

Ontology-Based Information Extraction: A machine learning approach



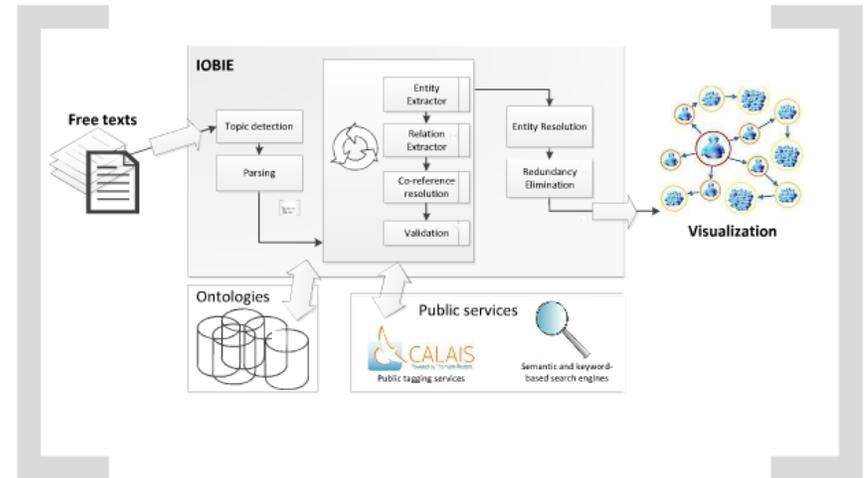
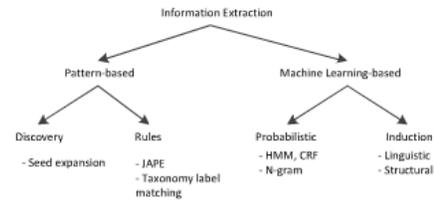
Slavko Žitnik
slavko.zitnik@fri.uni-lj.si

Ogilavo d.o.o.
Zupancičeva 8
5270 Ajdovščina

University of Ljubljana
Faculty of computer and information science
Tržaška cesta 25
1000 Ljubljana

Information extraction

- As a task: Filling slots in a database from text
- As a family of techniques:
segmentation+classification+association+clustering



Evaluation

- each component separately
- whole IOBIE system
 - MUC 1-7 (Message Understanding Conference)
 - RISE (Repository of Information Sources)
 - ACE (Automatic Content Extraction)
 - Enron Email Dataset
- Wikipedia with DBpedia

Ontology-Based Information Extraction: A machine learning approach



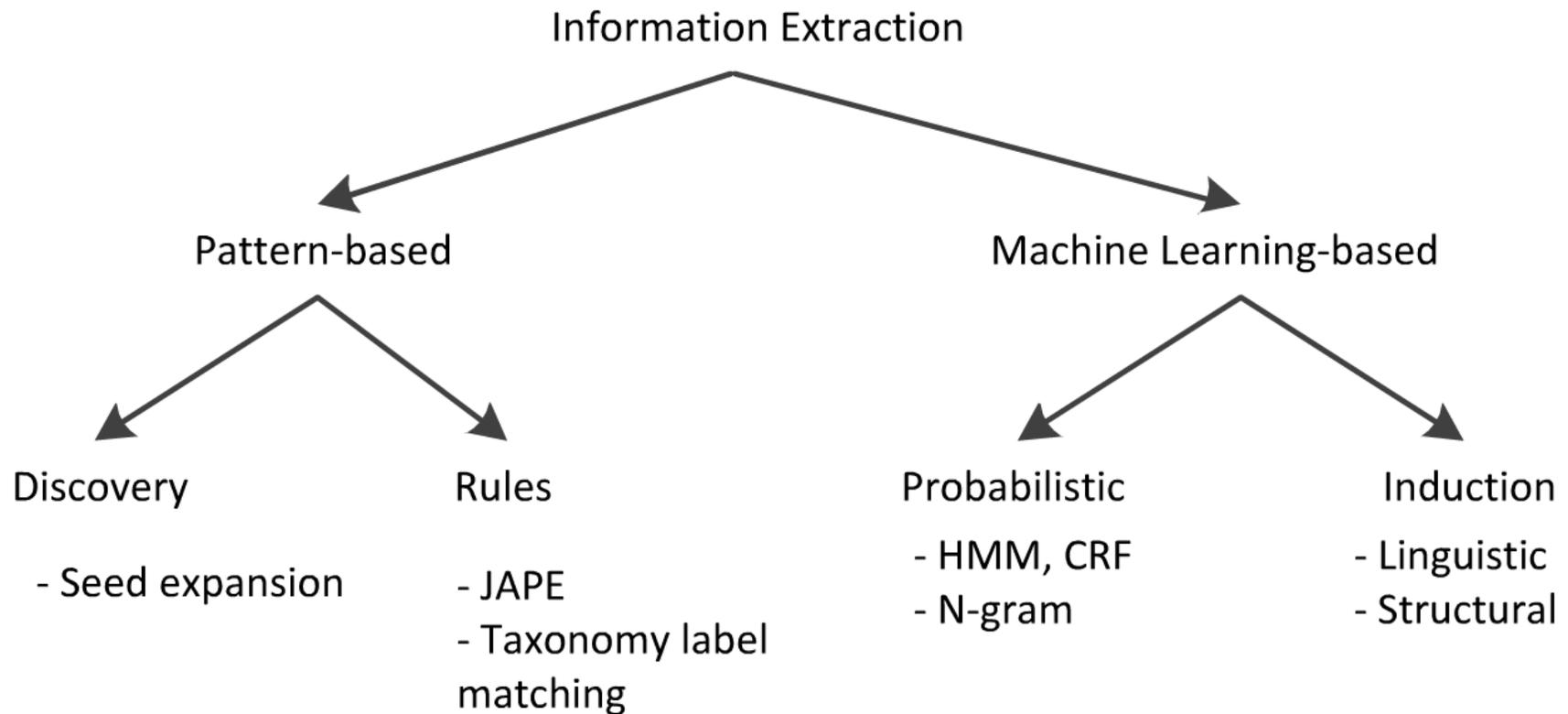
Optilab d.o.o.
Župančičeva 8
5270 Ajdovščina

