# Tool-Supported Approaches to Ontology Engineering
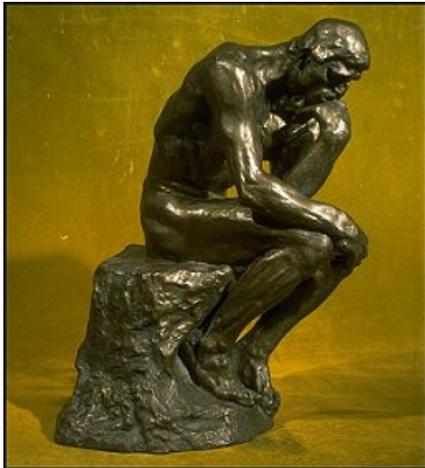
*Jon Atle Gulla*
*Professor of Information Systems*
*Norwegian University of Science and Technology, Trondheim*

Can we generate ontologies automatically?

Some examples of generated ontology structures

# Outline



♣
**Ontology Learning vs. Ontology engineering**
♣
**Principles of ontology learning**
♣
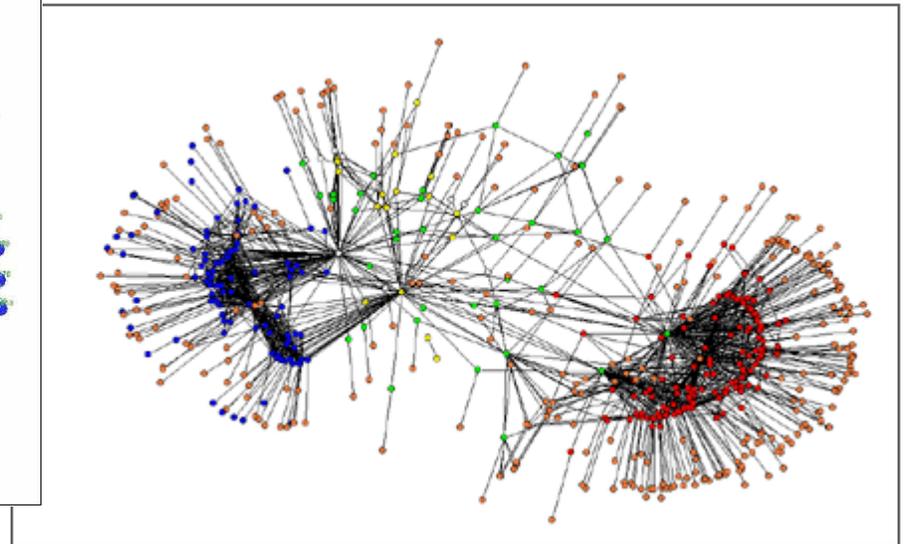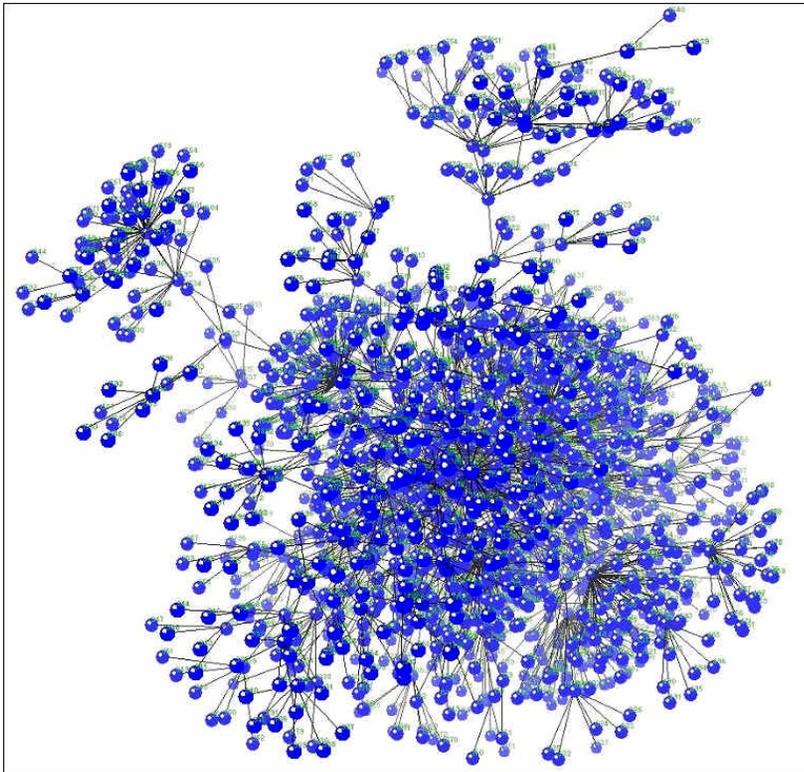**Ontology learning strategies**
♣
**Learning Classes, Individuals and Relationships in the movie domain**
♣
**Quality of ontology learning**
♣
**Conclusions**

