

Integration von Regel- und Statistik-basierten Methoden in der Maschinellen Übersetzung

Kurt Eberle
Lingenio GmbH

Köthen
4. Juli 2008

Überblick

- Das Übersetzungsproblem:
Mehrdeutigkeit, Translation mismatches
- Lösungsvorschläge:
1., 2., 3. Generation der MÜ
- Aktuelle Diskussion:
Hybride Systeme
- *translate*:
Hintergrund, Architekturen
- 'Hybridisierung' von *translate*:
Statistische Korpus-Information für Analyse, Transfer und Generierung

Mehrdeutigkeit

- lexikalisch
 - Kategorial
 - *Time_{N/V} flies_{N/V} like_{V/P} an arrow - Zeitfliegen lieben einen Pfeil*
 - Semantisch
 - *Sie stellt_{PHYS/SOC} den Drucker ein_{PERSON/TOOL} - She hires/adjusts the printer.*
- strukturell
 - Funktionale Ambiguität (Label-/Etiketten-Ambiguität)
 - Ein Bild Albrecht Dürers. - A picture of/by Albrecht Dürer*
 - Attachment-Ambiguität
 - Bilder von der Kanzlerin hinter dem Tresen. - Pictures of the chancellor behind the bar*
 - Skopus-Ambiguität
 - Viele Hunde jagen eine Katze - many dogs chase a cat*
- referentiell, pragmatisch
 - *They saw the Alps as they flew over Zurich.*

Mehrdeutigkeit und Übersetzung

- Filter für die Disambiguierung
 - Syntaktische Constraints
*time flies like ... - *Zeit fliegt mag/mögen ...*
 - Semantische Selektionsbeschränkungen
*Sie stellt den Drucker ein ... *She hires the print device*
- Nicht alle Mehrdeutigkeiten müssen aufgelöst werden!
 - *printer - Drucker*
 - *viel(x,hund,ein(y,katze,jagen(x,y))) / ein(y,katze,viel(x,hund,jagen(x,y)))*
- **Variable Analysetiefe**- Übersetzung als *negociator*
(Kay Gawron Norvig1994)

Translation mismatches

(Kameyama Ochitani Peters 1991)

- Lexikalische Divergenz

un novillo - ein Jungbulle - a young bull

Rappe, Schimmel, ...

Boden - soil (Durrell 1988)

- Thematische Divergenz und Scrambling
(Dorr 1990, Hutchins Somers 1992)

Mir gefällt die Aufführung - I like the performance

Il remet le bouquet à la femme - er überreicht der Frau den Strauß/ den Strauß der Frau

- Hinzufügen und Tilgen von Teilstrukturen

er durchschwimmt den Fluß - Il traverse la rivière en nageant

- Strukturumkehrung (Head switching)

Er raucht gerne - He likes to smoke

La plupart des gens aiment le foot - Die meisten Leute mögen Fußball

Translation mismatches und Repräsentationen

- Morphosyntaktische Repräsentation

Er schreibt an die Angestellten ↔ Il écrit aux employés

[Prep [DET_{def} N] NP]

- Funktional/Semantische Repräsentation

Peter würde den Wein nicht mögen. ↔ Peter n'aimerait pas le vin.

[PRED: "mögen((↑SUBJ) (↑OBJ))"
SUBJ: [PRED: "wein"
OBJ: [PRED: "peter"
NEG: +
TENSE: COND]]]

- Semantisch/Konzeptuelle Repräsentation

Peter raucht gerne ↔ Peter likes to smoke

s
s : ATT(peter, { <POS_DISP, λx.rauchen(x) > })

