

# Deceptive Planning for Resource Allocation

Shenghui Chen, Yagiz Savas, Mustafa O. Karabag, Brian M. Sadler, Ufuk Topcu

**Abstract**—We consider a team of autonomous agents that navigate in an adversarial environment and aim to achieve a task by allocating their resources over a set of target locations. An adversary in the environment observes the autonomous team’s behavior to infer their objective and responds against the team. In this setting, we propose strategies for controlling the density of the autonomous team so that they can deceive the adversary regarding their objective while achieving the desired final resource allocation. We first develop a prediction algorithm based on the principle of maximum entropy to express the team’s behavior expected by the adversary. Then, by measuring the deceptiveness via Kullback-Leibler divergence, we devise convex optimization-based planning algorithms that deceive the adversary by either exaggerating the behavior towards a decoy allocation strategy or creating ambiguity regarding the final allocation strategy. A user study with 320 participants demonstrates that the proposed algorithms are effective for deception and reveal the inherent biases of participants towards proximate goals.

## I. INTRODUCTION

In many scenarios, a team of autonomous agents needs to accomplish a task in an adversarial environment. Consider a swarm of autonomous drones tasked with securing a region and conducting surveillance missions amidst potential intruders [1], or autonomous military robots aiming to control strategic locations from opposing parties on battlefields [2]. Operating in such hostile environments often leads to the team inadvertently leaking critical information, enabling the adversary to devise counter-strategies that thwart task completion.

In this paper, we present a systematic approach for a team of agents to complete their task while managing information leakage through deliberate deception. We consider a setting in which a team consisting of a large number of autonomous agents distribute their resources, i.e., team members, to certain goal locations in an environment. Knowing only that the true goal distribution is among a limited set of distributions, an adversary observes the team’s behavior in the environment to deduce the team’s goal distribution and respond against the team. In this setting, we develop a swarm control strategy for the autonomous team to allocate their resources to desired locations in a way to deceive the adversary regarding the true distribution. The approach is summarized in Fig. 1.

We model the prior prediction of the adversary via the maximum entropy principle [3], [4]. Specifically, inspired by the experimental studies from the psychology literature [5],

we assume that the adversary expects the autonomous team to reach their final allocation in the environment through the shortest paths with a certain degree of inefficiency. We generate the expected behavior by solving a constrained optimization problem that combines a cost minimization objective with an entropy regularization.

We model the behavior of the team in the environment as a Markov decision process (MDP). MDPs model sequential decision-making problems under uncertainty and have been widely used to control the high-level behavior of autonomous agents in various applications [6]–[10]. We utilize MDPs to synthesize a strategy that controls the density of the team members in the environment while they progress toward achieving the desired final resource allocation. Specifically, we synthesize a strategy that maximizes the deceptiveness of the transient behavior while guaranteeing the attainment of the intended final allocation.

We quantify the deceptiveness of the team’s behavior as a function of the statistical distance between the observed behavior and the behavior expected to achieve the true objective. In particular, we consider two types of deception, namely, exaggeration and ambiguity, and show how Kullback-Leibler divergence between certain distributions can be used to develop deception metrics.

This paper has three main contributions. First, we show that an entropy-regularized cost minimization problem in MDPs subject to multiple probabilistic constraints can be formulated as a convex optimization problem and solved efficiently via off-the-shelf solvers. Unlike the existing literature on deception that typically focuses on a single agent with a single reachability objective, this work enables the modeling of adversary predictions in scenarios that involve a swarm of agents with multiple reachability objectives. Second, we introduce novel metrics to quantify the deceptiveness of the team’s behavior and present efficient convex optimization-based algorithms to synthesize strategies that control the density of the team and yield globally optimal deceptive behaviors while satisfying multiple reachability constraints. Third, we validate the deceptiveness of the synthesized strategies via a user study with 320 participants. Our results show that the proposed deceptive algorithms are effective and reveal the inherent biases of participants towards goals closer to the starting point.

**Related work:** Several lines of work are related to the deception problem considered in this paper. The most closely related ones are the authors’ previous work on supervisory control [11] and deception under uncertainty [12]. The former studies how to deceive a supervisor who provides a reference policy for the agent to follow, inspiring our use

S. Chen, Y. Savas, M. O. Karabag, and U. Topcu are with the University of Texas at Austin, TX, USA. E-mails: {shenghui.chen, yagiz.savas, karabag, utopcu}@utexas.edu

B. M. Sadler is with the U.S. Army Research Laboratory, MD, USA. E-mail: brian.m.sadler6.civ@army.mil

of hypothesis testing theory to formulate different deception techniques. The latter considers a deception problem in a single-agent setting, presenting a maximum-entropy-based algorithm for prediction and a linear-programming-based algorithm to optimally reach a single goal. Unlike [12], we develop a convex optimization-based prediction algorithm that incorporates reachability objectives for multiple goals. Rather than defining a cost function based on prediction probabilities, we employ Kullback-Leibler divergence to quantify deception.

There is a large body of literature on single-agent deception problems in which the agent aims to reach its goal while deceiving outside observers. The paper [13] presents a gradient-descent-based approach to synthesize locally optimal deceptive strategies for reaching a single goal in deterministic environments. They consider both exaggeration and ambiguity types of deception by quantifying deceptiveness as a function of prediction probabilities. Unlike [13], we quantify deceptiveness as a function of the statistical distances between density distributions. The paper [14] introduces the notion of the last deceptive point in an environment and presents heuristic approaches to synthesize deceptive policies based on this notion. Similarly, [15] develops deceptive strategies by modeling the observer predictions as a stochastic transition system over potential goals. Although the techniques presented in these works are quite insightful, they are designed for single-objective scenarios and are not applicable to situations requiring a team to allocate a certain fraction of their resources over multiple targets.

Game-theoretic approaches are also commonly used to develop deception strategies in various applications. In [16]–[18], the authors develop several algorithms to utilize decoys for deception in hypergames. Unlike the efficient algorithms presented in this paper, hypergame formulations, in general, yield computationally intractable solutions that can hardly be applied to large-scale systems. Other research [19] and [20] develop deceptive strategies for specific game types, focusing on finitely repeated and single-stage games, respectively. Our work differs from these papers as we consider a dynamic system model where an autonomous team needs to navigate in an environment to achieve an objective eventually.

Finally, the literature on goal recognition is also closely related to deception. In [21]–[23], the authors develop several algorithms for observers to infer an agent’s goal based on its past behavior. These algorithms typically focus on deterministic environments and assume that the agent aims to reach one of the finitely many goals. Since we model the team’s behavior as an MDP, the inference techniques presented in this paper also apply to stochastic environments. Moreover, unlike the existing algorithms on goal recognition, the proposed maximum-entropy-based approach can handle scenarios in which the team aims to reach multiple goals with associated probabilities.

