast as a stochastic shortest path (SSP) problem with dead ends and zero-cost loops. Specifically, the dead ends are the states \bar{S}_0 and zero-cost loops are formed by states \bar{S}_r . The existence of stationary policies for such SSP problems can be established by slightly modifying the statement of Theorem 1 in [13]. Since any stationary policy $\pi \in \Xi(\bar{\mathcal{M}}')$ is guaranteed to reach the set $\bar{B} \cup \bar{S}_0$ with probability 1 within finite number of stages, the result follows. \square

Theorem 1: For any given $\bar{\epsilon} > 0$, there exists $\epsilon > 0$ such that

$$\min_{\pi \in \Xi(\bar{\mathcal{M}}')} f_\epsilon(\pi) \leq \min_{\pi \in \Xi(\bar{\mathcal{M}}')} f_0(\pi) + \bar{\epsilon}.$$

Proof: For any policy $\pi \in \Xi(\bar{\mathcal{M}}')$ such that $\xi^\pi(s,a) < \infty$ for all $s \in \bar{S}_r$ and $a \in \mathcal{A}(s)$, we have

$$f_\epsilon(\pi) = f_0(\pi) + \sum_{s \in \bar{S}_r} \sum_{a \in \mathcal{A}(s)} \xi^\pi(s,a)\epsilon. \quad (11)$$

Now, for a given $\bar{\epsilon} > 0$, we evaluate both sides of the above equation at $\bar{\pi} \in \arg \min_{\pi \in \Xi(\overline{\mathcal{M}'})} f_0(\pi)$, which satisfies the condition $\xi^{\bar{\pi}}(s, a) < \infty$ due to Lemma 1. Choosing

$$\epsilon = \frac{\bar{\epsilon}}{\sum_{s \in \overline{S}_r} \sum_{a \in \mathcal{A}(s)} \xi^{\bar{\pi}}(s, a)} > 0$$

and taking the minimum of the left hand side of (11) over the set $\Xi(\overline{\mathcal{M}'})$, we conclude the result. \square

We conclude this section by noticing a remarkable property of the proposed incentive design. Specifically, to implement the proposed design (10), the principal should use only a simple switch mode which offers the same incentives until the agent reaches the set $B \cup S_0$ and shifts all incentives to zero after the agent either satisfies the principal's objective or fails to satisfy it.

V. COMPUTATION OF AN OPTIMAL INCENTIVE DESIGN

In the previous section, we developed a method to synthesize an $\bar{\epsilon}$ -optimal incentive design which require us to solve a constrained cost minimization problem given in (8). Specifically, to solve the incentive design problem (4a)-(4d), one should synthesize a *stationary deterministic* policy $\bar{\pi}$ such that

$$\bar{\pi} \in \arg \min_{\pi \in \Xi(\overline{\mathcal{M}'})} \mathbb{E}_{s_0}^{\pi} \left[\sum_{t=1}^{\infty} \phi^{\epsilon}(S_t, \mathcal{A}_t) \right] \quad (12)$$

In this section, we develop a method to solve the above optimization problem. For the ease of notation, we consider an scLTL formula of the form $\varphi = \diamond p$. The incentive design for general scLTL formulae is discussed in Section VI.

A. Construction of the feasible policy space

To solve the problem (12), we first represent the set $\Xi(\overline{\mathcal{M}'})$ of feasible policies as a set of policies that maximizes the expected total reward with respect to a specific reward function.

For a given MDP $\overline{\mathcal{M}}$, we partition the set of states into three disjoint sets \overline{B} , \overline{S}_0 , and \overline{S}_r as explained in Section IV, and make the states $s \in \overline{B} \cup \overline{S}_0$ absorbing to form the modified MDP $\overline{\mathcal{M}'}$. For the modified MDP, we define a reward function $r: \overline{S} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$ such that

$$r(s, a) = \begin{cases} \sum_{s' \in \overline{B}} \overline{\mathcal{P}}_{s, a, s'} & \text{if } s \in \overline{S}_r \\ 0 & \text{otherwise.} \end{cases}$$

