

# **Your Title Goes Here (It Can Be Really Really Really Really Long)**

Your Name Here

## **Abstract**

A nice abstract goes here.

## Acknowledgments

Some acknowledgments go here.

**Your Title Goes Here (It Can Be Really Really Really Really Long)**

Your Name Here

A departmental senior thesis submitted to the  
Department of Computer Science at Trinity University  
in partial fulfillment of the requirements for graduation  
with departmental honors.

April 1, 2005

---

Thesis Advisor

---

Department Chair

---

Associate Vice President  
for  
Academic Affairs

Student Copyright Declaration: the author has selected the following copyright provision:

☐ This thesis is licensed under the Creative Commons Attribution-NonCommercial-NoDerivs License, which allows some noncommercial copying and distribution of the thesis, given proper attribution. To view a copy of this license, visit <http://creativecommons.org/licenses/> or send a letter to Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

☐ This thesis is protected under the provisions of U.S. Code Title 17. Any copying of this work other than “fair use” (17 USC 107) is prohibited without the copyright holder’s permission.

☐ Other:

Distribution options for digital thesis:

☒ Open Access (full-text discoverable via search engines)

☐ Restricted to campus viewing only (allow access only on the Trinity University campus via [digitalcommons.trinity.edu](http://digitalcommons.trinity.edu))

**Your Title Goes Here (It Can Be  
Really Really Really Really Long)**

Your Name Here

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.0.1	Game-Theoretic Goal Recognition Problems in Security Domains . .	2
1.0.2	Game-Theoretic Goal Recognition Design Problems in Security Do- mains . . . . .	3
1.0.3	Related Work . . . . .	4
1.0.4	Our Contributions . . . . .	5
1.1	Preliminary: stochastic games . . . . .	6
<b>2</b>	<b>Game Model</b>	<b>8</b>
2.0.1	Game-theoretic goal recognition model . . . . .	8
2.0.2	Game-theoretic goal recognition design model . . . . .	10
<b>3</b>	<b>Computation</b>	<b>12</b>
3.0.1	Game-theoretic goal recognition model . . . . .	12
3.0.2	Game-theoretic goal recognition design model . . . . .	16
<b>4</b>	<b>Partially Observable Environments</b>	<b>19</b>
4.1	Introduction . . . . .	19

4.2	The Whale Method . . . . .	20
4.3	The Transmogrification Method . . . . .	22
4.4	Experiments . . . . .	25
<b>A</b>	<b>Example appendix</b>	<b>29</b>

# List of Tables

4.1	Performance results from simple environment. . . . .	27
4.2	Average method performance in complex environment. . . . .	28

# List of Figures

1.1	Example Problem (left) and with Blocked Actions in Red (right). . . . .	1
4.1	A partially observable graph. . . . .	20
4.2	A partially observable graph with two shadow sets. . . . .	20
4.3	A partially observable graph with a single state in place of shadow states. .	21
4.4	The partially observable environment from Figure 4.1 (left) transmogrified into an, fully observable graph. . . . .	23
4.5	The adversary's path to target $T1$ . . . . .	24
4.6	The adversary's path to $T1$ from the observer's perspective. . . . .	25
4.7	A simple scenario using NetworkX. . . . .	26
4.8	The transmogrified graph. . . . .	26
4.9	Complex testing environment. . . . .	28



# Chapter 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 [?, ?]. In AI, this problem is known as *goal recognition* (GR) or, more generally, *plan recognition* [?]. Plan and goal recognition problems have been used to model a number of applications ranging from software personal assistants [?, ?, ?]; robots that interact with humans in social settings such as homes, offices, and hospitals [?, ?]; intelligent tutoring systems that recognize sources of confusion or misunderstanding in students through their interactions with the system [?, ?, ?, ?]; and security applications that recognize the plan or goal of terrorists [?].

	1	2	3	4	5
A					G1
B	G1				
C					G2
D					
E					

Figure 1.1: Example Problem (left) and with Blocked Actions in Red (right).