## II. PRELIMINARIES

**Notation:** For a given set  $\mathcal{S}$ , we denote its cardinality by  $|\mathcal{S}|$ . We define  $\mathbb{N} := \{1, 2, 3, \dots\}$ ,  $\mathbb{R} := (-\infty, \infty)$ , and  $\mathbb{R}_{\geq 0} :=$

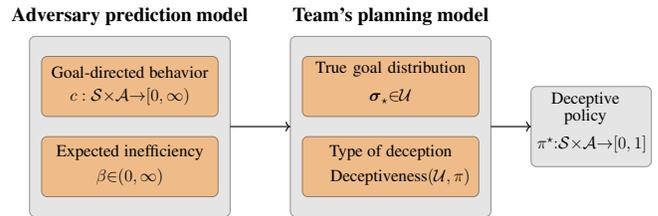


Fig. 1: The proposed deceptive resource allocation approach. For a goal distribution (allocation), the adversary prediction model describes how the adversary expects the autonomous team to achieve its final distribution. Based on the predicted policies, the team generates a deceptive policy either exaggerating the team’s behavior toward a decoy goal distribution or creating ambiguity regarding the true goal distribution while achieving the desired final allocation.

$[0, \infty)$ . For a matrix  $M \in \mathbb{R}^{n \times m}$ , we denote its  $(i, j)$ -th element by  $M_{i,j}$  and its transpose by  $M^T$ . Finally, for a constant  $K \in \mathbb{N}$ , we denote the set  $\{1, 2, \dots, K\}$  by  $[K]$ .

### A. Markov decision processes

We consider a team consisting of a large number of autonomous agents. We model the behavior of the team in a stochastic environment with a Markov decision process.

**Definition 1:** A *Markov decision process* (MDP) is a tuple  $\mathbb{M} = (\mathcal{S}, \alpha, \mathcal{A}, P)$  where  $\mathcal{S}$  is a finite set of states,  $\alpha : \mathcal{S} \rightarrow [0, 1]$  is an initial state distribution,  $\mathcal{A}$  is a finite set of actions, and  $P : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is a transition probability function such that  $\sum_{s' \in \mathcal{S}} P(s, a, s') = 1$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ , where  $\mathcal{A}(s) \subseteq \mathcal{A}$  is the set of available actions in state  $s$ .

For notational convenience, we denote the transition probability  $P(s, a, s')$  by  $P_{s,a,s'}$ . A state  $s \in \mathcal{S}$  is said to be *absorbing* if  $P_{s,a,s} = 1$  for all  $a \in \mathcal{A}(s)$ .

We control the temporal density of the team through a policy to be applied by every member of the team.

**Definition 2:** For an MDP  $\mathbb{M}$ , a *policy*  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is a mapping such that  $\sum_{a \in \mathcal{A}(s)} \pi(s, a) = 1$  for all  $s \in \mathcal{S}$ . We denote the set of all policies by  $\Pi(\mathbb{M})$ .

We note that a policy  $\pi$  is traditionally referred to as a *stationary* policy [7]. Although it is possible to consider more general policy classes, the set  $\Pi(\mathbb{M})$  is sufficient without loss of generality for the purposes of this paper.

A *path* is a sequence  $\varrho = s_1 a_1 s_2 a_2 s_3 \dots$  of states and actions which satisfies that  $\alpha(s_1) > 0$  and  $P_{s_t, a_t, s_{t+1}} > 0$  for all  $t \in \mathbb{N}$ . We define the set of all paths in  $\mathbb{M}$  with initial distribution  $\alpha$  generated under the policy  $\pi$  by  $Paths_{\mathbb{M}}^{\pi, \alpha}$  and use the standard probability measure over the set  $Paths_{\mathbb{M}}^{\pi, \alpha}$  [24]. Let  $\varrho[t] := s_t$  denote the state visited at the  $t$ -th step along  $\varrho$ . For a given state  $s \in \mathcal{S}$ , we define

$$\Pr^{\pi}(\alpha \models \diamond s) := \Pr \{ \varrho \in Paths_{\mathbb{M}}^{\pi, \alpha} : \exists t \in \mathbb{N}, \varrho[t] = s \}$$

as the probability with which the paths generated in  $\mathbb{M}$  with initial distribution  $\alpha$  under  $\pi$  reaches the state  $s \in \mathcal{S}$ .

### III. PROBLEM STATEMENT

We consider a team of autonomous agents that are distributed in a stochastic environment. The team aims to navigate through the environment and achieve a desired final distribution, expressing the optimal allocation of resources to certain goal locations. There is an adversary that observes the team's behavior and aims to predict their final distribution to respond against the team's allocation. We note that the distribution of the resources could be an outcome of an underlying game (e.g., a Colonel Blotto game [25], a zero-sum matrix game [26]) between the team and the adversary. We study the problem of generating a swarm-control policy for the team members so that they deceive the adversary regarding their final distribution for as long as possible while eventually achieving the desired final distribution.

Formally, let  $\mathcal{U} = \{\sigma_1, \dots, \sigma_N\}$  be a finite set of goal distributions where  $\sigma_*$  is the *true goal distribution*. If the adversary knew the true goal distribution, then deception would not be possible as the adversary would respond optimally regardless of the team's behavior in the environment. However, deception becomes possible when the adversary only knows that the true goal distribution belongs to the set  $\mathcal{U}$  of potential goal distributions.

For a given set  $\mathcal{U}$  of potential goal distribution and a policy  $\pi \in \Pi(\mathbb{M})$ , let  $Deceptiveness(\mathcal{U}, \pi)$  be a measure that quantifies the deceptiveness of the team's behavior. In this paper, we aim to synthesize a policy  $\pi^*$  such that

$$\pi^* \in \arg \max_{\pi \in \Pi(\mathbb{M})} Deceptiveness(\mathcal{U}, \pi) \quad (1a)$$

$$\text{subject to: } \Pr^\pi(\alpha \models \diamond g_i) = \sigma_*(g_i) \text{ for all } g_i \in \mathcal{G}. \quad (1b)$$

The above problem aims to enable the team to deceive the adversary regarding their true objective while guaranteeing that they achieve the final distribution  $\sigma_*$ . In what follows, we discuss how to formally define the deceptiveness of a team's induced behavior and develop algorithms to solve the problem in (1a)–(1b).

Throughout the paper, we assume that the problem in (1a)–(1b) is feasible. The validity of this assumption for a given problem instance can be verified efficiently by solving a linear program, as described in [27]. We note that the feasibility assumption holds in many practical settings. For example, it holds when the MDP model has deterministic transitions and there is at least one path from the initial state to each goal state. Some problem instances violate the assumption due to an unachievable goal state distribution. We do not consider these instances as deception is the main focus of this paper.

### IV. EXPRESSING PREDICTIONS THROUGH THE PRINCIPLE OF MAXIMUM ENTROPY

To deceive the adversary about the goal distribution, the team needs to know how they associate the observed behavior with a goal distribution. In this section, we introduce an inference model based on the principle of maximum entropy, which characterizes the adversary's predictions.