By making use of the known results, e.g., Theorem 10.100 in [10], it can be easily shown that for any $s \in \overline{S}$ and $\pi \in \Pi(\overline{\mathcal{M}'})$,

$$\mathbb{E}_s^{\pi} \left[\sum_{t=1}^{\infty} r(S_t, \mathcal{A}_t) \right] = \Pr^{\pi}(s \models \varphi)$$

where $\varphi = \diamond p$, $p \in \mathcal{AP}$, and $\{p\} \subseteq \mathcal{L}(s')$ if and only if $s' \in \overline{B}$. Let $x_s^* := \max_{\pi \in \Pi(\overline{\mathcal{M}'})} \Pr^{\pi}(s \models \varphi)$. Then, the problem (12) can be rewritten as

$$\min_{\pi \in \Pi(\overline{\mathcal{M}'})} \mathbb{E}_{s_0}^{\pi} \left[\sum_{t=1}^{\infty} \phi^{\epsilon}(S_t, \mathcal{A}_t) \right] \quad (13a)$$

$$\text{subject to: } \mathbb{E}_{s_0}^{\pi} \left[\sum_{t=1}^{\infty} r(S_t, \mathcal{A}_t) \right] = x_{s_0}^*. \quad (13b)$$

B. Synthesis of an optimal stationary deterministic policy

Using Lemma 1, one can formulate the problem (13a)-(13b) as a linear optimization problem and synthesize an optimal stationary policy. First, we compute the maximum probability of satisfying the specification φ , i.e., $x_{s_0}^* = \max_{\pi \in \Pi(\overline{\mathcal{M}'})} \Pr^{\pi}(s_0 \models \varphi)$, by solving a linear program (LP) [10] (see Chapter 10). Then we solve the following LP

$$\text{minimize}_{\lambda(s, a)} \sum_{s \in \overline{S}_r} \sum_{a \in \mathcal{A}} \lambda(s, a) \phi^{\epsilon}(s, a) \quad (14a)$$

$$\text{subject to: } \sum_{s \in \overline{S}_r} \sum_{a \in \mathcal{A}} \lambda(s, a) r(s, a) = x_{s_0}^* \quad (14b)$$

$$\forall s \in \overline{S}_r, \sum_{a \in \mathcal{A}(s)} \lambda(s, a) - \sum_{s' \in \overline{S}_r} \sum_{a \in \mathcal{A}(s')} \overline{\mathcal{P}}_{s', a, s} \lambda(s', a) = \alpha(s) \quad (14c)$$

$$\forall s \in \overline{S}_r, a \in \mathcal{A}(s), \lambda(s, a) \geq 0 \quad (14d)$$

where $\alpha: \overline{S} \rightarrow \{0, 1\}$ is a function such that $\alpha(s_0) = 1$ and $\alpha(s) = 0$ for all $s \in \overline{S} \setminus \{s_0\}$. The variable $\lambda(s, a)$ denotes the expected residence time in the state-action pair (s, a) [1], [14]. The constraint (14b) ensures that the probability of satisfying the specification φ is maximized, and the constraints (14c) represent the balance between the ‘‘inflow’’ to and ‘‘outflow’’ from states.

For each $s \in \overline{S}_r$ and $a \in \mathcal{A}(s)$, let $\lambda^*(s, a)$ be optimal decision variables in (14a)-(14d). An optimal stationary policy $\pi^* = \{d^*, d^*, \dots\}$ that solves the problem (13a)-(13b) is then given by

$$d^*(s, a) := \begin{cases} \frac{\lambda^*(s, a)}{\sum_{a \in \mathcal{A}(s)} \lambda^*(s, a)} & \text{if } \sum_{a \in \mathcal{A}(s)} \lambda^*(s, a) > 0 \\ \text{arbitrary} & \text{otherwise} \end{cases} \quad (15)$$

for $s \in \overline{S}_r$, and $d^*(s, a) = 1$ for an arbitrary $a \in \mathcal{A}(s)$ for $s \notin \overline{S}_r$.

We note that a policy constructed through (15) is randomized in general. One can argue that choosing one of the actions $a \in \mathcal{A}(s)$ such that $d^*(s, a) > 0$ deterministically yields an optimal stationary deterministic policy. However, the following example illustrates that such an approach may result in an infeasible policy for the problem (14a)-(14d).

