

**Explicit ethical agents: Representing ethics explicitly.**

# **AI for Privacy: A Multiagent Perspective**

# Cybersecurity

Systems

Networks

Programs

Data

Privacy

“the claim of individuals, groups, or institutions to determine for themselves when, how, and to what extent information about them is communicated to others.”  
(...)

- Alan Westin

# Preserving Privacy in an Online World

How to represent the actual privacy preferences of users?

How to elicit the privacy preferences from users?

How to advise the users to take actions to preserve their privacy?

How to agree on how a co-owned content will be shared?

How to explain privacy decisions?



# Real Life Scenarios\*

OMG I HATE MY JOB!! My b  
always making me do shit stuff just to p  
Yesterday at 18:03 · Comment · Like

Hi [redacted], i guess y  
me on here?  
Firstly, don't flatter yourself. Secondly  
months and didn't work out that i'm ga  
around the office like a queen, but it's  
Thirdly, that 'shit stuff' is called your 'j  
pay you to do. But the fact that you s

**Claudya** [redacted]  
Vodka Shots ♥  
Like · Comment · 22 hours ago via mobile

**Terry** [redacted] Hold up aren't you babysitting?????  
22 hours ago · Like · 1

**Claudya** [redacted] Yes  
16 hours ago · Like

## Celebrities' Photos, Videos May Reveal Location

July 16, 2010

By KI MAE HEUSSNER

Keeping tabs on your favorite celebrities might be easier than you think -- and much easier than they want. But they likely have no one to blame but themselves.

According to two teams of computer scientists, Hollywood stars could be unintentionally giving up the exact locations of their homes and private whereabouts through pictures uploaded to the Internet, leaving them wide open to attacks by tech-savvy thieves (not to mention unwanted visits by starstruck fans).

Our online survey with 330 participants shows that more than 90% of privacy violations occur through inference.