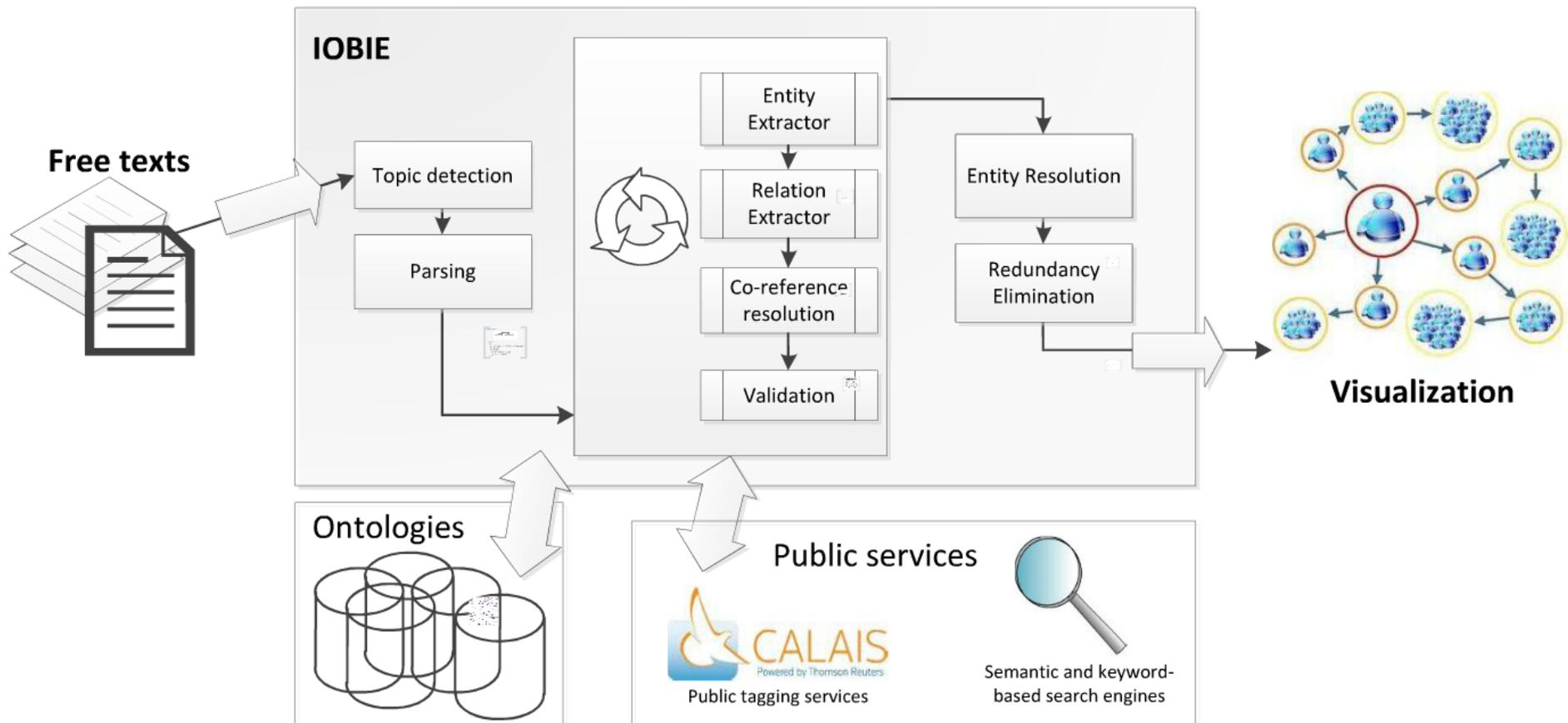
University of Ljubljana
Faculty of computer and information science
Tržaška cesta 25
1000 Ljubljana

