

Syntax natürlicher Sprachen

Vorlesung 10: Getypte Merkmalstrukturen und Unifikation

A. Wisiorek

Folien: Martin Schmitt

Centrum für Informations- und Sprachverarbeitung,
Ludwig-Maximilians-Universität München

10.01.2023

- Logik-Programmierung (z. B. Prolog)
- Pattern Matching (z. B. in funktionalen Programmiersprachen)
- Typinferenz (vor allem funktionale Programmiersprachen wie Haskell, Scala, etc. Eingeschränkt aber z. B. auch C#)
- **Merkmalstrukturen** (zur Beschreibung komplexer Objekte, z. B. grammatischer Merkmale)

Carpenter, Bob (1992). *The Logic of Typed Feature Structures*. Cambridge Tracts in Theoretical Computer Science. Cambridge University Press.

1. Formale Grundlagen

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - Bedingungen

- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

1.1. Merkmalstrukturen

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - Bedingungen
- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

Definition (Type, \sqsubseteq)

Sei **Type** eine endliche Menge von Typen mit **Vererbungshierarchie** \sqsubseteq .

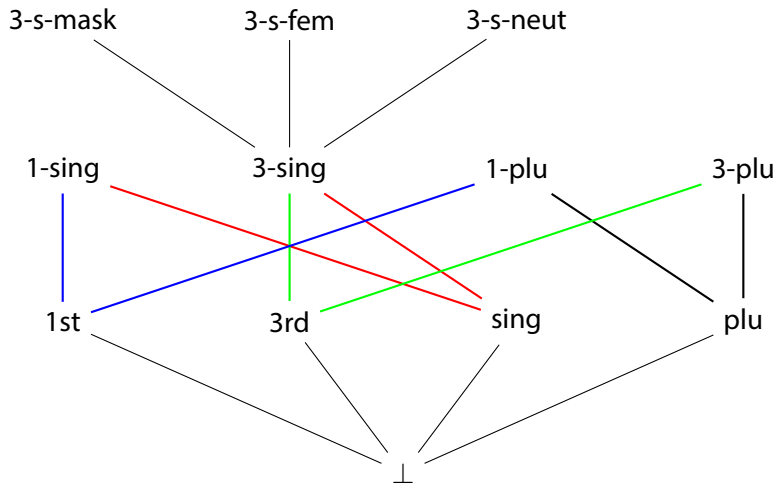
Wenn für $A, B \in \mathbf{Type}$ gilt, dass $A \sqsubseteq B$, dann

- *erbt* B Informationen von A .
- ist A Obertyp von B . (Alle A -Attribute sind auch B -Attribute.)
- A *subsumiert* B (B wird von A subsumiert).
- ist A „allgemeiner oder gleich“ B .
- ist B „spezieller oder gleich“ A .

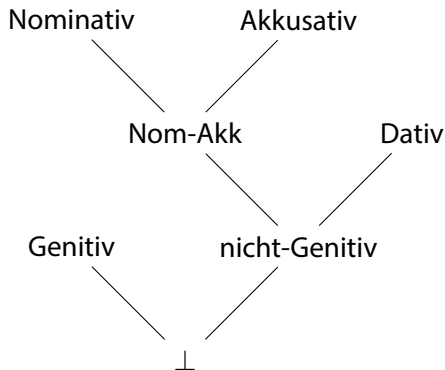
Eigenschaften von Type , \sqsubseteq

- wohldefinierte Unifikationsoperation
- transitiv ($\forall A, B, C \in \text{Type}. A \sqsubseteq B \wedge B \sqsubseteq C \implies A \sqsubseteq C$)
- reflexiv ($\forall A \in \text{Type}. A \sqsubseteq A$)
- antisymmetrisch ($\forall A, B \in \text{Type}. A \sqsubseteq B \wedge B \sqsubseteq A \implies A = B$) (keine Vererbungsschleifen)
 - \implies *partielle Ordnung* (d. h. nicht alle Elemente von **Type** müssen miteinander vergleichbar sein)
- Existenz eines eindeutigen allgemeinsten Typs ($\exists_1 A \in \text{Type}. \forall B \in \text{Type}. A \sqsubseteq B$)
 - $\implies \perp$ definiert als kleinstes Element von **Type** bzgl. \sqsubseteq

Beispiel: Typhierarchie



Noch ein Beispiel: Typhierarchie



Vgl. z. B. die Paradigmen:

der Hund, des **Hundes**, dem Hund, den Hund

das Buch, des **Buches**, dem Buch, **das** Buch

Definition (Feat)

Sei **Feat** eine *endliche* Menge von Merkmalen (engl. *features*).

