

Computerlinguistische Modellierung der deutschen Satzstruktur

Arkadiusz Jasiński

20. September 2006

formale Sprachbeschreibung – Muss das sein?

computerlinguistische Motivation:

formal definierter Objektbereich /
anwendungsorientierte CL

sprachwissenschaftliche Motivation:

präzise Modellierung der sprachlichen Kompetenz

*"I have always understood a generative grammar to be nothing more than an **explicit** grammar."* (CHOMSKY 1995:162)

formale Sprachen – ein Demonstrationsbeispiel

Annahme:

**Eine Mini-Sprache L mit einem Vokabular Σ_L
aus vier Elementen**

$$\Sigma_L = \{ \textit{Tobias}, \textit{gähnt}, \textit{kennt}, \textit{und} \}$$

Folglich:

**Die (unendliche) Menge Σ_L^+ beinhaltet alle
nicht-leeren Folgen von Elementen aus Σ_L**

formale Sprachen – ein Demonstrationsbeispiel

Länge 1

Tobias. gähnt. kennt. und.

Länge 2

Tobias Tobias.	Tobias gähnt.	Tobias kennt.	Tobias und.
gähnt Tobias.	gähnt gähnt.	gähnt kennt.	gähnt und.
kennt Tobias.	kennt gähnt.	kennt kennt.	kennt und.
und Tobias.	und gähnt.	und kennt.	und und.

Länge 3

Tobias Tobias Tobias.	Tobias Tobias gähnt.	Tobias Tobias kennt.	Tobias Tobias und.
Tobias gähnt Tobias.	Tobias gähnt gähnt.	Tobias gähnt kennt.	Tobias gähnt und.
Tobias kennt Tobias.	Tobias kennt gähnt.	Tobias kennt kennt.	Tobias kennt und.
Tobias und Tobias.	Tobias und gähnt.	Tobias und kennt.	Tobias und und.
gähnt Tobias Tobias.	gähnt Tobias gähnt.	gähnt Tobias kennt.	gähnt Tobias und.
gähnt gähnt Tobias.	gähnt gähnt gähnt.	gähnt gähnt kennt.	gähnt gähnt und.
gähnt kennt Tobias.	gähnt kennt gähnt.	gähnt kennt kennt.	gähnt kennt und.
gähnt und Tobias.	gähnt und gähnt.	gähnt und kennt.	gähnt und und.
kennt Tobias Tobias.	kennt Tobias gähnt.	kennt Tobias kennt.	kennt Tobias und.
kennt gähnt Tobias.	kennt gähnt gähnt.	kennt gähnt kennt.	kennt gähnt und.
kennt kennt Tobias.	kennt kennt gähnt.	kennt kennt kennt.	<i>usw.</i>

formale Sprachen – ein Demonstrationsbeispiel

Länge 1

Tobias.

gähnt.

kennt.

und.

Länge 2

Tobias Tobias.

gähnt Tobias.

kennt Tobias.

und Tobias.

Tobias gähnt.

gähnt gähnt.

kennt gähnt.

und gähnt.

Tobias kennt.

gähnt kennt.

kennt kennt.

und kennt.

Tobias und.

gähnt und.

kennt und.

und und.

Länge 3

Tobias Tobias Tobias.

Tobias gähnt Tobias.

Tobias kennt Tobias.

Tobias und Tobias.

gähnt Tobias Tobias.

gähnt gähnt Tobias.

gähnt kennt Tobias.

gähnt und Tobias.

kennt Tobias Tobias.

kennt gähnt Tobias.

kennt kennt Tobias.

Tobias Tobias gähnt.

Tobias gähnt gähnt.

Tobias kennt gähnt.

Tobias und gähnt.

gähnt Tobias gähnt.

gähnt gähnt gähnt.

gähnt kennt gähnt.

gähnt und gähnt.

kennt Tobias gähnt.

kennt gähnt gähnt.

kennt kennt gähnt.

Tobias Tobias kennt.

Tobias gähnt kennt.