# Ontology Engineering

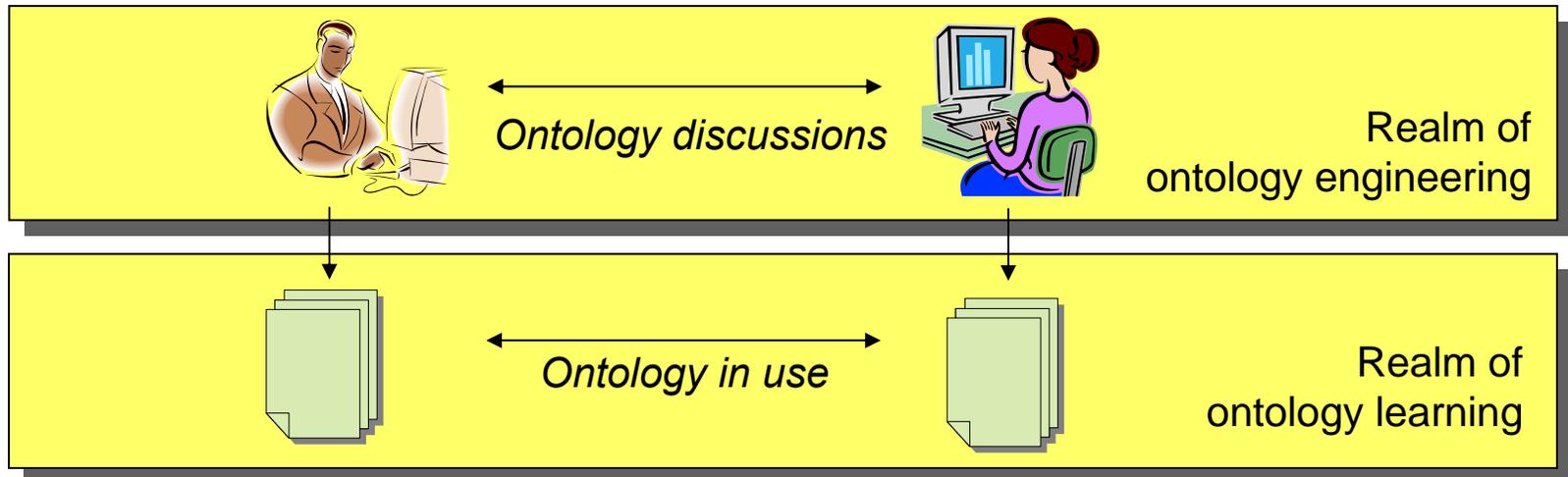- How to develop and maintain large complex ontologies?

# Ontology Modeling vs. Learning

- **Traditional ontology engineering approach**
  - *Project:*
    Form team of ontology and domain experts
  - *Ontology & domain experts:*
    Collaborative manual modeling process
  - *Domain experts:*
    Verify ontology against domain knowledge
  - *Ontology experts*:
    Verify ontology against syntactic and semantic quality measures

- **Expensive and time-consuming approach**

- **Ontology learning approach:**
  - *Domain experts:*
    Find representative domain text
  - *Tool:*
    Extract candidate classes, individuals and properties automatically from domain texts
  - *Ontology & domain experts:*
    Verify candidate structures and complete ontology