Experimental studies show that observers typically expect an agent's behavior to be goal-directed and efficient [5]. This expectation can be expressed through the principle of maximum entropy, which prescribes a probability distribution that is "maximally noncommittal with regard to missing information" [3]. In particular, suppose that the adversary believes that the true goal distribution of the team is  $\sigma_i$ . Then, the team's expected behavior for achieving the final distribution  $\sigma_i$  is described by a policy  $\bar{\pi}_i \in \Pi(\mathbb{M})$  such that

$$\bar{\pi}_i \in \arg \min_{\pi \in \Pi(\mathbb{M})} \mathbb{E}^\pi \left[ \sum_{t=1}^{\infty} \left( c(s_t, a_t) - \beta H(\pi(\cdot|s_t)) \right) \right] \quad (2a)$$

$$\text{subject to: } \Pr^\pi(\alpha \models \diamond g_i) = \sigma_i(g_i) \text{ for all } g_i \in \mathcal{G}. \quad (2b)$$

In (2),  $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_{\geq 0}$  is a cost function that specifies the cost incurred by the team while navigating in the environment. The term  $H(\pi(\cdot|s)) = \sum_{a \in \mathcal{A}(s)} \pi(s, a) \log \pi(s, a)$  denotes the entropy of the policy  $\pi$  in state  $s \in \mathcal{S}$ . Finally, the inefficiency parameter  $\beta \in [0, \infty)$  balances the incurred costs with the randomness of the policy followed by the team.

The objective in (2a) corresponds to minimizing the entropy-regularized total cost where the weight of the regularization is controlled by the inefficiency parameter  $\beta$ . The constraints in (2b) ensure that the resulting behavior of the team in the environment satisfies the expected final distribution  $\sigma_i$ . Note that as  $\beta \rightarrow 0$ , the team is expected to reach its final distribution only through optimal paths that minimize their total cost. On the other hand, as  $\beta \rightarrow \infty$ , the team is expected to be as random as possible while reaching their final distribution.

Let  $\mathcal{G} \cup \mathcal{S}_0 \cup \mathcal{S}_r$  be a partition of the set  $\mathcal{S}$  where  $\mathcal{G}$  is the set of goal states,  $\mathcal{S}_0$  is the set of states from which there is no path reaching the states in  $\mathcal{G}$ , and  $\mathcal{S}_r = \mathcal{S} \setminus \{\mathcal{G} \cup \mathcal{S}_0\}$ . These sets can be efficiently computed via simple graph search algorithms, e.g., breadth-first search. By slightly modifying the results presented in [28], it can be shown that the problem in (2a)–(2b) is equivalent to the following convex optimization problem:

$$\text{minimize}_{x(s,a) \geq 0} \sum_{s \in \mathcal{S}_r} \sum_{a \in \mathcal{A}(s)} x(s, a) \left[ c(s, a) + \beta \log \left( \frac{x(s, a)}{\nu(s)} \right) \right] \quad (3a)$$

subject to:

$$\nu(s) - \sum_{s' \in \mathcal{S}} \eta(s', s) = \alpha(s), \text{ for all } s \in \mathcal{S}_r, \quad (3b)$$

$$\sum_{s \in \mathcal{S}_r} \eta(s, g_i) = \sigma_i(g_i), \text{ for all } g_i \in \mathcal{G}, \quad (3c)$$

$$\eta(s, s') = \sum_{a \in \mathcal{A}(s)} x(s, a) P_{s,a,s'}, \text{ for all } s \in \mathcal{S}_r \text{ and } s' \in \mathcal{S}, \quad (3d)$$

$$\nu(s) = \sum_{a \in \mathcal{A}(s)} x(s, a), \text{ for all } s \in \mathcal{S}_r. \quad (3e)$$

In the above problem, the decision variables  $x(s, a)$  represent the density of the team members that occupy the state  $s$  and take the action  $a$ . These variables are traditionally referred to as occupancy measures [29]. The constraint in (3b) corresponds to balance equations, which express that the density entering a state should be equal to the density

leaving that state. Similarly, the constraint in (3c) ensures that the final distribution of the team satisfies the condition in (2b). Finally, the constraints in (3d)–(3e) are introduced just to simplify the notation.

The objective in (3a) is a convex function of  $x(s, a)$  which combines linear terms  $x(s, a)c(s, a)$  with the relative entropy of the distribution  $x(s, a)$  with  $\nu(s)$  [30]. Since the constraints are also linear functions of  $x(s, a)$ , the resulting convex optimization problem can be solved efficiently via off-the-shelf solvers. However, to ensure the existence of optimal solutions, we need to choose the cost function  $c(s, a)$  in a particular way, as described in the following proposition.

**Proposition 1:** If the problem in (1a)–(1b) is feasible and  $c(s, a) \geq \beta \log(|\mathcal{A}(s)|)$  for all  $s \in \mathcal{S}_r$  and  $a \in \mathcal{A}(s)$ , then the problem in (3a)–(3e) has a finite optimal solution.

**Proof:** We first note that if the problem in (1a)–(1b) is feasible, then the problem in (3a)–(3e) also has a feasible solution as shown in Lemma 1 in [27]. Additionally, this feasible solution has a finite value due to the conventions that  $0 \log 0 = 0$  and  $0 \log(0/0) = 0$  which are based on continuity arguments.

We now show that the optimal value is lower bounded by a finite constant. Suppose that  $c(s, a) \geq \beta \log(|\mathcal{A}(s)|)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ . For notational convenience, we drop the dependence of the set  $\mathcal{A}(s)$  on  $s$  in the following derivations. We express the objective function in (3a) as  $\sum_{s \in \mathcal{S}_r} \theta(s)$  where

$$\theta(s) := \sum_{a \in \mathcal{A}} x(s, a)c(s, a) + \beta \sum_{a \in \mathcal{A}} x(s, a) \log \left( \frac{x(s, a)}{\nu(s)} \right).$$

If  $\nu(s) = 0$ , then we have  $\theta(s) = 0$  by the convention  $0 \log(0/0) = 0$ . Additionally, for each  $s \in \mathcal{S}_r$  that satisfies  $\nu(s) > 0$ , we have

$$\begin{aligned} \theta(s) &\geq \beta \left[ \sum_{a \in \mathcal{A}} x(s, a) \log(|\mathcal{A}|) + \sum_{a \in \mathcal{A}} x(s, a) \log \left( \frac{x(s, a)}{\nu(s)} \right) \right] \\ &= \beta \left[ \nu(s) \log(|\mathcal{A}|) + \nu(s) \sum_{a \in \mathcal{A}} \frac{x(s, a)}{\nu(s)} \log \left( \frac{x(s, a)}{\nu(s)} \right) \right] \end{aligned} \quad (4b)$$

$$\geq \beta \nu(s) \left[ \log(|\mathcal{A}|) - \log(|\mathcal{A}|) \right] \geq 0. \quad (4c)$$

The inequality in (4a) follows from the fact that  $c(s, a) \geq \beta \log(|\mathcal{A}(s)|)$ . The equality in (4b) follows from the definition of  $\nu(s)$  in (3e) and the fact that  $\nu(s) > 0$ . Finally, the inequality in (4c) follows from the fact that the maximum entropy of a discrete probability distribution with a support size  $K \in \mathbb{N}$  is always less than or equal to  $\log(K)$ .