Tobias kennt kennt.

Tobias und kennt.

gähnt Tobias kennt.

gähnt gähnt kennt.

gähnt kennt kennt.

gähnt und kennt.

kennt Tobias kennt.

kennt gähnt kennt.

kennt kennt kennt.

Tobias Tobias und.

Tobias gähnt und.

Tobias kennt und.

Tobias und und.

gähnt Tobias und.

gähnt gähnt und.

gähnt kennt und.

gähnt und und.

kennt Tobias und.

kennt gähnt und.

usw.

pragmatische Fragestellungen → Performanzmodell

Länge 1

Tobias.

gähnt.

und.

Länge 2

gähnt Tobias.

Tobias gähnt.

Länge 3

Tobias kennt Tobias.

- *Wie heißt du?*
- *Tobias.*

kennt Tobias Tobias.

- *Was macht der Tobias gerade?*
- *Gähnt.*

- *Schau mal, der Tobias gähnt!*
- *Und?*

formale Sprachen – ein Demonstrationsbeispiel

Länge 2

gähnt Tobias.

Tobias gähnt.

Länge 3

Tobias kennt Tobias.

$$L = \{$$

Tobias gähnt.
 Gähnt Tobias?
 Tobias kennt Tobias.
 Kennt Tobias Tobias?
 ...
 Kennt Tobias Tobias und Tobias ... ?
 ...
 Tobias gähnt und Tobias kennt und Tobias ...
 ...

$$\}$$

kennt Tobias Tobias.

formale Grammatik

$$G = \langle V_N, V_T, R, S \rangle,$$

worin

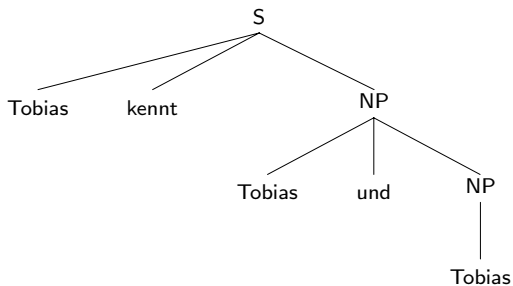
V_N – eine endliche Menge von Nichtterminalsymbolen,

V_T – eine endliche Menge von Terminalsymbolen,

R – eine endliche Menge von Regeln der Form $\alpha \longrightarrow \beta$,
wobei α und $\beta \in (V_N \cup V_T)^*$

S – das Startsymbol, $S \in V_N$.

kontextfreie Grammatik ...



$$G = \langle$$

$$\{ S, NP \},$$

$$\{ Tobias, gähnt, kennt, und \},$$

$$\{ S \rightarrow Tobias\ gähnt,$$

$$S \rightarrow gähnt\ Tobias,$$

$$S \rightarrow Tobias\ kennt\ NP,$$

$$S \rightarrow kennt\ Tobias\ NP,$$

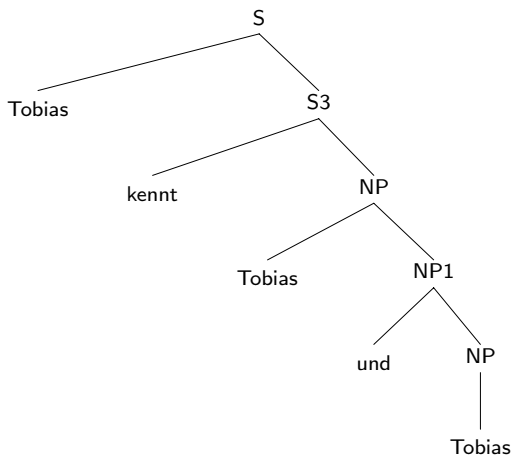
$$S \rightarrow S\ und\ S,$$

$$NP \rightarrow Tobias\ und\ NP,$$

$$NP \rightarrow Tobias \},$$

$$S \rangle$$

reguläre Grammatik ...



