

# Open Information Extraction

Dominik Both, Tonio Weidler

Proseminar *Text Mining*  
Andrea Zielinski

Institut für Computerlinguistik, Universität Heidelberg, 15.07.2016

# Strukturierung

- 1 Introduction
- 2 OIE - Principles
- 3 Example: LODifier
- 4 OIE Systems in Context
- 5 Conclusion

# Introduction to Information Extraction

# What is Information Extraction?

Goal of Information Extraction is automatically extracting information from unseen text Information: entities, relations, events...

To make the dough for a good pizza, we start with putting 1kg of flour into the mixing bowl. (1kg of flour, put into, mixing bowl)

○  
○  
○  
○○○○○○○○

○○○○  
○○○○○○○○  
○○○  
○

○  
○

○  
○

# Problems of Information Extraction

- Named Entity Recognition
- Relationship Extraction
- Coreference Resolution
- Comment Extraction
- many more

# OIE - Principles

# OIE - Principles

# Open Information Extraction

# Open Information Extraction

IE: Extractor for each target relation    Open: No pre-specified extractors  
Unsupervised learning of relation phrases    Extraction of information on every given domain



# Problems of Open Information Extraction

- Incoherent extractions:
- This guide contains dead links and omits sites -> contains omits
- Uninformative extractions:
- Faust made a deal with the devil -> (Faust, made, a deal)

# OIE - Principles

## Methods

# Text Runner and WOE

- 1. Label: Automatic sentence labeling by heuristics
- 2. Learn: A relation phrase extractor is learned
- 3. Extract: Identifying NP pairs and searching relations words between

# Problems

- Large number of labeled training examples required
- Alternative heuristic labeling leads to huge noise and stacked uncertainty
- Ignores both holistic and lexical aspects

# Lexical constraint

# Syntactic constraint

# Limitations of those constraints

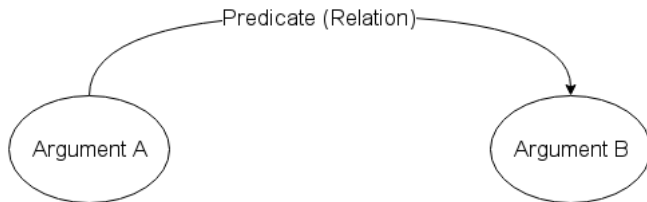
# ReVerb Extraction Algorithm



# OIE - Principles

## Data Representation

# Standard Patterns



**Argument A** is in a directed relation to **Argument B**.

# Unnormalized Annotation

(argument\_a, predicate\_x, argument\_b)  
(argument\_a, predicate\_y, argument\_c)  
(argument\_a, predicate\_y, argument\_d)

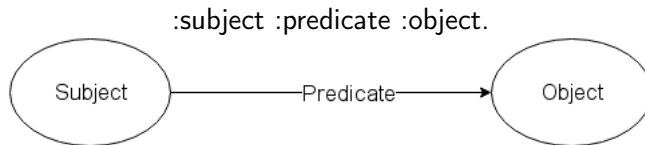
## Problems

- redundant
- unnormalized
- can only produce binary predicates

# RDF and Linked Data

## Resource Description Framework

Models propositions by constructing *triples* including **Subjects**, **Objects** and **Predicates**  
Generates a directed graph



# RDF Concepts and Notation

## ■ URIs

identifies ressources (S, R, O) distinctivly and references further informations (triples)

## ■ Conclusions

allows to draw conclusions using rules

## ■ Turtle

allows syntax abbreviations

## ■ Queries

can be searched by querying (eg SPARQL)

# Basic relations (built in)

Relation	Functionality
rdf:type (a)	x is of type y
owl:sameAs	x equals y

# RDF Syntax

```
dbr:Barack_Obama a foaf:person, :President;  
    dbo:spouse dbr:Michelle_Obama.  
dbr:Bernie_Sanders dbo:birthPlace dbr:New_York,  
    dbr:Brooklyn;  
dbr:Brooklyn dbo:isPartOf dbr:New_York
```

○  
○  
○  
○○○○○○○●

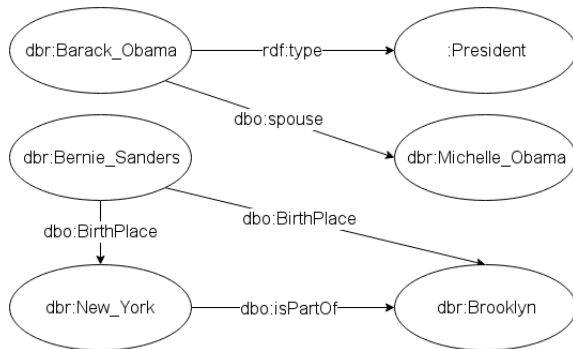
○○○○  
○○○○○○○○  
○○○  
○

○  
○

○  
○

## Data Representation

## ... as Graph





○  
○  
○  
○○○○○○○○

○○○○  
○○○○○○○○  
○○○  
○

○  
○

○  
○

# Example: LODifier



# LODifier: Generating Linked Data from Unstructured Text (Augenstein et al., 2012)

Generate an RDF Graph from unstructured Text

**Past Approaches:** Use Patterns to trade recall for precision

**LODifier:** Process the entire text

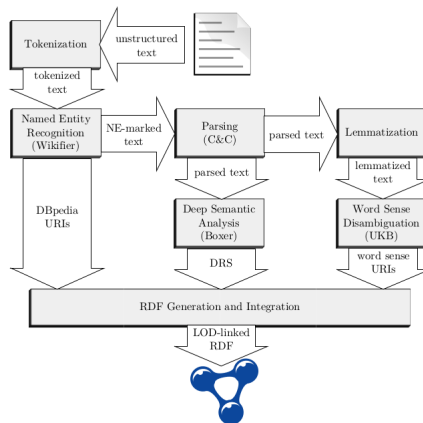


## Example: LODifier

# Architecture

## Architecture

## Architecture





# Approach

- 1 **Parse** the input text (POS, Treetagging, NER)
- 2 Apply **Deep Semantic Analysis** to get relations
- 3 Enrich NEs and words with **URIs** (DBpedia and WordNet)
- 4 Forge an **RDF Graph** of this information



# How does it happen?

Lets go through the process step-by-step!

