

# Privacy Preserving Plans in Partially Observable Environments

Using Goal Recognition Design for Improved Privacy  
IJCAI 16

Sarah Keren   Avigdor Gal   Erez Karpas



Faculty of Industrial Engineering and Management  
Technion — Israel Institute of Technology

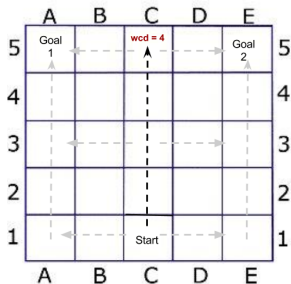
Haifa Security Research Seminar 11/2016

## Offline design as a way to facilitate online goal recognition



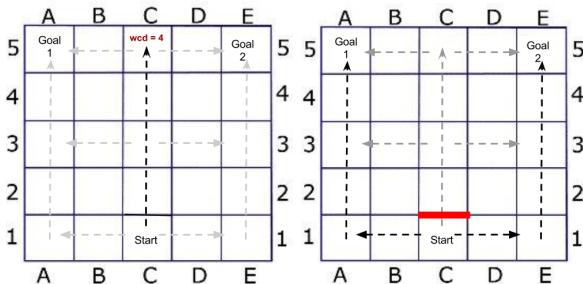
Worst case distinctiveness (wcd) as a measure of model quality

## Offline design as a way to facilitate online goal recognition



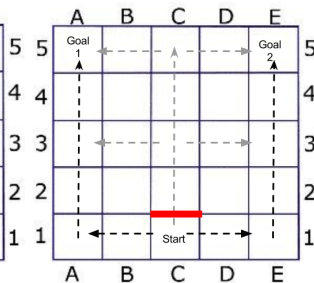
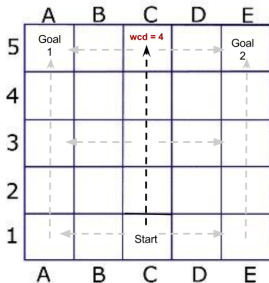
Worst case distinctiveness (wcd) as a measure of model quality

## Offline design as a way to facilitate online goal recognition



Worst case distinctiveness (wcd) as a measure of model quality

## Offline design as a way to facilitate **online** goal recognition



**Worst case distinctiveness (wcd) as a measure of model quality**

# Applications



(a) Intrusion Detection



(b) E-Commerce and Personalized Advertisement



(c) Human-Robot Teamwork



(d) Smart Home Design



(e) Virtual environments

## Deterministic Environment

- ▶ Optimal fully observable agents (ICAPS 2014)
- ▶ Sub-Optimal fully observable agents (AAAI 2015)
- ▶ Some Actions are Non-observable (AAAI 2016)
- ▶ Arbitrary sensor model (IJCAI 2016)
- ▶ Compilation to ASP (Son et. al., AAAI 2016)

## Stochastic Environment

- ▶ Solution using MDP (Wayllace et. al., IJCAI 2016)



# Extending the Goal Recognition Design Framework

## Deterministic Environment

- ▶ Optimal fully observable agents (ICAPS 2014)
- ▶ Sub-Optimal fully observable agents (AAAI 2015)
- ▶ Some Actions are Non-observable (AAAI 2016)
- ▶ Arbitrary sensor model (IJCAI 2016)
- ▶ Compilation to ASP (Son et. al., AAAI 2016)

## Stochastic Environment

- ▶ Solution using MDP (Wayllace et. al., IJCAI 2016)





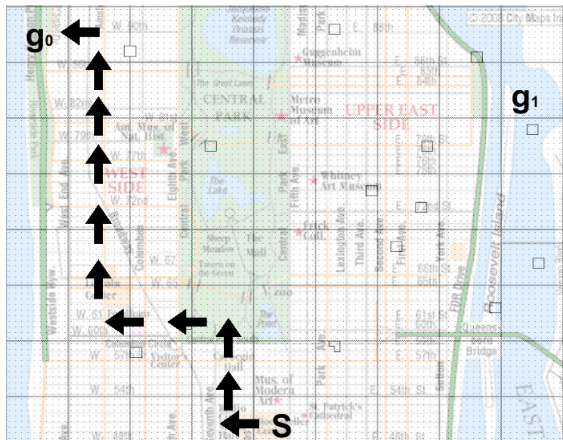
**Cloaking : How long can an agent keep his goal ambiguous ?**