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 ?? (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.

### 1.0.1 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 ?. Consider the security scenario in Figure ? (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.

### 1.0.2 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 ??, 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 [?, ?] 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 ?? (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 GTGRD problem where the observer can modify the underlying environment (i.e., adding  $K$  roadblocks) as to restrict the actions of the agent.

### 1.0.3 Related Work

GR and its more general forms, plan recognition and intent recognition, have been extensively studied [?] since their inception almost 40 years ago [?]. Researchers have made significant progress within the last decade through synergistic integrations of techniques ranging from natural language processing [?, ?] to classical planning [?, ?, ?] and deep learning [?]. 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 [?], a model where the game is over attack graphs [?], and an extension that allows for stochastic action outcomes [?]. 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 [?], 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 [?], actions of the agents can be stochastic [?], and actions of the agents can be only partially observable [?]. Further, aside from all the decision-theoretic approaches above, researchers have also modeled and solved the original GRD problem using answer set programming [?]. 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.

#### 1.0.4 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.

## 1.1 Preliminary: stochastic games

In our two-player zero-sum single-controller stochastic game  $G$ , (a) we have a finite set  $S$  of states, and an initial state  $s_0 \in S$ , (b) given a state  $s \in S$ , a finite action set  $J_s$  and  $I = I_s$  for the first player and for the second player, respectively, (c) given a state  $s \in S$  and  $j \in J_s$ , a single-controller transition function  $\chi(s, j)$  that deterministically maps state and action to a state, and (d) given a state  $s \in S$ ,  $j \in J_s$ , and  $i \in I$ , a reward function  $r(s, i, j, \theta) \in \mathbb{R}$ . Since this is a zero-sum game, without loss of generality, we define  $r$  to be the reward for player 2 and the reward of player 1 is the negative reward of player 2. We consider two-player zero-sum single-controller stochastic game where player 2 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  $G'_\theta$ s 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 first player is informed of  $\theta$  while the second player does not know  $\theta$ . At each stage of game  $t$  with current state  $s_t \in S$ , the first player selects  $j_t \in J_s$  and the second player selects  $i_t \in I$ , and  $s_{t+1}$  is reached according to  $\chi(s_t, j_t)$ . However, we assume that player 1 does not know  $\theta$ , and both of the players do not know  $r^\theta(s_t, i_t, j_t)$ . Note that player 2 can infer the action of player 1 given the new state since our transition function is deterministic. Hence, player 2 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 finite timestep 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 1.1.1.** 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 1.1.2.** 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 fixed initial state  $s_0$  under  $P$ ,  $\sigma_1$ , and  $\sigma_2$ , respectively.

**Definition 1.1.3.** 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 [?]. 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.

## Chapter 2

# Game Model

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

### 2.0.1 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 game is sequential and is played over several timesteps where both of the players move simultaneously. At each timestep, the observer selects a potential target in  $B$  to protect, and the agent moves from its current position 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 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 1.1. 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

$$r(s, i, j, \theta) = \begin{cases} d & j \neq \theta \text{ \& } i \neq \theta \\ d + q & j \neq \theta \text{ \& } i = \theta \\ d - u^\theta & j = \theta \text{ \& } i \neq \theta \\ d + q - u^\theta & j = \theta \text{ \& } i = \theta. \end{cases} \quad (2.1)$$

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 1.1 follow immediately for our games.

### 2.0.2 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 blocking 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 2.0.1.

## Chapter 3

# Computation

### 3.0.1 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 1.1 and Definition 1.1.1.

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 are strategies that depend only on the current state (for the adversary, also on  $\theta$ ) and not on the history of observations (see Definition 1.1.2). While for stochastic games with complete information, it is known that there always exist an optimal solution that consists of stationary strategies [?], 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 approximately optimal solutions: The state (i.e., adversary's location) already capture a large amount of information about the strategic intention of the adversary.