Slavko Žitnik
slavko.zitnik@fri.uni-lj.si

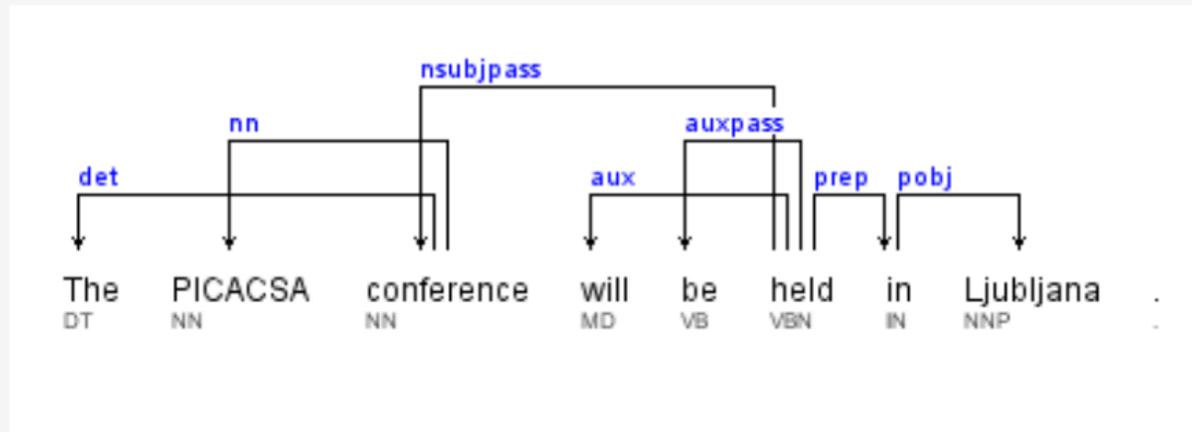
Information extraction

- As a task: Filling slots in a database from text
- As a family of techniques:
segmentation+classification+association+clustering





Stanford parser



(ROOT

(S

(NP (DT The) (NN PICACSA) (NN Conference))

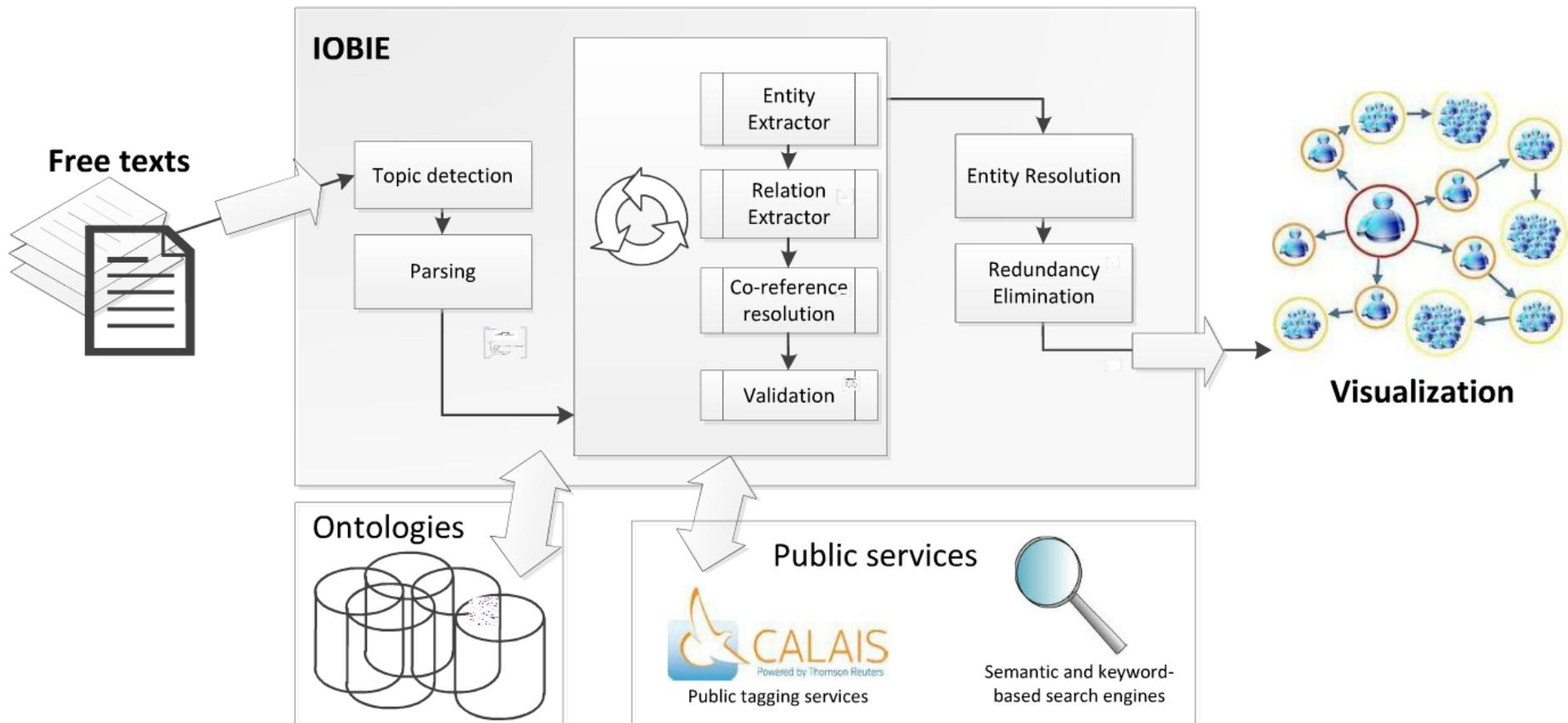
(VP (MD will)