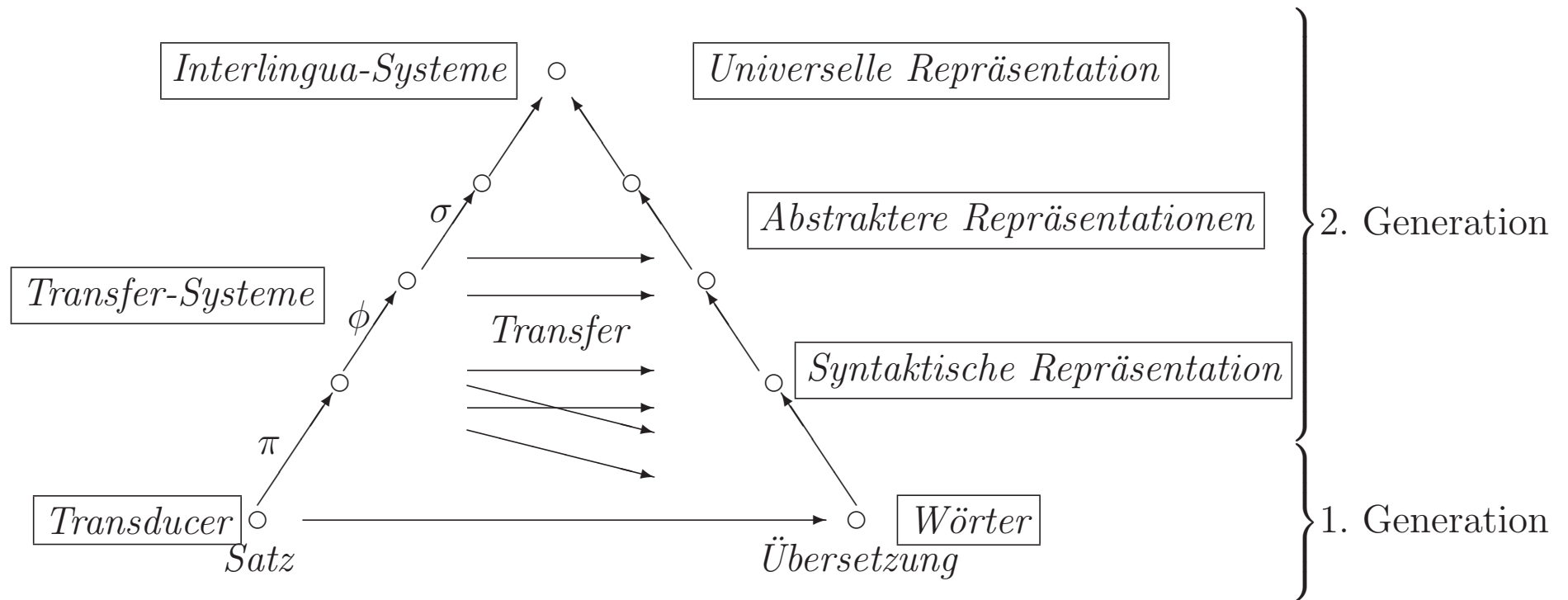
Mehrdeutigkeit, Mismatches und Analysen

- Analysen filtern!
- je abstrakter die Analyse desto geringer der Source-Target-Unterschied
- je abstrakter die Analyse desto mehr Disambiguierungsaufwand

→ **Optimierungsproblem !**

Traditionelle Architekturen

(Vauquois 1975)



Regel-basierte Architekturen

Lösungsversuche - 3. Generation

- Statistik-basierte Systeme (SMT)
Source-Channel-Modell (Brown et al 1990):

$$\hat{e}_1^I = \operatorname{argmax}_{e_1^I} \{P(e_1^I|f_1^J)\} = \operatorname{argmax}_{e_1^I} \{P(e_1^I) \times P(f_1^J|e_1^I)\}$$

(mit Lexikon-Sprach-Alignment-Modellen)

- Beispiel-basierte Systeme (EBMT)
(Sumita et al 1990, Maruyama Watanabe 1992)
Finde optimale Überdeckung aus Beispielen, verknüpfe Übersetzungen der Beispiele

Souvent le charpentier travaille le bois - Oft bearbeitet der Zimmermann das Holz

{le charpentier travaille} – {der Zimmermann arbeitet}
{travaille le bois} – {bearbeitet das Holz}
{le charpentier} – {der Zimmermann}

{Souvent } ∪ {le charpentier } ∪ {travaille le bois}

Aktuell - Hybride Systeme

- Maximum-Entropie-Modell und linguistische Features (Och Ney 2002)

$$\hat{e}_1^I = \operatorname{argmax}_{e_1^I} \left\{ \sum_{m=1}^M \lambda_m h_m(e_1^I, f_1^J) \right\}$$

h_m : POS-Unterschiede, Differenz VP-Knoten, ...

- Linguistische Vor- und Nachbearbeitung: (Quirk Menezes Cherry 2005): Dependency treelet translation

Text → *PRE* → *Suche* {← *Lexikon-*, *Alignment-*, *Sprach-Modell + Features*} → *POST*

- Syntaktisch motivierte Beispiele, rekursive Beispielkomposition (Chiang 2006): HIERO

$\langle (1)_{NP1} \text{ 's } (2)_{NP2}, \text{DET } (2)_{NP2} \text{ de } (1)_{NP1} \rangle$

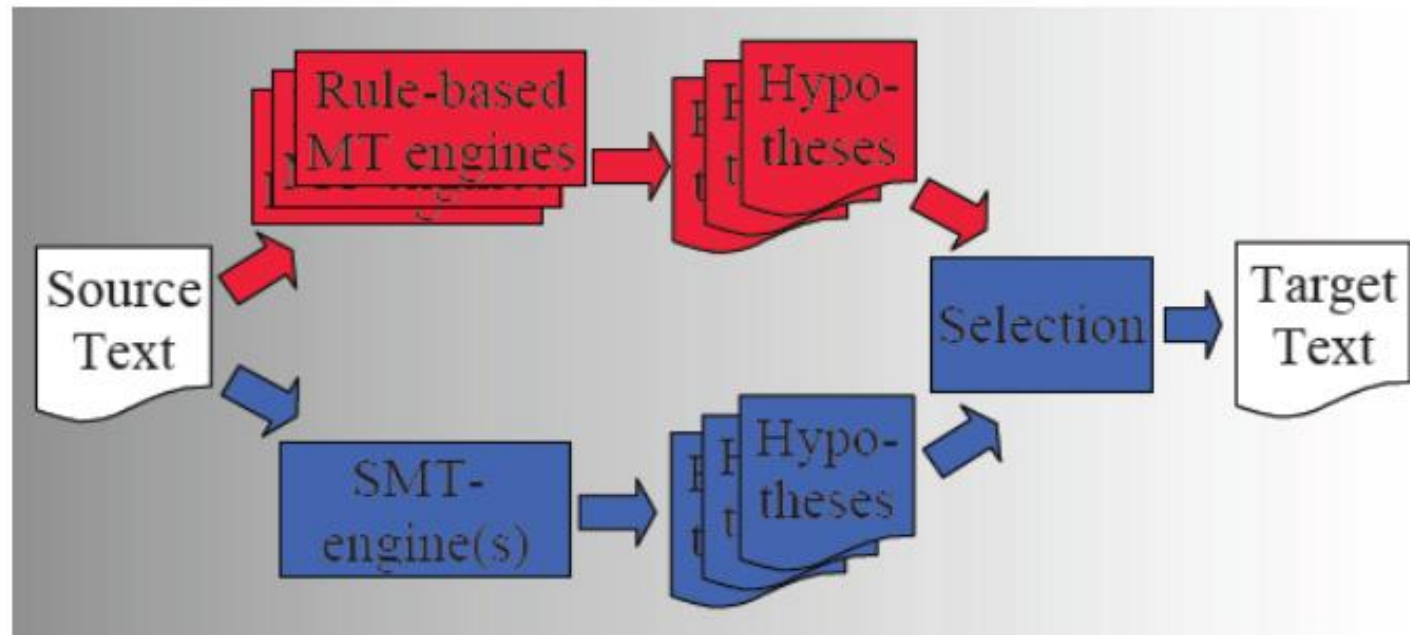
- Multisystem: Integriere Konkurrenzergbnisse

