

A tool for Corpus Analysis using partial Disambiguation and Bootstrapping of the Lexicon

Kurt Eberle, Ulrich Heid, Manuel Kountz, Kerstin Eckart

SFB-732, B3, Universität Stuttgart
Institut für maschinelle Sprachverarbeitung
– Computerlinguistik –
Azenbergstr. 12
D 70174 Stuttgart

KONVENS 2008, Berlin

Overview

- Motivation and Objectives
- Objectives – State of the Art – Requirements
- Sample data and descriptive problems
- Tool Architecture – sample Representations
- Tool Functions and Applications
- Conclusions and future work

Motivation and Objectives

Ambiguity and data extraction from corpora

- Observation: structural ambiguity
Er löst schwierige Aufgaben und Probleme
 - Task: extract verb+object pairs:
 - * *Aufgabe + lösen*: ok
 - * *Problem + lösen*: ok
 - ⇒ No problem for extraction
 - Task: extract adjective+noun pairs:
 - * *schwierig + Aufgabe*: ok
 - * *schwierig + Problem*: ??
 - ⇒ Structural ambiguity with impact on data extraction
- Objective 1:
Ambiguity awareness in corpus exploration:
not all ambiguities affect extraction

Motivation and Objectives

Data extraction: from targeted phenomenon to extracted data (1/2)

- Trivial observation:

Data extraction involves several steps:

- Knowledge of phenomenon to be extracted
- Query formulation, application to (preprocessed) text
- Query result, including local analysis
- Interpretation

- Even on preprocessed corpora:

Some phenomena are hard to express in terms of queries

- Objective 2:

Make relation visible and reproducible between:

- corpus sentence
- analysis of preprocessing + query
- parameters used for disambiguation

⇒ Integrated representation

Motivation and Objectives

Data extraction: from targeted phenomenon to extracted data (2/2)

- Observation:

To test descriptive and extraction hypotheses:

- Need to find analogous cases in corpus, with respect to
 - * analysis of preprocessing
 - * query
 - * parameters used for disambiguation

- Objective 3:

Integrated representation

should be fully searchable:

all aspects separately – together – any partial combination

Motivation and Objectives

Using data extraction to test descriptive hypotheses

- Observation:
To test descriptive hypotheses, often
a few contextual parameters are modified,
all others being kept the same
- Objective 4:
Keep track of variable depth of analysis
(cf. Kay/Gawron/Norvig 1994):
possibility to vary the degree of abstraction of analysis representations

Motivation and Objectives

Bootstrapping in corpus-based data extraction

- Observation:
Data extraction is similar to hypothesis testing:
improvement through results obtained in intermediate steps:
search → evaluation of results → refined search → better results
- Objective 5:
Enable bootstrapping of data extraction
by making intermediate results
available to the search tool

Objectives vs. the state of the art

Could we use existing tools and resources?

Objectives: reminder	Corpus tools and resources	Parsing
1 Ambiguity awareness		Deep parsing: formal grammars
2 Integrated representation: sentence/analysis/ disamb. parameters	treebanks, parsed corpora	
3 Querable integrated representation	treebank query	Search in parse forests
4 Variable depth	multilayered corpora	grammar
5 Bootstrapping	machine learning?	Feedback to parser

Objectives vs. the state of the art

Using existing technology (1/2)

Corpus tools and resources

- Treebanks are typically disambiguated
- Size of treebanks: rather small
- Multilayered corpora:
Possible in principle – tools and resources: prototypes
- Machine learning:
Parameters not clear enough in many cases

Objectives vs. the state of the art

Using existing technology (2/2)

Parsing

- Deep parsing: appropriate;
We prefer dependency-style grammars,
where variable depth of analysis
is easier to realize
- Search in parse forests is hard;
We prefer underspecified representations
- Feedback to parser:
via flexible lexicon interface

From objectives to requirements

Outline for an architecture

Integrated tool for analysis and data extraction:

- (1) ambiguity awareness:
dependency-based parsing
and underspecified representation of analysis results
- (2) integrated representation:
sentence – analysis/analyses – disambiguation parameters
- (3) integrated representation should be searchable
- (4) variable depth of analysis:
underspecification and partial disambiguation
- (5) bootstrapping:
via flexible interface to the lexicon:
easy way to insert extraction results into lexicon

Sample data – descriptive problems

Sortal ambiguity of nominalizations as an extraction problem

German *-ung*-nominalizations are sortally ambiguous:

- *Teilung*: event (e) or result state (s)
- *Messung*: event (e) or object (data,o)
- *Abdeckung*: event (e), state (s) or object (o)

Sortal disambiguation depends on ...

- ... verb semantics (Roßdeutscher et al. 2007): e.g.
 - deadjectival, property-denoting: e,s (*trocknen*, *lähmeln*)
 - effected object: e,o (*sammeln*, *bilden*)
 - denominal, affected object: e,s,o (*pflastern*)
- ... contextual clues

Sample data – descriptive problems

Contextual clues for sortal disambiguation: indicators

Context partners of nominalizations as possible sort indicators

- Modifiers, e.g.
 - Attributive adjectives: *kontinuierliche Messungen* (e)
 - Temporal/local PPs: *Messungen im Januar* (e)
 - Complement-indicating genitives:
die Fälschung von antiken Münzen (e)
- Selectors, e.g.
 - Embedding verbs (whose objects the nominals are):
Spaltung überwinden (s), *Messung durchführen* (e), *Fälschung verkaufen* (o)

⇒ Task for data extraction:

Identify sortal reading(s) by use of contextual data,

But: not all context partners are indicators:

wir brauchen neue Messungen

Sample data – descriptive problems

Structural ambiguity and its impact on extraction

Indicator status may depend on structural reading:

- *Messungen in Gomel*: PP-loc \Rightarrow e-Indicator
- *Messungen des Radiologischen Instituts in Gomel*:
PP-attachment ambiguous:
 - (a) [Messungen [des R.I._{NP-gen}][in G._{PP-loc}]] \Rightarrow e
 - (b) [Messungen [des R.I. [in G.]_{NP-gen}]] \Rightarrow ? \Rightarrow no indicator
- *Peterson plante stundenweise Absperrungen mit verschiedenen Materialien*
 - (a) *stundenweise Absperrungen* \Rightarrow e
 - (b) *Absperrungen mit verschiedenen Materialien* \Rightarrow o
stundenweise planen “absorbs” potential e-Indicator

Sample data – descriptive problems

Unexpected readings of *-ung*-nominals: facts

Examples:

- *die Verschärfung der Gesetze hat bewirkt, daß ...*
 - *die Messungen zeigen/beweisen, daß ...*
 - *die Überzeugung, wonach ...*
 - Fact-readings:
 - “the fact that the laws have been made more restrictive ...”
 - “the fact that [x] has been measured ...”
- ⇒ Task for data extraction:
Find verbs which license fact-readings of *-ung*-nominals as subjects:
same as verbs taking subject clauses?

Tool Architecture

Requirements – Specifications

- Flexible interface to the lexicon
 - “Messungen (e,o) zeigen subj: (f,o,e)...”
(clash: e & f / no clash: o & o)
- Underspecified representations
 - search for sentences which can show investigated phenomena:
 - “Messungen des Instituts in Gomel zeigen”
- Partial disambiguations
 - concentrate on considered phenomena
 - avoid irrelevant detailing of readings
 - tractability, efficiency

Tool Architecture

Realization – Implementation

- Adaptation of the MT tool *translate* (research version)
 - cooperation with the developer *Lingenio GmbH*
 - dependency based unification grammar
slot grammar (McCord 89/91)
 - underspecified semantic representation
Flat underspecified discourse representation theory (FUDRT)
(Eberle 97/04)
(Extension of UDRT (Reyle 93) for broad coverage representations,
implemented during SFB 340)
 - Large grammars, very large dictionaries, good performance
 - IBM's Logic based Machine Translation project (LMT)
 - Continuous exposition to market for 15 years
 - corresponding continuous improvement and extension