(VP (VB be)

(VP (VBN held) (PP (IN in)

(NP (NNP Ljubljana))))))

(. .))

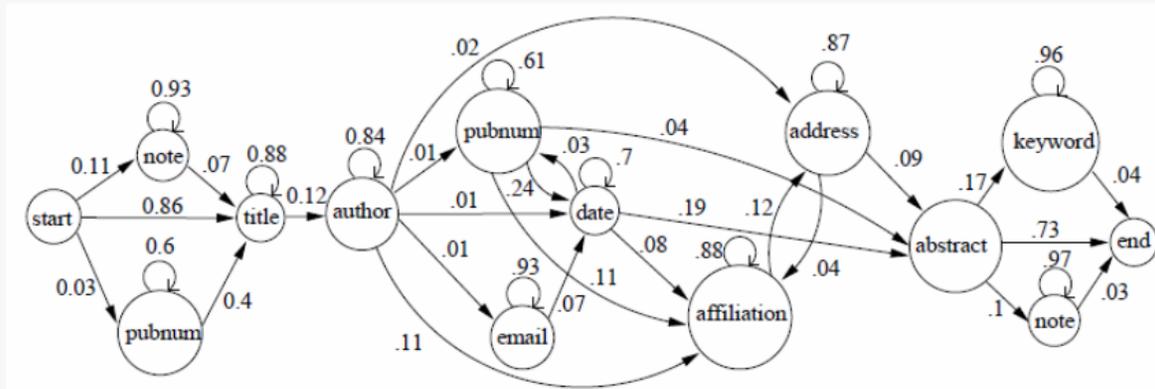


Entity Extraction

- Example

<faculty>FRI</faculty>,
located in
<location>Ljubljana</location>
will host
<conference>PICACSA</conference>
conference on
<date>June 3-4 2011</date>.