Multisysteme

- Schwache Integration
Wähle aus Ergebnissen / nutze zum Ressourcenaufbau
- Starke Integration
'interleaved architecture' (Reaktion auf jeweilige Teilergebnisse)

Multi-Engine MT (from Eisele, 2007)

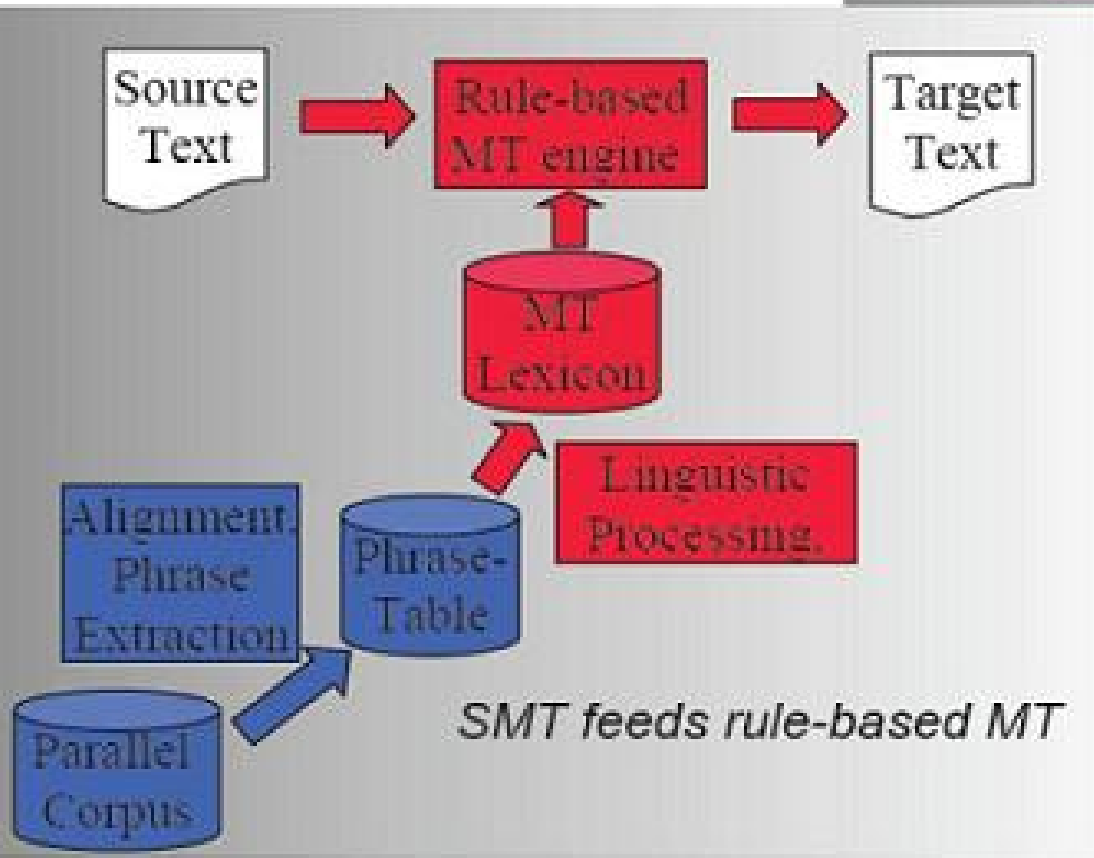
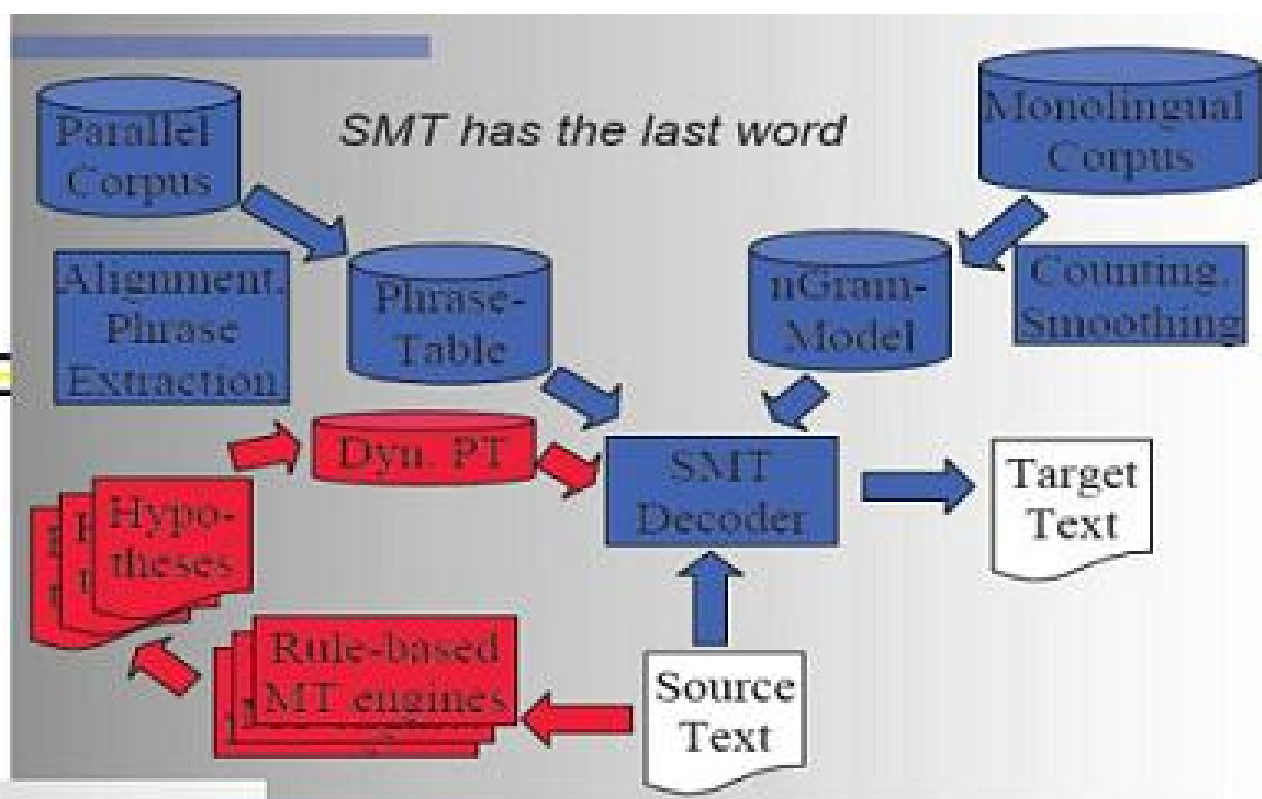
- Pangloss (1996)
- Moore (2008)



