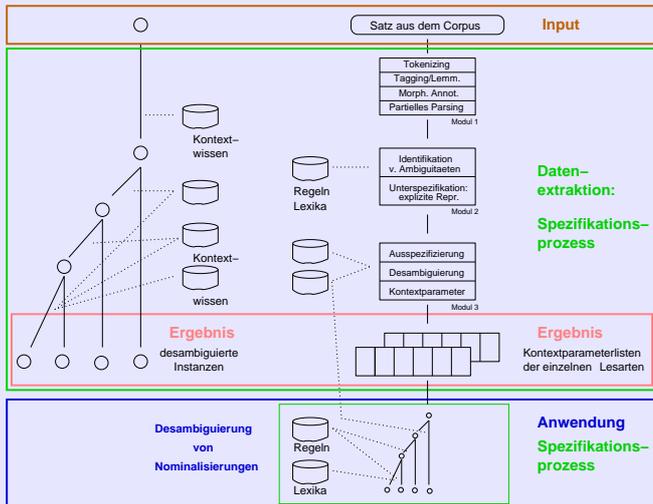


Unterspezifizierte Repräsentation und Disambiguierung sortal ambiger Nominalisierungen auf -ung

Kurt Eberle, Manuel Kountz, Ulrich Heid, Kerstin Eckart
Universität Stuttgart, Institut für Maschinelle Sprachverarbeitung, Azenbergstraße 12, D-70174 Stuttgart

Aufgabenstellung

- Automatische Datenextraktion als *inkrementeller Spezifizierungsprozess*: Berücksichtigung von Mehrdeutigkeiten bei der Extraktion von linguistischen Daten aus Textkorpora
 - **Ziele:**
 - * Identifikation mehrdeutiger Ausdrücke eines vorgegebenen Typs in Texten
 - * Identifikation von disambiguierungsrelevanter Information ("Indikatoren"): Welche Typen gibt es, welche Konstrukte finden sich im Kontext?
 - * Verfahren zur Markierung von Ergebnissen der Korpusuche: Sind sie von Mehrdeutigkeiten "betroffen", d.h. nicht "sicher"?
 - * Verfahren zur semi-automatischen Identifikation solcher Phänomene
 - **Ansatz:**
 - * Unterspezifizierte Repräsentation von Mehrdeutigkeiten in Texten
 - * Partielle Disambiguierung, wo dies für eine Extraktionsaufgabe nötig ist: Bootstrapping (s. "Methode")
- **Beispielfall:** sortale Ambiguität von -ung-Nominalisierungen:
 - *Absperrung*: Ereignis – Zustand – Objekt
vgl. aber im FR.: *barrière*_{OBJEKT} vs. *barrage*
 - Indikatoren aus dem Kontext:
 - z.B. Adjektive: *rot-weiße Absperrung*_{OBJEKT}
 - z.B. Verb vs. Subjekt: *die Absperrung*_{EREIGNIS} *ist in 5 Minuten erfolgt*
 - z.B. Verb vs. Objekt: *die Absperrung*_{ZUSTAND} *wird aufgehoben*
- Halbautomatische Extraktion *neuer disambiguierungsrelevanter Informationen*: Erweiterung des Indikatoren-Inventars um zusätzliche lexikalische und strukturelle Eigenschaften



Methode: Deterministisches Parsing mit "Ambiguity Awareness"

- Eindeutige **unterspezifizierte Satzanalyse** (Syntax und Semantik): vermeidet "forced guessing" (vgl. KOUNTZ ET AL. 2007)
 - Ergebnis:**
 - Sichere, eindeutige syntaktische Dependenzanalyse mit integrierter Ambiguitätsbeschreibung
 - Eindeutige unterspezifizierte semantische Repräsentation: Flache unterspezifizierte Diskursrepräsentationsstrukturen (FUDRS, EBERLE 2004)
- Aufgabenbezogene deterministische **partielle Disambiguierung** mittels Kontextwissen, z.B. mittels Selektionsrestriktionen (vgl. SPRANGER/HEID 2007)
 - Ergebnis:**
 - Menge aufgabenbezogen disambiguiertes unterspezifizierter Repräsentationen
- **Bootstrapping** zur Erweiterung des Indikatoren-Inventars
 1. Korpusuche: unterspezifizierte Suchanfragen
 2. Extraktion von Indikatoren-Kandidaten: partielle Disambiguierung, wo nötig
 3. Manuelle Sichtung, Erfassung von Indikatoren und Analysen der -ung-Nominalisierungen in Lexikon und Grammatik
 4. Verbesserte Extraktion