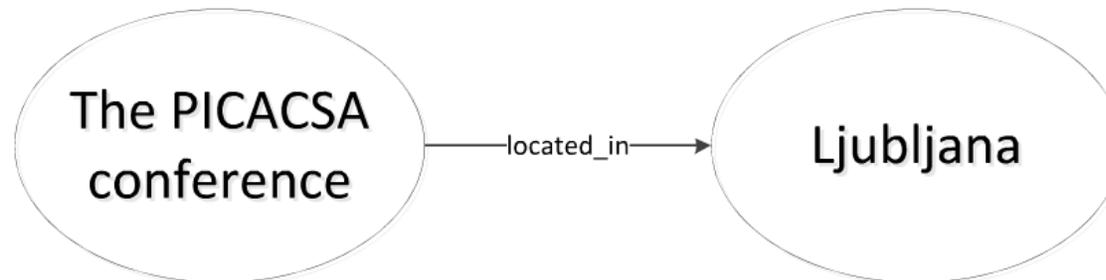
- Hidden Markov Models



Relation extraction

- Example

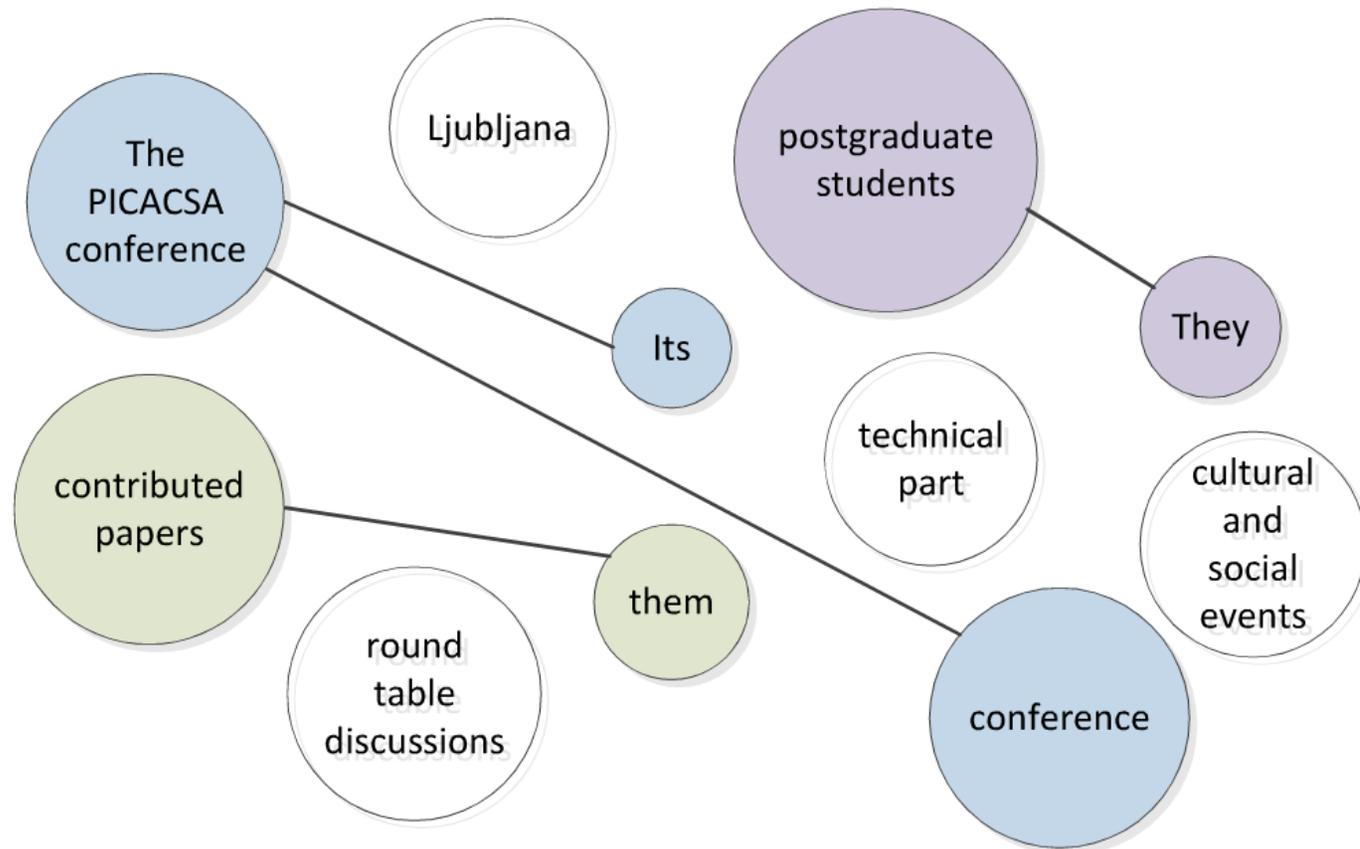
<conference>The PICACSA conference</conference>
will be held in
<location>Ljubljana.</location>