(Ohne weitere Anforderungen an Struktur oder Eigenschaften)

Beispiel

Feat = {GEN, CASE, NUM, AGR, PER, MOOD, CAT, TENSE}

Definition

Eine Merkmalstruktur über **Type** und **Feat** ist definiert als Tupel

$F = (Q, \bar{q}, \theta, \delta)$ mit:

- Q : endliche Menge von Knoten (Einträge)
- $\bar{q} \in Q$: Wurzelknoten
- $\theta: Q \rightarrow \mathbf{Type}$: totale Typisierungsfunktion
- $\delta: \mathbf{Feat} \times Q \rightarrow Q$: partielle Merkmal-Wert-Funktion

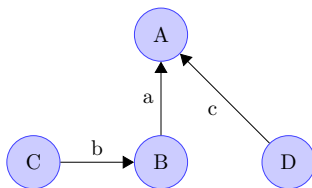
Sei \mathcal{F} die Menge aller Merkmalstrukturen.

$$\left[\begin{array}{cc} \text{CAT} & N \\ \text{AGR} & \left[\begin{array}{cc} \text{CASE} & \textit{Nom} \\ \text{NUM} & \textit{Pl} \end{array} \right] \end{array} \right]$$

Beschrifteter Graph

Ein *beschrifteter* Graph ist definiert als Tupel $G = (V, E, l_V, l_E, L_V, L_E)$ mit

- V : Menge der Knoten (engl. *vertices*)
- $E \subseteq V \times V$: Menge der Kanten (engl. *edges*)
- $l_V: V \rightarrow L_V$: Beschriftungsfunktion für Knoten (engl. *label*)
- $l_E: E \rightarrow L_E$: Beschriftungsfunktion für Kanten
- L_X : Menge von Beschriftungen für X



Visualisierung

Der Graph zu einer Merkmalstruktur $F = (Q, \bar{q}, \theta, \delta)$ ist gegeben durch:

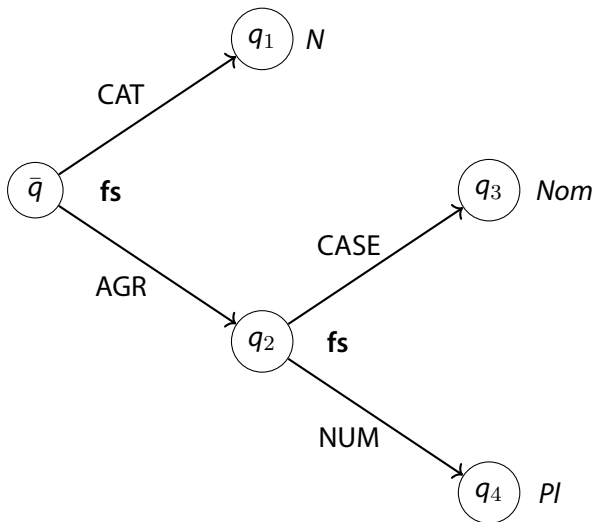
- $V := Q$
- $E := \{(q_1, q_2) \mid \exists f. \delta(f, q_1) = q_2\}$
- $L_V := \mathbf{Type}; l_V := \theta$
- $L_E := \mathbf{Feat}; l_E(q_1, q_2) := \{f \mid \delta(f, q_1) = q_2\}$

Anmerkung

Zur Vereinfachung werden einelementige Mengen ohne Mengenklammern geschrieben.

Also a statt $\{a\}$.

Beispiel: Graphdarstellung



Variablen

- **Var** sei eine abzählbar unendliche Menge von Variablen.
- Häufig wird $\mathbf{Var} = \mathbb{N}$ benutzt.
- Es gibt aber auch andere Möglichkeiten;
z. B. im NLTK: ASCII-Identifer ($?x, ?y, \dots$)

Definition (Zuweisungsfunktion, Valuation)

Eine Zuweisung $\alpha : \mathbf{Var} \rightarrow \mathcal{F}$ ist eine totale Funktion, die alle Variablen an Merkmalstrukturen (Knoten, Einträge) bindet.