Finally, since the problem in (3a)–(3e) has a feasible solution and its optimal value is lower bounded by zero, we conclude that it has a finite optimal solution.  $\square$

We note that if  $\beta$  is large, i.e., the agent's behavior is highly inefficient, then the  $x(s, a)$  may be unbounded. In this case, the agents are expected to spend infinite time in the environment thereby making the goal-directedness ineffective. The condition  $c(s, a) \geq \beta \log(|\mathcal{A}(s)|)$  given

Proposition 1 ensures that this pathological case does not happen, and inefficiency and goal-directedness are balanced.

Using the condition given in Proposition 1, we can choose a cost function  $c$  with sufficiently high values so that the problem in (3a)–(3e) has a finite solution. Let  $\{x^*(s, a) : s \in \mathcal{S}, a \in \mathcal{A}(s)\}$  be a set of optimal decision variables for the problem in (3a)–(3e). We can obtain the policy  $\bar{\pi}$ , which describes the expected behavior for the team to achieve the final distribution  $\sigma_i$ , by the rule

$$\bar{\pi}_i(s, a) = \begin{cases} \frac{x^*(s, a)}{\sum_{a \in \mathcal{A}(s)} x^*(s, a)} & \text{if } \sum_{a \in \mathcal{A}(s)} x^*(s, a) > 0 \\ \frac{1}{|\mathcal{A}(s)|} & \text{otherwise.} \end{cases} \quad (5)$$

In the following section, we will show how to utilize the policies  $\bar{\pi}_i$  for quantifying deception and generating behaviors that manipulate the predictions of the adversary.

## V. GENERATING DECEPTIVE BEHAVIOR

In this section, we introduce several measures to quantify the deceptiveness of the team's behavior and present efficient algorithms to synthesize deceptive policies.

### A. Quantifying Deceptiveness through Statistical Distance

We propose to quantify deception through the statistical distance between the team's behavior and the behavior expected by the adversary. Specifically, we utilize the Kullback-Leibler (KL) divergence to formally define deception.

**Definition 3:** Let  $Q_1$  and  $Q_2$  be discrete probability distributions with a countable support  $\mathcal{X}$ . The Kullback-Leibler (KL) divergence between  $Q_1$  and  $Q_2$  is defined as

$$\text{KL}(Q_1 || Q_2) := \sum_{x \in \mathcal{X}} Q_1(x) \log \left( \frac{Q_1(x)}{Q_2(x)} \right).$$

The KL divergence  $\text{KL}(Q_1 || Q_2)$  measures the deviation of the distribution  $Q_1$  from the distribution  $Q_2$ . As we will discuss shortly, in the context of deception, it provides us with a method to quantify the statistical deviation of the team's observed behavior from the behavior expected by the adversary.

We consider two different types of deception, namely, exaggeration and ambiguity. Before providing formal definitions for these deceptive behaviors, we first introduce some notation. For an arbitrary policy  $\pi \in \Pi(\mathbb{M})$ , let  $\Gamma^\pi$  be the distribution of paths in  $\mathbb{M}$  generated under  $\pi$ . Note that the support of the distribution  $\Gamma^\pi$  is the set  $\text{Paths}_{\mathbb{M}}^{\pi, \alpha}$  of all paths, which may, in general, contain infinitely many elements. As we will shortly observe, for the purposes of deception, it is not necessary to explicitly construct this distribution.

**Exaggeration:** In this first type of deception, the team aims to exaggerate its behavior to convince the adversary that they allocate their resources with respect to a *decoy* goal distribution  $\sigma_i \in \mathcal{U} \setminus \{\sigma_*\}$ . Without loss of generality, let  $\sigma_1$  be the true goal distribution, i.e.,  $\sigma_* = \sigma_1$ . Then, for

a given policy  $\pi \in \Pi(\mathbb{M})$ , we quantify the exaggeration of the team's resulting behavior through the following formula

$$\text{Deceptiveness}(\mathcal{U}, \pi) = \max_{i \in [N]} \left[ \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_1}) - \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i}) \right]. \quad (6)$$

The term  $\text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i})$  quantifies the KL-divergence between the path distributions induced by the policies  $\pi$  and  $\bar{\pi}_i$ . Therefore, the above deceptiveness metric measures the relative statistical distance of the paths induced by  $\pi$  to the true policy  $\bar{\pi}_1$  and decoy policy  $\bar{\pi}_i$ .

The intuition behind (6) comes from the likelihood-ratio test, which is the most powerful hypothesis testing method for a given significance level [31]. Recall that, for each  $i \in [N]$ , the adversary expects the team to follow a policy  $\bar{\pi}_i$  to achieve the final distribution  $\sigma_i$ . Now, suppose that the adversary runs the likelihood-ratio test to decide whether the team follows the policy  $\bar{\pi}_i$  or  $\bar{\pi}_j$ . Let  $\varrho_1, \dots, \varrho_n$  be the paths followed by  $n$  members of the team under the policy  $\pi$ . Moreover, let  $\Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_i)$  and  $\Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_j)$  be the probabilities of  $\varrho_1, \dots, \varrho_n$  under  $\bar{\pi}_i$  and  $\bar{\pi}_j$ , respectively. By the likelihood-ratio test, for a given constant  $C \in \mathbb{R}_{\geq 0}$ , the adversary decides that the team aims to achieve the final distribution  $\sigma_{U^i}^*$  through the policy  $\bar{\pi}_i$  if

$$\log \left( \Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_i) \right) - \log \left( \Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_j) \right) \geq C.$$

To see how (6) is related to the likelihood-ratio test, note that,  $n \left[ \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_1}) - \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i}) \right]$  is equal to

$$\mathbb{E}^\pi \left[ \log \left( \Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_i) \right) \right] - \mathbb{E}^\pi \left[ \log \left( \Pr(\varrho_1, \dots, \varrho_n | \bar{\pi}_1) \right) \right].$$

Therefore, the term inside the parenthesis in (6) quantifies the expected log-likelihood of a goal distribution  $\sigma_i$  being the true goal distribution to the goal distribution  $\sigma_1$  being the true goal distribution when the team follows the policy  $\pi$ . Note that by taking the maximum over  $i \in [N]$  in  $\text{Deceptiveness}(\mathcal{U}, \pi)$ , we quantify deceptiveness with respect to the most likely decoy goal distribution. Consequently, the problem in (1a)–(1b) corresponds to synthesizing a policy  $\pi^*$  that maximizes the expected relative log-likelihood for a decoy goal distribution  $\sigma_i \in \mathcal{U} \setminus \{\sigma_*\}$  to be the true goal distribution while guaranteeing that the team's resulting behavior satisfies the final resource distribution  $\sigma_*$ .