Example 1: Consider the MDP given in Fig. 2, where the cost of control ϕ^{ϵ} is such that $\phi^{\epsilon}(s_1, a_2) = 1$ and $\phi^{\epsilon}(s, a) = 0$ otherwise. Suppose that the specification is $\varphi = \diamond s_2$, i.e., $r(s_1, a_2) = 1$ and $r(s, a) = 0$ otherwise. For the LP (14a)-(14d), a set of optimal decision variables is given by $\lambda^*(s_0, a_1) = 2$, $\lambda^*(s_1, a_1) = 1$, and $\lambda^*(s_1, a_2) = 1$. Therefore, an optimal policy synthesized through (15) is $d^*(s_0, a_1) = 1$, $d^*(s_1, a_1) = 1/2$, and $d^*(s_1, a_2) = 1/2$. Clearly, if we consider

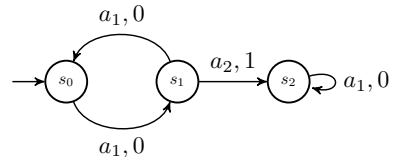


Fig. 2: An MDP example for which arbitrarily choosing one of the optimal actions and taking it deterministically yields an infeasible policy.

a deterministic policy such that $d(s_1, a_1)=1$, the probability of satisfying the specification φ under this policy is zero. Hence, choosing an arbitrary action $a \in \mathcal{A}(s)$ such that $d^*(s, a) > 0$ deterministically violates the constraint and yields an infeasible policy. \triangleleft

As Example 1 illustrates, a structured approach is required to synthesize an optimal deterministic policy from the solution of the LP (14a)-(14d). Let v^* be the optimal value of the LP in (14a)-(14d). To synthesize an optimal deterministic policy, we first solve the following LP,

$$\text{minimize}_{\lambda(s,a)} \sum_{s \in \bar{S}_r} \sum_{a \in \mathcal{A}} \lambda(s, a) \quad (16a)$$

$$\text{subject to: } \sum_{s \in \bar{S}_r} \sum_{a \in \mathcal{A}} \lambda(s, a) r(s, a) = x_{s_0}^* \quad (16b)$$

$$\sum_{s \in \bar{S}_r} \sum_{a \in \mathcal{A}} \lambda(s, a) \phi^\epsilon(s, a) = v^* \quad (16c)$$

$$\forall s \in \bar{S}_r, \sum_{a \in \mathcal{A}(s)} \lambda(s, a) - \sum_{s' \in \bar{S}_r} \sum_{a \in \mathcal{A}(s)} \bar{\mathcal{P}}_{s', a, s} \lambda(s', a) = \alpha(s) \quad (16d)$$

$$\forall s \in \bar{S}_r, a \in \mathcal{A}(s), \lambda(s, a) \geq 0. \quad (16e)$$

From the optimal decision variables $\lambda^*(s, a)$ of (16a)-(16e), an optimal policy $\pi^* = \{d^*, d^*, \dots\}$ can be generated as follows. Let $\mathcal{A}^*(s) := \{a \in \mathcal{A}(s) : \lambda^*(s, a) > 0\}$. If $\mathcal{A}^*(s) \neq \emptyset$, we choose $d^*(s, a) = 1$ for an arbitrary $a \in \mathcal{A}^*(s)$, and if $\mathcal{A}^*(s) = \emptyset$, we choose $d^*(s, a) = 1$ for an arbitrary $a \in \mathcal{A}(s)$.

Proposition 1: A stationary deterministic policy generated from the optimal decision variables $\lambda^*(s, a)$ of (16a)-(16e) is a solution to the problem (13a)-(13b).

A proof of Proposition 1 can be found in Appendix I. Intuitively, the LP in (16a)-(16e) computes the minimum expected time to reach the set \bar{B} with probability $x_{s_0}^*$ with the cost of v^* . Therefore, if $\lambda^*(s, a) > 0$, by taking the action $a \in \mathcal{A}(s)$, the agent has to “get closer” to the set \bar{B} with nonzero probability. Otherwise, the minimum expected time to reach the set \bar{B} would be strictly decreased. Consequently, by choosing an arbitrary action $a \in \mathcal{A}^*(s)$, the agent is guaranteed to reach the set \bar{B} with the desired probability.

VI. INCENTIVE DESIGN FOR GENERAL sCLTL SPECIFICATIONS

In previous sections, we have developed methods to synthesize $\bar{\epsilon}$ -optimal incentive designs for reachability specifications $\varphi = \diamond p$. For such specifications, the principal induces the desired agent behavior by sharing only the incentive sequences with the agent. In other words, the principal does not have to inform the agent explicitly about the specification. In this section, we show that for general sCLTL formulae, the problem (4a)-(4d) may not have a feasible solution, in which case the principal must share its objective with the agent to induce the desired behavior.

