

Aktuelle Architekturfragen in der Maschinellen  
Übersetzung  
semantischer Transfer und Integration statistischer  
Information in *translate*

Kurt Eberle  
Lingenio GmbH

26. Juni 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
  - Kategoriale Ambiguität von Lexemen und Gramemen
    - *Time<sub>N/V</sub> flies<sub>N/V</sub> like<sub>V/P</sub> an arrow - Zeitfliegen lieben einen Pfeil*
  - Semantische Ambiguität kategorial disambigierter Lexeme und Grameme
    - *Sie stellt<sub>PHYS/SOC</sub> den Drucker ein<sub>PERSON/TOOL</sub> - 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 showing the chancellor behind the bar/ pictures behind the bar showing the chancellor*
  - Skopus-Ambiguität
    - Few women like many men - Wenige Frauen mögen viele Männer*
- 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  
*\*She hires the print device*
- Nicht alle Mehrdeutigkeiten müssen aufgelöst werden!
  - *printer - Drucker*
  - *few(x, woman, many(y, man, like(x, y))) / many(y, man, few(x, woman, like(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<sub>def</sub> 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 swim*

s  
s : ATT(peter, { <POS\_DISP, λx.rauchen(x) > })

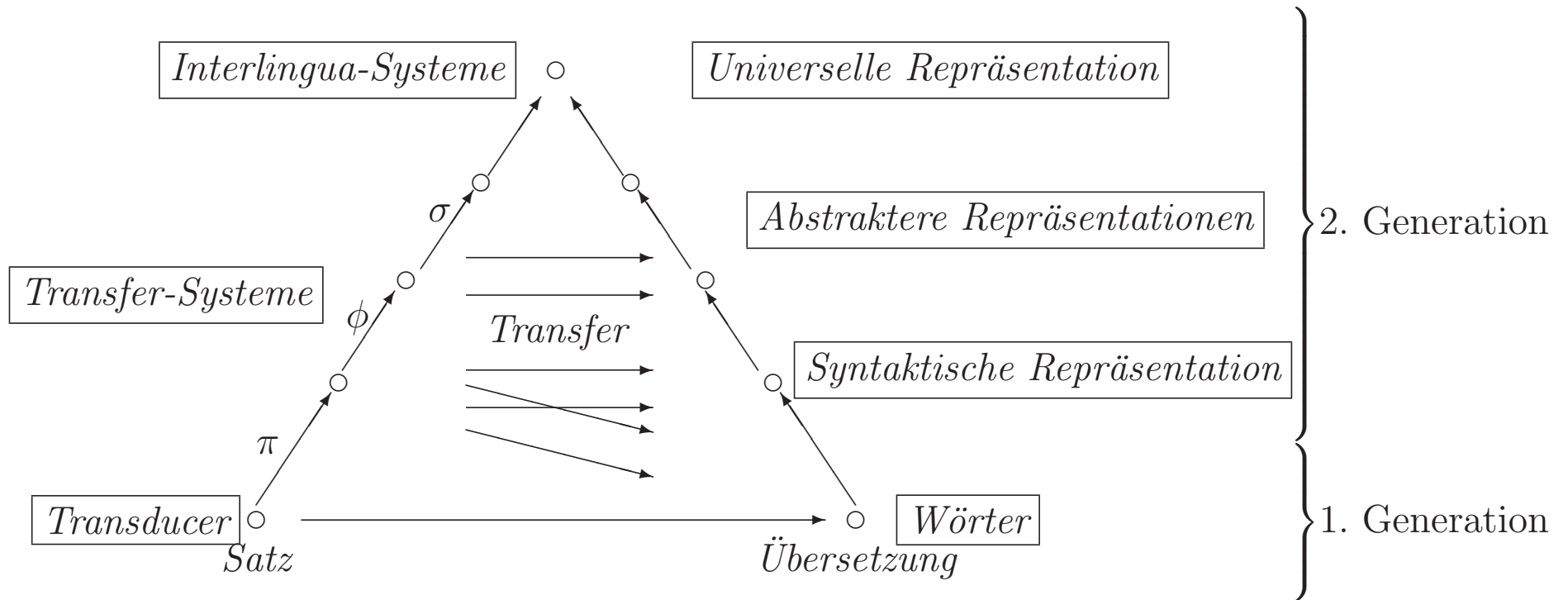
# Mehrdeutigkeit, Mismatches und Analysen