- **Can also be used to verify domain quality of existing ontology**
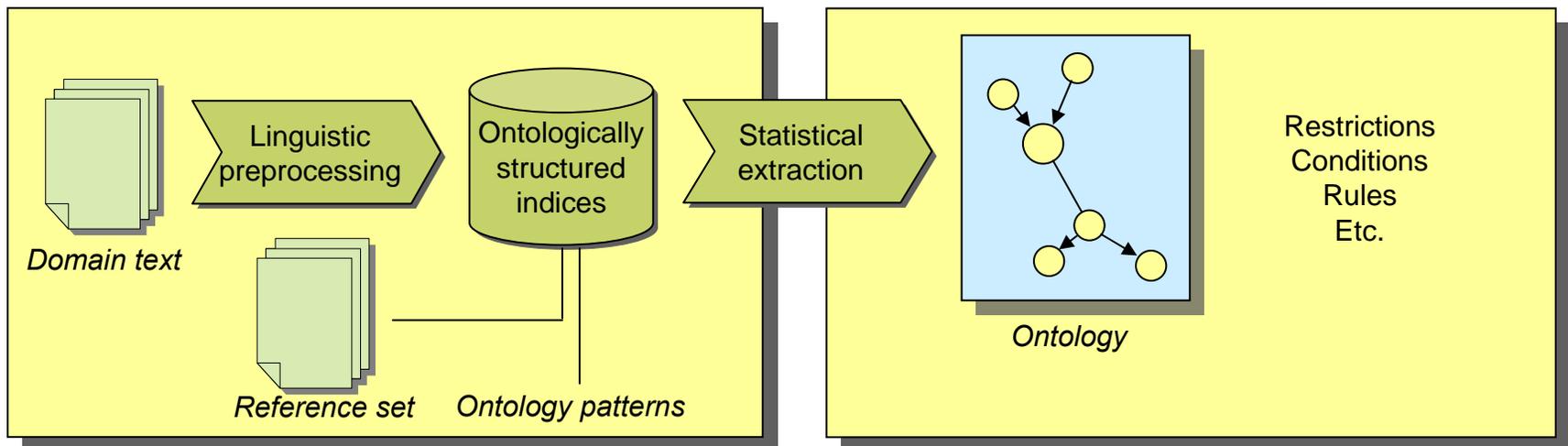- **Cost-effective approach**

# Ontology Learning Basis

- People communicate using domain-specific concepts
- People document using domain-specific concepts
- Ontology learning: *Extract ontology structures from written documentation*



- Requirements:
  - Documents representative for domain terminology
  - Documents cover all the terminology
  - Well-defined and consistent use of terminology in domain
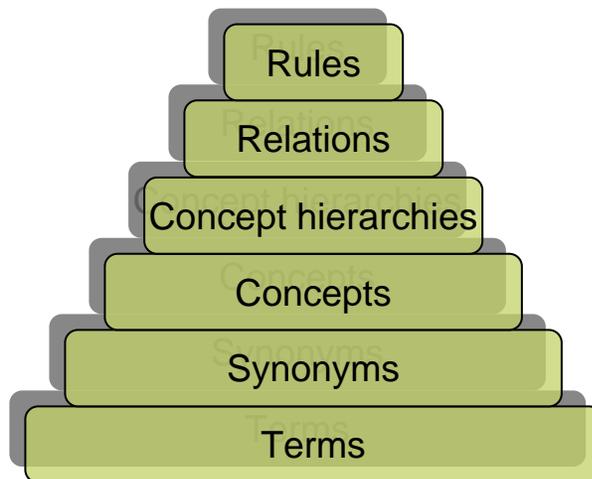
# Ontology Learning Process



Automatic extraction of ontology candidate structures

Manual verification of candidates and completion of ontology

# Levels of Ontology Learning

*Degree of difficulty*

Rules $\qquad\qquad\qquad$ $\forall x,y(manager(x,y) \rightarrow report(y,x))$

Relations $\qquad\qquad\quad$ *FINANCE(ag:SPONSOR, go: PROJECT)*

Concept hierarchies $\qquad$ *is_a(MANAGER, EMPLOYEE)*

Concepts $\qquad\qquad\quad$ *PROJECT*

Synonyms $\qquad\qquad\;$ *(leader, manager, lead)*

Terms $\qquad\qquad\qquad$ *sponsors, costs, charter*

# Ontology Learning Strategies

- Term extraction
  - *Linguistic analysis*
  - *Statistical analysis*
- Synonyms
  - *Classification-based techniques*
  - *Distribution-based techniques*
- Concept formation
  - *Structure recognition*
  - *Keyphrase generation*
  - *Instance learning*
- Concept hierarchy
  - *Clustering*
  - *Lexico-syntactic patterns*
  - *Head-modifier approaches*
  - *Subsumption approaches*
  - *Classification-based techniques*

