

# Game-theoretic Goal Recognition Models with Applications to Security Domains

Samuel Ang<sup>1</sup>, Hau Chan<sup>2</sup>, Albert Xin Jiang<sup>1</sup>, and William Yeoh<sup>3</sup>

<sup>1</sup> Department of Computer Science,  
Trinity University, San Antonio, TX 78212  
sang@trinity.edu, xjiang@trinity.edu

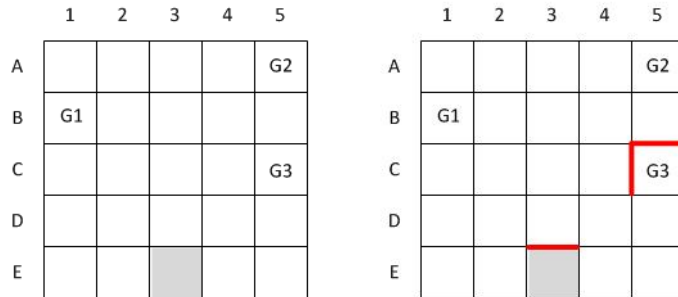
<sup>2</sup> Department of Computer Science and Engineering,  
University of Nebraska-Lincoln, NE, 68588  
hchan3@unl.edu

<sup>3</sup> Department of Computer Science and Engineering,  
Washington University in St. Louis, St. Louis, MO 63130  
wyeoh@wustl.edu

**Abstract.** Motivated by the goal recognition (GR) and goal recognition design (GRD) problems in the artificial intelligence (AI) planning domain, we introduce and study two natural variants of the GR and GRD problems with strategic agents, respectively. More specifically, we consider game-theoretic (GT) scenarios where a malicious adversary aims to damage some target in an (physical or virtual) environment monitored by a defender. The adversary must take a sequence of actions in order to attack the intended target. In the GTGR and GTGRD settings, the defender attempts to identify the adversary’s intended target while observing the adversary’s available actions so that he/she can strengthen the target’s defense against the attack. In addition, in the GTGRD setting, the defender can alter the environment (e.g., adding roadblocks) in order to better distinguish the goal/target of the adversary.

We propose to model GTGR and GTGRD settings as zero-sum stochastic games with incomplete information about the adversary’s intended target. The games are played on graphs where vertices represents states and edges are adversary’s actions. For the GTGR setting, we show that if the defender is restricted to playing only stationary strategies, the problem of computing optimal strategies (for both defender and adversary) can be formulated and represented compactly as a linear program. For the GTGRD setting, where the defender can choose  $K$  edges to block at the start of the game, we formulate the problem of computing optimal strategies as a mixed integer program, and present a heuristic algorithm based on LP duality and greedy methods. Experiments show that our heuristic algorithm achieves good performance (i.e., close to defender’s optimal value) with better scalability compared to the mixed-integer programming approach.

In contrast with our research, existing work, especially on GRD problems, has focused almost exclusively on decision-theoretic paradigms, where the adversary chooses its actions without taking into account the fact that they may be observed by the defender. As such an assumption



**Fig. 1.** Example Problem (left) and with Blocked Actions in Red (right).

is unrealistic in GT scenarios, our proposed models and algorithms fill a significant gap in the literature.

## 1 Introduction

Discovering the objective of an agent based on observations of its behavior is a problem that has interested both artificial intelligence (AI) and psychology researchers for many years [23, 7]. In AI, this problem is known as *goal recognition* (GR) or, more generally, *plan recognition* [25]. Plan and goal recognition problems have been used to model a number of applications ranging from software personal assistants [16–18]; robots that interact with humans in social settings such as homes, offices, and hospitals [26, 8]; intelligent tutoring systems that recognize sources of confusion or misunderstanding in students through their interactions with the system [14, 6, 12, 15]; and security applications that recognize the plan or goal of terrorists [5].

One can broadly summarize the existing research in GR as one that primarily focuses on developing better and more efficient techniques to recognize the plan or the goal of the user given a sequence of observations of the user’s actions. For example, imagine a scenario shown in Figure 1 (left), where an agent is at cell  $E3$ , it can move in any of the four cardinal directions, and its goal is one of three possible goals  $G1$  (in cell  $B1$ ),  $G2$  (in cell  $A5$ ), and  $G3$  (in cell  $C5$ ). Additionally, assume that it will move along a shortest path to its goal. Then, if it moves left to cell  $E2$ , then we can deduce that its goal is  $G1$ . Similarly, if it moves right to cell  $E4$ , then its goal is either  $G2$  or  $G3$ .

Existing research has focused on agent GR models that are non-strategic or partially strategic: The agent’s objective is to reach its goal with minimum cost, and the agent does not explicitly reason about its interaction with the observer. However, when the observer’s recognition of the agent’s goal affects the agent in some way, then it is in the agent’s best interest to be *fully strategic* – to

*explicitly* reason about how the agent’s choice affects the observer’s recognition. As a result, the observer will need to take into account the agent’s strategic reasoning when making decisions.

### 1.1 Game-Theoretic Goal Recognition Problems in Security Domains

Naturally, GT settings with strategic agents are common in many real-world (physical and cyber) security scenarios between an adversary and a defender. The adversary has a set of targets of interests and would be equally happy in attacking one of them. In physical security domains, the adversary must make a sequence of physical movements to reach a target; in cyber security domains, this could be a sequence of actions achieving necessary subgoals to carry out the attack. In any case, the defender is trying to recognize the adversary’s goal/target. We coined this the game-theoretic goal recognition (GTGR) problem.

Let us describe the security games of interests using Figure 1. Consider the security scenario in Figure 1 (left), where an agent (i.e., terrorist) wants to reach its intended target and carry out an attack, while we, the observer (the defender) try to recognize the agent’s goal as early as possible. Suppose once we recognize the agent’s goal, we will strengthen the agent’s target to defend against the attack. The more time we have between recognition and the actual attack, the less successful the attack will be. In this scenario, it is no longer optimal for the agent to simply choose a shortest path to its goal, as that could allow the observer to quickly identify its goal. On the other hand, the agent still wants to reach its goal in a reasonably short time, as a very long path could allow the observer time to strengthen all the targets. So, an optimal agent would need to explicitly reason about the tradeoffs between the cost of its path (e.g., path length) and the cost of being discovered early.