To solve the problem (4a)-(4d) for general sCLTL formulae, one needs to utilize the techniques from automata theory [10]. In particular, we use the fact that for any sCLTL formula

φ built up from \mathcal{AP} , we can construct a *deterministic finite automata* (DFA) $A_\varphi = (\mathcal{Q}, q_0, 2^{\mathcal{AP}}, \delta_\varphi, \mathcal{F})$ where \mathcal{Q} is a finite set of memory states, $2^{\mathcal{AP}}$ is the alphabet, $\delta_\varphi: \mathcal{Q} \times 2^{\mathcal{AP}} \rightarrow \mathcal{Q}$ is a transition function and $\mathcal{F} \subseteq \mathcal{Q}$ is the set of accepting states [11]. Then, after forming the expanded MDP $\bar{\mathcal{M}}$ for a given MDP \mathcal{M} and a decision horizon N as explained in Section IV-B, one can construct the product MDP which is defined as follows.

Definition 7: Let $\bar{\mathcal{M}} = (\bar{\mathcal{S}}, \bar{s}_0, \mathcal{A}, \bar{\mathcal{P}}, \mathcal{AP}, \bar{\mathcal{L}})$ be an expanded MDP and $A_\varphi = (\mathcal{Q}, q_0, 2^{\mathcal{AP}}, \delta_\varphi, \mathcal{F})$ be a DFA. The *product MDP* $\mathcal{M}_p = (\mathcal{S}_p, s_{0_p}, \mathcal{A}, \mathbb{P}, \mathcal{AP}, \mathcal{L}_p, \mathcal{F}_p)$ is a tuple where

- $\mathcal{S}_p = \bar{\mathcal{S}} \times \mathcal{Q}$,
- $s_{0_p} = (\bar{s}_0, q)$ such that $q = \delta(q_0, \bar{\mathcal{L}}(\bar{s}_0))$,
- $\mathbb{P}((s, q), a, (s', q')) = \begin{cases} \bar{\mathcal{P}}_{s, a, s'} & \text{if } q' = \delta(q, \bar{\mathcal{L}}(s')) \\ 0 & \text{otherwise,} \end{cases}$
- $\mathcal{L}_p((s, q)) = \{q\}$,
- $\mathcal{F}_p = \bar{\mathcal{S}} \times \mathcal{F}$.

The incentive design problem (4a)-(4d) can now be solved on the product MDP \mathcal{M}_p in three steps. First, we partition the states of \mathcal{M}_p into three disjoint sets. Let $B := \mathcal{F}_p$, S_0 be the set of states that have zero probability of reaching the set B , and $S_r := \mathcal{S}_p \setminus B \cup S_0$. Second, we form the modified product MDP \mathcal{M}'_p by making all states $B \cup S_0$ absorbing. Finally, we apply the methods developed in Section IV to synthesize an $\bar{\epsilon}$ -optimal incentive sequence on \mathcal{M}'_p .

Note that the incentive sequence is designed on the product MDP \mathcal{M}_p . Therefore, the principal must share the DFA structure, i.e., its objective, with the agent to be able to use the computed design. However, for the existence of a solution to the problem (4a)-(4d), the incentive sequence should be designed on the MDP \mathcal{M} . The following example illustrates that the problem (4a)-(4d) may have no feasible solution, even though the existence of an $\bar{\epsilon}$ -optimal incentive sequence on \mathcal{M}_p is guaranteed.

Example 2: Consider the MDP given in Fig. 3, where the numbers next to actions a_i represent the transition probabilities, e.g., $\mathcal{P}_{s_0, a_1, s_1} = 0.4$, and the letters next to state numbers represent labels, e.g., $\mathcal{L}(s_0) = A$. Let the agent’s decision horizon be $N=3$, and the reward function \mathcal{R} be such that $\mathcal{R}(s_0, a_1) = 1$ and $\mathcal{R}(s, a) = 0$ otherwise. Additionally, let the principal’s objective be expressed by the sCLTL formula $\varphi = \diamond(B \wedge \diamond C)$, i.e., first visit state B and then state C . The maximum probability of satisfying φ is $x_0^* = 0.5$, which can be computed by solving an LP [10]. The value x_0^* is attainable if and only if the agent takes the action $a_2 \in \mathcal{A}(s_0)$ with probability 1 after visiting state s_2 .