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 exhibit 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 [?]. 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 target begin  $\theta$ . 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 an max-min

solution assuming both players are playing a stationary strategy.

$$\max_{V, \{f_i(s)\}_{i,s}} \sum_{\theta} P(\theta) V(\theta, s_o) \quad (3.1)$$

$$V(\theta, s) \leq \sum_{i \in B} r(s, i, j, \theta) f_i(s) + V(\theta, j) \quad \forall \theta \in B, \forall s \mid s \neq \theta, \forall j \in \nu(s) \quad (3.2)$$

$$V(\theta, s) = 0 \quad \text{when } s = \theta \quad (3.3)$$

$$\sum_i f_i(s) = 1 \quad \forall s \quad (3.4)$$

$$f_i(s) \geq 0 \quad \forall s, i \quad (3.5)$$

In the above linear program, (3.1) 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 that the probability distribution is well defined at each state of the games, (3.4) and (3.5) 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 (3.2) 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, (3.3) 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

$$\min \sum_s t_s \quad (3.6)$$

$$t_s \geq \sum_{\theta, j} \lambda_{s,j}^\theta r(s, i, j, \theta) \quad \forall s, i \quad (3.7)$$

$$I_{s=s_0} P(\theta) + \sum_{s' \neq \theta: s \in \nu(s')} \lambda_{s',s}^\theta = \sum_{j \in \nu(s)} \lambda_{s,j}^\theta \quad \forall \theta \in B, \forall s \neq \theta \quad (3.8)$$

$$\lambda_{s,j}^\theta \geq 0 \quad \forall \theta, s, j \quad (3.9)$$

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 (3.8): 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 (3.7)).

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

*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.  $\square$

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 (3.1), the adversary cannot do better than the value of the LP by deviating to a non-stationary strategy.

**Corollary 3.0.1.** *If the defender plays the stationary strategy  $f_i(s)$  given by the solutions of LP (3.1), the adversary's stationary strategy  $\pi$  as prescribed by LP (3.6) 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 strategy that is stationary, 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 (3.1) do not overestimate the value of the game.

### 3.0.2 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 (3.2) with

$$V(\theta, s) \leq \sum_{i \in B} r(s, i, j, \theta) f_i(s) + V(\theta, j) + M z(s, j) \quad (3.10)$$

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 (3.2). This dual has the standard interpretation as the *shadow price*: it is the rate of change to the objective if we infinitesimally relax constraint (3.2).

Looking at the MIP, in particular constraint (3.10), we see that by blocking off an action from  $s$  to  $j$  we are effectively relaxing the corresponding LP constraints (3.2) 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 (3.1), updated with the current blocked edges. If edge  $(s, j)$  blocked, the corresponding constraint (3.2) 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.

---

<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.

## Chapter 4

# Partially Observable Environments

### 4.1 Introduction

Until now, both the GTGR and GTGRD models have given the observer full knowledge of the adversary's state for the entirety of the game. In real-world environments, observers may not have perfect information regarding the states and actions of an adversary.

To accommodate for scenarios with incomplete information for the adversary, we introduce a partially observable variant of the GTGR scenario. In partially observable scenarios, the rules of the game remain largely unchanged, except for addition of shadow states.” The observer can not discern the current state of the adversary, while the adversary occupies a shadow state. When the adversary enters an observable portion of the graph, the observer will become aware of the adversary's position once more.

Figure 4.1 illustrates a partially observable environment. Visible states, in which the observer can see the adversary white. Shadow states, in which the adversary is hidden from the observer, are black. The agent starts the game in state  $S$ . When the adversary moves to states 4, 5, 6, or 7, the observer is unable to determine their position until the adversary

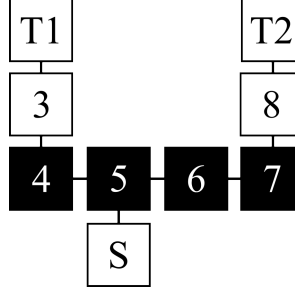


Figure 4.1: A partially observable graph.

re-enters a visible portion of the graph. We will examine two solutions to the partially observable model, both of which involving linear programming.