### 1.2 Game-Theoretic Goal Recognition Design Problems in Security Domains

So far we have been discussing the defender’s task on recognizing goals. However, the task could become very difficult in general. For instance, going back to our security example in Figure 1, if the agent moves up to  $D3$ , the observer cannot make any informed deductions. In fact, if the agent moves along any one of its shortest paths to goal  $G3$ , throughout its entire path, which is of length 4, we cannot deduce whether its goal is either  $G2$  or  $G3$ ! This illustrates one of the challenges with this approach, that is, there are often a large number of ambiguous observations that can be a result of a large number of goals. As such, it is difficult to uniquely determine the goal of the agent until a long sequence of observations is observed.

The work of [9, 10] proposed an orthogonal approach to *modify the underlying environment of the agent*, in such a way that *the agent is forced to reveal its goal as early as possible*. They call this problem the goal recognition design (GRD) problem. For example, if we block the actions ( $E3, up$ ), ( $C4, right$ ), ( $C5, up$ ) in

our example problem, where we use tuples  $(s, a)$  to denote that action  $a$  is blocked from cell  $s$ , then the agent can make at most 2 actions (i.e., right to E4 then up to D4) before its goal is conclusively revealed. Figure 1 (right) shows the blocked actions. This problem finds itself relevant in many of the same applications of GR because, typically, the underlying environment can be easily modified.

As such, in addition to studying the GTGR problem, we consider the GT-GRD problem where the observer can modify the underlying environment (i.e., adding  $K$  roadblocks) as to restrict the actions of the agent.

### 1.3 Related Work

GR and its more general forms, plan recognition and intent recognition, have been extensively studied [25] since their inception almost 40 years ago [23]. Researchers have made significant progress within the last decade through synergistic integrations of techniques ranging from natural language processing [27, 3] to classical planning [20–22] and deep learning [15]. The closest body of work to ours is the one that uses game-theoretic formulations, including an adversarial plan recognition model that is defined as an imperfect information two-player zero-sum game in extensive form [13], a model where the game is over attack graphs [1], and an extension that allows for stochastic action outcomes [4]. The main difference between these works and ours is that ours focuses on goal recognition instead of plan recognition.

While GR has a long history and extensive literature, the field of GRD is relatively new. Keren *et al.* introduced the problem in their seminal paper [9], where they proposed a decision-theoretic STRIPS-based formulation of the problem. In the original GRD problem, the authors make several simplifying assumptions: (1) the observed agent is assumed to execute an optimal (i.e., cost-minimal) plan to its goal; (2) the actions of the agent are deterministic; and (3) the actions of the agent are fully observable. Since then, these assumptions have been independently relaxed, where agents can now execute boundedly-suboptimal plans [10], actions of the agents can be stochastic [28], and actions of the agents can be only partially observable [11]. Further, aside from all the decision-theoretic approaches above, researchers have also modeled and solved the original GRD problem using answer set programming [24]. The key difference between these works and ours is that ours introduced a game-theoretic formulation that can more accurately capture interactions between the observed agent and the observer in security applications.

### 1.4 Our Contributions

As a result of the strategic interaction in the GTGR and GTGRD scenarios, the concept of cost-minimal plan (the solution concept in GR problem) and worst-case distinctiveness (the solution concept in GRD problem) are no longer a suitable solution concept since it does not reflect the behavior of strategic agents. Instead, our objective here is to formulate game-theoretic models of the

agent’s and observer’s interactions under GR and GRD settings. More specifically, we propose to model GTGR and GRGRD settings as zero-sum stochastic games with incomplete information where the adversary’s target is unknown to the observer. For the GTGR setting, we show that if the defender is restricted to playing only stationary strategies, the problem of computing optimal strategies (for both defender and adversary) can be formulated and represented compactly as a linear program. For the GTGRD setting, where the defender can choose  $K$  edges to block at the start of the game, we formulate the problem of computing optimal strategies as a mixed integer program, and present a heuristic algorithm based on LP duality and greedy methods. We perform experiments to show that our heuristic algorithm achieves good performance (i.e., close to defender’s optimal value) with better scalability compared to the mixed-integer programming approach.

## 2 Preliminary: stochastic games

In our two-player zero-sum single-controller stochastic game  $G$ , we have a finite set  $S$  of states, and an initial state  $s_0 \in S$ . The first player acts as an adversary attempting to reach some target within the environment, while second player acts as the observer of the environment. Given a state  $s \in S$ , there exist finite action sets  $J_s$  and  $I_s$  for the adversary and the observer respectively. Given a state  $s \in S$  and  $j \in J_s$ , a single-controller transition function  $\chi(s, j)$  deterministically maps state and action to a new state. Given a state  $s \in S$ ,  $j \in J_s$ ,  $i \in I$ , and intended target of the adversary  $\theta$ , we define a reward function  $r(s, i, j, \theta) \in \mathbb{R}$ . Since this is a zero-sum game, without loss of generality, we define  $r$  as the reward for the observer and the additive inverse of the reward for the adversary. We consider a two-player zero-sum single-controller stochastic game where observer has incomplete information. In particular, the game consists of a collection of zero-sum single-controller stochastic games  $\{G_\theta\}_{\theta \in B}$  and a probability distribution  $P \in \Delta(B)$  over  $B$ . For our setting, we assume that each stochastic game  $G_\theta$  could have different reward function  $r^\theta$ , but all of the games  $G_\theta$  have the same sets of states, actions, and transition rules. The game is played in stages over some finite time. First, a game  $G_\theta$  is drawn according to  $P$ . The adversary is informed of  $\theta$  while the observer does not know  $\theta$ , but rather a set of states  $B$  of which  $\theta$  is a part of. At each stage of game  $t$  with current state  $s_t \in S$ , the adversary selects  $j_t \in J_s$  and the observer selects  $i_t \in I$ , and  $s_{t+1}$  is reached according to  $\chi(s_t, j_t)$ . However, we assume that the adversary does not know  $i_t$ , and both of the players do not know  $r^\theta(s_t, i_t, j_t)$ . Note that observer can infer the action of the adversary given the new state since our transition function is deterministic. Hence, the observer knows  $j_t$ ,  $i_t$ , and  $s_{t+1}$ .