- Relations
  - *Association rules*
  - *Concept vectors*
- Rules
  - *Structure recognition for meta-property recognition*
  - *Dependency trees and path similarities*

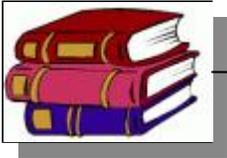# *Examples: Learning Classes, Individuals and Relationships*

## Core techniques for ontology learning

*Domain: Movie industry*

*Web data sources: IMDB, Videoload, Wikipedia, etc.*

*Resulting ontology: Semantic search application*

# Keyphrase Extraction for Learning Classes

Scope planning is the process of progressively elaborating and documenting the project work (project scope) that produces the product of the project.

*POS tagging*

Scope/NNP planning/NN is/VBZ the/DT process/NN of/IN progressively/RB elaborating/VBG and/CC documenting/VBG the/DT project/NN work/NN (/( project/NN scope/NN )/) that/WDT produces/VBZ the/DT product/NN of/IN the/DT project/NN ./.

*Stopword removal (571 words)*

Scope planning **is the** process **of** progressively elaborating **and** documenting **the** project work (project scope) **that** produces **the** product **of the** project.

*Lemmatization/stemming (POS tags not shown)*

Scope plan process progress elaborate document project work project scope produce product project

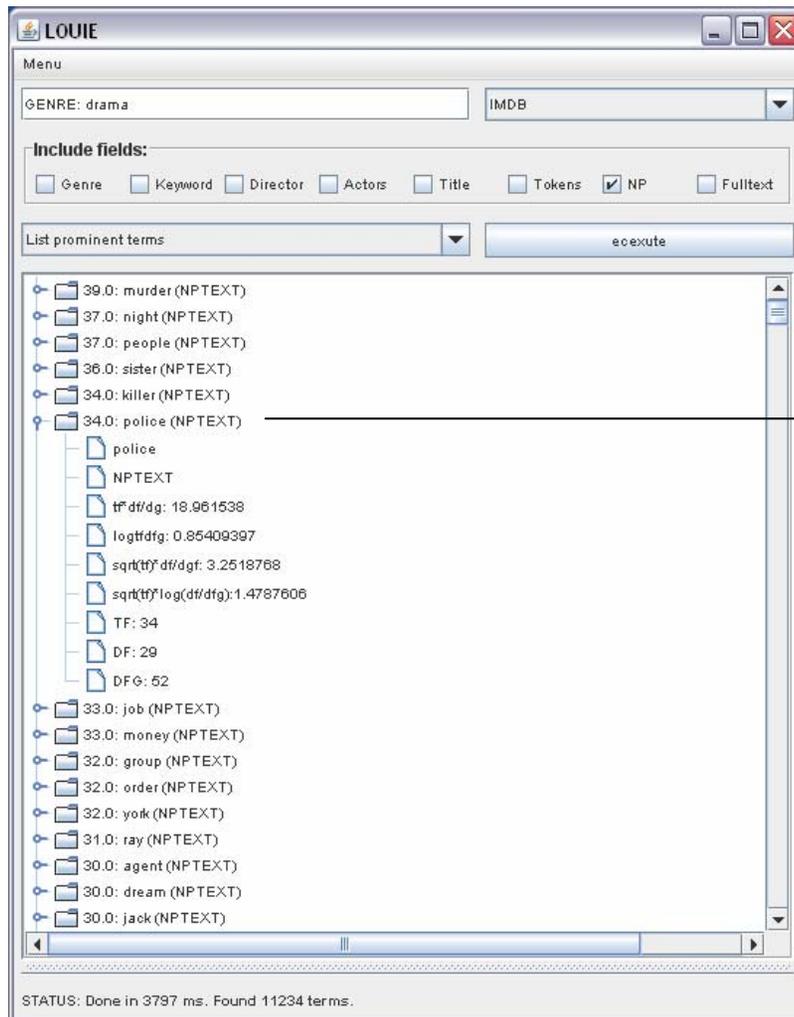*Select consecutive nouns as candidate phrases*