- Crucial: selection/combination
- complex
- monolingual

Architectures for weak

Integration (Eisele 2007)



translate

Hintergrund

Lingenio

- Wissenschaftliches Zentrum der IBM
- 1999 Ausgründung unter dem Namen *linguatec* Entwicklung & Services
- Zusammenarbeit mit *linguatec* Sprachtechnologien
- seit 2004 unter dem Namen *lingenio*
- Zusammenarbeit mit *digital publishing*

translate

Produktentwicklung

- Kommerzielles Maschinelles Übersetzungssystem
- *Logic based Machine Translation (LMT)*
(McCord 1989), IBM
- 1996 *Personal Translator* (IBM, linguattec)
- 1999 *Personal Translator Französisch* (linguateg)
- 2004 *translate* (Lingenio)
- 2005 *office* Wörterbcher

Produkte



version 8

deutschenglisch | englischdeutsch
deutschfranzösisch | französischdeutsch



translate

Die neuen Übersetzer
für Ihre Texte, Internet
und E-Mails

t translate pro - Deutsch - Französisch

☰ Datei Bearbeiten Ansicht Format Wörterbuch Übersetzen Satzarchiv Sprachausgabe Hilfe

Deutsch - Französisch

Arial 10 A B K U

Unbekannte Wörter

- Atomdrohungen (Kompositum) h18
- Opposition
- Kuhn
- Atomdrohung (Kompositum)
- Kraftmeierei
- Blumentopf (Kompositum)
- Topthemen (Kompositum)

Status

Quelltext - (unbenannt)

Topthemen

Merkel soll Pariser Atomdrohungen missbilligen

23/01/2006 08h18

Die Opposition hat Bundeskanzlerin Angela Merkel aufgefordert, sich bei ihrem he...

Frankreichs Präsident Jacques Chirac von den umstrittenen französischen Nuklea...

"Ereignis" Merkel muss endlich klar sagen, dass die französische Atomdrohung in Deu...