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

→ Optimierungsproblem

# Lösungsversuche

(Vauquois 19975)



Regel-basierte Architekturen



# Lösungsversuche - 3. Generation

- Statistik-basierte Systeme (SMT)
- Beispiel-basierte Systeme (EBMT)

# Statistik-basierte Systeme (SMT)

(Brown et al 1990), CANDIDE

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

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

*Lexikonmodell*      *travaille*  $\rightarrow$  {*bearbeitet* ( $p_{l,1}$ ), *arbeitet* ( $p_{l,2}$ ), ...} ,       $P(f_1^J | e_1^I)$   
                             *le*  $\rightarrow$  {*der, die, das, ...*}

*Sprachmodell*      *das Holz* ( $p_{s,1}$ ), *der Holz* ( $p_{s,2}$ )       $P(e_1^I)$

*Alignment-Modell*      Pos. 4 (*travaille*)  $\rightarrow$  Pos. 2 ( $p_{a,1}$ ), ...

# Beispiel-basierte Systeme (EBMT)

(Sumita et al 1990, Maruyama Watanabe 1992)

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

<i>{le charpentier travaille}</i>	–	<i>{der Zimmermann arbeitet}</i>
<i>{travaille le bois}</i>	–	<i>{bearbeitet das Holz}</i>
<i>{le charpentier}</i>	–	<i>{der Zimmermann}</i>

Finde (berechne) optimale Überdeckung:

*{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*  $\rightarrow$  *VB*  $\rightarrow$  *Suche*  $\{ \leftarrow$  *Lexikon-*, *Alignment-*, *Sprach-Modell + Features*  $\} \rightarrow$  *NB*

- 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

*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* (linguattec)
- 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 F 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 heftigen Gespräch mit dem französischen Präsidenten Jacques Chirac von den umstrittenen französischen Nuklear-Atomdrohungen zu distanzieren. "Frau Merkel muss endlich klar sagen, dass die französische Atomdrohung in Deutschland nicht akzeptiert wird", forderte die Opposition.

Zieltext - (unbenannt)