{scope planning, process, project work, project scope, product, project}
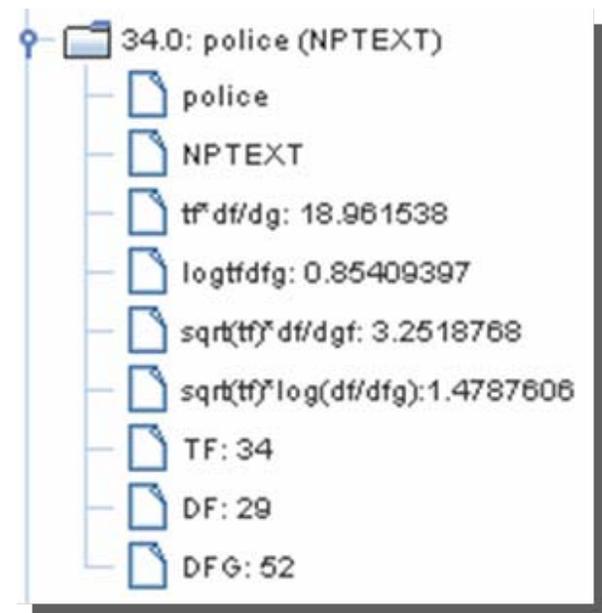
*Calculate tf.idf score for phrases*

{(scope planning, 0.0097), (project scope, 0.0047), (product, 0.0043), (project work, 0.0008), (project, 0.0001), (process, 0.0000)}

$$\text{tf} = \frac{n_i}{\sum_k n_k} \qquad \text{tfidf} = \text{tf} \cdot \log\left(\frac{|D|}{|(d_j \supset t_i)|}\right)$$
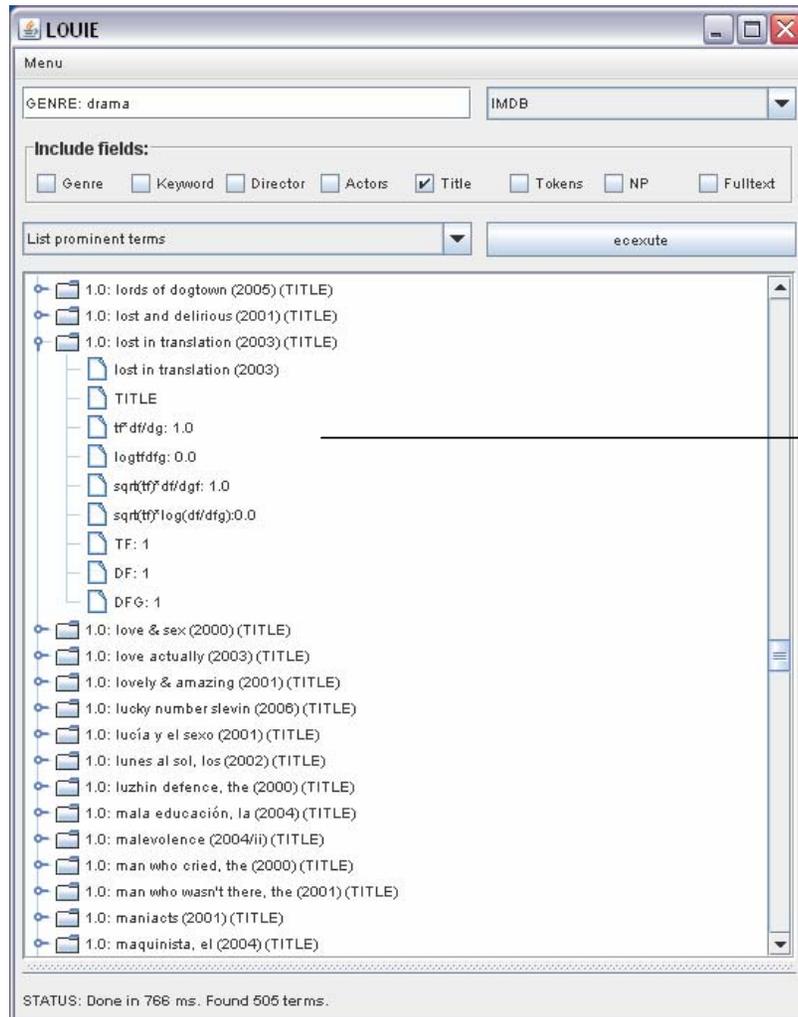
# Classes Relevant to the Drama Genre



