# **Automated Text Analysis for eHumanities**

#### Silvana Hartmann













#### **UKP Lab: Snapshot of January 2011**





- I Junior Prof.
- 3 junior research group leaders
- 2 PostDocs
- 14 PhD students







12.02.2011 | Computer Science Department | UKP Lab - Prof. Dr. Iryna Gurevych | Silvana Hartmann | 3



### **UKP Lab: Software and Data**



 Darmstadt Knowledge Processing (DKPro) Repository

- Wikipedia API
- Wiktionary API
- OmegaWiki API



Wiktionary ['wɪkʃənrɪ] n., a wiki-based Open Content dictionary



http://www.ukp.tu-darmstadt.de/software/





12.02.2011 | Computer Science Department | UKP Lab - Prof. Dr. Iryna Gurevych | Silvana Hartmann | 5

는 UKP

# UKP Lab: Language Technology for eHumanities









#### 12.02.2011 | Computer Science Department | UKP Lab - Prof. Dr. Iryna Gurevych | Silvana Hartmann | 7

# DARIAH-DE

- DARIAH-EU network
- Digital Research Infrastructures for the Arts and Humanities
- Sister project: CLARIN-EU





Digital Research Infrastructure

for the Arts and Humanities

Bundesministerium für Bildung und Forschung



### **DARIAH-DE@UKP/TU** Darmstadt

#### WP 1 E-Infrastructure

- Develop services for arts/humanities
  - User-oriented demonstrators
  - interactive community platforms
- Focus on interoperability of data and services

#### WP 2 Research and Education

- Guide paradigm shift to information technology based research methods
- Develop
  - Concepts for curricula for the eHumanities
  - Discipline-specific virtual research environments







### LOEWE research area "Digital Humanities"



🔆 LOEWE

Structure of research area "Digital Humanities"









### Language Technology (LT) for eHumanities







### LT for eHumanities: Tools and Methods

#### Tools: DKPro component repository

- Java-based toolkit using the Apache UIMA framework
- Process "flawed text", e.g. not orthographically normalized text

#### Method: Mining Multiword Expressions from Wikipedia

- Multiword Expressions (MWEs):
  - Significant part of natural language
  - Many types are productive
  - Important for corpus linguistics and natural language processing of modern language
  - Relation to formulaic expressions







# DKPro I

Darmstadt Knowledge Processing Repository

#### Highlights

- Conventiently compose natural language processing (NLP) pipelines
- Reusable software components for NLP
- General-purpose components for preprocessing
- Special-purpose components for specific domains and problems
- Based on Apache UIMA (Unstructured Information Management Architecture)



#### Applications at UKP

- Semantic text similarity
- Key-phrase extraction
- Opinion mining
- Information Retrieval
- Self organizing Wikis

http://www.ukp.tu-darmstadt.de/research/projects/dkpro/



# DKPro II

Darmstadt Knowledge Processing Repository

#### DKPro Core

- Linguistic Preprocessing
  - Tokenization, Compound splitting, POS-tagging, Parsing
- Support for various file formats and data sources
  - XML, PDF, WSDL, Wikipedia
- Available as open source: <u>http://dkpro-core-asl.googlecode.com</u>

#### DKPro IR

- Text indexing and retrieval
- Support for classical and semantic methods

#### DKPro Semantics

 Keyphrase extraction, text segmentation, summarization, semantic relatedness, link discovery





# DKPro III

Darmstadt Knowledge Processing Repository

#### Use-case: user generated discourse

- NLP tools choke on
  - Emoticons
  - Shorthand
  - Spelling errors
  - Spurious formatting/markup
  - Incorrect encoding conversion artifacts
  - ...

#### Data Cleansing

- Spell Checking, Spelling Correction
- Dictionary Annotator
- Regular Expression Annotator/Stemmer
- Apply Changes Annotator







# **DKPro IV**







- Removed emoticons
- Expanded shorthand
- Corrected spelling errors

Analogy: Similar "flaws" in manuscripts



### LT for eHumanities



Application Scenario: Mining Multiword Expressions from Wikipedia



### LT for eHumanities: Multiword Expressions



- Example: Mining English Multiword Expressions (MWEs) from Wikipedia
  - MWEs: "Idiosyncratic interpretations that cross word boundaries" (Sag et al. 2002)
  - Statistically, syntactically or semantically irregular

#### Motivation

- MWEs should not by split by parsers, during indexing
- MWEs occur frequently, are productive
- Application in Information Retrieval, Question Answering, Machine Translation, Ontology Creation,...