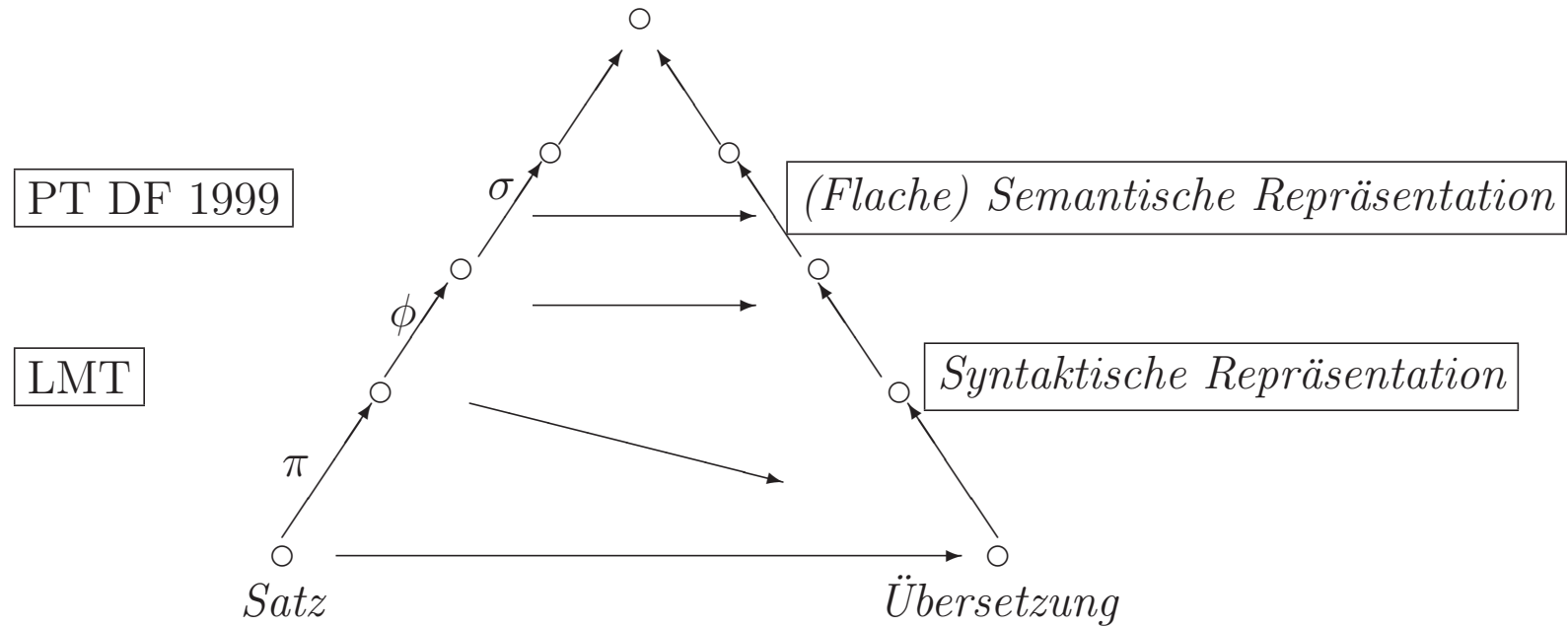
Zielltext - (unbenannt)

Des sujets de haut niveau

Merkel doit désapprouver des menaces d'atome parisiennes

23/01/2006 08 h18

translate-Architekturen



translate - Flache Semantik

Ziel

Flat underspecified discourse representation theory (FUDRT) (Eberle 1997)
(Erweiterung/Modifikation der UDRT (Reyle 1993))

- Lexikon:
Semantische Repräsentationen sind **Funktionen**
(die bei Bedarf schrittweise durch den Kontext ausgewertet werden)

$\underline{\text{drucker}}(x@ \neg \text{ARTEFACT}) := l_{x@ \text{PROF}}:$

x
druck_arbeiter(x)

$\underline{\text{drucker}}(x@ \neg \text{HUMAN}) := l_{x@ \text{EGERAET}}:$

x
druck_geraet(x)

- Satzrepräsentationen:
Mengen von DRSen und **DRS-Modifikatoren** ...
und Aussagen zu Ordnung und **Art der Applikation**

translate - Flache Semantik

Satzrepräsentationen

bilder(X) { ngen: kanzlerin(y), **xprep**(hinter): tresen(z) }

translate - **Flache Semantik**

implementiert

- Abhängigkeits-Strukturen (Prädikat-(gram.) Argument-Strukturen)
- Semantische Abstraktionen bei Koordinationen und bestimmten Angaben
- Informationsstruktur (Fokus-Hintergrundbestimmung bei Fokus-Adverbien)
- Skopusauflösung bei Bedarf
- Akzessibilitätsinformation für die Pronomenauflösung
- Transfer auf FUDRSen
- (variable Analysetiefe)

Ambiguitätserhaltende Übersetzung

Skopus-Ambiguität

Viele Hunde jagen eine Katze.

Repräsentation:

e jagen(e) subj(e,x) obj(e,y)	{	subj: <u>viele Hunde(x)</u> obj: <u>eine Katze(y)</u> ⋮	}
----------------------------------------	---	---------------------------------------------------------------	---

Übersetzung ohne Skopus-Disambiguierung:

Many dogs chase a cat.

Default-Transferalgorithmus

$$\tau(\text{jagen} \left\{ \begin{array}{l} \text{subj: } \underline{\text{viele Hunde}}(x) \\ \text{obj: } \underline{\text{eine Katze}}(y) \\ \vdots \end{array} \right\}) := \tau_n(\text{jagen} \left\{ \begin{array}{l} \tau_s(\text{subj}): \tau(\underline{\text{viele Hunde}})(x) \\ \tau_s(\text{obj}): \tau(\underline{\text{eine Katze}})(y) \\ \vdots \end{array} \right\})$$