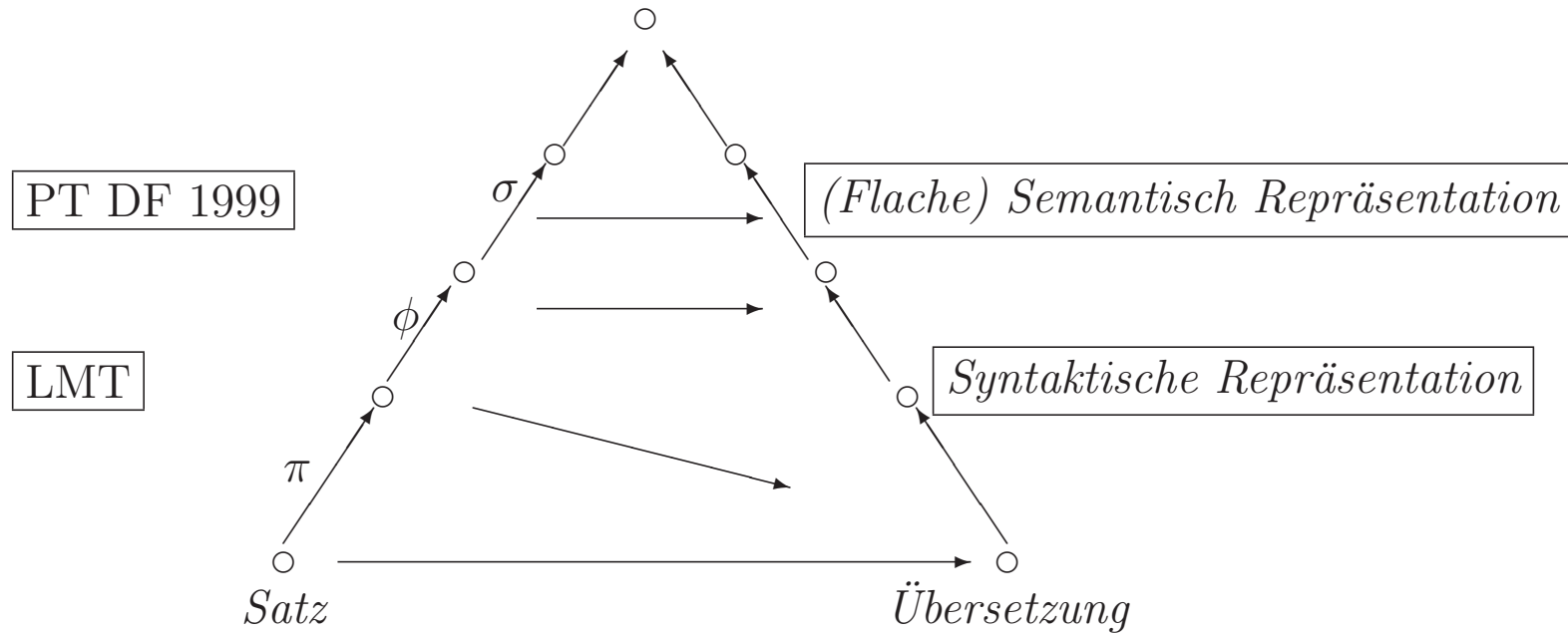
### Des sujets de haut niveau

#### Merkel doit désapprouver des menaces d'atome parisiennes

23/01/2006 08 h18



# *translate*-Architekturen



## *translate* - **Flache Semantik**

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

Disambiguierung nach Bedarf!

## *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)

# Dependenz-Strukturen

## Eingabeaufforderung (2) - prolog

::  
::  
::  
::  
::  
::  
::  
::

! : Kurz nachdem der vermeintliche Mörder verhaftet worden war, wurde der richtige gefasst.

Dependence tree.

■	top	s(fass,fass)	ntv(ind:del:nwh,tf(past,0,X1),p):[[app,fass]]	
├──	vsubconj	s(nachden,488819)	subconj(1):[[nachden,temp_subconj]]	
│	scadv	s(kurz,417117)	adv(p,1):[[dur_adv,kurz,l_adv,lnj,tmj]]	
├──	scomp	s(verhaft,744589)	ntv(dep:del:nwh,tf(past,1,X2),p):[[verhaft]]	
│	obj(n)	s(mörder,483526)	noun(cn,acc,pers3-sg-n,1):[[human,male,mörder]]	
│	ndet	s(der,d)	det(acc,pers3-sg-n,[def]):[[d,der]]	
│	└─	nadj	s(verneintlich,749211)	adj(p,acc,pers3-sg-n,1):[[verneintlich]]
│	obj(n)	s(richtig,578668)	adj(p,X3,pers3-sg-n,[[non,sg,n,usa!X4]]):[[adv_pred,jug_adj,richtig]]	
│	adet	s(der,d)	det(X5,pers3-sg-n,[def]):[[d,der]]	

Brièvement après que l'assassin présumé avait été arrêté, le juste fut pris.

::  
::  
::  
::  
::  
::  
::  
::

## ”Auffaltung” koordinierter Strukturen (bei Bedarf)

### ☞ Eingabeaufforderung (2) - prolog

Eingabe:

! : Zum Schluss war er Alkohol-abhängig und sogar süchtig nach Heroin.

Dependence tree.

■	top	s(sei,2)	ntv(ind:dcl:nwh,tf(past,0,X1),a):[[lse,sei,v_adj]]
	vprep	s(zu,814900)	prep([zu dat],[nwh]):[[dirp,zu]]
	objprep(dat)	s(schluss,schluss)	noun(cn,dat,pers3-sg-n,[]):[[accompl,result0,schluss]]
	ndet	s(dem,d)	det(dat,pers3-sg-n,[def]):[[d,dem]]
	lconj	s(sei,2)	ntv(ind:dcl:nwh,tf(past,0,X2),a):[[lse,sei,v_adj]]
	subj(n)	s(er,206261)	noun(pron(pers3),nom,pers3-sg-n,[]):[[er,perspron]]
	pred(a)	s(abhängig,4350)	adj(p,X3,X4,[]):[[abhängig]]
	acomound	s(alkohol,20600)	noun(cn,[nom,dat,acc],pers3-sg-m,[]):[[alkohol,beverage,dr]]
	rconj	s(sei,2)	ntv(ind:dcl:nwh,tf(past,0,X5),a):[[lse,sei,v_adj]]
	vadv	s(sogar,644400)	adv(p,[]):[[dsa,focus_adv,ppadv,shade,sogar]]
	pred(a)	s(süchtig,856507)	adj(p,X6,X7,[]):[[noadv,süchtig]]
	acomp(p([nach dat]))	s(heroin,33220)	noun(cn,dat,pers3-sg-nt,[]):[[droge0,heroin,mass]]
	subj(n)	empty	coref(5)

À la fin il était dépendant de l'alcool et il s'adonnait à l'héroïne  
même.