#### Approach

- Corpus-based approach
  - Identification of candidates
  - Ranking (statistical association metrics)
  - Selection
- Exploit structural properties of Wikipedia text ("wisdom of the crowds")



### **Mining Multiword Expressions**



- MWEs: "Idiosyncratic interpretations that cross word boundaries"
- Statistically, syntactically or semantically irregular
- Example classification:





### **Mining Multiword Expressions**





WIKIPEDIA Die freie Enzyklopädie

- MWEs: "Idiosyncratic interpretations that cross word boundaries"
- Statistically, syntactically or semantically irregular
- From Wikipedia (WikiMwe):





### MWE from Wikipedia (WikiMwe)



TECHNISCHE UNIVERSITÄT DARMSTADT

- Candidate sources:
  - Marked-up sequences
  - Meaningful units
  - Benefit from user effort





# JWPL





- High performance access to Wikipedia content
- Parser for the WikiMedia syntax
- Articles, discussion pages and categories as Java objects
- Access to information nuggets
  - · Redirects, links, l
  - ink anchors, interlanguage links,
  - sections, first paragraph, etc.
- Supports all Wikipedia language editions

Iow resource languages

transfer from languages with better resources

Available open source: http://jwpl.googlecode.com

#### Reference

Torsten Zesch, Christof Müller and Iryna Gurevych (2008).

Extracting Lexical Semantic Knowledge from Wikipedia and Wiktionary. In: Proceedings of the

Conference on Language Resources and Evaluation (LREC'08), electronic proceedings.



TECHNISCHE

DARMSTADT

# JWPL

Java-based Wikipedia Programming Library

#### JWPL TimeMachine

- Tool for reconstructing past states of Wikipedia
  - Define a point in time and revert Wikipedia to that date
  - Makes reproduction of Wikipedia-based research results easy

#### JWPL RevisionMachine

- Tool for easy access to article revisions
- Iterate over all changes to any article in Wikipedia
- Dedicated diff format for efficient storage of revision data
  - Complete uncompressed dump of a current English Wikipedia: ~5600 GB
  - Complete JWPL database of English Wikipedia with revisions: ~ 66 GB
  - study text development
  - study collaborative writing processes, interaction

#### Soon to be available open source





### WikiMwe Candidate Extraction



• MWE candidate extraction pipeline:



- Segmentation
- POS-tagging
- Named Entity Tagging
- Markup Annotation
- Output: candidate list (> 5 million candidates, 1.6 from titles, 4.3 from markup)
  - Later restrictions on candidates consisting of 2 to 4 constituent words



### WikiMwe Candidate Ranking and Selection



- Extraction of frequencies for candidates from Wikipedia text corpus
- Candidate ranking using statistical association metrics
- Candidate selection based on cutoff value
- Additional Filtering:
  - Wikipedia-typical patterns ("List of Countries")
  - Separation of Named Entities (NEs) based on
    - Stanford NE-tags
    - POS-sequences
- Numbers:
  - > 500,000 WikiMwe-NEs
  - > 350,000 WikiMwes



#### WikiMwe Manual Evaluation



UKP

- Annotation study on a random sample of 2500 WikiMwes
- 2 Raters (linguist background), disambiguation to Gold Standard by expert
- Distribution of classes in the Gold Standard:
  - > 55% "true MWES"



### WikiMwe Examples



#### Noncompositional MWE

- Upper triangular matrix
- Oxygen mask
- Choir screen

#### Collocation

- Low-cost airline
- Burning sensations
- First-person narrator

#### Regular Phrase

- Collaboration with Israel
- Protest against the war
- Camera models

#### Named Entities

- American Screenwriters Association
- Breslau Seminary
- Treatise on rhetoric



### WikiMwe Evaluation: Comparison



#### WikiMwe resource:

- ca 350,000 (noisy) MWEs of size 2-4
- Estimate based on annotated sample: >190,000 "true MWEs"
- Other MWE resources:
  - MWE data sets: Small (often manually crafted) evaluation sets
    - Noun Compound by C. Bannard: ~ 500
    - Noun Compounds by Tratz & Hovy 2010: 17,000
  - Specialized dictionaries: idioms, etc.
  - WordNet: >68,000 MWEs, 63,000 noun compounds of size 2-4
    - covered > 70% by WikiMwe
  - Wiktionary: > 40,000 MWEs

(...and this is for English..)



### Summary WikiMwe



#### Conclusion

- Identification of multiword expressions is still a big problem for language processing
- We leverage Wikipedia annotations to help with this task
- Result:
  - WikiMwe resource: large, noisy set of MWEs and NEs
  - Mainly nouns
  - Strong in technical terminology
  - Domain-transcending

#### Potential applications for WikiMwe:

- IR/QA (few harm by noise)
- Ontology creation (technical terminology)

#### Open questions:

- How to disambiguate MWEs (type vs token): plain dress, red carpet, don't drink the water?
- Wikipedia as a corpus for collocation extraction (compared to newswire)
- Transfer to other languages



### LT for eHumanities: Summary



#### Flexible Tools:

- DKPro processing pipeline
- Data cleansing
- Additional analysis components as needed
- Example application: Mining MWEs from Wikipedia
  - Exploit structural properties of available data
  - Efficiently process textual data
  - Apply to languages with few resources

Analogy: User generated "flawed text" and historical documents
Future goal: Tools adapted to eHumanities (DKPro eHum)



#### Thank you for your attention!





http://www.ukp.tu-darmstadt.de