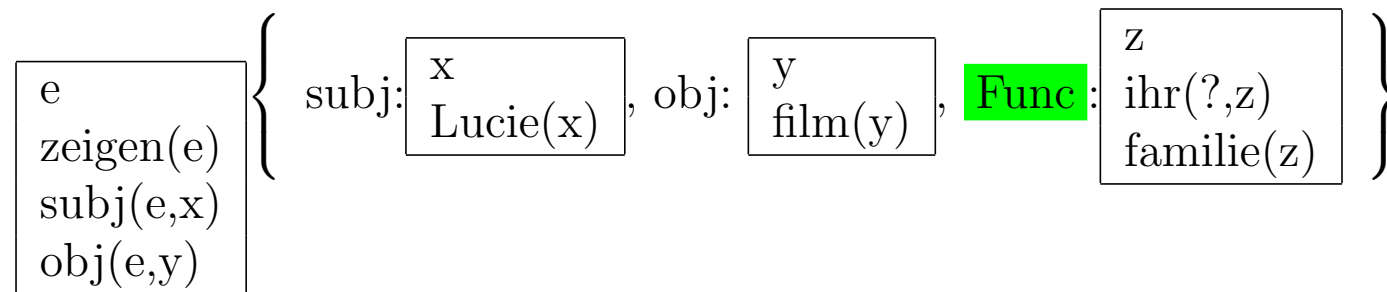
$$\boxed{\begin{array}{l} e \\ \text{chase}(e) \\ \text{subj}(e,x) \\ \text{obj}(e,y) \end{array}} \left\{ \begin{array}{l} \text{subj: } \underline{\text{many dogs}}(x) \\ \text{obj: } \underline{\text{a cat}}(y) \\ \vdots \end{array} \right\}$$

Partielle Disambiguierung

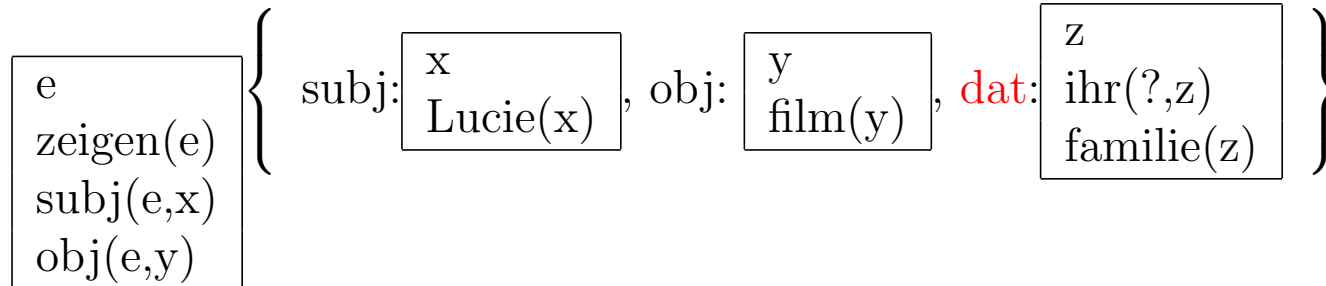
Funktionale Ambiguität

Lucie zeigte den Film ihrer Familie.