! :

4

# Implementierung im Test

## Funktionale Ambiguität, Attachment-Ambiguität

Zusammenarbeit mit IMS Stuttgart, SFB 732 Projekt B3

```
C:\ Eingabeaufforderung (2) - prolog
::
::
::
::
:: Peterson plante stundenweise Absperrungen mit verschiedensten 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]]
├── obj(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]]
-----
::
```

# Beispiele für die Verwendung von flachen Repräsentationen

## Skopus-Ambiguität - Darstellung

*In wenigstens fünf europäische Urlaubsgebiete flogen wenigstens drei Reiseveranstalter ihre Kunden in den letzten Jahren mehr als 20 Mal.*

*To at least five European tourist areas, at least three operators flew their clients more than 20 times within the last years.*

# Beispiele für die Verwendung von flachen Repräsentationen

## Skopus-Ambiguität - Darstellung

*In wenigstens fünf europäische Urlaubsgebiete flogen wenigstens drei Reiseveranstalter ihre Kunden in den letzten Jahren mehr als 20 Mal.*

*To at least five European tourist areas, at least three operators flew their clients more than 20 times within the last years.*

Repräsentation:

e fliegen(e) subj(e,x) obj(e,γ)	{	subj: <u>w 3 Veranstalter(x)</u>	}
		obj: <u>ihre_kunden(γ)</u>	
		vprep: <u>in w 5 Gebiete</u>	
		⋮	



# Default-Transferalgorithmus

$$\tau(\underline{\text{fliegen}} \left\{ \begin{array}{l} \text{subj: } \underline{\text{w 3 Veranstalter}}(x) \\ \text{obj: } \underline{\text{ihre\_kunden}}(\gamma) \\ \text{vprep: } \underline{\text{in w 5 Gebiete}} \\ \vdots \end{array} \right\}) := \tau_n(\underline{\text{fliegen}} \left\{ \begin{array}{l} \tau_s(\text{subj}): \tau(\underline{\text{w 3 Veranstalter}})(x) \\ \tau_s(\text{obj}): \tau(\underline{\text{ihre\_kunden}})(\gamma) \\ \tau_s(\text{vprep}): \tau(\underline{\text{in w 5 Gebiete}}) \\ \vdots \end{array} \right\})$$

e
fly(e)
subj(e,x)
obj(e,γ)

$$\left\{ \begin{array}{l} \text{subj: } \underline{\text{al 3 operators}}(x) \\ \text{obj: } \underline{\text{their clients}}(\gamma) \\ \text{vprep: } \underline{\text{to al 5 areas}} \\ \vdots \end{array} \right\}$$

# Partielle Disambiguierung

## Funktionale Ambiguität

*Lucie zeigte den Film ihrer Familie.*

a)                          V                          dat Obj  
*Lucie **présenta** le film **à** sa famille.*

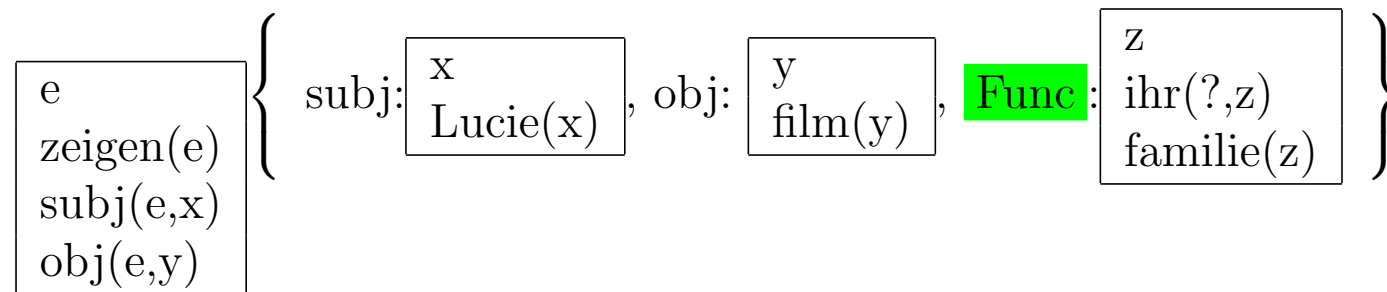
b)    N                          gen Obj  
*Lucie **passa** le film **de** sa famille.*