Tool Architecture

Semantic representations (FUDRSs) – tool functions

- Lexical representations are **functions**
(to be evaluated gradually by context)
- Sentence representations are sets
of DRSs and **DRS-modifiers and modifiers of ...**
and statements about order and **type of application**
- Carry out partial disambiguation and compute effects
- Inspect underspecified representations
- Use FUDRSs in search queries

Tool Architecture - Lexical Representation

The example of *Absperrung* (e,s,o)

Absperrung defines a mapping from DRFs in labeled DRSs (FUDRSs)
evaluated (partly/gradually)
depending on contextual information acting as a trigger.

absperrung(ref) \Rightarrow nsem_l

absperrung(e @ event) : I_{e@event}:

e (s, o)
absperrn(e)
abgesperrt(s)
meets(e,s)
absperrung(s,o)

absperrung(s @ state) : I_{s@state}:

s (e, o)
abgesperrt(s)
absperrn(e)
meets(e,s)
absperrung(s,o)

absperrung(x @ object) : I_{x@object}:

x
absperrung(x)

e @ event/s @ state... conditions in the sense of *wait-statements*:
wait((absperrung(SORT: event) →_))

Tool Architecture - Lexical Representation

Alternative lexical entries

Make use of derivation rules...

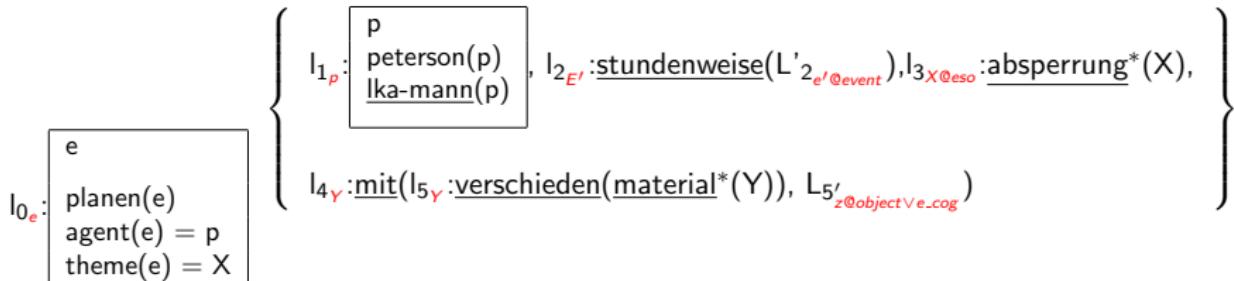
absperlung(x) := |_{x@eso}:**deriv(|_{e@event}:absperren(e))**

Tool Architecture - Sentence Semantics

Structured sets of partial representations

LKA-Mann Peterson plante stundenweise Absperrungen mit verschiedenen Materialien.

- (i) Peterson planned hour-long blockings using all kinds of materials.
- (ii) Repeatedly for hours, Peterson planned barriers consisting of all kinds of materials.



Tool Architecture - Sentence Semantics

Output of the *translate* implementation

```
Eingabeaufforderung (2) - prolog
:::
:::
:::
:::
:: Peterson plante stundenweise Absperrungen mit verschiedenen Materialien.

Dependence tree.

■ top      s(plan,534364)      mtv(ind:dcl:nwh,tf(past,0,X1),a):[[cogv,creationv,nocmp,plan]]
└ subj(n)  s(Peterson,1876456) noun(prop,nom,pers3-sg-X2,[]):[[Peterson,human,lastname]]
   └ xmod   s(stundenweise,683010) adv(p,[]):[[stundenweise]]
      └ ob(n)  s(absperrung,1580731) noun(cn,acc,pers3-pl-f,[]):[[absperrung,nonliv,stat]]
         └ xprep  s(mit,477132)      prep([mit|dat],[nwh]):[[mit]]
            └ objprep(dat)  s(material,461366) noun(cn,dat,pers3-pl-nt,[]):[[mat,material]]
               └ nadj   s(verschieden,751701) adj(s,dat,pers3-pl-nt,[nwh]):[[verschieden]]
```

