

# Einführung in die Computerlinguistik

## Merkmalsstrukturen (Feature Structures)

Laura Kallmeyer  
Heinrich-Heine-Universität Düsseldorf  
Sommersemester 2013

---

Merkmalsstrukturen 1 Sommersemester 2013

### Überblick

1. Einleitung
2. Merkmalsstrukturen als Graphen
3. Subsumption und Unifikation
4. Getypte Merkmalsstrukturen
5. Erweiterungen

---

Merkmalsstrukturen 2 Sommersemester 2013

### Einleitung (1)

Die in CFGs verwendeten Nichtterminalen sind in der Regel nicht ausreichend, um linguistische Generalisierungen auszudrücken.

Bsp. Kongruenz

$$S \rightarrow \text{NPSg VPSg} \quad S \rightarrow \text{NPPI VPPi}$$

Besser:

$$S \rightarrow \text{NP VP} \quad \text{Bedingung: gleicher Numerus in NP und VP}$$


---

Merkmalsstrukturen 3 Sommersemester 2013

### Einleitung (2)

Um diese Generalisierung auszudrücken, faktorisieren wir die Nichtterminalen:

- Ein Nichtterminales ist nicht mehr atomar, sondern hat eine Struktur.
- Der Inhalt des Nichtterminalen wird durch **Attribute** beschrieben, die bestimmte **Werte** haben können.

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \text{PI} \end{bmatrix}$$

- Man kann sich an verschiedenen Stellen auf denselben Attributwert beziehen (**Structure Sharing**)

$$\begin{bmatrix} \text{CAT} & \text{S} \end{bmatrix} \rightarrow \begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \boxed{1} \end{bmatrix} \begin{bmatrix} \text{CAT} & \text{VP} \\ \text{NUM} & \boxed{1} \end{bmatrix}$$

Die Variable  $\boxed{1}$  denotiert immer den gleichen Wert.

---

Merkmalsstrukturen 4 Sommersemester 2013

**Einleitung (3)**

**Unterspezifikation:** Es sind nicht immer alle Merkmale bekannt. Anstatt die verschiedenen Möglichkeiten explizit aufzulisten, gibt man nur die Merkmale an, die feststehen.

$$\begin{bmatrix} \text{CAT} & \text{N} \\ \text{NUM} & \text{Sg} \\ \text{GEN} & \text{m} \end{bmatrix} \rightarrow \text{man} \quad \begin{bmatrix} \text{CAT} & \text{N} \\ \text{GEN} & \text{n} \end{bmatrix} \rightarrow \text{fish}$$

$$\begin{bmatrix} \text{CAT} & \text{Det} \\ \text{NUM} & \text{Sg} \end{bmatrix} \rightarrow \text{a} \quad \begin{bmatrix} \text{CAT} & \text{Det} \end{bmatrix} \rightarrow \text{the}$$

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{NUM} & \boxed{1} \\ \text{GEN} & \boxed{2} \\ \text{PERS} & 3 \end{bmatrix} \rightarrow \begin{bmatrix} \text{CAT} & \text{Det} \\ \text{NUM} & \boxed{1} \end{bmatrix} \begin{bmatrix} \text{CAT} & \text{N} \\ \text{NUM} & \boxed{1} \\ \text{GEN} & \boxed{2} \end{bmatrix}$$

**Einleitung (4)**

Attribute müssen nicht notwendig atomare Werte haben. Der Wert eines Attributs kann wiederum eine Merkmalsstruktur sein:

$$\begin{bmatrix} \text{CAT} & \text{N} \\ \text{AGR} & \begin{bmatrix} \text{GEN} & \text{n} \end{bmatrix} \end{bmatrix} \rightarrow \text{fish}$$

$$\begin{bmatrix} \text{CAT} & \text{Det} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \end{bmatrix} \end{bmatrix} \rightarrow \text{a}$$

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{AGR} & \boxed{1} \begin{bmatrix} \text{PERS} & 3 \end{bmatrix} \end{bmatrix} \rightarrow \begin{bmatrix} \text{CAT} & \text{Det} \\ \text{AGR} & \boxed{1} \end{bmatrix} \begin{bmatrix} \text{CAT} & \text{N} \\ \text{AGR} & \boxed{1} \end{bmatrix}$$

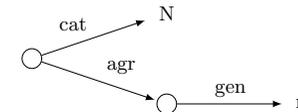
**Merkmalsstrukturen als Graphen (1)**

Merkmalsstrukturen werden in der Regel als gerichtete Graphen formalisiert.