## 4.2 The Whale Method

The first method of solving partially observable environments, which we will call the "Whale Method," will utilize disjoint sets of shadow states. We call these sets of shadow states "shadow sets." We say that two shadow states belong to the same shadow set, if the adversary can travel between the two without entering an observable state.

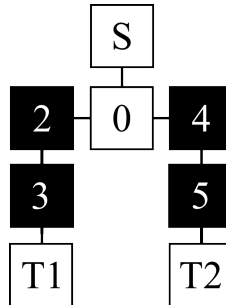


Figure 4.2: A partially observable graph with two shadow sets.

The graph in Figure 4.2 has two disjoint shadow sets, one composed of states 2 and 3,

the other composed of states 4 and 5. In using the Whale Method, the observer treats each shadow set a single state, which we will call a "whale state." We can identify the single shadow set in Figure 4.1, composed of states 4, 5, 6, and 7. In using the Whale method, the observer treats each state in the shadow set as the same state. While the adversary may require several turns to travel among states 4, 5, 6, and 7, the adversary will act as if the has decided to remain stationary in the newly created state  $W$  seen in Figure 4.3

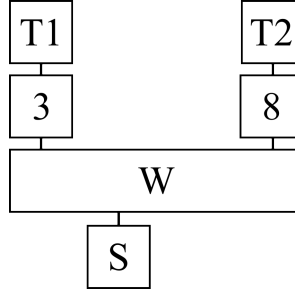


Figure 4.3: A partially observable graph with a single state in place of shadow states.

We can add the following to the mixed integer program to accommodate for partially observable environments when using the Whale method.

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

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

$$\sum_i f_i(s) = 1 \quad \forall s \quad (3)$$

$$\sum_w f_w(s) = 1 \quad \forall s \quad (3)$$

$$f_i(s) \geq 0 \quad \forall s, i \quad (4.1)$$

$$f_w(s) \geq 0 \quad \forall s, w \quad (4.2)$$

We let  $H$  denote the set of all shadow states, and  $w$  denote the whale state the observer knows the adversary to be occupying. An observer action for any whale state  $w$  is written as  $f_w(s)$ . These changes to the linear program require the observer to take the same action for each turn the adversary spends in a particular shadow set. The performance of the Whale method will be examined in a later section.

### 4.3 The Transmogrification Method

When using the Whale method, the observer ignores some of the information available to them. The Whale method does not account for where the adversary entered a shadow set, or how long the adversary has remained hidden. The "Transmogrification Method" takes both of these pieces of information into account, by generating a fully observable environment from a partially observable environment. To do this, the observer must make some basic assumptions about the adversary's strategy. The following lemmas and corollary assume the agent is playing optimally against the observer's stationary strategy.

**Lemma 2.** *If the adversary's target does not lie within a shadow set, the adversary will eventually exit the shadow set.*

*Proof (Sketch).* As mentioned previously, the game could theoretically go on forever.

But because of the potential per-timestep cost of  $d$  and the observer’s predictions, any sufficiently long path for the adversary would be dominated by the strategy of taking the shortest path to their target  $\theta$ . If the adversary never leaves the set of shadow states, then the game will go on forever. Thus, the adversary will eventually leave a shadow set.

**Corollary 2.** *An optimal agent occupying a state in a shadow set will take a shortest path to the exit state of their choosing.*

*Proof (Sketch).* We established that an optimal adversary in a shadow set must exit that shadow set. Each unnecessary turn the adversary spends in a shadow state invites the observer to guess their intended target. Thus, it is in the observer’s interest to reach their chosen exit as quickly as possible.

With Corollary 2, the observer can generate a new graph to play on.