Dependence tree diagram:

- top node: plan (534364)
 - subj(n): Peterson (1876456)
 - xmod: stundenweise (683010)
 - ob(n): absperrung (1580731)
 - xprep: mit (477132)
 - objprep(dat): material (461366)
 - nadj: verschieden (751701)

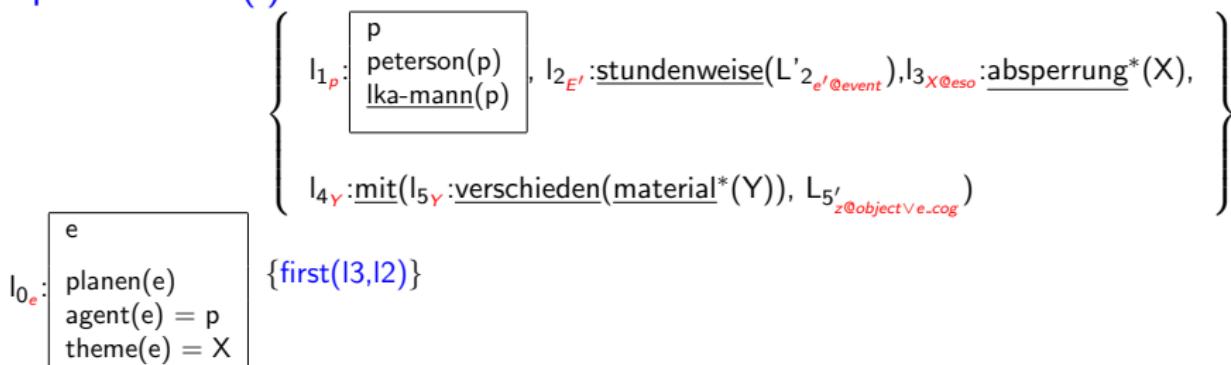
Tool Architecture - Partial Disambiguation

Specifications: constrain interpretations

Peterson plante *stundenweise Absperrungen* mit verschiedenen Materialien.

- (i) Peterson planned *hour-long blockings*, using all sorts of material.
- (ii) Repeatedly for hours, Peterson planned barriers consisting of all sorts of material.

Specification (i)



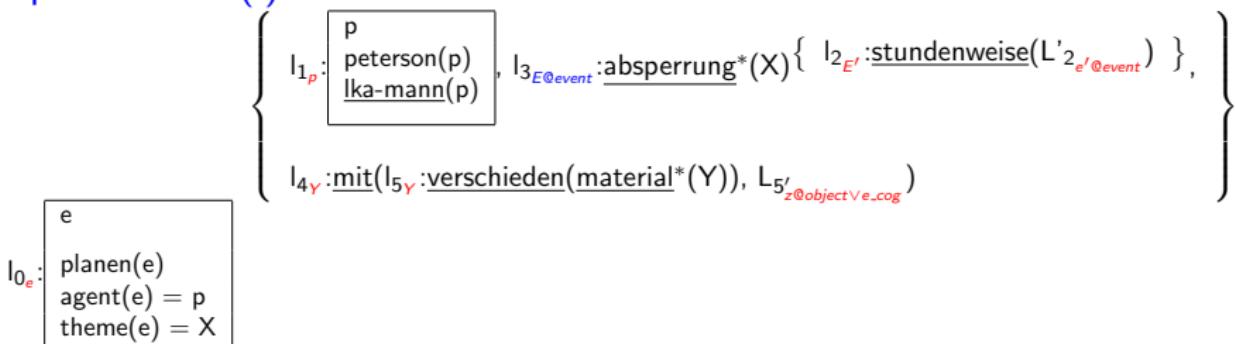
Tool Architecture - Partial Disambiguation

Specifications: consequences

Peterson plante **stundenweise Absperrungen** mit verschiedenen Materialien.