The principal should decide on which actions to incen-

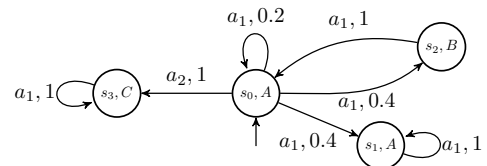


Fig. 3: An MDP example for which there exists no feasible incentive design for the sCLTL specification $\varphi = \diamond(B \wedge \diamond C)$.

one to the target state. Specifically, the shortest path would take 8 stages to reach the target state and cost 12 UR to the principal, whereas the lowest cost path takes 10 stages to reach the target and cost 9 UR. Quantitatively, the proposed incentive design allows the principal to save 25% of the resources that would be paid to the agent if it was to follow the shortest path.

B. Incentives for general scLTL specifications

In this example, we consider the same grid world environment introduced in the previous example with different state labels. The agent's decision horizon is $N=4$, and its objective is to reach the state labeled as C in Fig. 5. The principal's objective is to induce an agent policy that satisfies the scLTL specification $\varphi=\diamond(A \wedge \diamond(B \wedge \diamond C))$, i.e., the agent should first visit state A , then B , and then C , with probability 1.

The agent receives the reward of 2 for transitioning to the top left state and the reward of 5 for transitioning to the top right state. Its optimal path in the absence of incentives is shown in Fig. 5 with blue arrows (top path). We synthesize an optimal incentive sequence under which the agent's optimal path is shown in Fig. 5 with red arrows (bottom path).

The total cost of the incentives to the principal is computed as $2+13\epsilon$ units of resources. Specifically, the principal provides $2+\epsilon$ incentives for the *right* action in the start state and then ϵ incentives at each stage for desired actions. An interesting property of the incentivized (red) path is that the agent stays in the same state in third stage by taking *down* action. This is due to the fact that the state s on the left of the state labeled as B has value $V_n(s)=0$ for all n . Therefore, the principal wants that state to be the agent's initial state when it computes its second 4-stage policy. By doing so, the principal ensures that the states s' occupied by the agent in the next 4 stages will always have a value zero, i.e., $V_n(s')=0$ if $\sum_{a \in \mathcal{A}(s')} \mu_{4+n}^\pi(s', a) > 0$, and therefore the cost of control will only be ϵ .

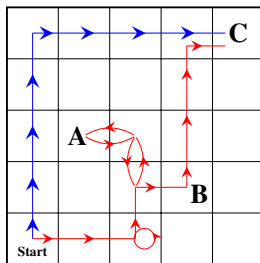


Fig. 5: The motion of an agent on a grid world. The agent's decision horizon is $N=4$, and it starts from the *Start* state. The principal's objective is to induce an agent policy that satisfies the scLTL specification $\varphi=\diamond(A \wedge \diamond(B \wedge \diamond C))$, i.e., first visit A , then B , and then C . The optimal path of the agent in the absence of incentives is shown by blue arrows (top path). Red arrows indicate the agent's optimal path under the provided incentives (bottom path).

VIII. CONCLUSIONS AND FUTURE DIRECTIONS

We considered a principal-agent model and studied the problem of designing an optimal sequence of incentives that

the principal should offer to the agent in order to induce a desired agent behavior expressed as a syntactically co-safe linear temporal logic (scLTL) formula. For reachability objectives, we presented a polynomial-time algorithm to synthesize an incentive design that minimizes the cost to the principal. By providing an example scenario, we showed that a feasible incentive design may not exist for general scLTL formulae, and the principal may need to share its objective with the agent to induce the desired behavior. Furthermore, we provided sufficient conditions under which the principal can induce the desired behavior without sharing the scLTL formula with the agent.

The results that we present in this paper are obtained under the assumptions that the agent's reward function and the length of its decision horizon are known by the principal. An interesting future direction may be to develop methods to infer the length of the agent's decision horizon through perfect/imperfect observations, or to design an incentive sequence that does not require the knowledge of the length of the decision horizon.