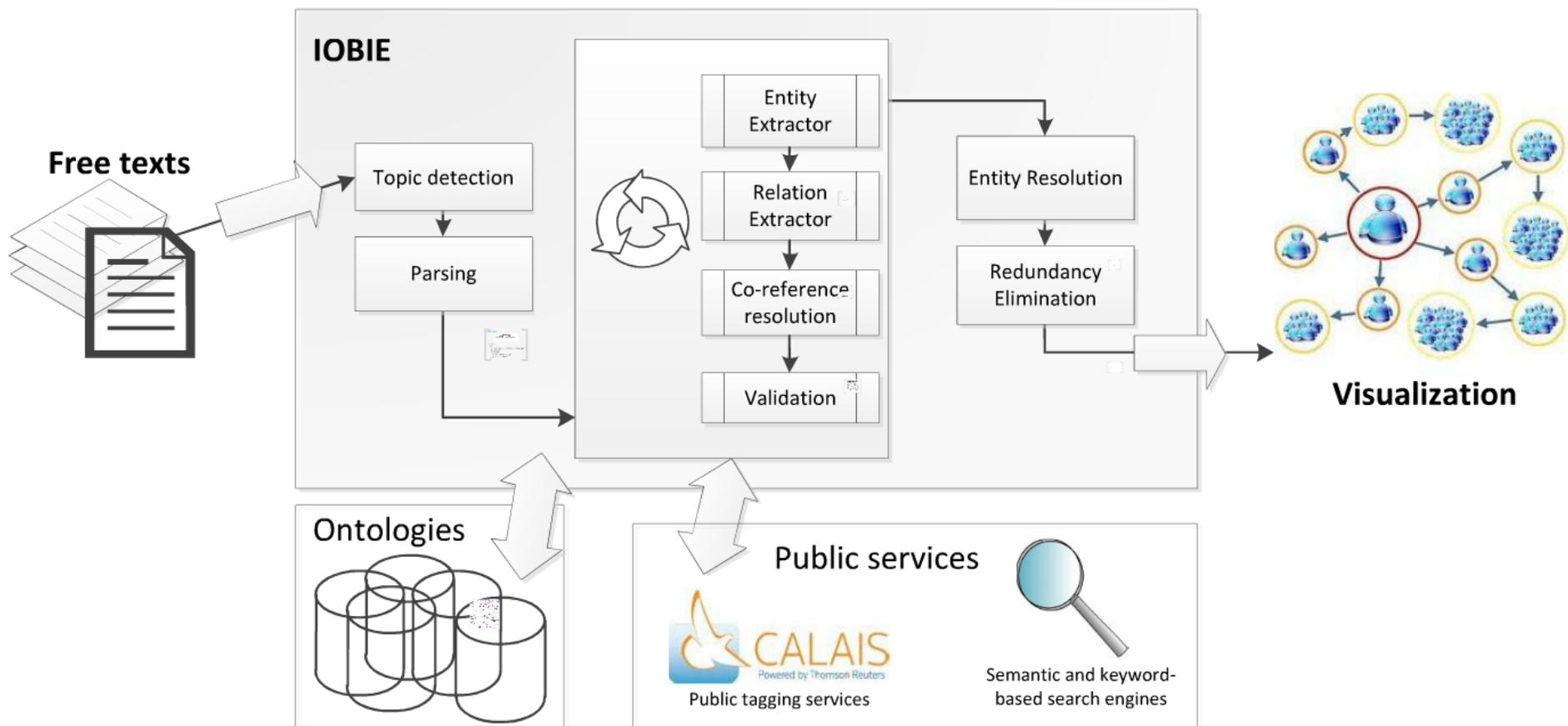
- Feature-based methods

Co-reference Resolution

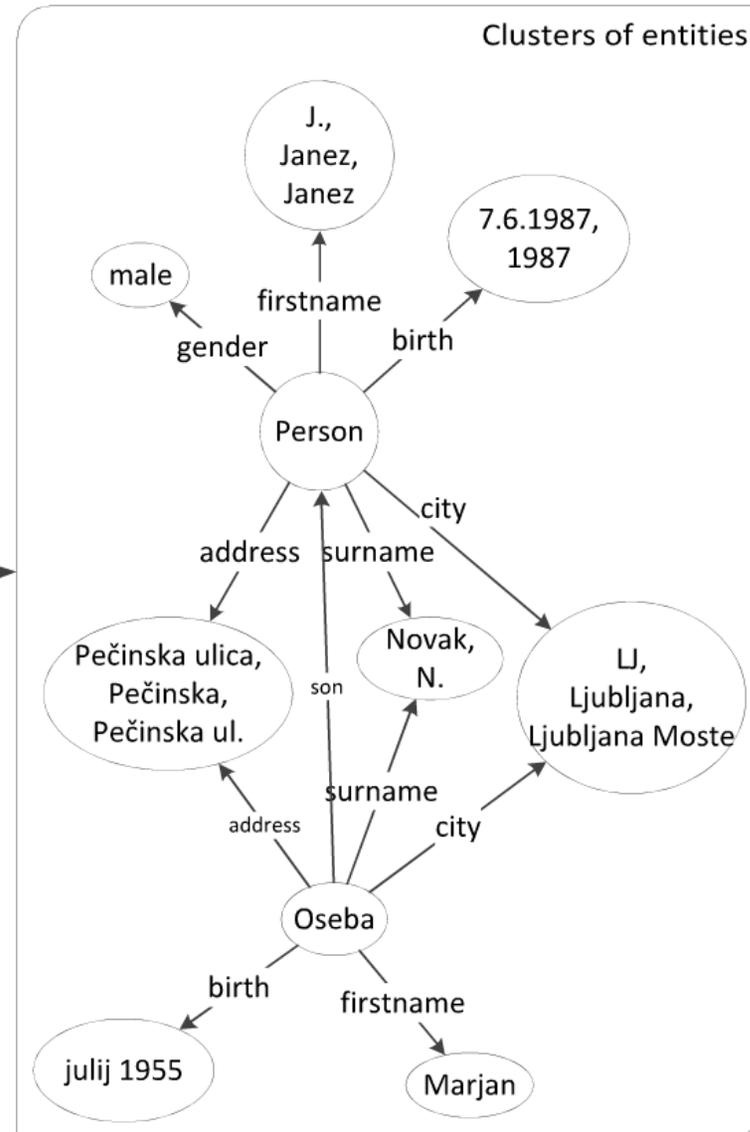
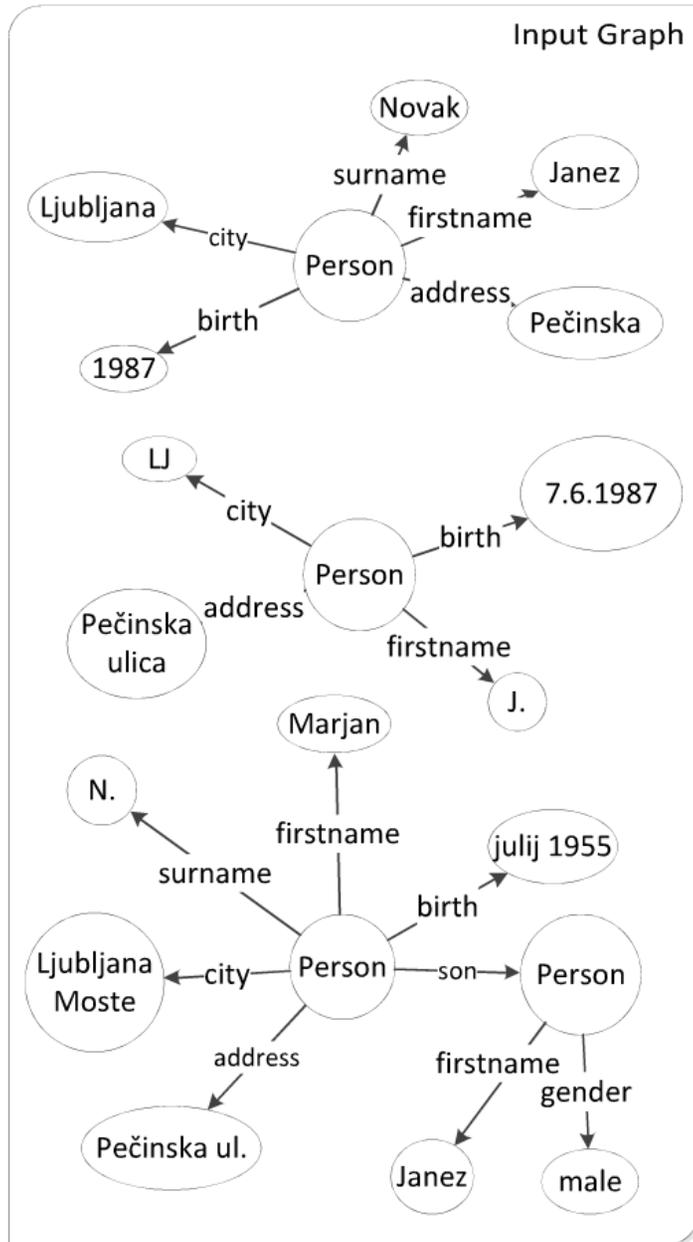
The PICACSA Conference¹ will be held in Ljubljana². Its¹ purpose is to bring together a panel of postgraduate students⁴. They⁴ will present contributed papers⁵ and have round table discussions⁶ about them⁵. The technical part⁷ of the conference¹ will be accompanied by cultural and social events⁸.