**Ambiguity:** In this second type of deception, the team aims to behave in a way to make its true goal distribution  $\sigma_*$  ambiguous to the adversary. Specifically, for a given policy  $\pi \in \Pi(\mathbb{M})$ , we quantify the ambiguity of the team's behavior through the following formula

$$\text{Deceptiveness}(\mathcal{U}, \pi) = - \max_{i \in [N]} \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i}). \quad (7)$$

Similar to the exaggeration behavior, the intuition behind the equation in (7) comes from the likelihood-ratio test. Specifically, in (7), we measure the deceptiveness of a policy as the minimum expected log-likelihood of any utility matrix  $\sigma_i \in \mathcal{U}$ . As a result, the problem in (1a)–(1b) corresponds to synthesizing a policy  $\pi^*$  that *minimizes* the maximum

log-likelihood for any goal distribution to be the true goal distribution while guaranteeing that the team's resulting behavior satisfies the final resource distribution  $\sigma_* \in \mathcal{U}$ .

## B. Synthesis of Policies through Convex Optimization

In this section, we present algorithms to solve the problem in (1a)–(1b) when  $\text{Deceptiveness}(\mathcal{U}, \pi)$  is defined as in (6) and in (7).

Although the problem in (1a)–(1b) is feasible, it is, in general, possible that the optimal value is not bounded below when  $\text{Deceptiveness}(\mathcal{U}, \pi)$  is defined as in (6) and in (7). This is due to the fact that, for given  $U^i$  and  $U^j$ , the support of the final distributions  $\sigma_{1,U^i}^*$  and  $\sigma_{1,U^j}^*$  may be different. As a result, the KL divergence between the path distributions  $\Gamma^{\bar{\pi}_1}$  and  $\Gamma^{\bar{\pi}_i}$  may be infinite.

To ensure the finiteness of the optimal value in (1a)–(1b), we propose to divide the team's behavior into two phases, namely, *deceptive* and *goal-directed phases*. During the deceptive phase, the team aims to deceive the adversary regarding its goal distribution by optimizing its behavior with respect to the measures in (6) or in (7). Let  $T \in \mathbb{N}$  be a critical decision stage at which the team switches from the deceptive phase to the goal-directed phase. After  $T$ , the team aims to reach its final distribution  $\sigma_*$  through the shortest path.

We utilize extended MDPs to compactly represent the deceptive and goal-directed phases in a single decision model. Formally, let  $\bar{\mathbb{M}}_T = (\bar{\mathcal{S}}, \bar{\alpha}, \mathcal{A}, \bar{P})$  denote an *extended MDP* where  $\bar{\mathcal{S}} = \mathcal{S} \times [T + 1]$  is a finite set of states,  $\bar{\alpha} : \bar{\mathcal{S}} \rightarrow [0, 1]$  is an initial distribution such that, for each  $\langle s, t \rangle \in \bar{\mathcal{S}}$ ,  $\bar{\alpha}(\langle s, t \rangle) = \alpha(s)$  if  $t = 1$  and  $\bar{\alpha}(\langle s, t \rangle) = 0$  otherwise, and  $\bar{P} : \bar{\mathcal{S}} \times \mathcal{A} \times \bar{\mathcal{S}} \rightarrow [0, 1]$  is a transition function such that

$$\bar{P}_{\langle s, t \rangle, a, \langle s', t' \rangle} = \begin{cases} P_{s, a, s'} & \text{if } t \leq T \text{ and } t' = t + 1 \\ P_{s, a, s'} & \text{if } t = T + 1 \text{ and } t' = t \\ 0 & \text{otherwise.} \end{cases}$$

In an extended MDP, we can clearly distinguish the deceptive and goal-directed phases by defining the objectives separately for the states  $\mathcal{S} \times [T]$  and  $\mathcal{S} \times \{T + 1\}$  as will be discussed shortly.

In the above construction,  $T$  is a design variable that can be used to tune the duration of the team's deceptive behavior. One practical approach is to set  $T$  as a function of the shortest path. Specifically, let  $T_{\min}$  be the minimum expected time for the team to reach their final distribution  $\sigma_{1,U^i}^*$ .  $T_{\min}$  can be computed by replacing the objective function in (3a) with  $\sum_{s \in \mathcal{S}_r} \sum_{a \in \mathcal{A}(s)} x(s, a)$ . Then, we can simply set  $T = k \lceil T_{\min} \rceil$  where  $k \in \mathbb{N}$  determines the balance between suboptimality of the behavior and deception effort.

**Exaggeration:** We achieve the exaggeration behavior for deception by solving  $N$  separate linear programs (LPs). Let  $\mathbf{x}^\pi$  be the vector of occupancy measures that correspond to the policy  $\pi$  constructed through the formula in (5). By simple algebraic manipulations, it can be shown [11] that, for each  $i \in [N]$ , we have

$$\text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_1}) - \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i}) = \sum_{s \in \mathcal{S}_r} \sum_{a \in \mathcal{A}} x^\pi(s, a) \log \left( \frac{\bar{\pi}_i(s, a)}{\bar{\pi}_1(s, a)} \right).$$

Note in the above equation that the logarithmic term is a constant and corresponds to a virtual reward that quantifies the statistical likelihood of the decoy goal distribution  $\sigma_i$  with respect to the true goal distribution  $\sigma_1$ . The virtual reward may be infinite when there is a support mismatch between the policies  $\bar{\pi}_1$  and  $\bar{\pi}_i$ . To ensure the finiteness of the virtual reward and avoid computational issues, we add a small constant  $\epsilon$  to both the numerator and the denominator in the logarithmic term. Accordingly, for each  $i \in [N]$ , we consider the following LP:

$$\begin{aligned} \text{maximize}_{x(\langle s, t \rangle, a) \geq 0} & \sum_{\langle s, t \rangle \in \mathcal{S} \times [T]} \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a) \log \left( \frac{\bar{\pi}_i(s, a) + \epsilon}{\bar{\pi}_1(s, a) + \epsilon} \right) \\ & - \sum_{\langle s, t \rangle \in \mathcal{S} \times \{T+1\}} \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a) \end{aligned} \quad (8a)$$

subject to:

$$\nu(\langle s, t \rangle) - \sum_{\langle s', t' \rangle \in \mathcal{S} \times [T]} \eta(\langle s', t' \rangle, \langle s, t \rangle) = \bar{\alpha}(\langle s, t \rangle), \quad \text{for all } \langle s, t \rangle \in \mathcal{S}_r \times [T+1] \quad (8b)$$

$$\sum_{\langle s, t \rangle \in \mathcal{S}_r \times [T]} \sum_{t' \in [T]} \eta(\langle s, t \rangle, \langle g_i, t' \rangle) = \sigma_{1, U^i}^*(g_i), \quad \text{for all } g_i \in \mathcal{G} \quad (8c)$$

$$\eta(\langle s, t \rangle, \langle s', t' \rangle) = \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a) \bar{P}_{\langle s, t \rangle, a, \langle s', t' \rangle}, \quad \text{for all } \langle s, t \rangle \in \mathcal{S}_r \times [T] \text{ and } \langle s', t' \rangle \in \mathcal{S} \times [T] \quad (8d)$$

$$\nu(\langle s, t \rangle) = \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a), \quad \text{for all } \langle s, t \rangle \in \mathcal{S}_r \times [T] \quad (8e)$$