## Example Text:

The New York Times reported that John McCarthy died. He invented the programming language LISP.

example taken from Augenstein et al., 2012

# Example: LODifier

## Preprocessing



# Named Entity Recognition - Wikifier

## Wikifier

Recognizes NE and replaces them with the Wikipedia Page Link  
Disambiguates by comparing links between pages.

### Example Text Output:

[The New York Times] reported that [John McCarthy (computer scientist)|John McCarthy] died. He invented the [Programming language|programming language] [Lisp (programming language)|LISP].





# Parsing Syntax - C&C

## C&C Parser

Syntactical Parser that tags POS and builds Parse Trees (CCG).



# Parsing - Output

```
ccg(1, rp(s:dcl,
  ba(s:dcl,
    lx(np, n,
      t(n, 'The_New_York_Times', 'The_New_York_Times', 'NNS', 'I-NP', '0')),
    fa(s:dcl\np,
      t((s:dcl\np)/s:em, 'reported', 'report', 'VBD', 'I-VP', '0'),
      fa(s:em,
        t(s:em/s:dcl, 'that', 'that', 'IN', 'I-SBAR', '0'),
        ba(s:dcl,
          lx(np, n,
            t(n, 'John_McCarthy', 'John_McCarthy', 'NNP', 'I-NP', 'I-PER')),
            t(s:dcl\np, 'died', 'die', 'VBD', 'I-VP', '0')))),
      t(period, '.', '.', '.', '0', '0'))).
ccg(2, rp(s:dcl,
  ba(s:dcl,
    t(np, 'He', 'he', 'PRP', 'I-NP', '0'),
    fa(s:dcl\np,
      t((s:dcl\np)/np, 'invented', 'invent', 'VBD', 'I-VP', '0'),
      fa(np:nb,
        t(np:nb/n, 'the', 'the', 'DT', 'I-NP', '0'),
        fa(n,
          t(n/n, 'programming_language', 'programming_language', 'NN', 'I-NP', '0'),
          t(n, 'LISP', 'LISP', 'NNP', 'I-NP', '0')))),
      t(period, '.', '.', '.', '0', '0'))).
```



# Find Relations - Boxer

## Boxer

Creates DRSs from C&C Output



# Find Relations - Boxer

## Boxer

Creates DRSs from C&C Output

## Discours Representation Structure (DRS)

Represents the discourse via *relations* between *entities*

Allows referencing over the entire discourse

# Find Relations - Boxer

## Boxer

Creates DRSs from C&C Output

## Discours Representation Structure (DRS)

Represents the discourse via *relations* between *entities*

Allows referencing over the entire discourse

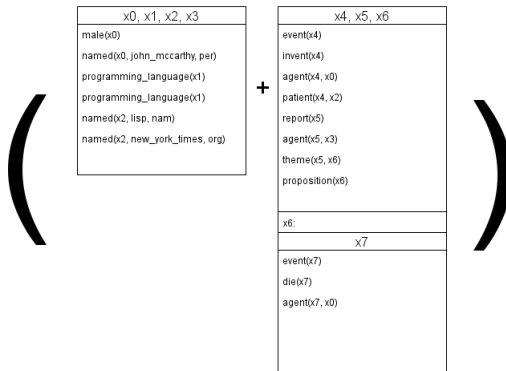
## Boxers DRS Relations (Conditions):

- **Unary Relations (Classes):** eg. *topic*, *person*, *event*, *male*, ...  
+ all verbs
- **Binary Relations:** agent, patient, ... (semantic roles)



## Preprocessing

## Boxer Output





# Assign WordNet URIs

## RDF WordNet

**WN:** Lexicography containing senses linked by semantic relations

**RDF WN:** LD Representation of WN providing URIs for words

### Steps:

- 1 Lemmatization
- 2 WSD (UKB)
- 3 Assign RDF WN URIs to word senses



# Preprocessing Result

**We now have ...**

- URIs for all NEs
- URIs for all (disambiguated) words
- Relations between entities (those URIs)



# Example: LODifier

## RDF Construction



# What now?

Let's now construct the RDF Graph from this information!



# Namespaces/Vocabularies

LODifier creates several namespaces:

- drsclass:
- class:
- drsrel:
- ne:
- reify:

And uses standard namespaces:

- rdf:
- owl:

As well as the two ontologies:

- wn30:
- dbpedia:

# Example: LODifier

## Conclusions

○  
○  
○  
○○○○○○○○

○○○○  
○○○○○○○○  
○○○  
○

○  
○

○  
○

# OIE Systems in Context

# OIE Systems in Context

## Comparison



# OIE Systems in Context

## Evaluating the Approaches

# Conclusion



○  
○  
○  
○○○○○○○○

○○○○  
○○○○○○○○  
○○○  
○

○  
○

●  
○

# Conclusion

## Problems and Obstacles

# Conclusion

# Future Opportunities