The strategies of the players can be based on their own history of the previous states and strategies. In addition, player 1 can condition his strategies based on  $\theta$ . We consider a finite timestep to be at most  $T$ . Let  $h_t^1 = (s_0, j_0, s_1, j_1, \dots, j_{t-1}, s_t)$  and  $h_t^2 = (s_0, j_0, i_0, s_1, \dots, j_{t-1}, i_{t-1}, s_t)$  to denote a possible history of length  $t$  of player 1 and player 2 where  $j_k \in J_{s_k}$  and  $i_k \in I$  for  $k = 1, \dots, t$ . Let  $H_{s_t}^1$

and  $H_{s_t}^2$  be the set of all possible histories of length  $t$  ended up at state  $s_t$ . Then, the sets of deterministic strategies for player 1 and player 2 are therefore  $\prod_{t=0 \leq t \leq T, s_t \in S, h_{s_t}^1 \in H_{s_t}^1} J_{s_t}$  and  $\prod_{t=0 \leq t \leq T, s_t \in S, h_{s_t}^2 \in H_{s_t}^2} I$ , respectively. Indeed, for each possible history, the players need to select some actions. Naturally, the players mixed strategies are distributions over the deterministic strategies.

**Definition 1.** Given  $\theta \in B$ ,  $0 \leq t \leq T$ ,  $s_t \in S$ ,  $h_{s_t}^1 \in H_{s_t}^1$ , player 1's behavioral strategy  $\sigma_1(\theta, h_{s_t}^1, j_{s_t})$  returns the probability of playing  $j_{s_t} \in J_{s_t}$  such that  $\sum_{j_{s_t} \in J_{s_t}} \sigma_1(\theta, h_{s_t}^1, j_{s_t}) = 1$ . (Player 2's behavioral strategy  $\sigma_2$  is defined similarly and does not depend on  $\theta$ ).

**Definition 2.** A behavioral strategy  $\sigma$  is stationary if and only if it is independent of any timestep  $t$  and depends only on the current state (i.e.,  $\sigma_1(\theta, h_s^1, j_s) = \sigma_1(\theta, \bar{h}_s^1, j_s)$  such that  $h_s^1$  and  $\bar{h}_s^1$  have the same last state and  $\sigma_2$  can be defined similarly).

Given a sequence  $\{(s_t, i_t, j_t)\}_{t=1}^T$  of actions and states, the total reward for player 2 is  $r_T = \sum_{t=1}^T r^\theta(s_t, i_t, j_t)$ . Thus, the expected reward  $\gamma_T(P, s_0, \sigma_1, \sigma_2) = \mathbf{E}_{P, s_0, \sigma_1, \sigma_2}[r_T]$  is the expectation of  $r_T$  over the set of stochastic games  $\{G_\theta\}_{\theta \in B}$  given the the fixed initial state  $s_0$  under  $P$ ,  $\sigma_1$ , and  $\sigma_2$ , respectively.

**Definition 3.** The behavioral strategy  $\sigma_2$  is a best response to  $\sigma_1$  if and only if for all  $\sigma'_2$ ,  $\gamma_T(P, s_0, \sigma_1, \sigma_2) \geq \gamma_T(P, s_0, \sigma_1, \sigma'_2)$ . The behavioral strategy  $\sigma_1$  is a best response to  $\sigma_2$  if and only if for all  $\sigma'_1$ ,  $\gamma_T(P, s_0, \sigma_1, \sigma_2) \leq \gamma_T(P, s_0, \sigma'_1, \sigma_2)$ .

For two-player zero-sum games, the standard solution concept is the max-min solution:  $\max_{\sigma_2} \min_{\sigma_1} \gamma_T(P, s_0, \sigma_1, \sigma_2)$ . One can also define min-max solution  $\min_{\sigma_1} \max_{\sigma_2} \gamma_T(P, s_0, \sigma_1, \sigma_2)$ . For zero-sum games, the max-min value, min-max value, and Nash equilibrium values all coincide [2]. For simultaneous-move games this can usually be solved by formulating a linear program. In this work, we will be focusing on computing the max-min solution.

### 3 Game Model

We begin by describing our settings and introducing the GTGR and GTGRD models.

#### 3.1 Game-theoretic goal recognition model

Consider a deterministic environment such as the one in the introduction. We can model the environment with a graph in which the nodes correspond to the states and the edges connect neighboring states. Given the environment and the graph, as in many standard GR problems, the agent wants to plan out a sequence of moves (i.e., determining a path) to reach its target location of the graph. The target location is unknown to the observer, and the observer's goals

are to identify the target location based on the observed sequence of moves and to make preventive measure to protect the target location.

We model this scenario as a two-player zero-sum game, between the agent/adversary and the observer. Given the graph  $G = (L, E)$  of the environment, the adversary is interested in a set of potential targets  $B \subseteq L$  and has a starting position  $s_0 \in L \setminus B$ . The adversary’s aim is to attack a specific target  $\theta \in B$ , which is chosen at random according to some prior probability distribution  $P$ . The observer does not know the target  $\theta$ , and only the adversary knows its target  $\theta$ . However, the observer knows the set of possible targets  $B$  and the adversary’s starting position  $s_0$ . For any  $s \in L$ , we let  $\nu(s)$  is the set of neighbors of  $s$  in the graph  $G$ .