- (i) Peterson planned **hour-long blockings**, using all sorts of material.
- (ii) Repeatedly for hours, Peterson planned barriers consisting of all sorts of material.

Specification (i)



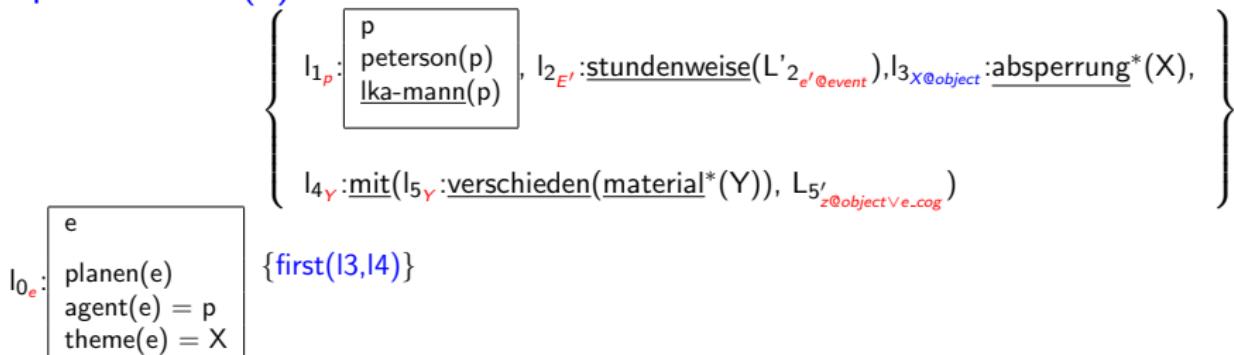
Tool Architecture - Partial Disambiguation

Specifications: constraints and consequences

Peterson plante stundenweise Absperrungen mit verschiedenen Materialien.

- (i) Peterson planned hour-long blockings, using all sorts of material.
- (ii) Repeatedly for hours, Peterson planned barriers consisting of all sorts of material.

Specification (ii)



Tool Architecture - Partial Disambiguation

Constraints and consequences: tool output

```
File: Eingabeaufforderung (2) · prolog
```

Partly specified.

```
top          s(plan,534364)      mtv(ind:dcl:nvh,tf{past,0,X1},a):[[cogv,creationv,nocmp,plan]]  
subj(n)     s(Peterson,1876456) noun(prop,non,pers3-sg-X2,[1]):[[Peterson,human,lastname]]  
xmod        s(stundenweise,683010) adv(p,[1]):[[stundenweise]]  
obj(n)      s(absperrung,1580731) noun(cn,acc,pers3-pl-f,[1]):[e,s,o][absperrung,nonliv,stat]]  
objp       s(mit,477132)      prep([init:dat],[nvh]):[[mit]]  
xprep       s(material,461366) noun(cn,dat,pers3-pl-nt,[1]):[[mat,material]]  
objprep(dat)s(verschieden,751701) adj(s,dat,pers3-pl-nt,[nvh]):[[verschieden]]
```

Partly specified.

```
top          s(plan,534364)      mtv(ind:dcl:nvh,tf{past,0,X1},a):[[cogv,creationv,nocmp,plan]]  
subj(n)     s(Peterson,1876456) noun(prop,non,pers3-sg-X2,[1]):[[Peterson,human,lastname]]  
obj(n)      s(absperrung,1580731) noun(cn,acc,pers3-pl-f,[1]):[e,geq(3,4),first(4,3)][absperrung,nonliv,stat]]  
xmod        s(stundenweise,683010) adv(p,[1]):[[stundenweise]]  
xprep       s(mit,477132)      prep([init:dat],[nvh]):[[mit]]  
objp       s(material,461366) noun(cn,dat,pers3-pl-nt,[1]):[[mat,material]]  
nadj        s(verschieden,751701) adj(s,dat,pers3-pl-nt,[nvh]):[[verschieden]]
```