Es gibt zwei Möglichkeiten: Man kann eine Attribut-Wert-Matrix

wie  $\begin{bmatrix} \text{CAT} & \text{N} \\ \text{AGR} & \begin{bmatrix} \text{GEN} & \text{n} \end{bmatrix} \end{bmatrix}$  als

1. einen gerichteten Graphen auffassen

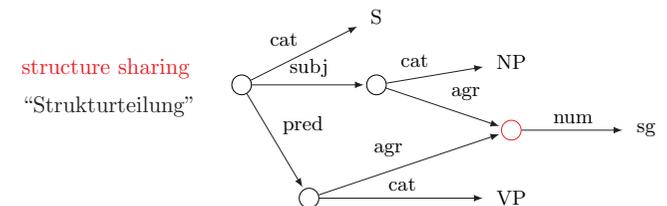


2. oder als eine Beschreibung derartiger Graphen, die im Prinzip von einer unendlichen Menge von Graphen erfüllt wird.

**Merkmalsstrukturen als Graphen (2)**

Merkmalsstrukturen als Graphen sind nicht notwendig Bäume:

$$\begin{bmatrix} \text{CAT} & \text{S} \\ \text{SUBJ} & \begin{bmatrix} \text{CAT} & \text{NP} \\ \text{AGR} & \boxed{1} \begin{bmatrix} \text{NUM} & \text{sg} \end{bmatrix} \end{bmatrix} \\ \text{PRED} & \begin{bmatrix} \text{CAT} & \text{VP} \\ \text{AGR} & \boxed{1} \end{bmatrix} \end{bmatrix}$$



### Merkmalsstrukturen als Graphen (3)

Wichtig: Die Attribute sind funktional, d.h., sie müssen, wenn vorhanden, eindeutige Werte haben.

⇒ Bedingung auf die Graphen: für jeden Knoten  $v$  und jedes Attribut  $a$  gibt es nur höchstens eine von  $v$  ausgehende Kante, die das Label  $a$  hat.

Mit anderen Worten: Jedes Attribut  $a$  kann aufgefasst werden als eine partielle Funktion auf der Menge der Knoten des Graphen einer Merkmalsstruktur.

### Subsumption und Unifikation (1)

Subsumption: Relation auf Merkmalsstrukturen:  $S_1$  subsumiert  $S_2$  ( $S_1 \sqsubseteq S_2$ ), wenn in  $S_2$  mindestens die Information aus  $S_1$  enthalten ist.

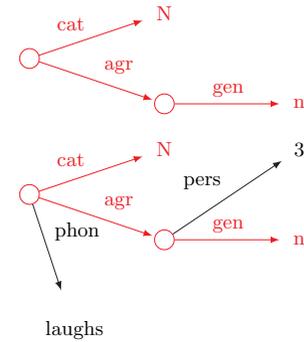
$$\text{Bsp. } S_1: \begin{bmatrix} \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \end{bmatrix} \end{bmatrix} \quad S_2: \begin{bmatrix} \text{ORTH} & \text{laughs} \\ \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{PERS} & 3 \\ \text{NUM} & \text{Sg} \end{bmatrix} \end{bmatrix}$$

$$S_1 \sqsubseteq S_2$$

Mit anderen Worten: Es gibt eine injektive Abbildung der Knoten aus  $S_1$  in die Knoten aus  $S_2$ , die kanten- und labelerhaltend ist, also einen Homomorphismus.

### Subsumption und Unifikation (2)

$S_1$  als Graph und sein homomorphes Bild in  $S_2$ :



### Subsumption und Unifikation (3)

Subsumption ist eine **partielle Ordnung**, d.h. sie ist

1. reflexiv: Jede Struktur subsumiert sich selbst  $S \sqsubseteq S$  für alle  $S$ ;
2. transitiv: aus  $S_1 \sqsubseteq S_2$  und  $S_2 \sqsubseteq S_3$  folgt  $S_1 \sqsubseteq S_3$  für alle  $S_1, S_2, S_3$ ;
3. antisymmetrisch: aus  $S_1 \sqsubseteq S_2$  und  $S_2 \sqsubseteq S_1$  folgt  $S_1 = S_2$ .