# Partielle Disambiguierung

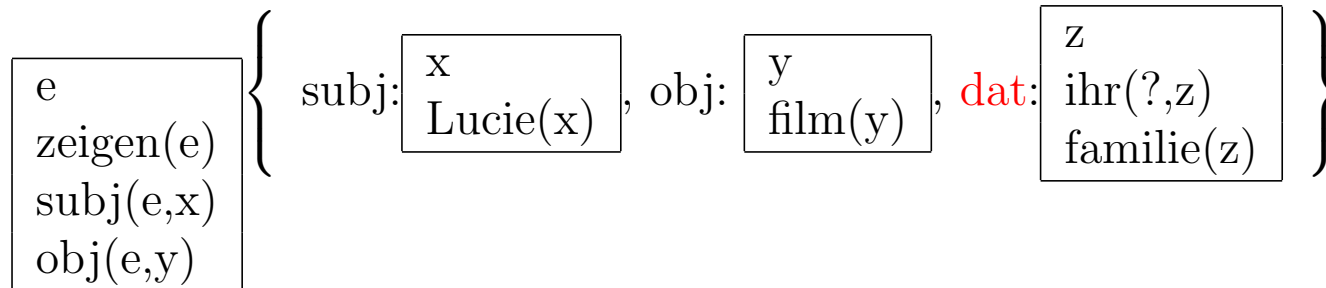
## Funktionale Ambiguität

*Lucie zeigte den Film ihrer Familie.*

Ausgangssituation



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

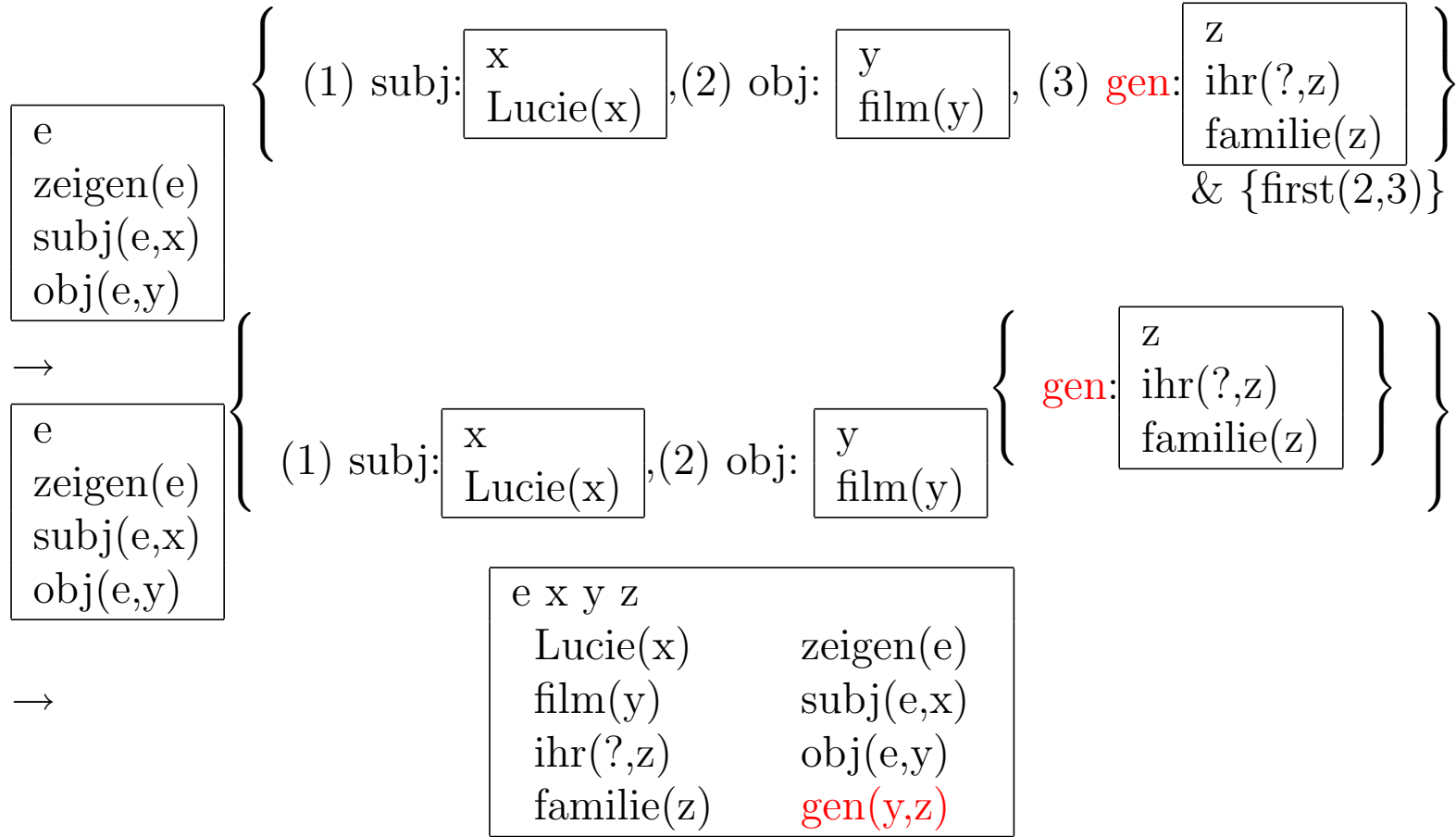


→

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

## Funktionale Ambiguität

Wodurch werden tiefere Auswertungen ausgelöst?

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

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

# Interaktionen

## Pronomen - Nomenambiguität

Pronomen haben nicht immer dieselbe Übersetzung:

elle	→	sie:	elle est intelligente	←	la femme
	→	er:	elle est visible	←	la lune
	→	es:	elle est bien élevée	←	la jeune fille

Der Bezug entscheidet!

Koreferenz-Entscheidungen haben komplexe Konsequenzen:

*Enfin il trouva la balle. Elle était dégonflée.*

elle	→	la balle	→	die Kugel	→	sie
				der Ball	→	er

Pronomenauflösung und Antezedent-Übersetzung interagieren!

... la *balle<sub>u</sub>*... *Elle<sub>x</sub>* était *dégonflée*.

s s: dégonflé(x)	$\vdash_D$	hohlkörper(x) norm_mit_gas_gefüllt(x)
---------------------	------------	--

u balle_kugel(u)	$\vdash x=u :$	Widerspruch!
---------------------	----------------	--------------

versus

u balle_ball(u)	$\vdash x=u :$	kein Widerspruch!
--------------------	----------------	-------------------



# Schlussfolgerungen

*Enfin il trouva la balle. Elle était dégonflée.*

- a) Der Bezug *balle* - *elle* ist möglich.
- b) Wenn der Bezug gewählt wird, muss *balle* im Sinne von *Ball* verstanden werden!
- c) Dann ist *elle* mit *er* zu übersetzen.

und ...

Die Koreferenz-Information kann weitere Effekte in immer größeren Kontexten zeitigen!

*Il tire la balle. Puis, tout d'un coup, elle se dégonfle.*

- Achte auf sinnvolle Abfolge von Übersetzungs- und Auswertungsschritten

# Interaktionen

## Pronomen - Strukturambiguität

### Pronomenbezug und Teilrepräsentationen interagieren!

*Der Polizist zeigte jedes der sichergestellten Videos einer Verdächtigen.  
Er versuchte, sie dadurch in Widersprüche zu verwickeln.*

- Die Videos zeigen alle dieselbe Verdächtige (N gen )