$S \rightarrow Tobias\ S1,$
 $S1 \rightarrow gähnt,$
 $S1 \rightarrow gähnt\ S4,$
 $S1 \rightarrow kennt\ NP,$
 $S \rightarrow gähnt\ S2,$
 $S2 \rightarrow Tobias,$
 $S2 \rightarrow Tobias\ S4,$
 $S \rightarrow kennt\ S3,$
 $S3 \rightarrow Tobias\ NP,$
 $NP \rightarrow Tobias,$
 $NP \rightarrow Tobias\ S4,$
 $NP \rightarrow Tobias\ NP1,$
 $NP1 \rightarrow und\ NP,$
 $S4 \rightarrow und\ S$

Chomsky's Diktum 1

*"[...] if you invent a computer language, it doesn't really matter which rules you pick to characterize its expressions; it's the **expressions** that are the language, not the specific computational system that characterizes them.*

*That's not the way natural language works. In natural language there is something in the head, which IS the computational system. The **generative system** is something real, as real as the liver; the utterances generated are like an epiphenomenon. This is the opposite point of view."*

(CHOMSKY 1999:8)

Anforderungen an ein Grammatikmodell

Beobachtungsadäquatheit

Korrekte Erfassung der sprachlichen Daten

Beschreibungsadäquatheit

Korrekte Erfassung der Regularitäten eines Sprachsystems

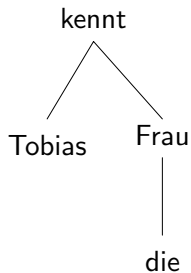
Erklärungsadäquatheit

Hypothese über die universelle Sprachausstattung

dependenzorientierte Ansätze ...

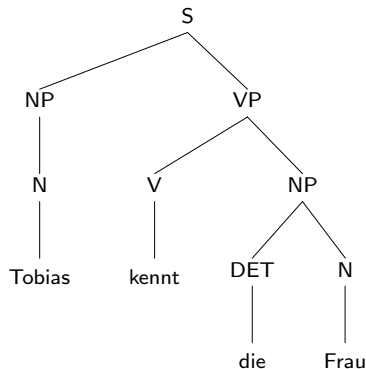
Relationen zwischen Wörtern (beginnend mit Tesni'ere:1980)

Tobias kennt die Frau

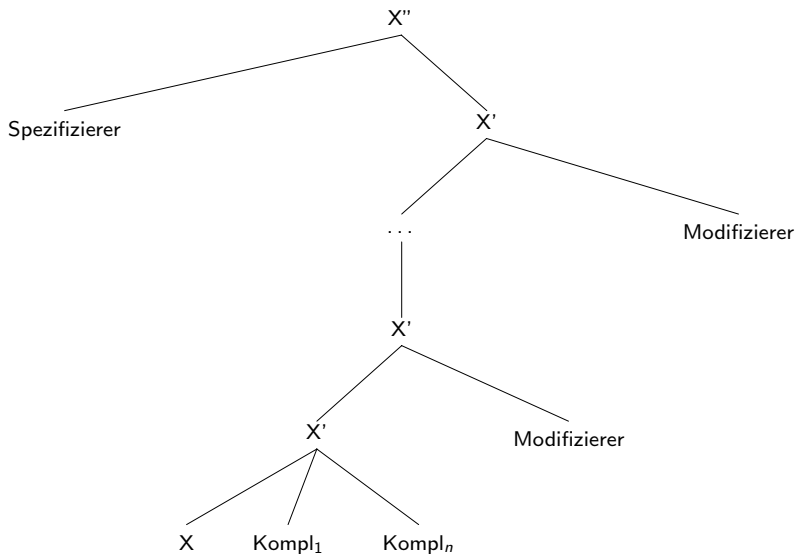


konstituentenorientierte Ansätze ...