Die leere Merkmalsstruktur  $[\ ]$  subsumiert alle anderen Merkmalsstrukturen.

### Subsumtion und Unifikation (4)

Eine Merkmalsstruktur  $S$  heißt **Unifikation** von  $S_1$  und  $S_2$  ( $S_1 \sqcup S_2$ ), wenn  $S$  sowohl von  $S_1$  als auch von  $S_2$  subsumiert wird und wenn außerdem  $S$  alle anderen Merkmalsstrukturen subsumiert, die ebenfalls von  $S_1$  und von  $S_2$  subsumiert werden.

$$\begin{bmatrix} \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \end{bmatrix} \end{bmatrix} \sqcup \begin{bmatrix} \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{PERS} & 3 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \\ \text{PERS} & 3 \end{bmatrix} \end{bmatrix}$$

Damit  $\sqcup$  immer definiert ist, führen wir ein Symbol  $\perp$  für die inkonsistente Merkmalsstruktur ein, die von allen Merkmalsstrukturen subsumiert wird.

$$\begin{bmatrix} \text{CAT} & \text{NP} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \end{bmatrix} \end{bmatrix} \sqcup \begin{bmatrix} \text{CAT} & \text{V} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & \text{Sg} \\ \text{PERS} & 3 \end{bmatrix} \end{bmatrix} = \perp$$

### Subsumtion und Unifikation (5)

Die Menge der Merkmalsstrukturen, versehen mit der Relation  $\sqsubseteq$  bildet dann einen **Verband** (lattice):  $\sqsubseteq$  ist eine partielle Ordnung und zusätzlich gilt für alle Merkmalsstrukturen  $S_1, S_2$ :

- (sup) Es gibt eine Merkmalsstruktur  $S$ , so dass  $S_1 \sqsubseteq S$  und  $S_2 \sqsubseteq S$  und so dass  $S$  alle anderen Merkmalsstrukturen subsumiert, die ebenfalls von  $S_1$  und  $S_2$  subsumiert werden.  $S$  heißt dann das **Supremum** von  $\{S_1, S_2\}$ .
- (inf) Es gibt eine Merkmalsstruktur  $S$ , so dass  $S \sqsubseteq S_1$  und  $S \sqsubseteq S_2$  und so dass  $S$  von allen anderen Merkmalsstrukturen subsumiert wird, die ebenfalls  $S_1$  und  $S_2$  subsumieren.  $S$  heißt dann das **Infimum** von  $\{S_1, S_2\}$ .

Das bezüglich  $\sqsubseteq$  kleinste Element ist  $[\ ]$ , das größte Element ist  $\perp$ .

### Getypte Merkmalsstrukturen (1)

Die obigen Beispiele sind implizit davon ausgegangen, dass CAT die syntaktische Kategorie angibt und AGR für Kongruenz zuständig ist. D.h., folgende Merkmalsstrukturen sollten nicht möglich sein:

$$\begin{bmatrix} \text{CAT} & \text{Sg} \\ \text{AGR} & \begin{bmatrix} \text{NUM} & 3 \\ \text{PERS} & \text{V} \end{bmatrix} \end{bmatrix} \quad \begin{bmatrix} \text{CAT} & \begin{bmatrix} \text{AGR} & \begin{bmatrix} \text{PERS} & 3 \end{bmatrix} \end{bmatrix} \end{bmatrix}$$

Bisher haben wir diese Merkmalsstrukturen aber nicht ausgeschlossen, da wir die entsprechende Generalisierung nicht definieren können.

Ziel: Einschränkungen formulieren von der Art "eine Agreement Merkmalsstruktur kann höchstens die Attribute NUM, PERS und GEN haben".