The objective function in the above LP consists of two terms that enable the team to perform a deception phase followed by a goal-directed phase. Specifically, the first sum in the objective corresponds to  $\text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_1}) - \text{KL}(\Gamma^\pi || \Gamma^{\bar{\pi}_i})$  on the extended state-space  $\mathcal{S} \times [T]$ . On the other hand, the second term ensures that after the deception phase, i.e.,  $T+1$ , the team reaches its final distribution by minimizing their total residence time in the environment.

The constraints in (8b)–(8e) are the same as the constraints in (3b)–(3e) with a minor difference. Specifically, the constraints in (8b)–(8e) are now defined over the extended MDP  $\bar{\mathbb{M}}_T$  instead of the original MDP  $\mathbb{M}$ .

Now, for each  $i \in [N]$ , let  $v_i^*$  be the optimal value of the LP given in (8a)–(8e) and  $i^* \in \arg \max_{i \in [N]} v_i^*$ . Moreover, let  $\{x^*(\langle s, t \rangle, a) : \langle s, t \rangle \in \mathcal{S} \times [T+1], a \in \mathcal{A}\}$  be the set of optimal decision variables corresponding to the LP with the optimal value  $v_{i^*}^*$ . We obtain an optimal deceptive policy  $\pi^* \in \Pi(\bar{\mathbb{M}}_T)$  through the construction

$$\pi^*(\langle s, t \rangle, a) = \begin{cases} \frac{x^*(\langle s, t \rangle, a)}{\sum_{a \in \mathcal{A}} x^*(\langle s, t \rangle, a)} & \text{if } \sum_{a \in \mathcal{A}(\langle s, t \rangle)} x^*(\langle s, t \rangle, a) > 0 \\ \frac{1}{|\mathcal{A}(\langle s, t \rangle)|} & \text{otherwise.} \end{cases} \quad (9)$$

It follows from the standard results in the MDP theory, e.g., [7, Chapter 7], the policy  $\pi^*$  ensures that the team reaches its desired final distribution  $\sigma_1$ .

**Ambiguity:** We achieve ambiguous behavior for deception by solving a single convex optimization problem. Recall from (7) that the objective in this type of deception is to obtain a policy  $\pi$  that has the minimum statistical distance to each potential policy  $\bar{\pi}_i$ . Accordingly, using the derivations from the exaggeration behavior, we consider the following convex program:

$$\begin{aligned} \text{minimize}_{z, x(\langle s, t \rangle, a) \geq 0} & \sum_{\langle s, t \rangle \in \mathcal{S} \times \{T+1\}} \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a) + z \quad (10a) \\ \text{subject to:} & \quad (8b) \text{--} (8e), \end{aligned}$$

$$z \geq \sum_{\langle s, t \rangle \in \mathcal{S} \times [T]} \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a) \log \left( \frac{x(\langle s, t \rangle, a) / \sum_a x(\langle s, t \rangle, a)}{\bar{\pi}_i(s, a) + \epsilon} \right) \quad \text{for all } i \in [N] \quad (10b)$$

The objective in the above convex program consists of two terms similar to the exaggeration behavior. The first term corresponds to the goal-directed phase in which the team aims to reach its final distribution by minimizing their occupancy time in the environment. The second term, expressed by the scalar variable  $z$ , corresponds to the ambiguity of the behavior during the deception phase. In particular, through the constraint in (10b), this variable quantifies the maximum KL-divergence of the path distribution induced by  $\pi$  to the set of path distributions induced by  $\bar{\pi}_i$  where  $i \in [N]$ .

Finally, an optimal deceptive policy for achieving ambiguous behavior can be obtained from the optimal decision variables for the program in (10a)–(10b) through the construction introduced in (9).

## VI. SIMULATIONS

In this section, we present a numerical simulation to illustrate the proposed deception strategy in a motion planning example. We utilize the CVXPY interface [32] and the ECOS solver [33] to obtain solutions to the considered convex optimization problems.

We consider a scenario in which the team navigates in an environment represented by the  $10 \times 10$  grid-world shown in Fig. 2. Each grid cell represents a state, and there are four available actions  $\{up, down, right, left\}$  in each state under which the agents transition to the neighboring state in the corresponding direction. All team members start their motion from the state labeled with  $S$ . The team aims to allocate their resources over three goal states labeled with  $g_1$ ,  $g_2$ , and  $g_3$ . We consider two potential goal allocations  $\sigma_1 = [0, 0.5, 0.5]$  and  $\sigma_2 = [0.5, 0.5, 0]$ . Note that since we have  $\sigma_1 = [0, 0.5, 0.5]$ , the team aims to allocate half of the team members to  $g_2$  and the remaining half to  $g_3$ .

We synthesize a deceptive density control strategy for Team 1 using the proposed methods. For the synthesis of deceptive control strategies, we solve the optimization problem (8a)–(8e) for exaggeration behavior and the optimization problem (10a)–(10b) for ambiguity behavior.

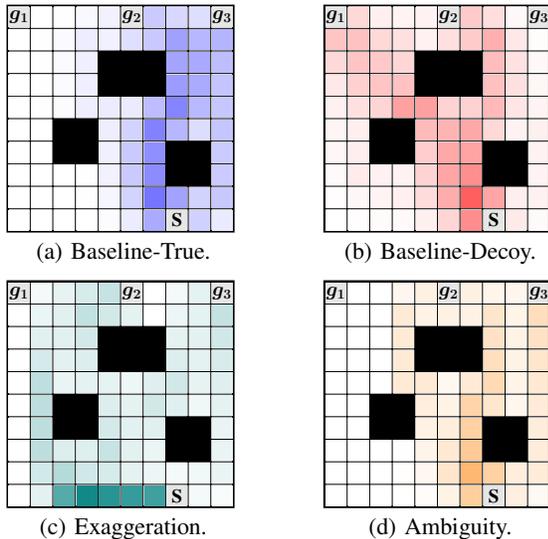


Fig. 2: Normalized density distributions of the autonomous team for  $\beta = 6$  and  $c(s, a) = 10$ . The true final distribution is  $[0.5, 0.5, 0]$  and the decoy final distribution is  $[0, 0.5, 0.5]$  (a) The baseline density distribution (blue) expected by the adversary for the true final distribution. (b) The baseline density distribution (red) expected by the adversary for the decoy final distribution. (c) The density distribution (teal) for the exaggeration behavior. (d) The density distribution (orange) for the ambiguity behavior.