Relationen zwischen Konstituenten (beginnend mit Chomsky:1957)



X-Bar-Theorie (Jackendoff, 1977)



Chomsky's Diktum 2

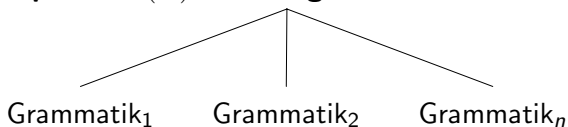
*"For formal languages, there is no "right" grammar; it's arbitrary, you pick any one you like. So by analogy, in language the linguist can pick any grammar, depending on one or another concern or interest; the only thing that is real is the **utterances**.*

*That's a false analogy to start with; human languages are biological objects. What is real – what is in the brain – is a **particular procedure** for characterizing information about sound, meaning and structural organization of linguistic expressions."*

(CHOMSKY 1999:16)

Chomsky's Diktum 2

Sprache $L(G)$ als Menge von Ausdrücken



die (mentale) **Grammatik**

natürliche Sprache (auch Menge von Äußerungen?)

Grammatikformalismen für das Kompetenzmodell

moderne Unifikationsgrammatiken

- Generalisierte Phrasenstrukturgrammatik (GPSG)
(Gazdar, Klein, Pullum und Sag, 1985)
- Lexikalisch Funktionale Grammatik (LFG)
(Bresnan, 1982a, 2001; Berman und Frank, 1996;
Berman, 2003)
- Head-Driven Phrase Structure Grammar (HPSG)
(Pollard und Sag, 1987, 1994; Muller, 1999a, 2002)

didaktische Motivation ...

Warum formal:

ein besseres "Gefühl" für sprachliches Regelwissen

durch Modellierung der zugrundeliegenden
Strukturen für Sprachausschnitte des Deutschen

klares Verständnis der Zusammenhänge in der Satzstruktur

durch selbständiges Entwerfen und Testen des
Computer-Codes

Behandlung mehr komplexer Fragestellungen

die sonst wegen der zu aufwendigen Mathematik
vermieden werden müssten

WERKZEUGE – Prolog und DCGs ...

Prolog:

- deklarative Umgebung
- ermöglicht Abstrahieren von prozeduraler Aspekte
- "eingebauter" Top-Down-Parser mit Backtracking
- eigentlich DIE Programmiersprache für Fragestellungen der Linguistik und der KI

WERKZEUGE – Prolog und DCGs ...

DCG:

- von Stuart Shieber (SHIEBER 1986) als **Werkzeugformalismus** (*tool formalism*) charakterisiert
- ermöglicht Abstrahieren von programminterner Aspekte
- ermöglicht direkte Modellierung von Phrasenstruktursyntaxen

np --> det, n.

- ermöglicht Erweiterung v. PSGn durch (komplexe) Merkmale

np(Num, Kas) --> det(Num, Kas, Gen), n(Num, Kas, Gen).

Modellentwicklung ...

... Generalisierung über Wortarten, Kategorien, lineare Abfolge und Hierarchiebeziehungen

n	\rightarrow	<i>Tobias</i>		S	\rightarrow	NP	VP	
v_{intr}	\rightarrow	<i>gähnt</i>		S	\rightarrow	VP	NP	
v_{tr}	\rightarrow	<i>kennt</i>		NP	\rightarrow	n		
$conj$	\rightarrow	<i>und</i>		VP	\rightarrow	v_{intr}		
				VP	\rightarrow	v_{tr}	NP	
				S	\rightarrow	S	$conj$	S
				NP	\rightarrow	NP	$conj$	NP

Merkmalsstrukturen

Modellierung komplexer Kategorien durch eine Menge von Merkmalspezifikationen
(Paare grammatischer **Merkmale** und zugehöriger **Werte**)

$$S = \begin{bmatrix} \text{MERKMAL}_1 & \text{Wert}_1 \\ \dots & \dots \\ \text{MERKMAL}_n & \text{Wert}_n \end{bmatrix}$$

z.B für Wortform *Hundes*: $N_{sg/mask/gen} = \begin{bmatrix} \text{CAT} & N \\ \text{KASUS} & gen \\ \text{NUM} & sg \\ \text{GENUS} & mask \end{bmatrix}$