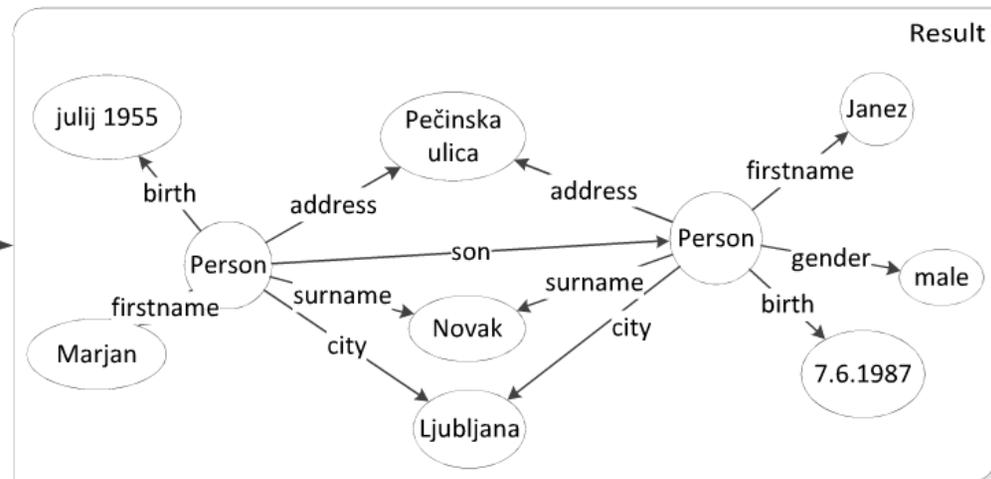
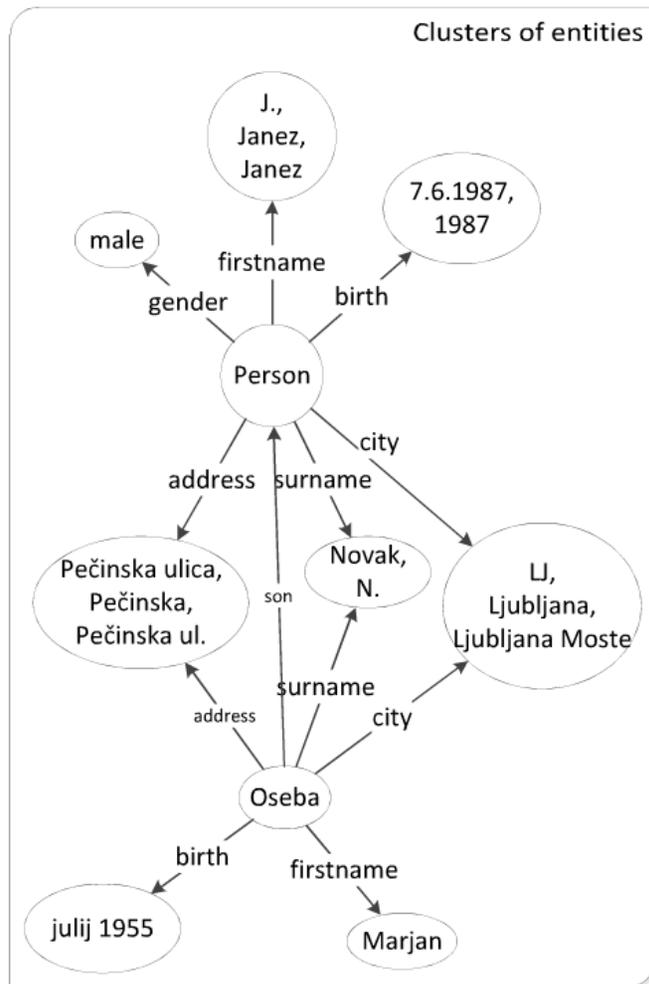
- Pairwise coreference model