The sequential game is played over several timesteps where both players move simultaneously. Each timestep, the observer selects a potential target in  $B$  to protect, and the agent moves to a neighboring node. We consider the zero-sum scenario: With each timestep, the adversary and the observer will lose and gain a value  $d$ , respectively. In addition, if the observer protects the correct target location  $\theta$ , an additional value of  $q$  will be added to the observer and subtracted from the adversary. The game ends when the attacker reaches its target  $\theta$ , a value of  $u^\theta$  will be added to the adversary’s overall score, and  $u^\theta$  will be subtracted from the observer’s overall score. Notice that during the play of the game, the adversary does not observe the observer’s action(s), and the players do not know of their current scores.

Because of the potentially stochastic nature of the adversary’s moves at each timestep, and the uncertainty of adversary’s target in the system, our setting is most naturally modeled as a *stochastic game with incomplete information* as defined in Section 2. More specifically, the set of states is  $L$  with an initial state  $s_0$ . Given a state  $s \in S$ ,  $\nu(s)$  is the action set for the adversary and  $B$  is the action set for the observer. Given a state  $s \in S$  and  $j \in \nu(s)$ , the single-controller transition function  $\chi(s, j) = j$ . Indeed, the transition between states are controlled by the adversary only and is deterministic: From state  $s$ , where  $s \neq \theta$ , given attacker action  $j \in \nu(s)$ , the next state is  $j$ . The state  $\theta$  is terminal: Once reached, the game ends. Given a state  $s \in S$ ,  $j \in \nu(s)$ , and  $i \in B$ , we define the reward function  $r^\theta(s, i, j) \equiv r(s, i, j, \theta)$  from the observer’s point of view as

$$r(s, i, j, \theta) = \begin{cases} d & j \neq \theta \ \& \ i \neq \theta \\ d + q & j \neq \theta \ \& \ i = \theta \\ d - u^\theta & j = \theta \ \& \ i \neq \theta \\ d + q - u^\theta & j = \theta \ \& \ i = \theta. \end{cases} \quad (1)$$

While, in theory, the game could go on forever if the adversary never reaches his target  $\theta$ , because of the per-timestep cost of  $d$ , any sufficiently long path for the adversary would be dominated by the strategy of taking the shortest path to  $\theta$ . Eliminating these dominated strategies allows us to set a finite bound for the duration of the game, which grows linearly in the shortest distance to the target that is furthest away. Even in games where the value of  $d$  is set to 0, the defender could potentially play a uniformly random strategy that imposes a

cost of  $\frac{q}{|B|}$  per timestep. Therefore, an adversary strategy taking forever would achieve a value of  $-\infty$  against the uniformly random defender strategy. In any Nash equilibrium the attacker will always reach their target in finite time.

We call this the game-theoretic goal recognition (GTGR) model. All of the definitions in Section 2 follow immediately for our games.

### 3.2 Game-theoretic goal recognition design model

As mentioned in the introduction, we also consider the game-theoretic goal recognition design (GTGRD) model. Formally, before the game starts, we allow the observer to block a subset of at most  $K$  actions from the game. In our model, that corresponds to blocking at most  $K$  edges from the graph. In one variant of the model, blocking an edge effectively removes that edge, i.e. the adversary can no longer take that action. In another variant, blocking an edge does not prevent the adversary from taking the action, but the adversary would incur a cost by taking that action. After placing the blocks, the game proceeds as described in Section 3.1.

## 4 Computation

### 4.1 Game-theoretic goal recognition model

With the game defined, we are interested in computing the solution of the game: What is the outcome of the game when both players behave rationally? Before defining rational behavior, we first need to discuss the set of strategies. In a sequential game, a pure strategy of a player is a deterministic mapping from the current state and the player’s observations/histories leading to the state, to an available action. For the adversary, such observations/histories include its own sequence of prior actions and its target  $\theta$ ; the observer’s observations/histories include the adversary’s sequence of actions and the observer’s sequence of actions. A mixed strategy is a randomized strategy, specified by a probability distribution over the set of pure strategies. The strategies are defined more formally in Section 2 and Definition 1.

As mentioned earlier, we are interested in computing the max-min solution, which is equivalent to the max-min value, min-max value, and Nash equilibrium value of the game. For simultaneous-move games this can usually be solved by formulating a linear program. However, for our sequential game, each pure strategy need to prescribe an action for each possible sequence of observations leading to that state and, as a result, the sets of pure strategies are exponential for both players.

To overcome this computational challenge, we focus on *stationary strategies*, which depend only on the current state (for the adversary, also on  $\theta$ ) and not on the history of observations (see Definition 2). While for stochastic games with complete information, it is known that there always exist an optimal solution that consists of stationary strategies [2], it is an open question whether the same



property holds for our setting, which is an incomplete-information game. Nevertheless, there are some heuristic reasons that stationary strategies are at least good approximations of optimal solutions: The state (i.e., adversary’s location) already captures a large amount of information about the strategic intention of the adversary.

An intuitively optimal non-stationary strategy in which the observer assigns resources to the target with maximal probability, determined through observing the actions of the adversary, presents additional challenges. An optimal strategy of this nature would require information regarding adversary’s strategy from the beginning of the game, so as to determine the likelihood of a given action assuming a particular target for the adversary. Making such assumptions is a straightforward process when restricting the observer to stationary strategies. Later in this paper we will demonstrate how given a stationary strategy for the observer, there exists a best response strategy for the adversary that is also stationary.

Restricting to stationary strategies, randomized strategies now correspond to a mapping from state to a distribution over actions. We have thus reduced the dimension of the solution space from exponential to polynomial in the size of the graph. Furthermore, our game exhibits the single-controller property: The state transitions are controlled by the adversary only. For complete information stochastic games with a single controller, a *linear programming* (LP) formulation is known [19]. We adapt this LP formulation to our incomplete information setting.