Merkmalsstrukturen

Werte der Merkmale:

- **atomar**

oder

- **komplex**

Hundes

CAT	<i>N</i>								
AGR	<table style="border-collapse: collapse; margin-left: auto; margin-right: auto;"> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px 5px 5px 5px;">KASUS</td> <td style="padding: 5px 5px 5px 5px;"><i>gen</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px 5px 5px 5px;">NUM</td> <td style="padding: 5px 5px 5px 5px;"><i>sg</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px 5px 5px 5px;">GENUS</td> <td style="padding: 5px 5px 5px 5px;"><i>mask</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px 5px 5px 5px;">PERS</td> <td style="padding: 5px 5px 5px 5px;"><i>3</i></td> </tr> </table>	KASUS	<i>gen</i>	NUM	<i>sg</i>	GENUS	<i>mask</i>	PERS	<i>3</i>
KASUS	<i>gen</i>								
NUM	<i>sg</i>								
GENUS	<i>mask</i>								
PERS	<i>3</i>								

Merkmalsstrukturen

Werte der Merkmale:

Hund

- **atomar**

oder

- **komplex**

CAT	<i>N</i>								
AGR	<table style="border-collapse: collapse; margin-left: auto; margin-right: auto;"> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px;">KASUS</td> <td style="padding: 5px;"><i>nom</i> ∨ <i>dat</i> ∨ <i>acc</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px;">NUM</td> <td style="padding: 5px;"><i>sg</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px;">GENUS</td> <td style="padding: 5px;"><i>mask</i></td> </tr> <tr> <td style="border-left: 1px solid black; border-right: 1px solid black; padding: 5px;">PERS</td> <td style="padding: 5px;"><i>3</i></td> </tr> </table>	KASUS	<i>nom</i> ∨ <i>dat</i> ∨ <i>acc</i>	NUM	<i>sg</i>	GENUS	<i>mask</i>	PERS	<i>3</i>
KASUS	<i>nom</i> ∨ <i>dat</i> ∨ <i>acc</i>								
NUM	<i>sg</i>								
GENUS	<i>mask</i>								
PERS	<i>3</i>								

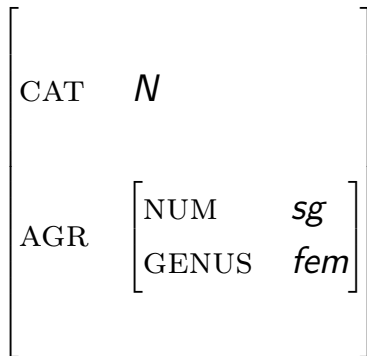
Merkmalsstrukturen

Unterspezifikation:

In einer Merkmalsstruktur können Merkmale, deren Werte noch nicht bekannt sind, unspezifiziert belassen werden

Spezifikation erfolgt nach Auswertung weiterer Information (Kontext)

Katze



- ... kein Info über Kasus (hier egal, wird vom DET bestimmt)
- ... Person (3) wird vom V in die komplette NP durchgereicht

Merkmalsstrukturen

Unifikation:

Lücken in partiellen Beschreibungen werden im Verlauf der Unifikation gefüllt

Reihenfolge der Merkmalsstrukturen, die unifiziert werden, spielt keine Rolle!

Katze

$$\left[\begin{array}{l} \text{CAT} \\ \text{AGR} \end{array} \right. \left[\begin{array}{l} \mathbf{N} \\ \left[\begin{array}{ll} \text{NUM} & \textit{sg} \\ \text{GENUS} & \textit{fem} \end{array} \right] \end{array} \right]$$

die

$$\left[\begin{array}{l} \text{CAT} \\ \text{AGR} \end{array} \right. \left[\begin{array}{l} \mathbf{DET} \\ \left[\begin{array}{ll} \text{KASUS} & \textit{nom} \\ \text{GENUS} & \textit{fem} \\ \text{NUM} & \textit{sg} \end{array} \right] \end{array} \right]$$

Wie geht das?