REFERENCES

- [1] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.
- [2] M. Kloetzer and C. Mahulea, "Multi-robot path planning for syntactically co-safe LTL specifications," in *International Workshop on Discrete Event Systems (WODES)*, 2016, pp. 452–458.
- [3] A. Armando, R. Carbone, and L. Compagna, "LTL model checking for security protocols," in *IEEE Computer Security Foundations Symposium*, 2007, pp. 385–396.
- [4] L. Tan, O. Sokolsky, and I. Lee, "Specification-based testing with linear temporal logic," in *IEEE International Conference on Information Reuse and Integration*, 2004, pp. 493–498.
- [5] H. Zhang and D. C. Parkes, "Value-based policy teaching with active indirect elicitation," in *AAAI Conference on Artificial Intelligence*, 2008, pp. 208–214.
- [6] H. Zhang, Y. Chen, and D. C. Parkes, "A general approach to environment design with one agent," in *International Joint Conference on Artificial Intelligence*, 2009, pp. 2002–2014.
- [7] H. Zhang, D. C. Parkes, and Y. Chen, "Policy teaching through reward function learning," in *ACM Conference on Electronic commerce*, 2009, pp. 295–304.
- [8] Y. Chen, J. Kung, D. C. Parkes, A. D. Procaccia, and H. Zhang, "Incentive design for adaptive agents," in *International Conference on Autonomous Agents and Multiagent Systems*, 2011, pp. 627–634.
- [9] Y.-C. Ho, P. B. Luh, and G. J. Olsder, "A control-theoretic view on incentives," *Automatica*, vol. 18, no. 2, pp. 167–179, 1982.
- [10] C. Baier and J.-P. Katoen, *Principles of Model Checking*. MIT Press, 2008.
- [11] C. Belta, B. Yordanov, and E. A. Gol, *Formal Methods for Discrete-Time Dynamical Systems*. Springer, 2017.
- [12] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [13] F. Teichteil-Königsbuch, "Stochastic safest and shortest path problems," in *AAAI Conference on Artificial Intelligence*, 2012.
- [14] K. Etessami, M. Kwiatkowska, M. Y. Vardi, and M. Yannakakis, "Multi-objective model checking of Markov decision processes," in *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, 2007, pp. 50–65.
- [15] L. Gurobi Optimization, "Gurobi optimizer reference manual," 2018. [Online]. Available: <http://www.gurobi.com>
- [16] M. ApS, *MOSEK Optimizer API for Python. Version 8.1.*, 2019. [Online]. Available: <https://docs.mosek.com/8.1/pythonapi/index.html>
- [17] S. Diamond and S. Boyd, "CVXPY: A Python-embedded modeling language for convex optimization," *Journal of Machine Learning Research*, vol. 17, no. 83, pp. 1–5, 2016.
- [18] R. Serfozo, *Basics of Applied Stochastic Processes*. Springer, 2009.

APPENDIX I
PROOF OF PROPOSITION 1

Note that any deterministic policy constructed from the optimal decision variables $\lambda^*(s, a)$ of (16a)-(16e) can only violate the reachability constraint (14b). In other words, the constructed policy is guaranteed to minimize the expected total cost.

We will need the following result to prove Proposition 1. For a given policy π , let $Reach^\pi(s, s')$ denote the probability of reaching s' from s under π . Note that $Reach^\pi(s, s) = \sum_{a \in \mathcal{A}(s)} d(s, a) Reach^\pi((s, a), s)$ where $Reach^\pi((s, a), s)$ is the probability of reaching state s from state action pair (s, a) under the policy π . Finally, let $\xi^\pi(s) := \sum_{a \in \mathcal{A}(s)} \xi^\pi(s, a)$, and note that for any $\xi^\pi(s) < \infty$, we have [18]

$$\xi^\pi(s) = \frac{Reach^\pi(s_0, s)}{1 - \sum_{a \in \mathcal{A}(s)} d(s, a) Reach^\pi((s, a), s)}. \quad (17)$$

To prove the claim of Proposition 1, we show that any policy that is constructed by choosing actions $a \in \mathcal{A}(s)$ such that $\lambda^*(s, a) > 0$ deterministically is optimal.

Let $\bar{\pi} = \{\bar{d}, \bar{d}, \dots\}$ be the stationary randomized policy constructed from $\lambda^*(s, a)$ of LP (16) through the formula (15). Additionally, let $\tilde{\pi} = \{\tilde{d}, \tilde{d}, \dots\}$ be a stationary randomized policy such that $\tilde{d}(s^*, a^*) = 1$, and $\tilde{d}(s) = \bar{d}(s)$ for all $s \in S \setminus \{s^*\}$. Informally, in state s^* , we choose one of the *active* actions deterministically and do not change the rest of the policy.