Ausgangssituation



a) Interpretiere **Func** als Dativ-Rolle der Verbrepräsentation

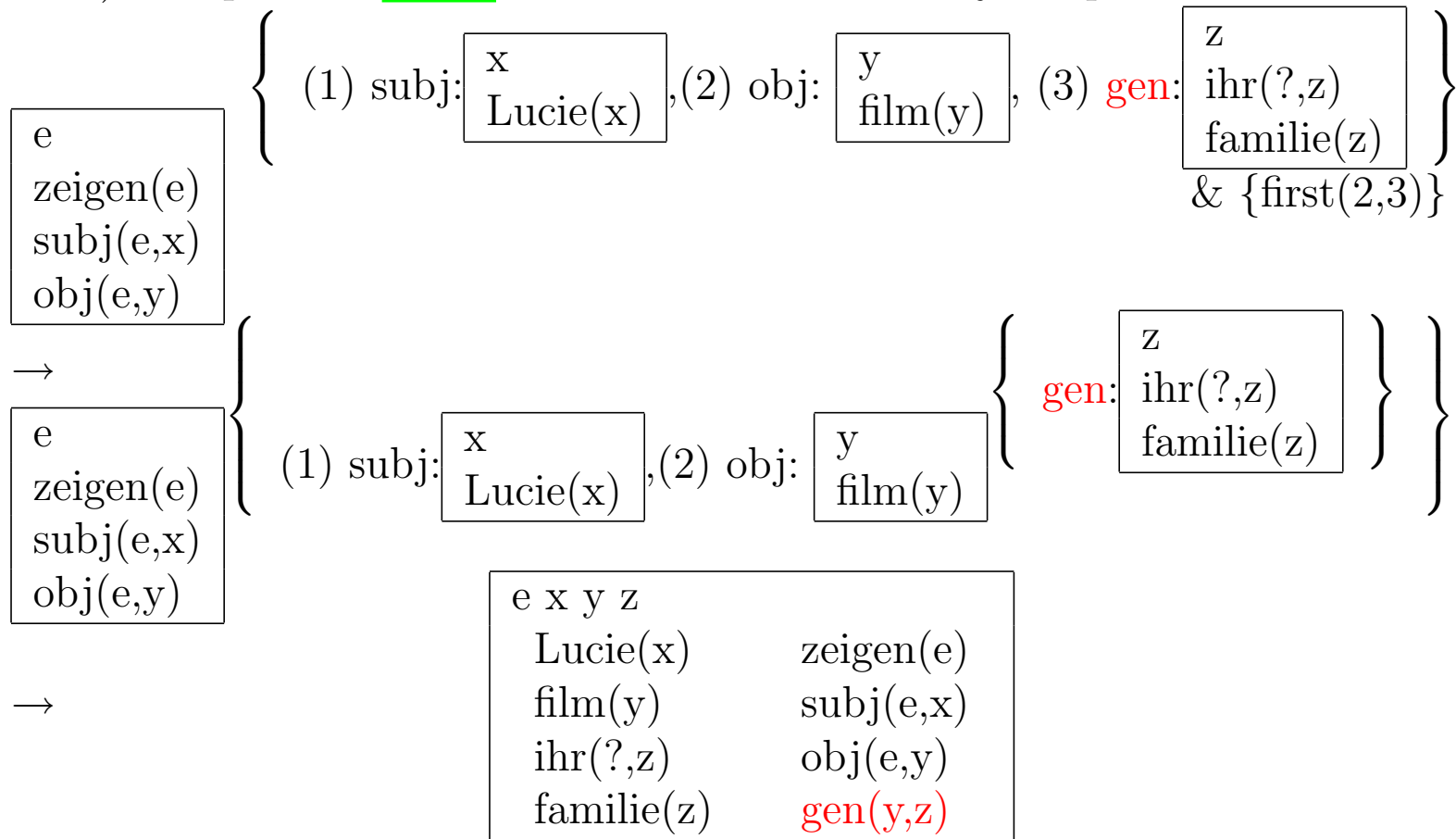


→

e	x	y	z
Lucie(x)		zeigen(e)	
film(y)		subj(e,x)	
ihr(?,z)		obj(e,y)	
familie(z)		dat(e,z)	

*Lucie **présente** le film **à** sa famille.*

b) Interpretiere **Func** als Genitiv-Rolle der Objektrepräsentation



*Lucie **pas**se le film **de** sa famille.*

Partielle Disambiguierung

Wodurch werden tiefere Auswertungen ausgelöst?

- Durch den Transfer
- nach Maßgabe im entsprechenden Eintrag im Lexikon

zeigen(e): Cond: $C \vdash_D \text{filled}(\text{obj}, \text{FILM}), \text{empty}(\text{dat})$ τ : passer

zeigen(e): Cond: $C \vdash_D \text{filled}(\text{obj}, \text{FILM}), \text{filled}(\text{dat})$ τ : préserter

→ Funktionale Ambiguität wird reduziert !

Partielle Disambiguierung

Wodurch werden tiefere Auswertungen ausgelöst?

- Durch Propagieren von Information im Kontext !

Beispiel: Pronomenresolution

Kürzlich erst hatte sie den Drucker eingestellt.

a) Jetzt kündigte er schon wieder.

b) Jetzt war er schon wieder defekt.

It was only recently that she had hired/adjusted the printer.

a) Now he already dismissed.

b) Now it already was defective again.

Pronomenauflösung und Antezedent-Übersetzung interagieren!

... den *Drucker_u*.... *Er_x* war *defekt*.

$\boxed{\begin{array}{l} s \\ s: \text{defekt}(x) \end{array}} \vdash_D \boxed{\text{MASCHINE}(x)}$

$\boxed{\begin{array}{l} u \\ \text{drucker_HUMAN}(u) \end{array}} \vdash x=u : \text{Widerspruch!}$

versus

$\boxed{\begin{array}{l} u \\ \text{drucker_MASCHINE}(u) \end{array}} \vdash x=u : \text{kein Widerspruch!}$

Verb- und Nomeninterpretation interagieren qua semantischer Selektionsbeschränkung!

einstellen(e,x,y@HUMAN) := $1_{e@SOC_ev}$:

e
einstellen_hire(e)
agens(e) = x
thema(e) = y

einstellen(e,x,y@MASCHINE) := $1_{e@PHYS_ev}$:

e
einstellen_adjust(e)
agens(e) = x
thema(e) = y

...

- einstellen [subj(n),obj(n)]

C: $\vdash_D SOC_ev$

τ : hire

C: $\vdash_D PHYS_ev$

τ : adjust

...

Partielle Disambiguierung

- Propagieren von Information
 - Selektionsbeschränkungen
 - Pronomen und Kennzeichnungen
 - ...
- Transfer-Notwendigkeit
 - Legitimation??

Integration regelbasierter und statistischer Methoden

- regelbasierter Kern
 - ökonomische Repräsentation (flach)
vermeide Entscheidungen die nicht durch ('harte') Regeln begründet sind
 - wenig Weltwissen (Typhierarchie ...)
 - variable Analysetiefe
Auswertung nur per *constraint propagation* und bei *Transfer-Bedarf*
 - mächtiger Lexikonformalismus
 - keine Transfergrammatik
 - modulares Design
unterscheide 'harte' Regeln von 'flexibler' Bewertung
- statistikbasierte Peripherie
 - Entscheidungshilfe

'Hybridisierung' von *translate*

Statistische Korpus-Information für

- Analyse
 - Lexem-Disambiguierung
 - strukturelle Disambiguierung
- Transfer
 - Lernen von Übersetzungsbeziehungen
 - Bewertung von Transfer-Äquivalenten
- Generierung
 - Wortstellung

Statistische Information für die Analyse

Funktional-/Attachment-Ambiguität

Er zeigt den Film der Familie (im Urlaub, ...)