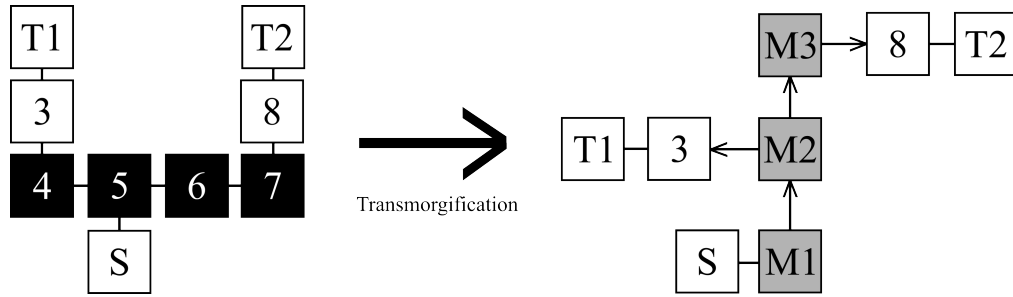


Figure 4.4: The partially observable environment from Figure 4.1 (left) transmogrified into an, fully observable graph.

Figure 4.4 illustrates a transmogrified version of the graph from Figure 4.1. The grey nodes in the transmogrified graph are introduced to represent the number of turns the adversary has remained hidden from the observer. For instance, if the adversary has been hidden for two turns, then the observer would see the adversary as occupying state  $M2$  in the transmogrified graph. Take note that the transmogrified graph now includes directed

edges, because we are (hopefully) not dealing with a time traveling agent. To generate these graphs, we examine each shadow set. For each shadow set, we examine each "entrance state" connected by an edge to the shadow state. We then compute the shortest path between every two entrance states for that shadow state. By Corollary 2, we assume the length of the longest shortest path will be the maximum number of turns an optimal agent will spend within the shadow set. We then create the same number of memory nodes with directed edges from one to the next (see states  $M1, M2, M3$  from state  $S$  in Figure 4.4). Then we create an edge from each entrance state, to the memory state corresponding to the length of the shortest path between the entrance state, and the original entrance state from which we spawned the memory states. For example, because a adversary would have to two shadow states on their journey from state  $S$  to state 3 in Figure 4.1, we connect state  $M2$  to state 3 in Figure 4.4. Note, that the post transmogrification graph in Figure 4.4 is actually incomplete, because memory states have only been spawned from the entrance state  $S$  and not from entrance states 3 or 8. Because an optimal agent has no reason to backtrack in this particular scenario, the memory nodes for 3 and 8 were omitted to keep the graph simple.

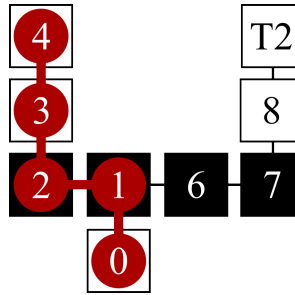


Figure 4.5: The adversary's path to target  $T1$ .

Next, let us examine how a game might play out from the adversary's perspective.



Using the environment from Figure 4.1, let us assume the adversary was assigned target  $T1$ . Because the graph is rather limiting, the adversary takes the shortest path from starting state  $S$  to the target  $T1$ . Figure 4.5 illustrates the adversary’s journey, each turn market with a red circle. Note that on turns 1 and 2, the observer would lose sight of the attacker until turn 3. Now, let us examine what the same scenario would look like from the observer’s perspective, when using the transmogrification method.

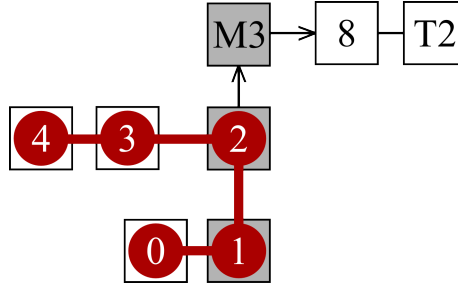


Figure 4.6: The adversary’s path to  $T1$  from the observer’s perspective.

Because the adversary takes two turns within the shadow set, the observer sees the adversary as moving from state  $M1$  to state  $M2$ , before exiting the shadow set on turn 3.