Partly specified.

```
top          s(plan,534364)      mtv(ind:dcl:nvh,tf{past,0,X1},a):[[cogv,creationv,nocmp,plan]]  
subj(n)     s(Peterson,1876456) noun(prop,non,pers3-sg-X2,[1]):[[Peterson,human,lastname]]  
xmod        s(stundenweise,683010) adv(p,[1]):[[stundenweise]]  
obj(n)      s(absperrung,1580731) noun(cn,acc,pers3-pl-f,[1]):[o,geq(5,4),first(4,5)][absperrung,nonliv,stat]]  
objp       s(mit,477132)      prep([init:dat],[nvh]):[[init]]  
xprep       s(material,461366) noun(cn,dat,pers3-pl-nt,[1]):[[mat,material]]  
objprep(dat)s(verschieden,751701) adj(s,dat,pers3-pl-nt,[nvh]):[[verschieden]]
```

Tool Functions and Applications

- Trigger partial disambiguation and compute effects
- Query representations
- Use underspecified representations in search queries

Tool Functions and Applications

Querying the integrated representation

Searching for modifiers and selectors of *-ung*-nominalizations

Tool output

```
Eingabeaufforderung (2) - prolog
:::
:::
:::
:::
::: +findmods.
::: Peterson plante selbst stundenweise Absperrungen mit verschiedenen Materialien.

Dependence tree.

■—— top      s(plan,534364)      mtv(ind:dcl:nwh,tf(past,0,X1),a):[[cogv,creationv,nocmp,plan]]
└—— subj(n)   s(Peterson,1876456) noun(prop,nom,pers3-sg-X2,[1]):[[Peterson,human,lastname]]
    └—— xmod     s(selbst,630799)  noun(pron(undef),[nom,gen,dat,acc],pers3-sg-X3,[coref(1))]:[[selbst]]
    └—— xmod     s(stundenweise,683010) adv(p,[1]):[o,geq(5,4),first(4,5),[stundenweise]]
    └—— obj(n)   s(absperrung,1580731) noun(cn,acc,pers3-pl-f,[],):[[absperrung,nonliv,stat]]
    └—— xprep    s(mit,477132)    prep(lmit,dat,[nwh]):[[mit]]
        └—— objprep(dat) s(material,461366) noun(cn,dat,pers3-pl-nt,[],):[[mat,material]]
        └—— nadj     s(verschieden,751701) adj(s,dat,pers3-pl-nt,[nwh]):[[verschieden]]

[[plan,obj(n)], [xmod,selbst], [xmod,stundenweise], [xprep,[mit|material]]]
```

Tool Functions and Applications

Searching for modifiers and selectors

Flexible output representations

Variable granularity: with/without morphosyntactic, semantic and contextual information...
without...

Selectors:

[feststell,vprep(bei),xprep(bei),xprep(nach)]
[gelt,iobj(p([bei|dat])),subj(n),xprep(von)]
...

Modifiers:

[nprep(in),frühjahr]
[xprep(in),grundwasser]
[nprep(in),monat]

...

[nadj,monatelang]
[xmod,stundenweise]

...

Tool Functions and Applications

Searching for modifiers and selectors – Flexible output representations

With morphosyntactic, semantic and contextual information...

```
[n(feststell,mtv(ind:dcl:nwh,tf(past,0,_),a), [cogv]), [xprep(ben), [166]], [xprep(nach), [293]]]  
[n(feststell,mtv(ind:dcl:nwh,tf(past,0,_),p), [cogv]), [vprep(ben), [40,297]]]
```

...

```
[nprep(in), [n('frühjahr',det(def,sg), [season0]), [40]]]  
[xprep(in), [n(grundwasser,det(def,sg), [liq]), [256]]]  
[nprep(in), [n(monat,det(def,pl), [intervall0,timeas]), [257]]]
```

...