We first show that $\Pr^{\tilde{\pi}}(s_0 | \varphi) < x_{s_0}^*$ implies $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$. Then by showing that $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$ cannot be true, we conclude that $\Pr^{\tilde{\pi}}(s_0 | \varphi) = x_{s_0}^*$.

As for the first claim, suppose for contradiction that $\Pr^{\tilde{\pi}}(s_0 | \varphi) < x_{s_0}^*$ and $Reach^{\tilde{\pi}}((s^*, a^*), s^*) < 1$. Note that if $Reach^{\tilde{\pi}}((s^*, a^*), s^*) < 1$, then $Reach^{\tilde{\pi}}(s^*, s^*) < 1$. Therefore, $Reach^{\tilde{\pi}}(s, s) < 1$ for all $s \in S_r$ satisfying $Reach^{\tilde{\pi}}(s^*, s) > 0$. Additionally, as $\bar{d}(s') = \tilde{d}(s')$ for all s' such that $Reach^{\tilde{\pi}}(s^*, s') = 0$, we have $Reach^{\tilde{\pi}}(s', s') < 1$. Consequently, probability of leaving the set S_r is 1. Since all actions that are chosen by policy $\tilde{\pi}$ satisfy $x_s^* = \mathcal{P}_{s, a, s'} x_{s'}^*$ where x_s^* is the maximum probability of reaching the set B from the state s (see e.g. Chapter 10 in [10]), probability of entering the set B must be equal to $x_{s_0}^*$. This raises a contradiction.

As for the second claim, suppose that $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$. Then, $Reach^{\bar{\pi}}((s^*, a^*), s^*) = 1$ since $\bar{\pi}$ differs from $\tilde{\pi}$ only in the state s^* . We now construct a policy $\hat{\pi}$ such that $\hat{d}(s) = \bar{d}(s)$ for all $s \in S \setminus \{s^*\}$, $\hat{d}(s^*, a^*) = 0$, and

$$\hat{d}(s^*, a_i) = \frac{\bar{d}(s^*, a_i)}{\sum_{a \in \mathcal{A}(s^*) \setminus \{a^*\}} \bar{d}(s^*, a_i)}. \quad (18)$$

Note that $\hat{\pi}$ satisfies $\Pr^{\hat{\pi}}(s_0 | \varphi) = x_{s_0}^*$. By showing that $\hat{\pi}$ attains an objective value in (16) that is strictly smaller than the policy $\bar{\pi}$, we will conclude that $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$ cannot be possible.

For the ease of notation, let $\bar{a}_i := \bar{d}(s^*, a_i)$, $\bar{R}_i := Reach^{\bar{\pi}}((s^*, a_i), s^*)$, and $\hat{a}_i := \hat{d}(s^*, a_i)$, and $\hat{R}_i := Reach^{\hat{\pi}}((s^*, a_i), s^*)$. Without loss of generality, we choose $a_1 = a^*$. By the construction of $\hat{\pi}$, it can be shown that

$$\xi^{\hat{\pi}}(s^*) = (1 - \bar{a}_1) \xi^{\bar{\pi}}(s^*) - \frac{\bar{a}_1(\bar{R}_1 - 1)(1 - \bar{a}_1)}{C} \quad (19)$$

where $C := (1 - \sum_{i=1}^n \bar{a}_i \bar{R}_i)(1 - a_1 - \sum_{i=2}^n \bar{a}_i \bar{R}_i) > 0$. Note that $\xi^{\hat{\pi}}(s^*, a_i)(1 - \bar{a}_1) = \xi^{\bar{\pi}}(s^*, a_i)$ due to (18). Then, since $\bar{a}_1 > 0$ and $\bar{R}_1 = 1$, we have

$$\xi^{\hat{\pi}}(s^*, a_i) \leq \xi^{\bar{\pi}}(s^*, a_i) \quad (20)$$

for all a_i $i=2, 3, \dots, n$ and $\xi^{\hat{\pi}}(s^*, a_1) < \xi^{\bar{\pi}}(s^*, a_1)$. Consequently, $\hat{\pi}$ attains an objective value in (16) that is strictly smaller than the policy $\bar{\pi}$.

Finally, since $Reach^{\bar{\pi}}((s^*, a^*), s^*) = 1$ cannot be true, $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$ cannot be true. If $Reach^{\tilde{\pi}}((s^*, a^*), s^*) = 1$ is not true, $\Pr^{\tilde{\pi}}(s_0 | \varphi) < x_{s_0}^*$ is not true. This concludes the proof. \square