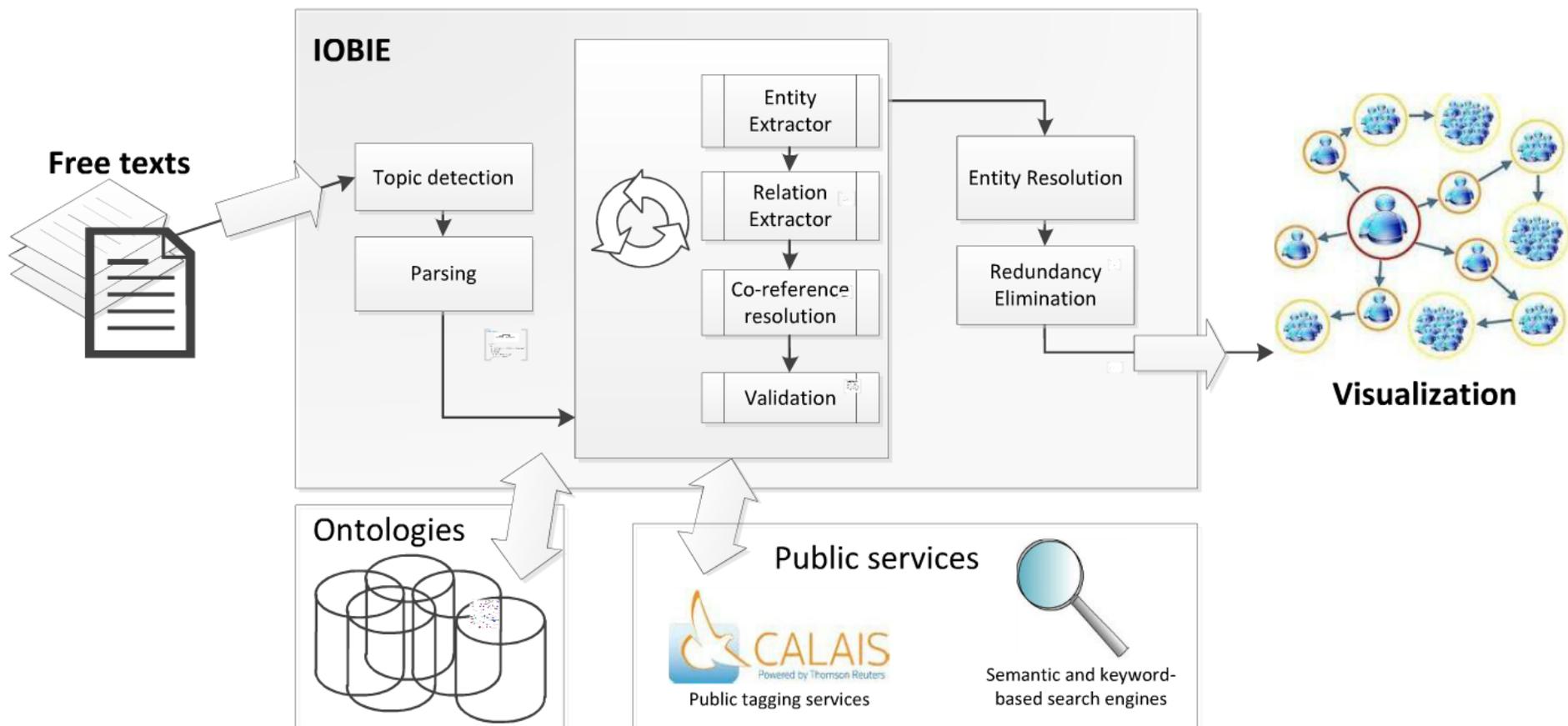


Entity Resolution



Redundancy elimination





Evaluation

- each component separately
- whole IOBIE system
 - MUC 1-7 (Message Understanding Conference)
 - RISE (Repository of Information Sources)
 - ACE (Automatic Content Extraction)
 - Enron Email Dataset
 - Wikipedia with DBPedia



Thanks