We define  $V(\theta, s)$  to be a variable that represents the expected payoff to the observer at state  $s$  and with adversary’s intended target  $\theta$ . We use  $P(\theta)$  to denote the prior probability of  $\theta \in B$  being the adversary’s target such that  $\sum_{\theta \in B} P(\theta) = 1$ . The observer’s objective is to find a (possibly randomized) strategy that maximizes his expected payoff given the prior distribution over the target set  $B$ , the moves of the adversary, and the adversary’s starting location. The following linear program computes the utility of the observer in a max-min solution assuming both players are playing a stationary strategy.

$$\max_{V, \{f_i(s)\}_{i,s}} \sum_{\theta} P(\theta) V(\theta, s_o) \quad (2)$$

$$V(\theta, s) \leq \sum_{i \in B} r(s, i, j, \theta) f_i(s) + V(\theta, j) \quad \forall \theta \in B, \forall s \mid s \neq \theta, \forall j \in \nu(s) \quad (3)$$

$$V(\theta, s) = 0 \quad \text{when } s = \theta \quad (4)$$

$$\sum_i f_i(s) = 1 \quad \forall s \quad (5)$$

$$f_i(s) \geq 0 \quad \forall s, i \quad (6)$$

In the above linear program, (2) is the objective of the observer. The  $f_i(s)$ ’s represent the probability of the observer taking an action  $i \in B$  given the state  $s$ . To ensure a well defined probability distribution for each state of the games, (5) and (6) impose the standard sum-equal-to-one and non-negative conditions

on the probability of playing each action  $i \in B$ . The Bellman-like inequality (3) bounds the expected value for any state using expected values of next states plus the expected current reward, assuming the adversary will choose the state transition that minimizes the observer's expected utility. Finally, (4) specifies the base condition when the adversary has reached their destination and the game ends. The size of the linear program is polynomial in the size of the graph.

The solution of this linear program prescribes a randomized stationary strategy  $f_i(s)$  for the observer and, from the dual solutions, one can compute a stationary strategy for the adversary. In more detail, the dual linear program is

$$\min \sum_s t_s \quad (7)$$

$$t_s \geq \sum_{\theta, j} \lambda_{s,j}^\theta r(s, i, j, \theta) \quad \forall s, i \quad (8)$$

$$I_{s=s_0} P(\theta) + \sum_{s' \neq \theta: s \in \nu(s')} \lambda_{s',s}^\theta = \sum_{j \in \nu(s)} \lambda_{s,j}^\theta \quad \forall \theta \in B, \forall s \neq \theta \quad (9)$$

$$\lambda_{s,j}^\theta \geq 0 \quad \forall \theta, s, j \quad (10)$$

where  $I_{s=s_0}$  is the indicator that equals 1 when  $s = s_0$  and 0 otherwise. The dual variables  $\lambda_{s,j}^\theta$  can be interpreted as the probability that adversary type  $\theta$  takes the edge from  $s$  to  $j$ . These probabilities satisfies the flow conservation constraints (9): given  $\theta$ , the total flow into  $s$  (the left hand side) is equal to the probability that type  $\theta$  visits  $s$ , which should equal the total flow out of  $s$  (the right hand side). The variables  $t_s$  can be interpreted as the contribution to defender's utility from state  $s$ , assuming that the defender is choosing an optimal action at each state (ensured by constraint (8)).

Given the dual solutions  $\lambda_{s,j}^\theta$ , we can compute a stationary strategy for the adversary: let  $\pi(j|\theta, s)$  be the probability that the adversary type  $\theta$  chooses  $j$  at state  $s$ . Then for all  $\theta \in B$  and  $s \neq \theta$ ,  $\pi(j|\theta, s) = \frac{\lambda_{s,j}^\theta}{\sum_{j' \in \nu(s)} \lambda_{s,j'}^\theta}$ . It is straightforward to verify that by playing the stationary strategy  $\pi$ , the adversary type  $\theta$  will visit each edge  $(s, j)$  with probability  $\lambda_{s,j}^\theta$ .

**Lemma 1.** *Given a stationary strategy for the defender, there exists a best response strategy for the adversary that is also a stationary strategy.*

*Proof (Sketch).* Given a stationary defender strategy  $f_i(s)$ , each adversary type  $\theta$  now faces a Markov Decision Process (MDP) problem, which admits a stationary strategy as its optimal solution.

More specifically, since the state transitions are deterministic and fully controlled by the adversary, each type  $\theta$  faces a problem of determining the shortest path from  $s_0$  to  $\theta$ , with the cost of each edge  $(s, j)$  as  $\sum_{i \in B} f_i(s) r(s, i, j, \theta)$ . Looking into the components of  $r(s, i, j, \theta)$ , since the adversary reward  $u^\theta$  for reaching target  $\theta$  occurs exactly once at the target  $\theta$ , it can be canceled out and the problem is equivalent to the shortest path problem from  $s_0$  to  $\theta$  with

edge cost  $d + f_\theta(s)q$ . Since edge costs are nonnegative the shortest paths will not involve cycles.

What this lemma implies is that if the defender plays the stationary strategy prescribed by the LP (2), the adversary cannot do better than the value of the LP by deviating to a non-stationary strategy.

**Corollary 1.** *If the defender plays the stationary strategy  $f_i(s)$  given by the solutions of LP (2), the adversary’s stationary strategy  $\pi$  as prescribed by LP (7) is a best response, i.e., no non-stationary strategies can achieve a better outcome for the adversary.*

While it is still an open question whether the defender has an optimal stationary strategy, we have shown that if we restrict to stationary strategies for the defender, it is in the best interest of the adversary to also stick to stationary strategies and our LP (2) does not overestimate the value of the game.