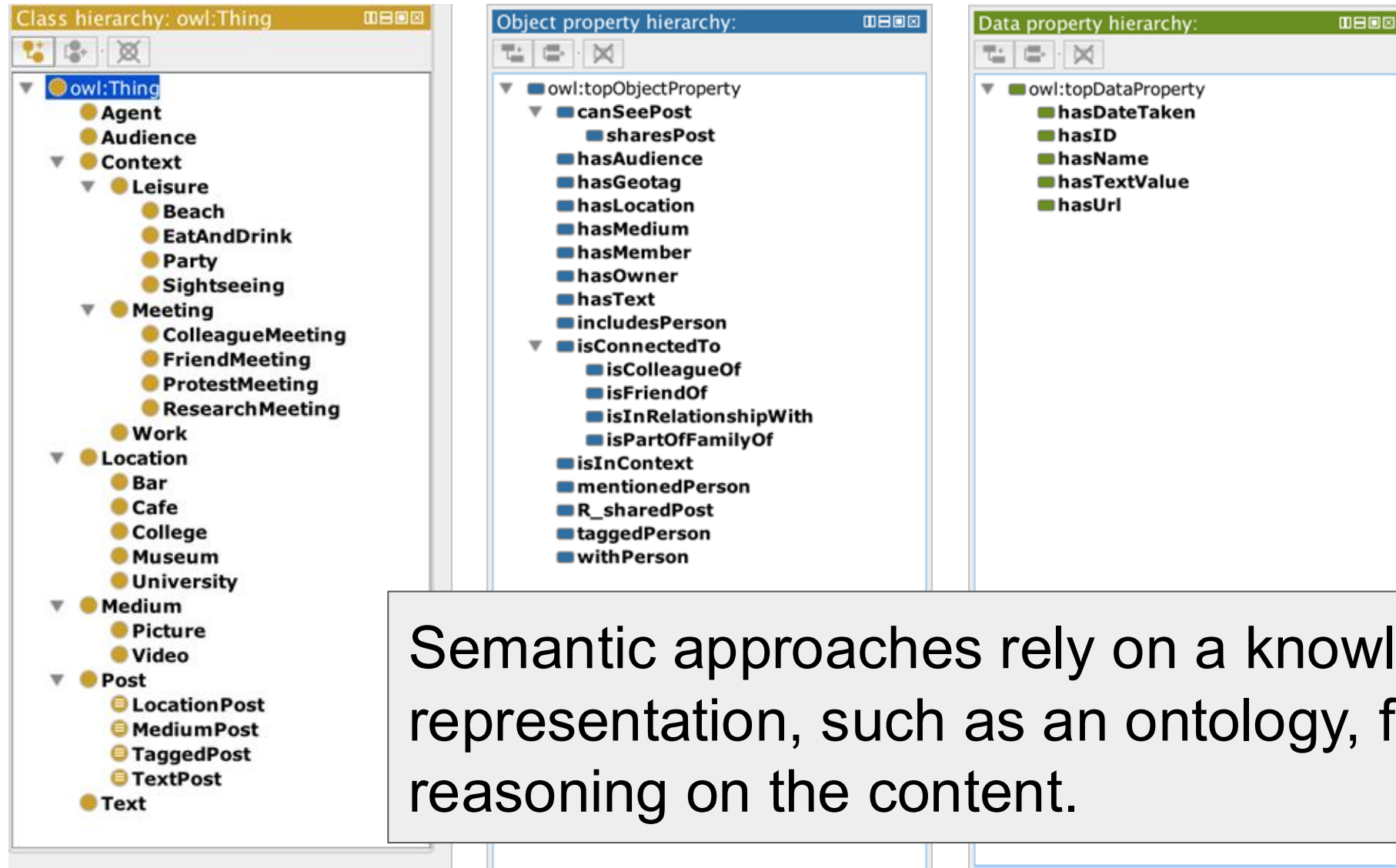
# Understanding privacy violations

## Privacy Concerns of Dennis

Dennis wants his friends to see his pictures but not his location.

	No inference	Inference
<b>User</b>	(i) Dennis checks in at a restaurant.	(iii) Dennis shares a picture without declaring his location. It turns out that his picture is geo-tagged.
<b>Others</b>	(ii) Charlie shares a picture with everyone. He tags Dennis in it as well.	(iv) Charlie checks in at a restaurant. At the same time, Dennis shares a picture of Charlie.

# Content Ontology



Semantic approaches rely on a knowledge representation, such as an ontology, for reasoning on the content.

# Privacy Preferences

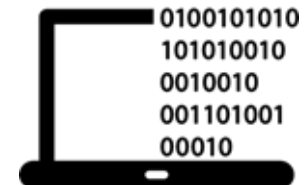


---

$P_{E_2}$ : *hasMedium*(?pr, ?m), *taggedPerson*(?m, :eve),  
*isInContext*(?pr, ?ctx), *Work*(?ctx)  $\rightarrow$  *rejects*(:eve, ?pr)  
[Eve rejects posts that are in work context.]

---

We can build software agents that can reason on users' privacy preferences.



# Detection of Privacy Violations: PriGuard

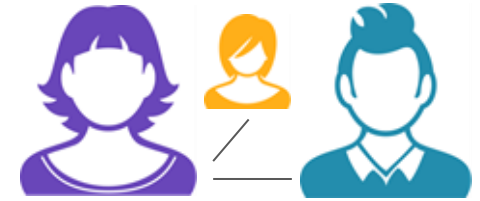


- We represent the social network as an *agent-based online social network* (ABSN).
- Agents know the privacy preferences of their users.
- We develop a sound and complete algorithm to detect privacy violations.
- We show the scalability of the approach on real-life social networks.

**PriGuard can detect privacy violations and notify the users to take actions.**



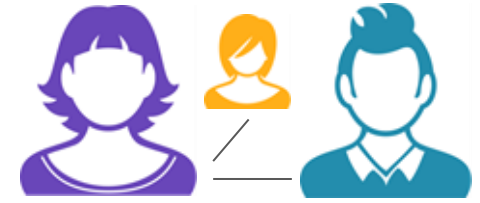
# Prevention of Privacy Violations: PriArg



- Agents discuss on a post *before* it is shared.
- We develop a framework that enables agents to carry out a dialogue with other agents.
- We adapt **computational argumentation** to enable privacy decision-making.