Projekthintergrund

- Sonderforschungsbereich 732
 - Thema: *Incremental Specification in Context*
 - Laufzeit: 1. Juli 2006 – 31. Juni 2010
- Projekt B3: *Disambiguierung von Nominalisierungen bei der Extraktion linguistischer Daten aus Corpora*

URLs

- SFB 732: <http://www.uni-stuttgart.de/linguistik/sfb732/>
- Lingenio: <http://www.lingenio.de>

Aktuelle Implementierung

Erweiterung der Analysekomponente des kommerziellen maschinellen Übersetzungssystems *translate* (LINGENIO (<http://www.lingenio.de>), MCCORD 1989)

- **Grammatik:** Umfangreiche Dependenzgrammatik (Deutsch, auch: Französisch, Englisch): unifikationsbasierter Formalismus der Slot Grammar (MCCORD 1991)
- **Unterspezifizierte Repräsentation:**
 - *Syntax:* Dependenzstrukturen
 - *Semantik:* FUDRSen: Flache unterspezifizierte Diskursrepräsentationsstrukturen, (EBERLE 2001, EBERLE 2004)
- **Ausgabe:** Für die jeweilige Extraktionsaufgabe relevante Lesarten: z.B. Nominalisierung plus sortale Ausdifferenzierung: *Absperrung*_{EREIGNIS} vs. *Absperrung*_{OBJEKT}- sowie effizient durchsuchbare Sammlung von Indikatoren aus Texten

Analyse- und Ausgabebeispiele

Peterson plante stundenweise Absperrungen mit verschiedensten Materialien.

- (I) STUNDENWEISES PLANEN VON ABSPERRUNGEN, DIE AUS VERSCHIEDENSTEN MATERIALIEN BESTEHEN (Absperrung als Objekt: *o*)
- (II) PLANEN VON STUNDENWEISEN ABSPERRUNGEN (Absperrung als Ereignis: *e*)
- (III) STUNDENWEISES, MIT VERSCHIEDENSTEN MATERIALIEN DURCHFÜHRTES PLANEN VON ABSPERRUNGEN (ohne Disambiguierung: *eso*)

Flache unterspezifizierte DRS (FUDRS)

mit Kern-DRS und Menge von Modifikatoren und (relevanten) DRF-'Dekorationen'

$$\begin{aligned}
 & \left\{ \begin{array}{l} 11 \text{ } p \text{ } \text{[Peterson(p)]} \\ 13 \text{ } E_i \text{ } \text{[stundenweise(L3')} \text{ } \{L3' \leq L3\} \\ 14 \text{ } X \text{ } \text{[eso]} \text{ } \emptyset \text{ } \text{[absperrung(x)]}, \\ 15 \text{ } y \text{ } \text{[mit(d6y'} \text{ } \emptyset \text{ } \text{[versch(material(y))]} \text{ } L4 \text{ } L4_{\text{[objective_instr]}} \end{array} \right\} \\
 & \left\{ \begin{array}{l} 1 \in i \\ \varepsilon \\ t < n \end{array} \right\} \\
 & 12 \text{ } E_{t, \text{Lakt}} \text{[heterogeneous]} \text{ } \left\{ \begin{array}{l} \text{planen}(e) \\ \text{agent}(e) = p \\ \text{theme}(e) = X \\ e \subset t \end{array} \right\}
 \end{aligned}$$