## 4.2 Game-theoretic goal recognition design model

One can solve this GTGRD problem by brute-force, i.e., try every subset of edges to block and then for each case solve the resulting LP. The time complexity of this approach grows exponentially in  $K$ . Instead, we can encode the choice of edge removal as integer variables added to the LP formulation, resulting in a mixed-integer program (MIP). For example, we could replace (3) with

$$V(\theta, s) \leq \sum_{i \in B} r(s, i, j, \theta) f_i(s) + V(\theta, j) + Mz(s, j) \quad (11)$$

where  $M$  is a positive number, and  $z(s, j)$  is a 0-1 integer variable indicating whether the action/edge from  $s$  to  $j$  is blocked.  $M$  thus represents the penalty that the attacker incurs if he nevertheless chooses to take the edge from  $s$  to  $j$  while it is blocked. By making  $M$  sufficiently large, we can make the actions of crossing a blocked edge dominated and therefore effectively removing the edges that we block. We also add the constraint  $\sum_{s,j} z(s, j) \leq K$ .

**Dual-based greedy heuristic.** The MIP approach scales exponentially in the worst case as the size of the graph and  $K$  grows. We propose a heuristic method for selecting edges to block. We first solve the LP for goal recognition and its dual. In particular, we look at the dual variable  $\lambda_{s,j}^\theta$  for the constraint (3). This dual has the standard interpretation as the *shadow price*: it is the rate of change to the objective if we infinitesimally relax constraint (3).

Looking at the MIP, in particular constraint (11), we see that by blocking off an action from  $s$  to  $j$  we are effectively relaxing the corresponding LP constraints (3) indexed by  $\theta, s, j$  for all  $\theta \in B$ . These are the adversary’s incentive constraints for going from  $s$  to  $j$ , for all adversary types  $\theta$ .

Utilizing the shadow price interpretation of the duals, the sum of the duals corresponding to the edge from  $s$  to  $j$ :  $\sum_{\theta \in B} \lambda_{s,j}^\theta$  gives the rate of change to the

objective (i.e. defender’s expected utility) if the edge  $(s, j)$  is blocked by an infinitesimal amount. Choosing the edge that maximizes this,  $\arg \max_{s,j} \sum_{\theta \in B} \lambda_{s,j}^\theta$  we get the maximum rate of increase of our utility. These rates of changes hold only when the amount of relaxation (i.e.,  $M$ ) is infinitesimal. However, in practice we can still use this as a heuristic for choosing edges to block.<sup>1</sup>

When  $K > 1$ , we could choose the  $K$  edges with the highest dual sums. Alternatively, we can use a greedy approach: pick one edge with the maximum dual sum, place a block on the edge and solve the updated LP for goal recognition, and pick the next edge using the updated duals, and repeat. In our experiments, the latter greedy approach consistently achieved significantly higher expected utilities than the former. Intuitively, by re-solving the LP after adding each edge, we get a more accurate picture of the adversary’s adaptations to the blocked edges. Whereas the rates of changes used by the former approach are only accurate when the adversary do not adapt at all to the blocked edges (see footnote 1). Our greedy heuristic is summarized as follows.

- for  $i = 1 \dots K$ :
  - Solve LP (2), updated with the current blocked edges. If edge  $(s, j)$  blocked, the corresponding constraint (3) indexed  $s, j, \theta$  for all  $\theta$  are modified so that  $M$  is added to the right hand side. Get the primal and dual solutions.
  - Take an edge  $(s^*, j^*) \in \arg \max_{s,j} \sum_{\theta \in B} \lambda_{s,j}^\theta$ , and add it to the set of blocked edges.
- return the set of blocked edges, and the primal solution of the final LP as the defender’s stationary strategy.

## 5 Experiments

Experiments were run on a machine using OSX Yosemite version 10.10.5, with 16 GB of ram and a 2.3 GHz Intel Core i7 processor, and were conducted on grid environments such as the one seen in Figure 2. In these environments, the adversary is allowed to move to adjacent nodes connected by an edge.  $S$  denotes the starting location of the adversary while  $T1$  and  $T2$  denote the locations of two potential targets.

In Figure 2, targets  $T1$  and  $T2$  each have a equal likelihood of being the adversary’s intended target. The adversary’s timestep penalty  $d$  and completion reward  $u^\theta$  are both set to 0. The defender’s reward for correctly guessing the adversary’s intended target  $q$  is set to 10. The attacker penalty value for crossing an edge penalized by the observer is set to 10. The observer is permitted to penalize 3 edges.

---

<sup>1</sup> Another perspective: from the previous section we see that  $\lambda_{s,j}^\theta$  is the probability that adversary type  $\theta$  traverses the edge  $s, j$ . Then if the adversary and defender do not change their strategies after the edge  $(s, j)$  is blocked, the defender would receive an additional utility of  $M \sum_{\theta \in B} \lambda_{s,j}^\theta$  from the adversary’s penalty for crossing that edge.

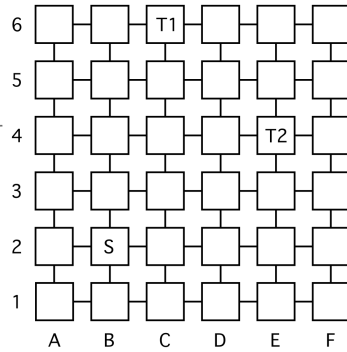


Fig. 2. An Instance of GTGR/GTGRD Games Used in Experiments.

### 5.1 A Comparison of MIP and Greedy Solutions

As seen in Figure 3 and Figure 4, the mixed integer program and greedy heuristic can yield different results. The mixed integer program yields an expected outcome of 43.3 for the observer, while utilizing the greedy heuristic yields an outcome of 40.0 for the observer. The default expected outcome for the observer (in which no edges are penalized) is 30.0. The following experiments averaged the results of similar grid problems.