## 4.4 Experiments

As with previous experiments, all tests 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. First, we will use the simple example from Figure 4.1 to compare the performance of the Whale, and Transmogrification methods. Additionally, we will compare the performance against a best case solution, in which we turn all shadow states into standard states. The best case solutions measures how the observer would perform if their strategy effectively negated the shadow states. Graphs were displayed using the Python NetworkX library. Arrowheads on directed

edges were added manually for clarity. The starting state is displayed in green, target states are displayed in red, shadow states are displayed in purple, and default states are displayed in cyan.

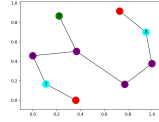


Figure 4.7: A simple scenario using NetworkX.

After transmogrifying the graph, the scenario becomes more complex, but allows the observer to play a fully observable game. There graph will no longer have any shadow states (purple). All states labeled with values greater than 1000 are nodes that have been added to represent turns hidden within the shadow set. Note that the graph in Figure 4.8 will have more nodes than the previously seen illustration of the transmogrified graph in Figure 4.4, because we spawn memory states for every entry state, and not just the starting state.

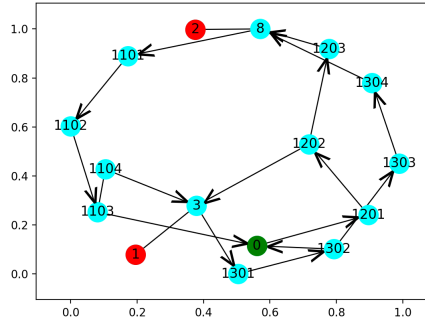


Figure 4.8: The transmogrified graph.

For this game, the adversary is not penalized for timesteps taken. The reward for

reaching the target is 0. The reward the observer receives for each correct guess was set to 1. The point value for each method at the end of the game equates to the number of correct guesses the observer can expect to make. At the start of the game, the adversary has a 75% chance of being assigned target 1, and a 25% of being assigned target 2.

<b>No Shadow States</b>	3.75 Correct Guesses
<b>Whale</b>	3.25 Correct Guesses
<b>Transmogrification</b>	3.50 Correct Guesses

Table 4.1: Performance results from simple environment.

If a strategy were to effectively remove the shadow states from the board, the observer could expect to make 3.75 correct guesses over the course of the game. Thus, we could not expect the whale or transmogrification methods to perform any better. Using the whale methods, in which all shadow states are mashed together into one whale of a state, the observer can expect to make 3.25 correct guesses. While the transmogrification method does not allow the observer to effectively see through shadow states, the method still outperforms the whale method with a score of 3.50. Next, we will test these methods against each other in a more complex scenario.

With the simple example out of the way, we move to a more complex environment which offers the adversary more freedom in approaching their target. For the next set of tests, we place the adversary starting state, and 3 potential targets on non-shadow states (marked in white) in Figure 4.9. We then create a random probability distribution, and solve the game with both methods. After ten-thousand iterations, the scores are averaged and compared.

Though the transmogrification method does not yield what the observer could score without shadow states, it still outperforms the whale method.

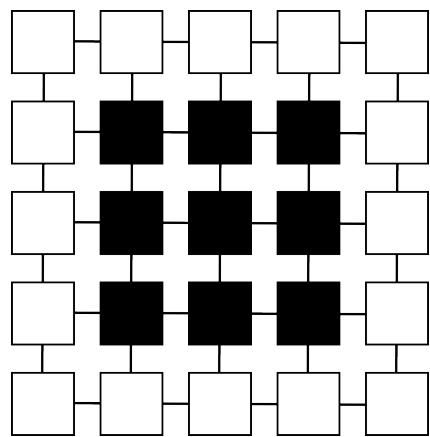


Figure 4.9: Complex testing environment.

No Shadow States	2.65 Correct Guesses
Whale	2.41 Correct Guesses
Transmogrification	2.59 Correct Guesses

Table 4.2: Average method performance in complex environment.

## Appendix A

### Example appendix

Here is my code.

```
#include <iostream>
int main(void) {
    cout << "Hello, world!\n";
    return 0;
}
```