- Keyphrase extraction technique
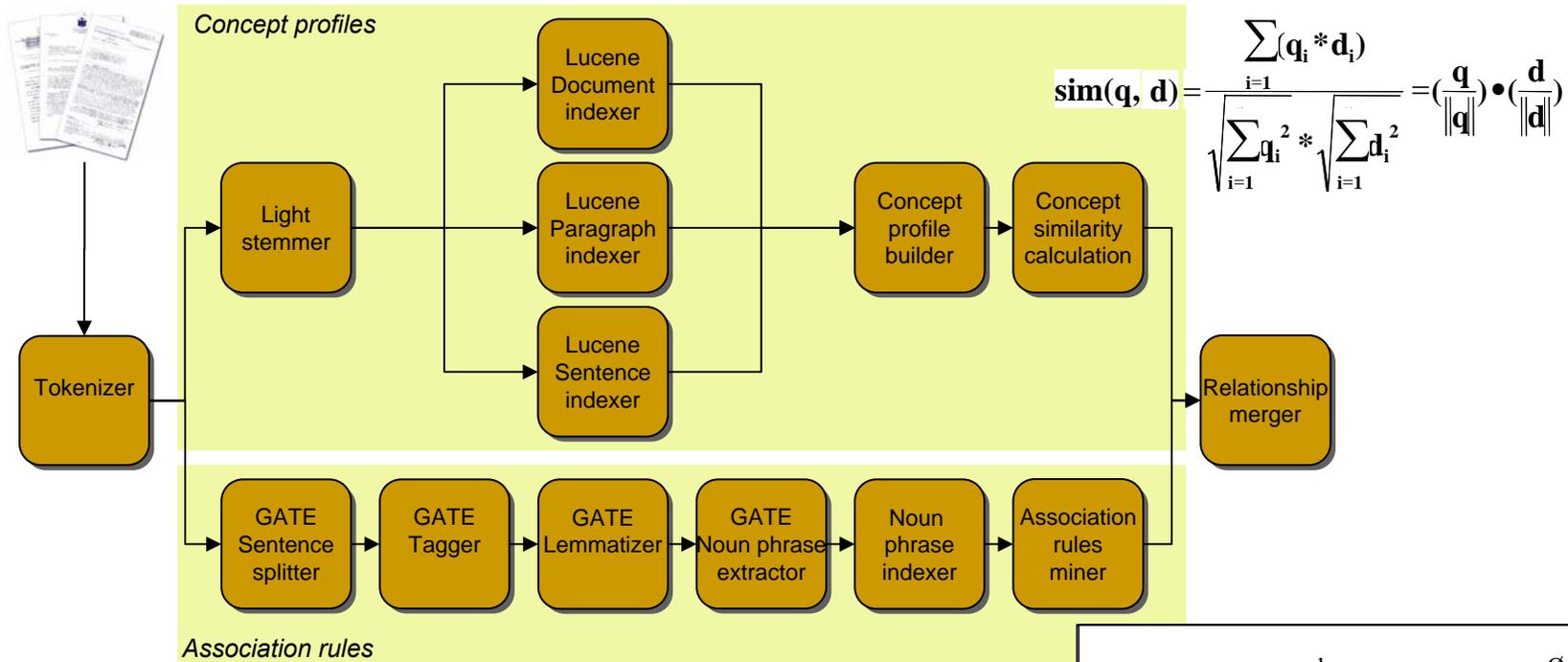- Noun phrases ranked according to various statistical measures