### Getypte Merkmalsstrukturen (2)

Daher führen wir **Typen** für Merkmalsstrukturen ein:

- Jede Merkmalsstruktur hat eine Typ  $\tau$ .
- Für jeden Typ  $\tau$  ist festgelegt, welche Attribute für ihn zugelassen sind und von welchem Typ wiederum die Werte dieser Attribute sind.
- Typen sind in einer Typenhierarchie organisiert, in der spezifischere Typen von allgemeineren erben können.
- Unifikation wird erweitert, so dass Typen berücksichtigt werden.

### Getypte Merkmalsstrukturen (3)

Die Festlegung von Typen und ihren möglichen Attributen setzt sich aus einer Typenhierarchie und Attributangaben für einzelne Typen zusammen.

$$\begin{array}{l} \left[ \begin{array}{ll} agr\text{-}structure & \\ AGR & agr \end{array} \right] \\ \left[ \begin{array}{ll} determiner & \\ QUANT & quant \end{array} \right] \end{array} \quad \begin{array}{l} \left[ \begin{array}{ll} agr & \\ NUM & num \\ GEN & gen \\ PERS & pers \end{array} \right] \\ \left[ \begin{array}{ll} noun & \\ CASE & case \end{array} \right] \\ \left[ \begin{array}{ll} syncat & \\ CAT & cat \end{array} \right] \end{array}$$

Typ *quant*: {every, most, some, none}, Typ *num*: {Sg, Pl}, Typ *gen*: {m,f,n}, Typ *pers*: {1, 2, 3}, Typ *case*: {nom, acc, dat}, Type *cat*: {N, V, NP, VP, S, ...}

### Getypte Merkmalsstrukturen (4)

$$\begin{array}{l} \left[ \begin{array}{ll} agr\text{-}structure & \\ AGR & agr \end{array} \right] \\ \left[ \begin{array}{ll} determiner & \\ QUANT & quant \end{array} \right] \end{array} \quad \begin{array}{l} \left[ \begin{array}{ll} syncat & \\ CAT & cat \end{array} \right] \\ \left[ \begin{array}{ll} noun & \\ CASE & case \end{array} \right] \end{array} \quad \begin{array}{l} \left[ \begin{array}{ll} agr & \\ NUM & num \\ GEN & gen \\ PERS & pers \end{array} \right] \end{array}$$

Typenhierarchie:

```

    agr-structure  syncat
      /           /
     /           /
    /             /
   /             /
  /             /
 /             /
/             /
determiner  noun
  
```

D.h., die Attribute von *noun* werden bestimmt durch die Klassen *agr-structure*, *syncat* und *noun*.

$$\left[ \begin{array}{ll} noun & \\ CAT & N \\ CASE & acc \\ AGR & \left[ \begin{array}{ll} agr & \\ NUM & Sg \\ GEN & m \\ PERS & 3 \end{array} \right] \end{array} \right]$$

### Erweiterungen (1)

Einige linguistische Theorien verwenden auch mengen- oder listenwertige Attribute.

Bsp.: Head-Driven Phrase Structure Grammar (HPSG) kodiert Syntaxbäume in Form von Merkmalsstrukturen, wobei die Töchter eines Knotens in Form einer Liste als Wert des Attributs DTRS ("daughters") formalisiert werden.

$$\left[ \begin{array}{ll} phrase & \\ DTRS & \left\langle \left[ \begin{array}{ll} CAT & PRO \\ ORTH & I \end{array} \right], \left[ \begin{array}{ll} CAT & VP \\ DTRS & \left\langle \left[ \begin{array}{ll} CAT & V \\ ORTH & love \end{array} \right], \left[ \begin{array}{ll} CAT & NP \\ ORTH & New\ York \end{array} \right] \right\rangle \end{array} \right\rangle \end{array} \right]$$

### Erweiterungen (2)

- Manche Systeme arbeiten direkt mit Merkmalsstrukturen, also mit Graphen.
- Andere verwenden Beschreibungen von Merkmalsstrukturen.

Vorteil von Beschreibungen: Je nach verwendeter Logik große Expressivität (die man sich allerdings in der Regel auch mit einer entsprechenden Komplexität erkaufte). Dinge, die nützlich sein könnten:

1. Disjunktion:  $CASE = acc \vee CASE = dat$
2. Negation:  $\neg(CASE = nom)$
3. Pfadungleichheiten:  $SUBJ [CASE] \neq OBJ [CASE]$