Unifikation:

Lücken in partiellen Beschreibungen werden im Verlauf der Unifikation gefüllt

Katze

$$\left[\begin{array}{ll} \text{CAT} & \mathbf{N} \\ \text{AGR} & \left[\begin{array}{ll} \text{NUM} & \textit{sg} \\ \text{GENUS} & \textit{fem} \end{array} \right] \end{array} \right]$$

NP →

$$\left[\begin{array}{ll} \text{CAT} & \mathbf{DET} \\ \text{AGR} & \left[\begin{array}{ll} \text{KASUS} & \alpha \\ \text{NUM} & \beta \\ \text{GENUS} & \gamma \end{array} \right] \end{array} \right]$$

$$\left[\begin{array}{ll} \text{CAT} & \mathbf{N} \\ \text{AGR} & \left[\begin{array}{ll} \text{KASUS} & \alpha \\ \text{NUM} & \beta \\ \text{GENUS} & \gamma \end{array} \right] \end{array} \right]$$

Von einfachen zu komplexen Kategorien ...

NP → **DET N**

$$\left[\begin{array}{l} \text{CAT} \quad \mathbf{NP} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \alpha \\ \text{NUM} \quad \beta \end{array} \right] \end{array} \right] \rightarrow \left[\begin{array}{l} \text{CAT} \quad \mathbf{DET} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \alpha \\ \text{NUM} \quad \beta \\ \text{GENUS} \quad \gamma \end{array} \right] \end{array} \right] \quad \left[\begin{array}{l} \text{CAT} \quad \mathbf{N} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \alpha \\ \text{NUM} \quad \beta \\ \text{GENUS} \quad \gamma \end{array} \right] \end{array} \right]$$

mit $\alpha \in \{\text{nom}, \text{gen}, \text{dat}, \text{akk}\}$, $\beta \in \{\text{sg}, \text{pl}\}$ und $\gamma \in \{\text{mask}, \text{fem}, \text{neut}\}$

$$\left[\begin{array}{l} \text{CAT} \quad \mathbf{NP} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \text{nom} \\ \text{NUM} \quad \text{sg} \end{array} \right] \end{array} \right] \rightarrow \left[\begin{array}{l} \text{CAT} \quad \mathbf{DET} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \text{nom} \\ \text{NUM} \quad \text{sg} \\ \text{GENUS} \quad \text{fem} \end{array} \right] \end{array} \right] \quad \left[\begin{array}{l} \text{CAT} \quad \mathbf{N} \\ \text{AGR} \quad \left[\begin{array}{l} \text{KASUS} \quad \text{nom} \\ \text{NUM} \quad \text{sg} \\ \text{GENUS} \quad \text{fem} \end{array} \right] \end{array} \right]$$

die

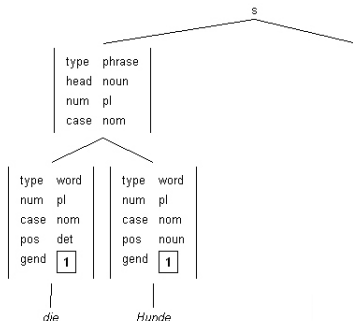
Katze

Beispiel 1: sysana

% Regel fuer Nominalphrase:

```
type:phrase  
..head:noun  
..num:NUMERUS  
..case:KASUS
```

```
--> type:word  
..pos:det  
..num:NUMERUS  
..case:KASUS  
..gend:GENUS  
,  
type:word  
..pos:noun  
..num:NUMERUS  
..case:KASUS  
..gend:GENUS  
.
```



% lexikalische Regeln:

```
type:word  
..pos:det  
..num:pl  
..case:nom --> "die".
```