# Pattern Matching for Learning Individuals



- Using structural information (headings, keywords, etc.) to recognize movie instances
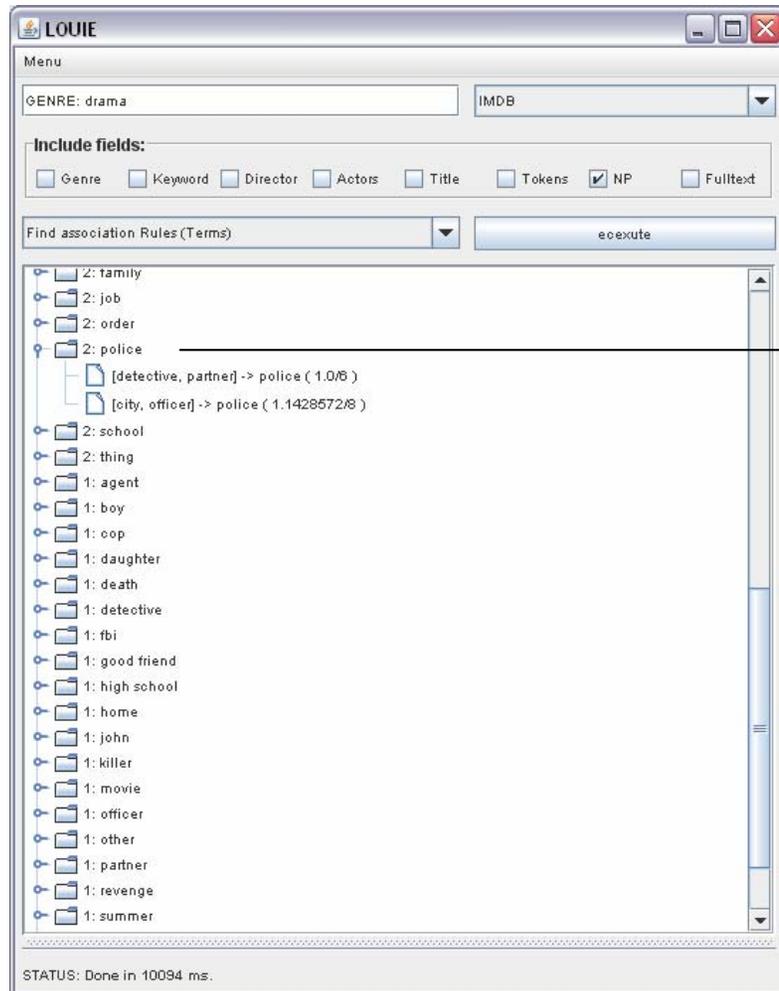- Instances ranked according to various statistical measures
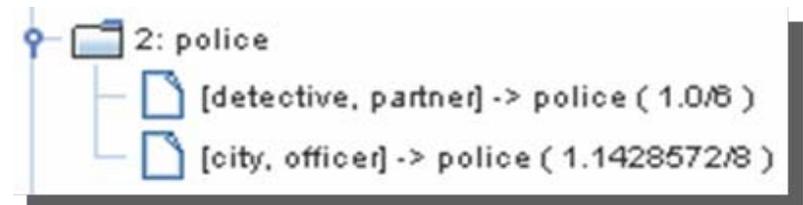
# Learning Relationships (Properties)



Concept profiles

Tokenizer → Light stemmer → Lucene Document indexer / Lucene Paragraph indexer / Lucene Sentence indexer → Concept profile builder → Concept similarity calculation → Relationship merger

Association rules

GATE Sentence splitter → GATE Tagger → GATE Lemmatizer → GATE Noun phrase extractor → Noun phrase indexer → Association rules miner

$$\mathbf{sim(q, d)} = \frac{\sum\limits_{i=1}(\mathbf{q_i * d_i})}{\sqrt{\sum\limits_{i=1}\mathbf{q_i}^2} * \sqrt{\sum\limits_{i=1}\mathbf{d_i}^2}} = (\frac{\mathbf{q}}{\|\mathbf{q}\|}) \bullet (\frac{\mathbf{d}}{\|\mathbf{d}\|})$$

$X \Rightarrow Y,$     where  $X \subset I, Y \subset I, X \cap Y = \varnothing$

A rule $X \Rightarrow Y$ holds in the transaction set D with *confidence c* if c% of the transactions in D that contain X also contain Y. The rule $X \Rightarrow Y$ has *support s* in the transaction set D if s% of the transactions in D contains $X \cup Y$.