In Fig. 2, we illustrate the density distributions that the adversary predicts, as well as the distribution that the team follows. The density of state  $s$  is equal to  $\sum_{t \in [T+1]} \sum_{a \in \mathcal{A}} x(\langle s, t \rangle, a)$  that is the expected time that the team members spend at state  $s$ .

For the synthesis of deceptive strategies, we set the parameter  $T = 5$  in the optimization problems, i.e., the team shows the deceptive behavior for 5 steps and then switches to the goal-directed behavior. In Fig. 2c, we observe that Team 1 exaggerates its behavior (shown in green color) and pretends to achieve the decoy distribution  $[0.5, 0.5, 0]$  during the initial deception phase. Then, during the goal-directed phase, they eventually reach their final distribution. For the ambiguity behavior shown in Fig. 2d, during the initial deception phase, the team follows a path (shown in orange color) that is the only significantly plausible path for both the true and decoy final distributions. This behavior preserves the ambiguity of the true distribution until the goal-directed phase.

## VII. USER STUDY

We now assess the effectiveness of the proposed deceptive strategies in altering real users' perceptions of final goal distribution by showcasing five paths sampled from each policy (see Fig. 2) at steps 4, 8, 12, 16, 20.

### A. Experiment Design

**Manipulated Variables:** We conducted two sets of experiments. In the first set, participants were informed of the deceptive purpose (YES), while in the second set, the true purpose was not disclosed until after the study (NO).

Within each set, we varied the policy types: a baseline policy for the true goal distribution (*Baseline-True*), a baseline policy for the decoy goal distribution (*Baseline-Decoy*), a deceptive policy that creates ambiguity regarding the true goal distribution (*Ambiguity*), or a deceptive policy that exaggerates towards a decoy goal distribution (*Exaggeration*). The baseline policies are generated using the principle of maximum entropy method discussed in Section IV, and the deceptive policies are generated using the methods given in V. We use  $\beta = 6$  for the maximum entropy method. In total, the experiment has 8 different conditions.

**Dependent Measures:** We have two dependent measures: correctness, measured as the ratio of correctly perceived goals, and confidence on a 5-point Likert scale. We combine the two in a **score**: the confidence if they get two goals correct, half the confidence if they get only one goal correct, and negative of the confidence if they are not correct at all [34]. This score captures that if one is incorrect, it is better not to be confident about it. The lower the score is, the more effective a deceptive policy is.

**Participants:** We used a between-subjects design, where each participant would only see paths from one condition, in order to avoid carryover and fatigue effects from having seen a different condition before. We ran this experiment on a total of 320 participants across the 8 conditions, recruited via Prolific [35]. We filtered out data from 3 (0.938%) participants who withdrew consent. The average age of the 317 consenting participants was 39.595 (SD = 13.939). The gender ratio was 0.497 female.

### B. Results

Across every condition studied, we observed increases in both correctness and confidence as the path advanced. This outcome is as expected, as longer path segments generally reduce interpretive uncertainty for participants. By the final step, the path endpoints in all conditions either reach or nearly reach the predetermined goals, leading to high average and low variance in scores.<sup>1</sup>

**H1: The proposed deceptive policies are effective:** Fig. 3 plots the mean scores for the Baseline-True, Ambiguity, and Exaggeration policies at each step of the path within both the NO and YES groups. Whether in the YES or NO group, participants consistently register lower mean scores for the Ambiguity and Exaggeration policies compared to the Baseline-True policy, a pattern particularly pronounced before Step 16. Upon closer examination of the two deceptive policies, we notice that Exaggeration yields even lower mean scores than Ambiguity. This not only supports H1 but further reveals exaggeration is more effective than ambiguity.

**H2: Informing participants of deception makes a policy less effective:** We made this hypothesis based on the presumption that users are more likely to second-guess themselves if they know they are being deceived, leading to reduced scores. Nevertheless, the outcomes presented in

<sup>1</sup>All data in the results is publicly available at [https://github.com/vivianchen98/deception\\_user\\_study\\_data](https://github.com/vivianchen98/deception_user_study_data).

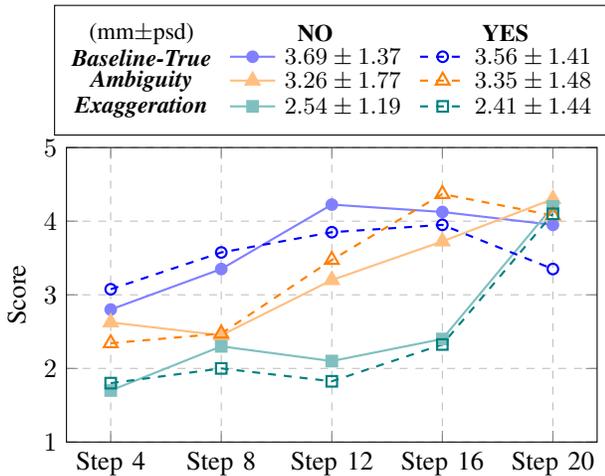


Fig. 3: Participant scores for paths sampled from the *Baseline-True*, *Ambiguity*, and *Exaggeration* policies, in both NO and YES groups (all with  $\beta = 6$ ). The legend delineates statistical details for each line, including the mean of means (mm) and the pooled standard deviation (psd) across path steps.

Fig. 3 from both the YES and NO groups for each policy type—depicted by solid and dashed lines of the same color—do not manifest a consistent trend, rendering the findings non-definitive with respect to this hypothesis.

**H3: Baseline-True and Baseline-Decoy have comparable scores:** Fig. 4 shows that the *Baseline-True* policy, with goals closer to the starting position on the right side of the interface, consistently attains higher scores than *Baseline-Decoy* (see solid lines). Two conceivable reasons may account for this trend: participants might either possess an inherent preference for goals located nearer to the starting point, or they could have a bias towards the right side of the interface, perhaps due to the ease of clicking. To discern the actual bias, we mirrored the experiment using horizontally-flipped paths, using data from 315 consenting participants. We limited the display to four steps, as the final step would not yield substantial insights. The results for the flipped experiment, represented by dashed lines in Fig. 4, show a similar trend, suggesting the higher scores for *Baseline-True* are not a consequence of the goals’ alignment on the interface. Rather, it reveals that participants have a consistent inclination towards proximate goals. This inclination can also be explained by the notion of last deceptive point (LDP) proposed [14], as the *Baseline-Decoy* paths inadvertently approach the LDP whereas the *Baseline-True* paths never come close to the LDP.