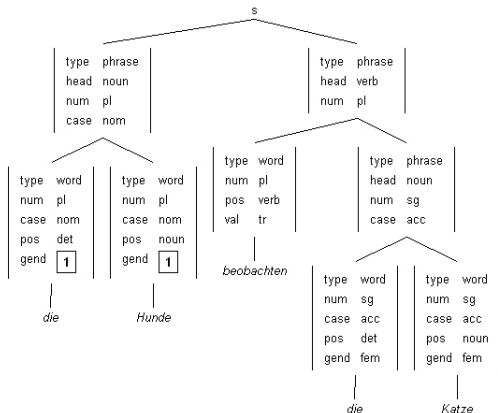
```
type:word  
..pos:noun  
..num:pl  
..case:nom --> "Hunde".
```

Beispiel 1: sysana

% Regel fuer Nominalphrase:

```
type:phrase  
..head:noun  
..num:NUMERUS  
..case:KASUS
```

```
--> type:word  
..pos:det  
..num:NUMERUS  
..case:KASUS  
..gend:GENUS  
,  
type:word  
..pos:noun  
..num:NUMERUS  
..case:KASUS  
..gend:GENUS  
.
```



% lexikalische Regeln:

```
type:word  
..pos:det  
..num:pl  
..case:nom --> "die".
```

```
type:word  
..pos:noun  
..num:pl  
..case:nom --> "Hunde".
```

Beispiel 2: trale

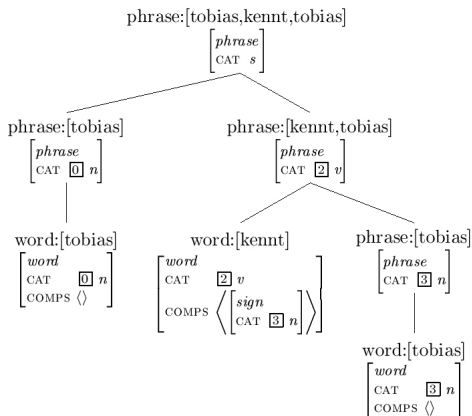
```

s_rule ## (phrase,
  dtrs: [D1,D2],
  cat:s) ==>
  cat> (D1,phrase,
    cat:n),
  cat> (D2,phrase,
    cat:v).

xp_rule1 ## (phrase,
  dtrs: [D1],
  cat:Cat) ==>
  cat> (D1,word,
    cat:Cat,
    comps: []).

xp_rule2 ## (phrase,
  dtrs: [D1,D2],
  cat:Cat) ==>
  cat> (D1,word,
    cat:Cat,
    comps: [(cat:Compl)]),
  cat> (D2,phrase,
    cat:Compl).

```



% lexikalische Regeln:

```

tobias --> (word,
  cat:n,
  comps: []).

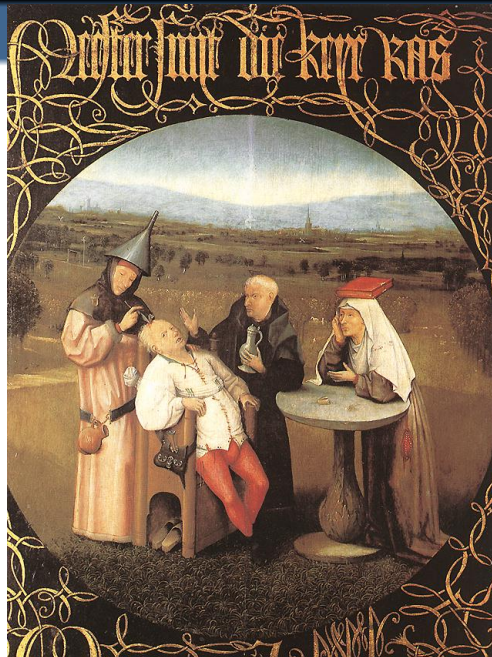
kennt --> (word,
  cat:v,
  comps: [(cat:n)]).

```


Vermittlung von Wissen an Lernende?

Vermittlung von Wissen an Lernende?

Das Steinschneiden
Hieronymus Bosch
15. Jahrhundert



**Vermittlung
von Wissen
an Lernende?**

oder gemeinsame
**Konstruktion
von Wissen**

Student

21. Jahrhundert