# Learning Class Relationships



- ■ Association rules on extracted concepts

# Extract from Police OWL Declaration



```
– <rdf:RDF xml:base="http://www.owl-ontologies.com/
    <owl:Ontology rdf:about=""/>
    <owl:Class rdf:ID="Police"/>
    <owl:Class rdf:ID="Partner"/>
    <owl:Class rdf:ID="Detective"/>
    <owl:Class rdf:ID="Officer"/>
    <owl:Class rdf:ID="City"/>
  – <owl:SymmetricProperty rdf:ID="related">
      <rdf:type rdf:resource="http://www.w3.org/2002
    – <rdfs:range>
      – <owl:Class>
        – <owl:unionOf rdf:parseType="Collection">
            <owl:Class rdf:about="#Detective"/>
            <owl:Class rdf:about="#Partner"/>
            <owl:Class rdf:about="#City"/>
            <owl:Class rdf:about="#Officer"/>
          </owl:unionOf>
        </owl:Class>
      </rdfs:range>
      <rdfs:domain rdf:resource="#Detective"/>
    </owl:SymmetricProperty>
  </rdf:RDF>
```

# Learning Relationships between Movies



Movies related to ”Lost in Translation” and confirmed by both methods:

”Far from heaven” (2002)
”Kaho naa... Pyaar hai” (2000)

Can choose how techniques are to be combined

*Concept vector similarities*              *Association rules*

# Extract from OWL Generation

# Quality of Class Learning

**Evaluation Procedure**

- Extracted candidates from project management domain (PMBOK):
    - 50,600 tokens (ca. 130 pages)
    - Generated candidates for each area (chapter)
- Constructed ontology from candidates (with help from STATOIL employee)
- Built an alternative ontology manually (with help from another STATOIL employee)
- Compared quality of two ontologies for domain representation
- (Compared quality of two ontologies in ontology-driven (semantic) search)

# Results for Class Learning Evaluation

- Domain representation:



*Semi-automatically constructed ontology for project management*

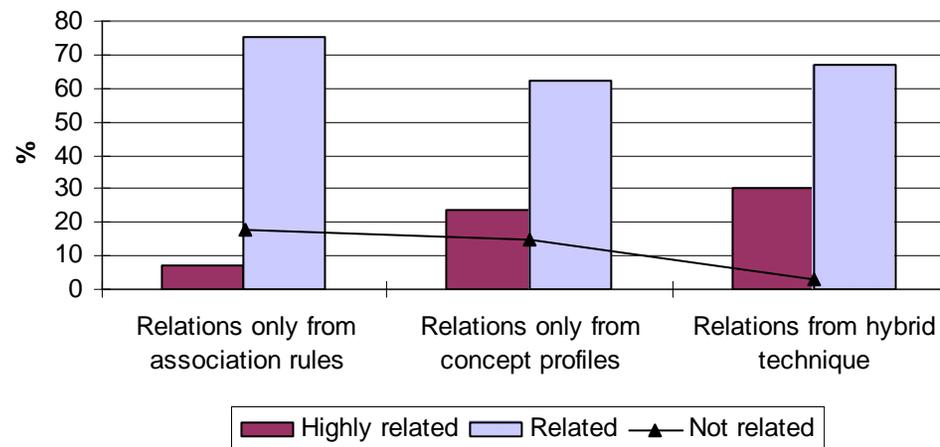|  | Classes | Hierarchical levels | Very good classes | Acceptable classes |
|---|---|---|---|---|
| Semi-automatic | 106 | 3 | 73 (79%) | 33 (21%) |
| Manual | 142 | 5 | 122 (86%) | 20 (14%) |

- ❑ 62 (58% of semi-automatic ontology) classes identical
- ❑ Tool-generated ontology:
  - Slightly smaller, with less abstraction levels
  - Almost as good as manually built ontology
  - Substantially faster to build
  - Easy to improve further

**Semi-automatic ontology construction very promising!**

# Quality of Relationship Learning

- Experiment with Statoil's project management standard (PMI)
  - Generated class relationships based on PMBOK
  - Quality of relationships verified by project management experts
  - Comparison between association rules and concept vector similarity

- Result of evaluation

# Conclusions

- Ontology Learning is the discipline of automatically or semi-automatically constructing ontologies

- Challenge to construct and maintain search ontologies

- Numerous learning strategies
    - Classes
    - Individuals
    - Relationships (properties)

- Ontology learning produces an intial fragmentary OWL model
    - Manual verification and correction
    - Manual completion of missing parts
    - But: Quality of techniques improving

- *Ontology learning a complement to traditional ontology engineering methodologies*