Disambiguierungsergebnisse der erweiterten *translate*-Analyse:

```

Partly specified.
top s(plan,534364) ntv(ind:dc1:nuh,tf(past,0,X1),a):[[cogw,creation,nocmp,plan]]
subj(n) s(Peterson,1876456) num(grp,non,per3-sg-A2,I1):[[Peterson,human,lastname]]
mod s(stundenweis,483810) adv(p,I1):[[stundenweise]]
obj(n) s(absperrung,1588731) num(cn,acc,per3-pl-f,I1):[e,seq(3,4),first(4,3)]absperrung,nonliv,stat]]
adv s(mit,47732) prep([init idat],lnh):[[mit]]
objprep(dat) s(material,461366) num(cn,dat,per3-pl-nt,I1):[[mat,material]]
adv s(verschieden,751701) adj(c,dat,per3-pl-nt,lnh):[[verschieden]]

Partly specified.
top s(plan,534364) ntv(ind:dc1:nuh,tf(past,0,X1),a):[[cogw,creation,nocmp,plan]]
subj(n) s(Peterson,1876456) num(grp,non,per3-sg-A2,I1):[[Peterson,human,lastname]]
mod s(stundenweis,483810) adv(p,I1):[[stundenweise]]
obj(n) s(absperrung,1588731) num(cn,acc,per3-pl-f,I1):[e,seq(3,4),first(4,3)]absperrung,nonliv,stat]]
adv s(mit,47732) prep([init idat],lnh):[[mit]]
objprep(dat) s(material,461366) num(cn,dat,per3-pl-nt,I1):[[mat,material]]
adv s(verschieden,751701) adj(c,dat,per3-pl-nt,lnh):[[verschieden]]

Partly specified.
top s(plan,534364) ntv(ind:dc1:nuh,tf(past,0,X1),a):[[cogw,creation,nocmp,plan]]
subj(n) s(Peterson,1876456) num(grp,non,per3-sg-A2,I1):[[Peterson,human,lastname]]
mod s(stundenweis,483810) adv(p,I1):[[stundenweise]]
obj(n) s(absperrung,1588731) num(cn,acc,per3-pl-f,I1):[e,seq(5,4),first(4,5)]absperrung,nonliv,stat]]
adv s(mit,47732) prep([init idat],lnh):[[mit]]
objprep(dat) s(material,461366) num(cn,dat,per3-pl-nt,I1):[[mat,material]]
adv s(verschieden,751701) adj(c,dat,per3-pl-nt,lnh):[[verschieden]]

```

Nächste Ziele

- **Meta-Analysewerkzeug:**
 - Einbeziehung und Anpassung verschiedener Analysewerkzeuge
 - Architektur für die Integration der Werkzeuge: Möglichkeit zur Kombination von Einzelresultaten, Ziel: Optimierung der Verlässlichkeit der Analysen
- **Datenbank** zur Verwaltung von Sätzen und ihren Analysen: Forschungswerkzeug: Suche in Analysen, z.B. für Bootstrapping
- **Abbildung** der unterspezifizierten Repräsentation auf *Korpusrepräsentations-Standards* (Auf Basis LAF/GrAF, z.B. KOUNTZ ET AL. 2008)
- **Ausweitung** des Ansatzes auf Daten aus anderen Sprachen (EN, FR), ggf. Nutzung von Übersetzungswissen (Paralleltext) zur Disambiguierung: *Die Bedienung ist schwierig* vs. *la serveuse / l'utilisation est difficile*

Literatur

- [Eberle2001] Kurt Eberle. 2001. FUDR-based MT, head switching and the lexicon. In: *Proceedings of the eighth Machine Translation Summit, Santiago de Compostela*.
- [Eberle2004] Kurt Eberle. 2004. Flat underspecified representation and its meaning for a fragment of German. *Habilitationsschrift*, Stuttgart.
- [McCord1989] Michael McCord. 1989. Design of LMT. In: *Computational Linguistics*, 15(1), pp33–52.
- [McCord1991] Michael McCord. 1991. The Slot Grammar System. In: Wedekind, J. & Rohrer, Ch. (eds.): *Unification in Grammar*, MIT-Press.
- [Kountz et al.2007] Manuel Kountz, Ulrich Heid, Kristina Spranger. 2007 Automatic sortal Interpretation of German Nominalisations with -ung: Towards using underspecified Representations in Corpora. In: *Proceedings of Corpus Linguistics 2007*
- [Kountz et al.2008] Manuel Kountz, Ulrich Heid, Kerstin Eckart. 2008 A LAF/GrAF based Encoding Scheme for underspecified Representations of Dependency Structures. erscheint in: *Proceedings of the Language Resources and Evaluation Conference (LREC)*, Marrakesch
- [Spranger/Heid2007] Kristina Spranger, Ulrich Heid. 2007 Applying Constraints derived from the Context in the Process of Incremental Sortal Specification. In: *Proceedings of the Workshop CSLP@Context 2007*, Roskilde.