**A user can choose a path that potentially maximizes its privacy**

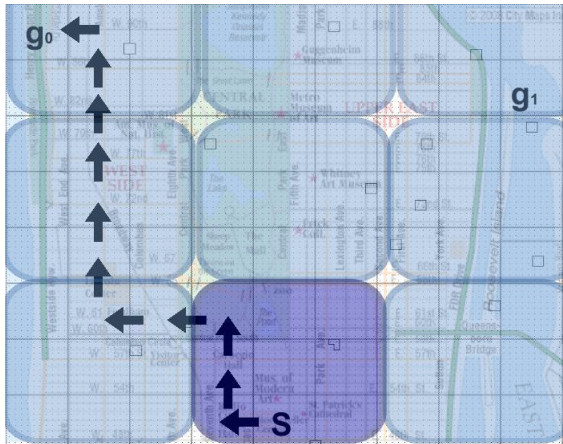
the *wcd*-path that allows him to stay ambiguous for at most *wcd* steps

# Example 1



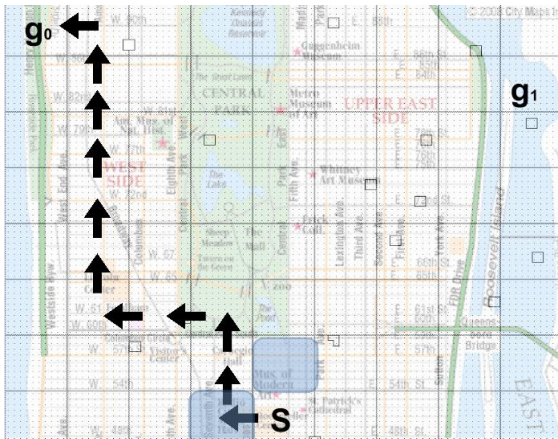
Full Observability

# Example 1



Coarse Sensors

# Example 1



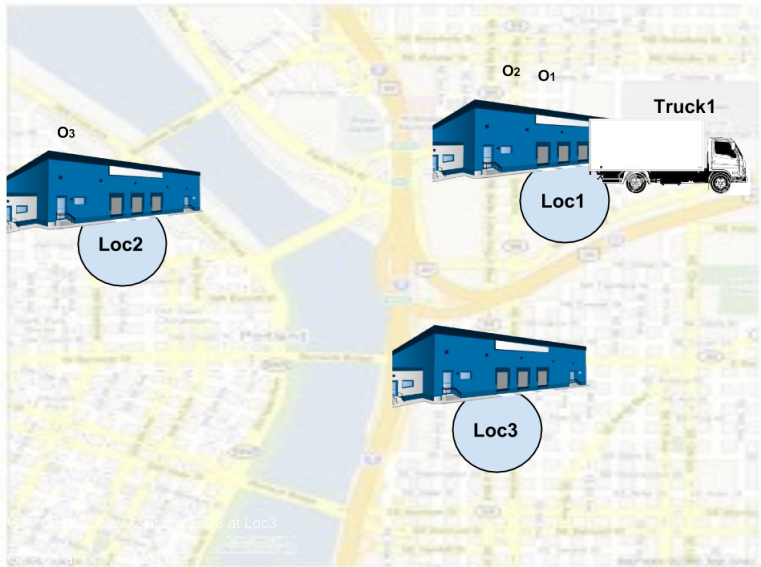
Noisy Sensors

## Example 2



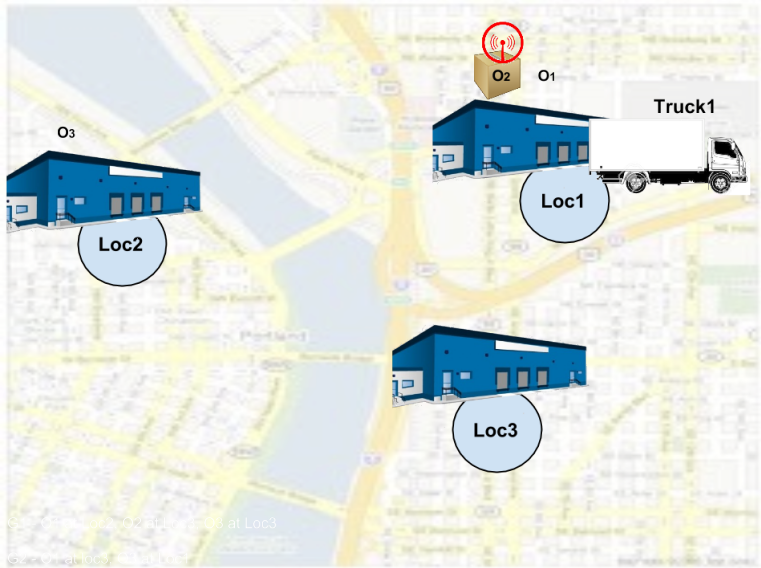
Full Observability

## Example 2



### Coarse Sensors

## Example 2



### Noisy Sensors

## Sensor Model

Maps each action to a set of possible observation tokens.  
The special token  $o_\emptyset$  denotes non-observable action.