**Argumentation serves as a useful technique to mimic how humans deal with privacy disputes.**

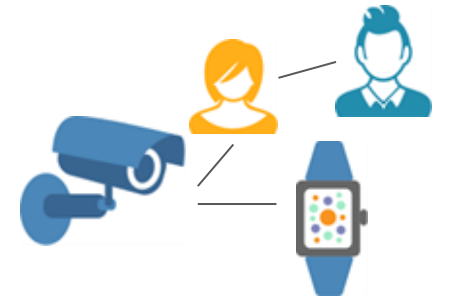
# Prevention of Privacy Violations: PriNego



- PriNego is a negotiation-based approach where agents negotiate with each other on their privacy preferences.
- Agents use different negotiation strategies to preserve their users' privacy.
- It exploits reciprocity as a heuristic (e.g., this time you help me, next time I help you).

**Agreement can be established over multiple posts.**

## Proactive Agents: an IoT example



- Each IoT entity follows **contextual norms** to calculate the appropriateness of sharing information.
- **Computational argumentation** enables the agent to reason on its knowledge and belief bases under **uncertainty**.
- To make inference based on others' information, a **trust model** needs to be in place.

**Agents can choose to violate privacy for a better outcome!**

# Preserving Privacy in an Online World

How to represent the actual privacy preferences of users?

How to elicit the privacy preferences from users?

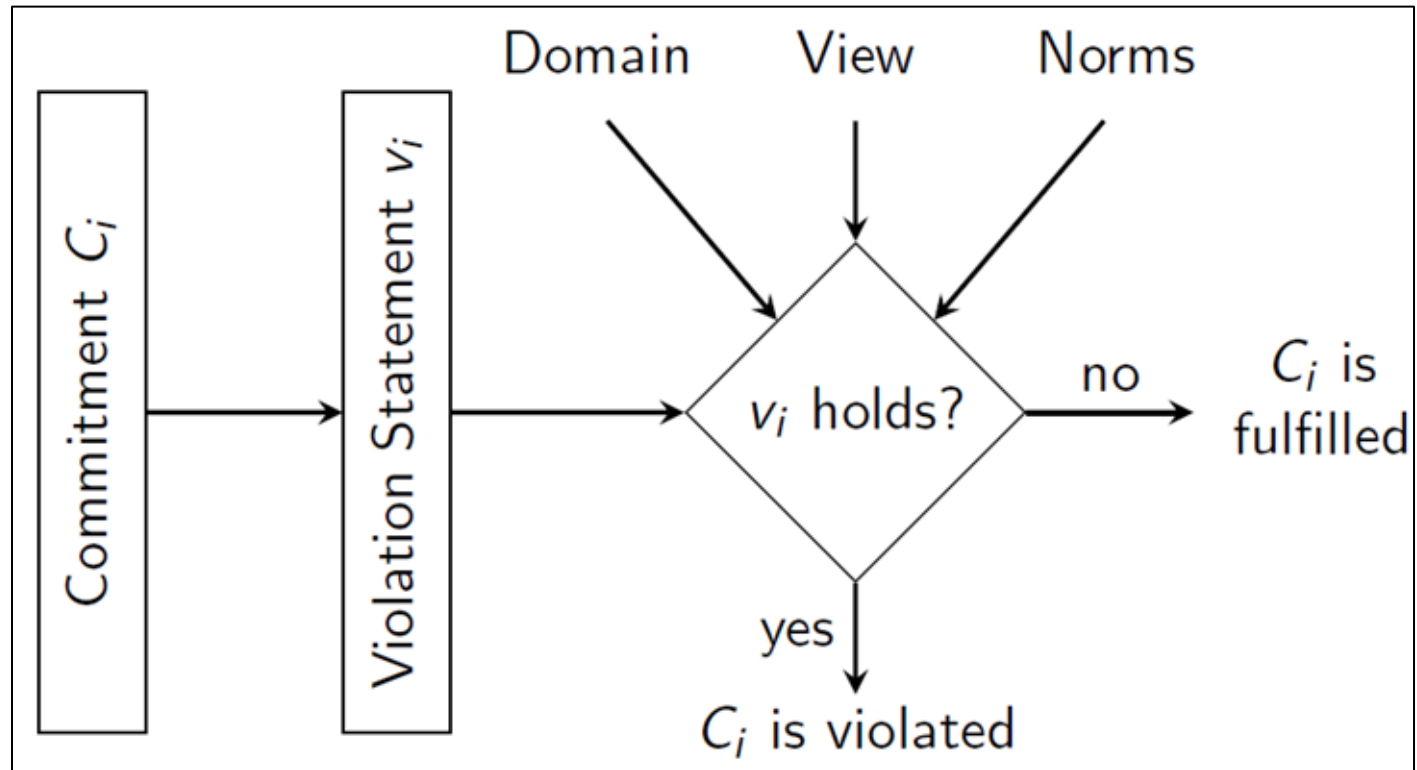
How to advise the users to take actions to preserve their privacy?