## VIII. CONCLUSIONS

We studied the problem of synthesizing deceptive resource allocation strategies for a team consisting of a large number of autonomous agents. We developed a prediction algorithm based on the principle of maximum entropy that models the predictions of adversarial observers regarding the autonomous team’s final allocation strategy over multiple goal

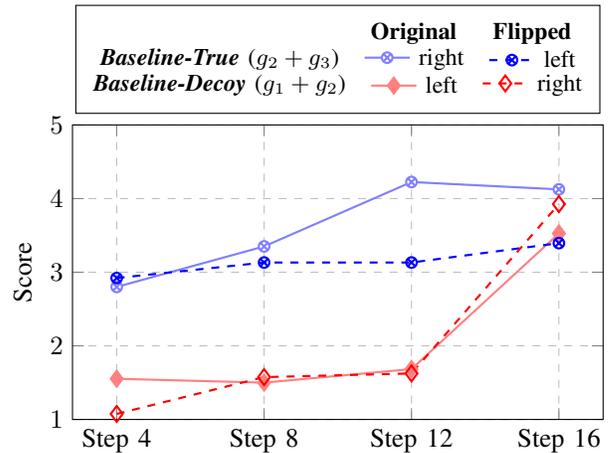


Fig. 4: Participant scores for paths sampled from the *Baseline-True* and *Baseline-Decoy* policies across path steps in the original and the flipped experiments (both with  $\beta = 6$  in NO group).

locations. By quantifying deceptiveness as a function of statistical distance between certain distributions, we then developed deceptive strategies, based on convex optimization, to control the density of the team members in the environment while they progress towards their final distribution. A user study validates the effectiveness of the proposed algorithms.

## REFERENCES

- [1] M. Saska, V. Vonásek, J. Chudoba, J. Thomas, G. Loianno, and V. Kumar, “Swarm distribution and deployment for cooperative surveillance by micro-aerial vehicles,” *Journal of Intelligent & Robotic Systems*, 2016.
- [2] P. Lin, G. Bekey, and K. Abney, “Autonomous military robotics: Risk, ethics, and design,” California Polytechnic State Univ San Luis Obispo, Tech. Rep., 2008.
- [3] E. T. Jaynes, “Information theory and statistical mechanics,” *Physical Review*, 1957.
- [4] B. D. Ziebart, A. L. Maas, J. A. Bagnell, A. K. Dey *et al.*, “Maximum entropy inverse reinforcement learning.” in *AAAI Conference on Artificial intelligence*, 2008.
- [5] G. Gergely, Z. Nádasdy, G. Csibra, and S. Bíró, “Taking the intentional stance at 12 months of age,” *Cognition*, 1995.
- [6] E. A. Feinberg and A. Shwartz, *Handbook of Markov Decision Processes: Methods and Applications*. Springer Science & Business Media, 2012.
- [7] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.
- [8] D. A. Dolgov and E. H. Durfee, “Resource allocation among agents with mdp-induced preferences,” *Journal of Artificial Intelligence Research*, 2006.
- [9] M. Hibbard, Y. Savas, Z. Xu, and U. Topcu, “Minimizing the information leakage regarding high-level task specifications,” *IFAC-PapersOnLine*, 2020.
- [10] S. Witwicki, J. C. Castillo, J. Messias, J. Capitan, F. S. Melo, P. U. Lima, and M. Veloso, “Autonomous surveillance robots: A decision-making framework for networked multiagent systems,” *IEEE Robotics & Automation Magazine*, 2017.
- [11] M. O. Karabag, M. Ornik, and U. Topcu, “Deception in supervisory control,” *IEEE Transactions on Automatic Control*, 2021.
- [12] Y. Savas, C. K. Verginis, and U. Topcu, “Deceptive decision-making under uncertainty,” in *AAAI Conference on Artificial intelligence*, 2022.
- [13] A. Dragan, R. Holladay, and S. Srinivasa, “Deceptive robot motion: synthesis, analysis and experiments,” *Autonomous Robots*, 2015.
- [14] P. Masters and S. Sardina, “Deceptive path-planning,” in *International Joint Conferences on Artificial Intelligence*, 2017.
- [15] M. Ornik and U. Topcu, “Deception in optimal control,” in *Allerton Conference on Communication, Control, and Computing*, 2018.

- [16] L. Li, H. Ma, A. N. Kulkarni, and J. Fu, "Dynamic hypergames for synthesis of deceptive strategies with temporal logic objectives," *IEEE Transactions on Automation Science and Engineering*, 2022.
- [17] A. N. Kulkarni, H. Luo, N. O. Leslie, C. A. Kamhoua, and J. Fu, "Deceptive labeling: hypergames on graphs for stealthy deception," *IEEE Control Systems Letters*, 2020.
- [18] A. N. Kulkarni, J. Fu, H. Luo, C. A. Kamhoua, and N. O. Leslie, "Decoy allocation games on graphs with temporal logic objectives," in *International Conference on Decision and Game Theory for Security*, 2020.
- [19] T. H. Nguyen, Y. Wang, A. Sinha, and M. P. Wellman, "Deception in finitely repeated security games," in *AAAI Conference on Artificial Intelligence*, 2019.
- [20] A. R. Wagner and R. C. Arkin, "Acting deceptively: Providing robots with the capacity for deception," *International Journal of Social Robotics*, 2011.
- [21] M. Ramírez and H. Geffner, "Probabilistic plan recognition using off-the-shelf classical planners," in *AAAI Conference on Artificial Intelligence*, 2010.
- [22] M. Ramírez and H. Geffner, "Goal recognition over pomdps: Inferring the intention of a pomdp agent," in *International Joint Conference on Artificial Intelligence*, 2011.
- [23] M. Shvo and S. A. McIlraith, "Active goal recognition," in *AAAI Conference on Artificial Intelligence*, 2020.
- [24] C. Baier and J.-P. Katoen, *Principles of Model Checking*. MIT Press, 2008.
- [25] B. Roberson, "The colonel blotto game," *Economic Theory*, 2006.
- [26] T. Raghavan, "Zero-sum two-person games," *Handbook of Game Theory with Economic Applications*, 1994.
- [27] 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.
- [28] Y. Savas, M. Ornik, M. Cubuktepe, M. O. Karabag, and U. Topcu, "Entropy maximization for markov decision processes under temporal logic constraints," *IEEE Transactions on Automatic Control*, 2019.
- [29] E. Altman, *Constrained Markov Decision Processes*. Chapman and Hall/CRC, 1999.
- [30] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [31] J. Neyman and E. S. Pearson, "Ix. on the problem of the most efficient tests of statistical hypotheses," *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, 1933.
- [32] S. Diamond and S. Boyd, "Cvxpy: A python-embedded modeling language for convex optimization," *The Journal of Machine Learning Research*, 2016.
- [33] A. Domahidi, E. Chu, and S. Boyd, "ECOS: An SOCP solver for embedded systems," in *European Control Conference*, 2013.
- [34] S. H. Huang, D. Held, P. Abbeel, and A. D. Dragan, "Enabling robots to communicate their objectives," *Autonomous Robots*, 2019.
- [35] S. Palan and C. Schitter, "Prolific. ac—a subject pool for online experiments," *Journal of Behavioral and Experimental Finance*, 2018.