*sie = diese eine Verdächtige  $\rightarrow_{\tau}$  elle*

*– Il **passait** chacune des vidéos **d'**une suspecte. ... elle*

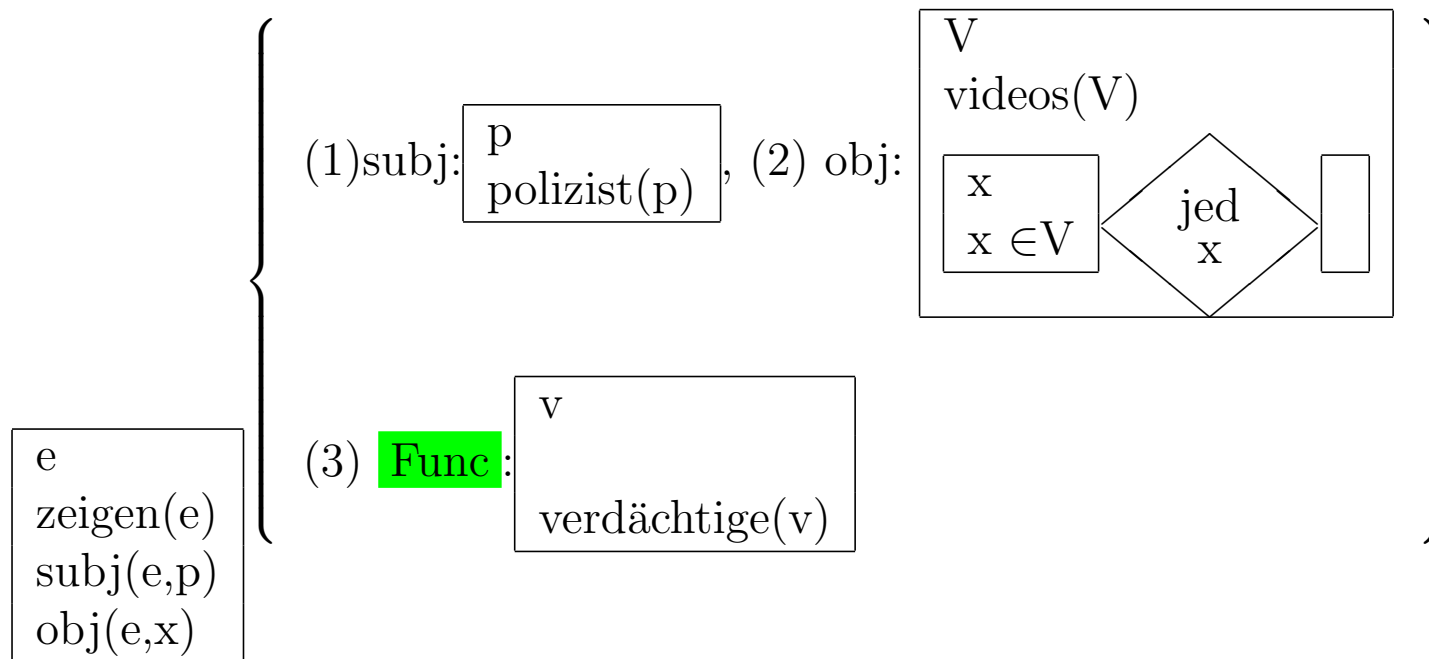
- Er zeigte jedes Video einer anderen Verdächtigen (V dat nach acc)

*sie = diese vielen Verdächtigen*

*– Il **présentait** chacune des vidéos **à** une suspecte (différente). ... elles*

- ...

# Interdependenz der Übersetzungen und der Bezüge der Pronomen und der Teilrepräsentationen!



# Zusammenhänge

$sie = \text{Pl\_Pro} \rightarrow_{\tau} \text{elles} \Rightarrow (2) > (3) \wedge \text{Func}=\text{dat}$   
 $\text{zeigen} \rightarrow_{\tau} \text{présenter à}$

$sie = \text{Sg\_Pro} \rightarrow_{\tau} \text{elle} \Rightarrow ( (3) > (2) \wedge \text{Func}=\text{dat} ) \vee ( \text{first}(2,3) \wedge \text{Func}=\text{gen} )$   
 $\text{présenter à} \qquad \qquad \qquad \text{passer}$

Pronomenresolution in *translate*:

- morpho-syntaktischer Filter (Lappin McCord 90)
  - semantische Bedingungen des gegebenen Kontexts
  - pragmatische Bewertungen (Lappin Leass 94)
  - Entscheidung !
- neue semantische Bedingungen
- abschnittsweise Übersetzung (aktuelle Entwicklung)

# Schwierige Mismatches

## Head switching

(Sadler Thompson 1991, Butt 1994)


*Peter schwimmt gerne.*

*Peter likes to swim.*

### Strukturelle Umkehrung:

gerne (Adjunkt) → like (syntaktischer Kopf)

schwimmen (Kopf) → swim (Komplement)

$$\tau((1) \text{ schwimmen } \left\{ \begin{array}{l} \text{subj: } (2) \text{ Peter} \\ \text{vadv: } (3) \text{ gerne} \end{array} \right\}) = \text{like} \left\{ \begin{array}{l} \text{subj: } \tau(2) \\ \text{comp: } \tau(1') \end{array} \right\} ?$$


# Lösung

- vermeide Zirkularität
- durch partielle Skopusfestlegung
- Anweisung im Lexikon!