- identify regularities from syntactico-semantic information
Example: [xprep(in), intervall0] as event indicator (?)
- incorporate indicator property into lexicon
- test assumption
- revise / adopt
- bootstrapping of the lexicon

Tool Functions and Applications

Querying underspecified representations – Search with underspecified representations

Example: potential fact interpretations of *-ung*-nominalizations

Scenario:

- Hypothesis: subject of *zeigen* is interpreted as fact, if its representation is of type e (*event*)
- Test hypothesis:
 - Extract *zeigen*-sentences where subject **can** be event
 - Reject *zeigen*-sentences where it cannot or where *o* reading of *-ung*-subject seems preferred (because of missing additional hints)
 - * “Messungen des Instituts . . . von . . . in Gomel . . .” can be e
 - * “Messungen des Instituts in Gomel, die dort 99 veröffentlicht wurden, . . .” cannot
 - * “Messungen des Instituts in Gomel, das erfahrener ist als das in Celle, . . .” gives no hint
- use corresponding underspecified representation in query!

Tool Functions and Applications

Querying underspecified representations – Search with underspecified representations

Example: potential fact interpretations of *-ung*-nominalizations

Case: Sentences whose VP possibly characterizes facts and whose subj

- -*ung*-nominalization derived from e-root
(has e-reading, Rossdeutscher 2007)
- subject of VP which may characterize facts
(described by potential sentential complement)
- e-reading of nominalization possibly supported by additional e-modifier

FUDRS representation:

$$\{ \text{I}_{1_{u@eso}}:\text{ung_deriv}(\text{I}'_{1_{e'@event}}:\text{L}'_1), \quad \text{I}_{2_u}:\text{mod}(\text{L}'_{2_{x@eso}}), \quad \text{I}_{3_e}:\text{fin}(e'') \}$$
$$\{\diamond \text{first}(\text{I}_1, \text{I}_2) \}$$

$$\text{I}_{0_e}:\text{L}_0 \quad \& \quad \boxed{\begin{array}{c} e \\ \text{subj}(e,u) \end{array}}$$

Tool Functions and Applications

Querying underspecified representations – Search with underspecified representations

Excerpt from results (newspaper corpus *Wortschatz* University of Leipzig);

Die Verschärfung der Korruptionsgesetze im Jahr 1997 hat bewirkt, dass selbst das so genannte "Anfüttern" von Beamten mit Geschenken ohne Gegenleistung als strafbar gilt.

Und weiter heißt es in der vom braunen Zeitgeist geprägten Broschüre: "Dass die nationalsozialistische Gedankenwelt Allgemeingut der ganzen Kolonie geworden ist, zeigt die Beteiligung aller Volksgenossen am Winterhilfswerk und an gemeinsamen Eintopfessen."

Der Arbeitsalltag der Hausdame ist bestimmt von der gelebten Überzeugung, wonach eine Frau in ihrem Fach entweder "mit einem Mann oder einem Hotel verheiratet ist".

- potential e modification: *im Jahr 1997, am Winterhilfswerk*
- potential sentential complement *wonach ...*
- 'deep structure' information *bestimmt von der gelebten Überzeugung*

Conclusions

A tool for the syntactico-semantic analysis of corpus text

- generates and uses underspecified representation in corpus analysis and search
- based on slot grammar, FUDRT
- implemented as adaptation of the research version of *translate*
- provides functions for
 - partial disambiguation
 - extraction of potential sortal indicators
 - testing of sortal theories

Conclusions

A tool for the syntactico-semantic analysis of corpus text

- has modular, flexible architecture
 - can use lexical information of variable granularity
 - supports bootstrapping of the lexicon
- currently used to test theories about sortal properties of
-ung-nominalizations:
detailed analysis of cases of sortal paradoxes, fact readings, . . .
Basis: fictional text (6000 analyses), newspaper corpus (30000 analyses)
- extensible to the investigation of other phenomena

Future work

Phenomena, tools and documentation

- extend investigation to other phenomena:
subcategorization behaviour and nominalization of German particle
and prefix verbs,
extraction of collocations, ...
- bootstrap the system by corresponding results
- compare the results to results of other tools (Schiehlen 2003), ...
- implement database for storing and relating
 - corpus data
 - analyses of different tools,
ranging over different levels of linguistic description
 - metadata
- stepwise improvement of reliability of corpus analyses and of analysis
tools