How to agree on how a co-owned content will be shared?

How to explain privacy decisions?



# PriGuard: Detection of Privacy Violations



# An Example

Dennis wants his friends to see his pictures but not his location. He posts a picture without declaring his location. However, it turns out that his picture is geotagged.

$C_1(:osn, :dennis, \text{isFriendOf}(:dennis, X), \text{isAbout}(P, :dennis), \text{LocationPost}(P), \text{not}(\text{canSeePost}(X, P)))$

$V_1 - :osn, :dennis, \text{isFriendOf}(:dennis, X), \text{isAbout}(P, :dennis), \text{LocationPost}(P), \text{canSeePost}(X, P))$

```
SELECT ?x ?p WHERE {  
  ?x osn:isFriendOf osn:dennis .  
  ?p osn:isAbout osn:dennis .  
  ?p rdf:type osn:LocationPost .  
  FILTER EXISTS (?x osn:canSeePost ?p) }
```

# The Social Network Domain

Agent, Post, Audience, Context, Content $\sqsubseteq \top$	Leisure, Meeting, Work $\sqsubseteq$ Context
Beach, EatAndDrink, Party, Sightseeing $\sqsubseteq$ Leisure	Bar, Cafe, College, Museum, University $\sqsubseteq$ Location
Picture, Video $\sqsubseteq$ Medium	Medium, Text, Location $\sqsubseteq$ Content
$\text{Post} \sqcap \exists \text{sharesPost}^-. \text{Agent} \equiv \exists R_{\text{sharedPost}}. \text{Self}$	$\text{LocationPost} \equiv \exists R_{\text{locationPost}}. \text{Self}$
$\text{LocationPost} \equiv \text{Post} \sqcap \exists \text{hasLocation}. \text{Location}$	$\text{MediumPost} \equiv \text{Post} \sqcap \exists \text{hasMedium}. \text{Medium}$
$\text{TaggedPost} \equiv \text{Post} \sqcap \exists \text{isAbout}. \text{Agent}$	$\text{TextPost} \equiv \text{Post} \sqcap \exists \text{hasText}. \text{Text}$

# Norms

---

$N_1:$   $sharesPost(X,P) \rightarrow canSeePost(X,P)$

[Agent can see the posts that it shares.]

---

$N_2:$   $sharesPost(X,P) \wedge hasAudience(P,A) \wedge hasMember(A,M) \rightarrow canSeePost(M,P)$

[Audience of a post can see the post.]

---

$N_3:$   $hasMedium(P,M) \wedge taggedPerson(M,X) \rightarrow isAbout(P,X)$

[Post is about agents tagged in a medium.]

---

$N_4:$   $Post(P) \wedge hasMedium(P,M) \wedge hasGeotag(M,T) \rightarrow LocationPost(P)$

[Geotagged medium gives away the location.]

---



# A Facebook Application: PriGuardTool