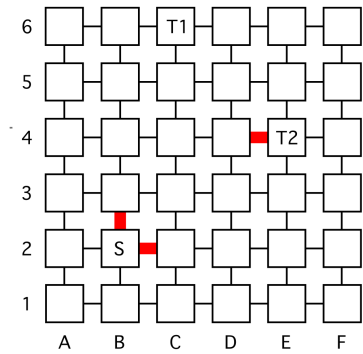


Fig. 3. MIP Solution

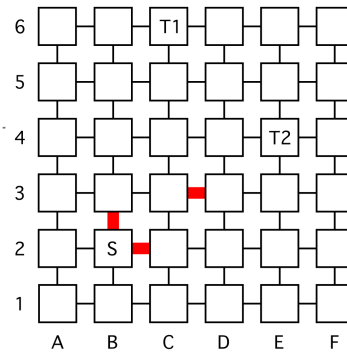


Fig. 4. Greedy Solution

### 5.2 Running Time and Solution Quality

Results from the following experiments were averaged over 1000 grid environments. For each experiment, the adversary's timestep penalty  $d$  and completion

reward  $u^\theta$  were set to 0. For each environment, the starting location of the adversary and all targets are placed randomly on separate nodes. Additionally, each target  $\theta$  is assigned a random probability  $P(\theta)$  such that  $\sum_{\theta \in B} P(\theta) = 1$ . In all of our figures below, the greedy heuristic for the GTGRD is graphed in orange, the MIP is graphed in blue, and the default method (LP) for GTGR is graphed in grey, in which the game is solved with no penalized edges. The defenders reward for correctly guessing the adversary's intended target  $q$  was set to 10. The attacker penalty value for crossing an edge penalized by the observer was set to 10. Each game, the observer was permitted to penalize 2 edges.

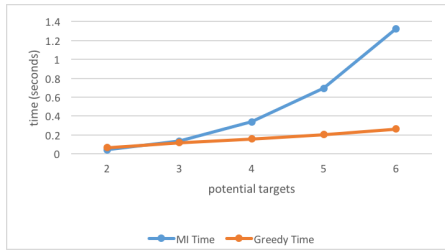


Fig. 5. Average time given targets.

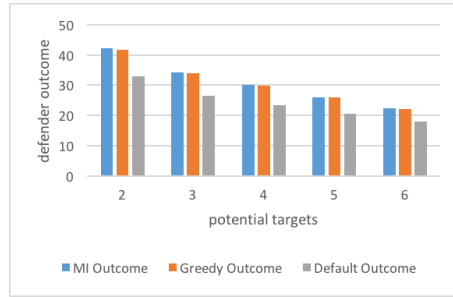


Fig. 6. Average outcome given targets.

**Various Potential Target Sizes** In this set of experiments, we want to investigate the effect of different potential target sizes (i.e.,  $|B|$ ) to the running time (Figure 5) and solution quality (Figure 6) of our algorithms. The results are averaged over 1000 simulations of 6 by 6 grids. Each game, the observer was permitted to penalize 2 edges.

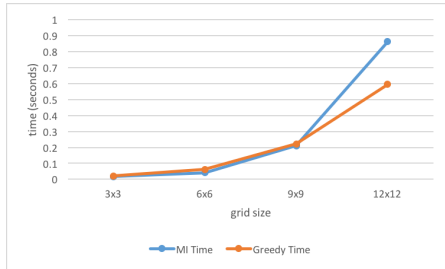


Fig. 7. Average time given size.

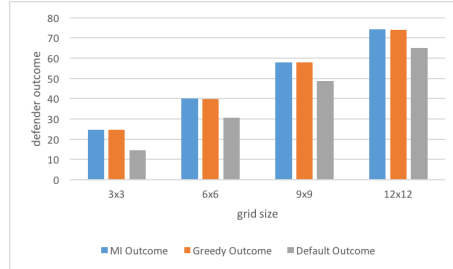
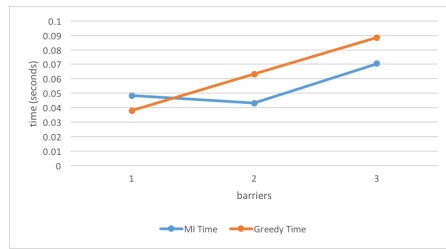


Fig. 8. Average outcome given size.

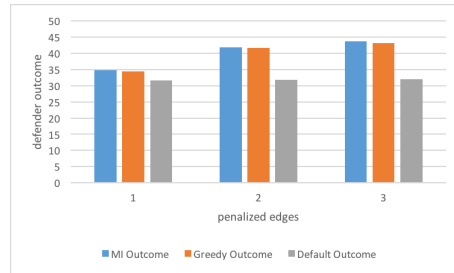
As indicated in Figure 5, the MIP running time increases exponentially while the greedy heuristic running time remains sublinear as we increase the number of potential targets. Moreover, the solution quality (measured by defender’s utility) as seen in Figure 6 suggests that MIP’s solution is closely aligned with our greedy heuristics. This gives evidence that our greedy heuristic provides good solution quality while achieving high efficiency. It is no surprise that the defender’s utility is higher in the GTGRD setting compared to those of GTGR.

**Various Instance Sizes** In this set of experiments, we investigate the effect of different instance sizes (i.e., grids) to the running time (Figure 7) and solution quality (Figure 8) of our algorithms.

Unlike our earlier observations on various target sizes, the average running times for both the MIP and our greedy heuristic increase significantly as we increase the instance sizes (see Figure 7). This is not surprising as now we have more variables and constraints in the integer programs. Despite this, the defender’s utilities generated by greedy heuristic are relatively similar to those generated using MIP (see Figure 8).



**Fig. 9.** Average time given penalized edges.

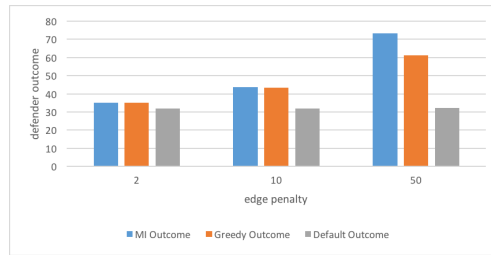


**Fig. 10.** Average outcome given penalized edges.

**Various Number of Barriers/Blocks** In this set of experiments, we want to investigate the effect of different number of barriers (i.e.,  $K$ ) to the running time (Figure 5) and solution quality (Figure 6) of our algorithms in the GTGRD models. The results are averaged over 1000 simulations of 6 by 6 grids.