*gerne*

> *adv* []

>  $\tau$ :

(u-cat(verb)  $\rightarrow$  aimer [item(subj(X'), $\tau$ (u-d(subj(X))))],item(obj(bin<sub>f</sub>), $\tau$ (u <sup>$\cancel{id}$</sup> )))]

...stößt eventuell weitere Disambiguierungen an...

## Worauf bezieht sich *gerne* ?

*Wie Franzi würde Peter gerne schwimmen.*

a) *gerne* [*wie Franzi schwimmen*] (können)

← *wie Franzi* befindet sich im Skopus von *gerne*

*Pierre aimerait (savoir) nager comme (/de la façon de) Franzi.*

b) *wie Franzi gerne* [*schwimmen*] (gehen)

← *wie Franzi* befindet sich nicht im Skopus von *gerne*

*Comme Franzi, Pierre aimerait (aller) nager.*

- ein semantisches Skopusphänomen!
- andere Lösungsvorschläge:  
Restriktionsoperator (Kaplan Wedekind 1993) , nicht ko-deskriptiv (Verbmobil)

# Partielle Skopusfestlegung

*Wie Franzi würde Peter gerne schwimmen.*

$$\boxed{\begin{array}{l} e \\ \text{schwimmen}(e) \\ \text{subj}(e,x) \end{array}} \left\{ \text{subj: } \boxed{\begin{array}{l} x \\ \text{Peter}(x) \end{array}}, \text{adv:gerne, pp: wiePetra} \right\}$$

(A)

$$\left( (1) \left( \boxed{\begin{array}{l} e \\ \text{schwimmen}(e) \\ \text{subj}(e,x) \end{array}} \left\{ \text{pp: (4)wieFranzi} \right\} \right) \left\{ \text{adv:(3)gerne} \right\} \right) \left\{ (2) \text{subj: } \boxed{\begin{array}{l} x \\ \text{Peter}(x) \end{array}} \right\}$$

gerne [schwimmen wie Franzi]  $\rightarrow_{\tau}$  aimer savoir nager comme Franzi



(B)

$$\left( \left( \left( \begin{array}{l} e \\ (1) \text{schwimmen}(e) \\ \text{subj}(e,x) \end{array} \right) \left\{ \text{adv: (3)gerne} \right\} \right) \left\{ \text{pp: (4)wieFranzi} \right\} \right) \left\{ \text{subj:(2)} \begin{array}{l} x \\ \text{Peter}(x) \end{array} \right\}$$

wie Franzi gerne [schwimmen]  $\rightarrow_{\tau}$  comme Franzi, aimer aller nager

# Repräsentationen, Formalismen - Fazit

- ökonomische Repräsentation (flach)
- Auswertung nach Bedarf
- mächtiger Lexikonformalismus
- keine Transfergrammatik
- modulares Design

# 'Hybridisierung' von *translate*

Statistische Korpus-Information für Analyse, Transfer und Generierung

## Analyse

- Lexem-Disambiguierung

*Little John was looking for his toy box. Finally, he found it.*

*The box was in the pen. John was very happy.*

(Bar-Hillel 1959)

Weltwissen → (pen = playpen/Laufstall)

oder → *Word sense disambiguation* (WSD)

$\text{pen}(x) \rightarrow \underline{\text{pen}}(x)|_C$      $C \vdash_d \text{sit}_x \in \text{BABY\_FRAME}$      $\rightarrow \text{Laufstall}(x)$   
 $C \vdash_D \text{sit}_x \in \text{WRITING\_FRAME}$      $\rightarrow \text{Stift}(x)$

- Funktional-/Attachment-Ambiguität

Beispiel: V NP Prep NP (Hindle/Rooth)

→ Bootstrapping: flache Analyse, Evaluation

# 'Hybridisierung' von *translate*

## Transfer

- *translate 11*:

Bootstrapping: Analyse und bilinguales Lexikonwissen

→ Kontextbewertungen für Übersetzungsalternativen:

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

- Lernen von (analytischen) Übersetzungsbeziehungen  
Bootstrapping: Lexikonwissen, Source-Targetanalysen

# 'Hybridisierung' von *translate*

## Generierung

- Wortstellung

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

→ Frequenzanalysen

Flache Analysen zu Korpusätzen 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