## Sensor Model

Maps each action to a set of possible observation tokens.  
The special token  $o_{\emptyset}$  denotes non-observable action.

## Observable Projection

The **observable projection** of a path is a **set** of possible observation sequences, determined by the sensor model.

## Sensor Model

Maps each action to a set of possible observation tokens.  
The special token  $o_{\emptyset}$  denotes non-observable action.

## Observable Projection

The **observable projection** of a path is a **set** of possible observation sequences, determined by the sensor model.

## Non-distinctive Path

A path is **non-distinctive** if it has an observable projection, which is also the observable projection of a path leading to a different goal.

# Goal Recognition Design with Arbitrary Sensor Models

## Sensor Model

Maps each action to a set of possible observation tokens.  
The special token  $o_{\emptyset}$  denotes non-observable action.

## Observable Projection

The **observable projection** of a path is a **set** of possible observation sequences, determined by the sensor model.

## Non-distinctive Path

A path is **non-distinctive** if it has an observable projection, which is also the observable projection of a path leading to a different goal.

## Worst Case Distinctiveness

The **worst case distinctiveness** (*wcd*) is the maximal non-distinctive path .

## Our language :

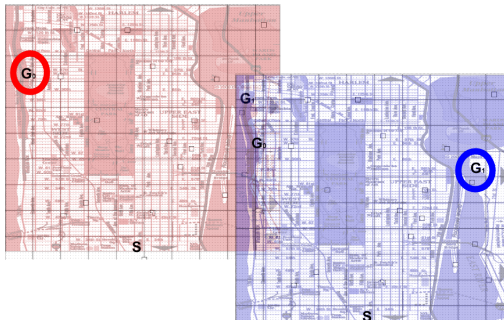
- ▶ STRIPS-like model:
  - ▶ Fluents  $F$
  - ▶ Actions  $A$  with  $a = \langle pre(a), add(a), del(a) \rangle$
  - ▶ Initial state  $s_0 \subseteq F$
  - ▶ Set of possible goals  $\mathcal{G}$
  - ▶ (Optional) sensor model which maps actions  $A$  to observation tokens

## Our tools:

- ▶ Off-the-shelf solvers (optimal and approximate)

## Calculating wcd: Compilation to Classical Planning

- ▶ We **compile** a goal recognition design problem with two goals as a planning problem with two agents each aiming at a separate goal
- ▶ Actions divided into
  - ▶ 'real' actions: change the state of the world
  - ▶ 'declare' actions: declare the observation token a 'real' action emits
- ▶ As long as both agents have declared the same observation sequence, they can get a discount when they declare the same observation token



# Empirical Evaluation : wcd

	LOGISTICS					BLOCKS WORLD					GRID-NAVIGATION			
	FULL	NO	POD-Obj	POD-Ac	POND	FULL	NO	POD-Obj	POD-Ac	POND	FULL	NO	POD	POND
<i>wcd</i>	1	1.2	1.2	1.3	1.3	5.3	6.1	6.1	8.5	8.5	2.8	3.02	3.09	3.18
time(LS)	2.85	—	—	—	—	4.9	—	—	—	—	0.3	—	—	—
time(LE)	35.1	83.75	—	—	—	72.4	74.1	—	—	—	0.3	0.24	—	—
time(CD)	263.8	107.1	94.7	117.3	397.3	82	103.3	96.1	113.2	373.5	0.63	0.64	0.48	1.33
% CD	0.8	0.9	0.9	0.85	0.7	1.0	1.0	1.0	1.0	0.75	1.0	1.0	1.0	1.0

Table 1: *wcd* Values, Running Time, and Coverage Ratio

- ▶ Measure effect non-deterministic partially observable sensor models have on the *wcd* value of a model and the efficiency of *wcd* calculation using the compilation.
- ▶ For each setting we manually created 5 sensor models : Fully observable (FULL), Non observable actions (NO), two versions of Partially observable deterministic (POD) and Partially observable non-deterministic (POND)
- ▶ For all domains, *wcd* increases with the decrease of observability and increase of uncertainty

## Summary and future work

### We have :

- ▶ Extended Goal Recognition Design to handle **arbitrary sensor models**
- ▶ Allows us to find plans for **privacy** preserving agents

### We plan to :

- ▶ Handle partial knowledge **of the agent**
- ▶ Apply Goal Recognition Design to new applications (e.g. pentesting)



Code and benchmarks available on our website:  
<http://ie.technion.ac.il/~sarahn/grd>