It turns out that as we increase the number of barriers, the running times of our greedy heuristic are longer than the MIP as shown in Figure 9. Nonetheless, as in the earlier experiments, both algorithms have similar solution quality as shown in Figure 10.

**Various Edge Penalties** Finally, consider the effect of different edge penalties to the solution quality of our greedy heuristic. The results are averaged over 1000



**Fig. 11.** Average outcome given penalty value.

simulations of 6 by 6 grids. As indicated in Figure 11, the solution gap between the MIP and greedy heuristic as we increase the edge penalty.

## References

1. Braynov, S.: Adversarial planning and plan recognition: Two sides of the same coin. In: Proceedings of the Secure Knowledge Management Workshop (2006)
2. Fudenberg, D., Tirole, J.: Game theory (1991)
3. Geib, C., Steedman, M.: On natural language processing and plan recognition. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJ-CAI). pp. 1612–1617 (2007)
4. Guillaume, N.L., Mouaddib, A.I., Lerouvreur, X., Gatepaille, S.: A generative game-theoretic framework for adversarial plan recognition. In: Proceedings of the Workshop on Distributed and Multi-Agent Planning (2015)
5. Jarvis, P., Lunt, T., Myers, K.: Identifying terrorist activity with AI plan recognition technology. *AI Magazine* 26(3), 73–81 (2005)
6. Johnson, W.L.: Serious use of a serious game for language learning. *International Journal of Artificial Intelligence in Education* 20(2), 175–195 (2010)
7. Kautz, H.A.: A Formal Theory of Plan Recognition. Ph.D. thesis, Bell Laboratories (1987)
8. Kelley, R., Wigand, L., Hamilton, B., Browne, K., Nicolescu, M., Nicolescu, M.: Deep networks for predicting human intent with respect to objects. In: Proceedings of the International Conference on Human-Robot Interaction (HRI). pp. 171–172 (2012)
9. Keren, S., Gal, A., Karpas, E.: Goal recognition design. In: Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS). pp. 154–162 (2014)
10. Keren, S., Gal, A., Karpas, E.: Goal recognition design for non-optimal agents. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 3298–3304 (2015)
11. Keren, S., Gal, A., Karpas, E.: Goal recognition design with non-observable actions. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 3152–3158 (2016)
12. Lee, S., Mott, B., Lester, J.: Real-time narrative-centered tutorial planning for story-based learning. In: Proceedings of the International Conference on Intelligent Tutoring Systems (ITS). pp. 476–481 (2012)



13. Lisý, V., Píbil, R., Stiborek, J., Bosanský, B., Pechoucek, M.: Game-theoretic approach to adversarial plan recognition. In: Proceedings of the European Conference on Artificial Intelligence (ECAI). pp. 546–551 (2012)
14. McQuiggan, S., Rowe, J., Lee, S., Lester, J.: Story-based learning: The impact of narrative on learning experiences and outcomes. In: Proceedings of the International Conference on Intelligent Tutoring Systems (ITS). pp. 530–539 (2008)
15. Min, W., Ha, E., Rowe, J., Mott, B., Lester, J.: Deep learning-based goal recognition in open-ended digital games. In: Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment (AIIDE) (2014)
16. Oh, J., Meneguzzi, F., Sycara, K., Norman, T.: ANTIPA: An agent architecture for intelligent information assistance. In: Proceedings of the European Conference on Artificial Intelligence (ECAI). pp. 1055–1056 (2010)
17. Oh, J., Meneguzzi, F., Sycara, K., Norman, T.: An agent architecture for prognostic reasoning assistance. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI). pp. 2513–2518 (2011)
18. Oh, J., Meneguzzi, F., Sycara, K., Norman, T.: Probabilistic plan recognition for intelligent information agents: Towards proactive software assistant agents. In: Proceedings of the International Conference on Agents and Artificial Intelligence (ICAART). pp. 281–287 (2011)
19. Raghavan, T.E.S.: Finite-step algorithms for single-controller and perfect information stochastic games. In: Neyman, A., Sorin, S. (eds.) *Stochastic Games and Applications*, pp. 227–251. Springer Netherlands (2003)
20. Ramírez, M., Geffner, H.: Plan recognition as planning. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI). pp. 1778–1783 (2009)
21. Ramírez, M., Geffner, H.: Probabilistic plan recognition using off-the-shelf classical planners. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI). pp. 1121–1126 (2010)
22. Ramírez, M., Geffner, H.: Goal recognition over POMDPs: Inferring the intention of a POMDP agent. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI). pp. 2009–2014 (2011)
23. Schmidt, C., Sridharan, N., Goodson, J.: The plan recognition problem: An intersection of psychology and artificial intelligence. *Artificial Intelligence* 11(1-2), 45–83 (1978)
24. Son, T.C., Sabuncu, O., Schulz-Hanke, C., Schaub, T., Yeoh, W.: Solving goal recognition design using asp. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI) (2016)
25. Sukthankar, G., Geib, C., Bui, H.H., Pynadath, D., Goldman, R.P.: *Plan, activity, and intent recognition: Theory and practice*. Newnes (2014)
26. Tavakkoli, A., Kelley, R., King, C., Nicolescu, M., Nicolescu, M., Bebis, G.: A vision-based architecture for intent recognition. In: Proceedings of the International Symposium on Advances in Visual Computing. pp. 173–182 (2007)
27. Vilain, M.: Getting serious about parsing plans: A grammatical analysis of plan recognition. In: Proceedings of the National Conference on Artificial Intelligence (AAAI). pp. 190–197 (1990)
28. Wayllace, C., Hou, P., Yeoh, W., Son, T.C.: Goal recognition design with stochastic agent action outcomes. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) (2016)