Bootstrapping (verwende/verbessere Analyse)

- Flache Analyse hat folgende Charakteristik:

$$\underline{V}(e @ e_TYP) \left\{ \begin{array}{l} \text{subj:}\underline{NP}, \text{obj:}\underline{NP}(y @ \text{obj_TYP}), \\ \text{Func:}\underline{NP}(z @ \text{func_TYP}), \text{xprep(in):}\underline{NP}(p @ \text{prep_TYP}) \end{array} \right\}$$

- Korpus-Suche mit entsprechender Query
(Zusammenarbeit SFB 732, Projekt B3, Prof. Heid)
- Evaluation
signifikante Muster für verschiedene Lesarten
- (Hindle/Rooth): V NP P NP - Strukturen, ...

Statistische Information für den Transfer

Bewertung von Transfer-Äquivalenten

- *translate 11*:

Bootstrapping (Analyse und bilinguales Lexikonwissen)

- Transfer-Vorschläge des Lexikons
- Korpus-Test (Europarl)
- signifikante Verteilungen

→ Kontextbewertungen für Übersetzungsalternativen:

Er stellt das ein - He adjusts {adjust 0.35, stop 0.3, suspend 0.15} it.

Statistische Information für den Transfer

Bewertung von Transfer-Äquivalenten

The screenshot shows the 'translate pro' software interface. The title bar reads 'translate pro - Deutsch - Englisch'. The menu bar includes 'Datei', 'Bearbeiten', 'Ansicht', 'Format', 'Wörterbuch', 'Übersetzen', 'Satzarchiv', 'Sprachausgabe', and 'Hilfe'. The toolbar contains various icons for file operations and translation functions. The main window is divided into three panes:

- Assistent (Assistant):** Shows the translation status as 'Übersetzung beendet' (Translation completed) and 'Bereit' (Ready). It asks 'Was möchten Sie als Nächstes tun?' (What do you want to do next?) and provides options: 'Übersetzung nachbearbeiten' (Edit translation), 'Quelltext speichern' (Save source text), 'Zieltext speichern' (Save target text), and 'Als Übersetzungsprojekt speichern' (Save as translation project).
- Quelltext (Source Text):** Displays the German text: 'Er stellt mit Freude den Engländer ein.' and 'Er stellt mit Freude den Franzosen auf die passende Norm ein.'
- Zieltext (Target Text):** Displays the English translation with statistical annotations: 'He hires {accommodates, engages, stops} the Englishman {monkey wrench} with joy {pleasure, delight}.' and 'He adjusts {accommodates, engages, puts in} the Frenchman to the suitable {right} norm {standard, specification} with joy {pleasure, delight}. to the suitable {right} norm {standard specification} with joy {pleasure, delight}.'

The bottom left pane shows 'Einstellungen' (Settings) and 'Szenario-Manager' (Scenario Manager) options.

Statistische/EBMT- Information für den Transfer

Lernen von (analytischen) Übersetzungsbeziehungen

Aus bestimmten Gründen stellten die beiden Fraktionen ihre Feindseligkeiten vorübergehend durch einen Waffenstillstand ein und vereinbarten ...

For some reason, a temporary cease-fire in the hostilities between the two factions was established and (Datei ep-96-09-18.al, Zeile 1318)

Bootstrapping (Lexikonwissen, Source-Targetanalysen)

- Source-Analyse
- Target-Analyse
- Zuordnung von Teilen mittels bilingualen *Cognates* aus dem Lexikon
- Generalisierung von Bedingungen
- Signifikanz-Test

(Konvens 2008)

Aus bestimmten Gründen stellten die beiden Fraktionen ihre Feindseligkeiten vorübergehend durch einen Waffenstillstand ein und vereinbarten ...

For some reason, a temporary cease-fire in the hostilities between the two factions was established and (Datei ep-96-09-18.al, Zeile 1318)

- l_0 :einstellen [subj(n),obj(n)]

- C: $d(\text{vadv}):l_1$:vorübergehend & $d(\text{subj}):l_2$:fraktion
& $d(\text{obj}):l_3$:feindseligkeit & $d(\text{prep}(\text{durch})):l_4$:waffenstillstand

- τ : establish [\emptyset ,obj(n): $\tau(l_4)$]
& $\tau(d-l_1)=\tau(l_0)-d(\text{obj})-d(\text{nadj})$
& $\tau(d-l_3)=\tau(l_0)-d(\text{obj})-d(\text{prep}(\text{in}))$
& $\tau(d-l_2)=\tau(l_0)-d(\text{obj})-d(\text{prep}(\text{in}))-d(\text{prep}(\text{between}))$

Verallgemeinerungen ...

jmd./ stellt **STATE/** durch **EVENT/** ein - **s.o./** establishes **EVENT/** in **STATE**

Statistische Information für die Generierung

Lernen von Wortstellungsregeln

- Wortstellung

Poirot remet la lettre à la femme → *den Brief der Frau/ der Frau den Brief*

→ Frequenzanalysen

Flache Analysen zu Korpussätzen mit Typinformationen für NPs
(Zusammenarbeit SFB 732, Projekt D2, Prof. Rohrer)

Zusammenfassung

- Aktuell: 'Hybridisierung'
- Gängig: Integration linguistischen Wissens in SMT
- *translate*: Integration von SMT- und EBMT-Methoden in RBMT
 - Flache Repräsentationen
 - variable Analysetiefe, Auswertungstrigger
 - Modularisierung:
Harte Regeln (deklarative Grammatiken) - *Weiche* Entscheidungskriterien (Disambiguierung)
 - bei Analyse, Transfer, Generierung
- Signifikante Qualitätsverbesserung
- ökonomisch, effizient
- autonome, Human-lesbare Wörterbücher