Reentrance (dt. *Wiedereintritt*)

Durch das Aufstellen von Bedingungen (s. später) können Variablen an verschiedene Teile von Merkmalstrukturen gebunden werden. *Diese müssen gleich sein.*

Beispiel

ORTH	<i>folgt</i>	
SYN	SBJ	1
	OBJ	2
SEM	AGT	1
	PAT	2

ORTH	<i>folgt</i>	
SYN	SBJ	1 <i>Hund</i>
	OBJ	2 <i>Katze</i>
SEM	AGT	1
	PAT	2

1.2. Subsumption

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - **Subsumption**
 - Unifikation
 - Bedingungen
- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

Erweiterung auf Merkmalstrukturen

$F = (Q, \bar{q}, \theta, \delta)$ subsumiert $F' = (Q', \bar{q}', \theta', \delta')$, genau dann wenn es eine totale Funktion $h : Q \rightarrow Q'$ gibt, sodass:

- $h(\bar{q}) = \bar{q}'$
- $\theta(q) \sqsubseteq \theta'(h(q))$ für alle $q \in Q$
- $h(\delta(f, q)) = \delta'(f, h(q))$ für alle $q \in Q$ und $f \in \mathbf{Feat}$,
für die $\delta(f, q)$ definiert ist

Beispiel

$$\bar{q} \begin{bmatrix} q_1 \text{ CAT} & N \end{bmatrix} \sqsubseteq \bar{q}' \begin{bmatrix} q'_1 \text{ CAT} & N \\ q'_2 \text{ GEN} & \text{mask} \end{bmatrix}$$

$$h(\bar{q}) = \bar{q}' \quad h(q_1) = q'_1$$

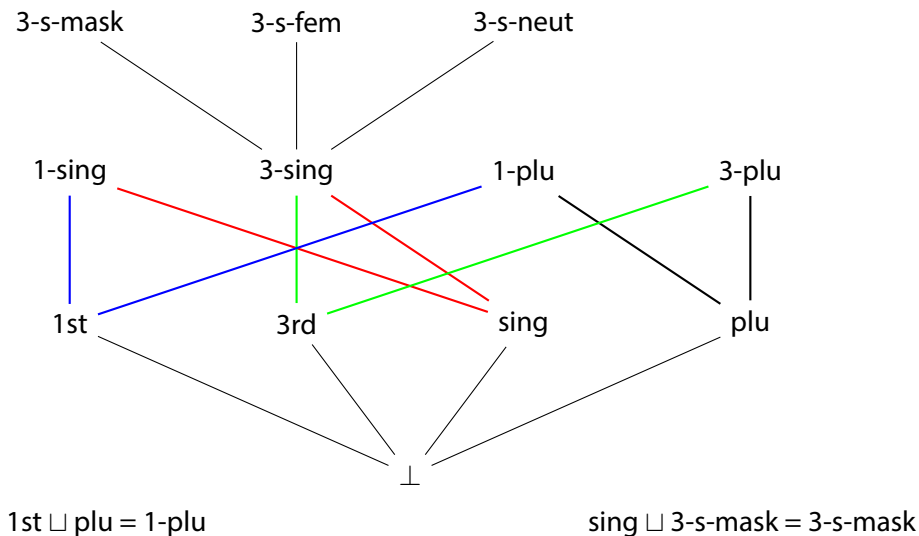
1.3. Unifikation

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - **Unifikation**
 - Bedingungen
- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

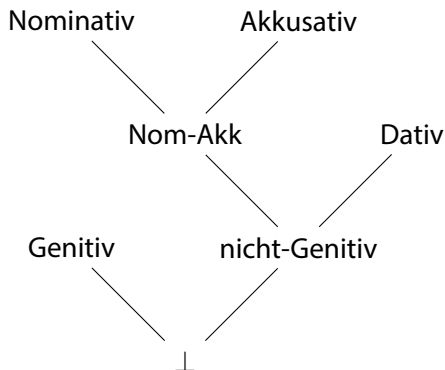
Unifikation (\sqcup) für Typen

- Das Ergebnis der Unifikation zweier Typen $A, B \in \mathbf{Type}$ ist ihre kleinste obere Schranke in \mathbf{Type} bzgl. \sqsubseteq .
- Diese kann auch undefiniert sein (Typen unifizieren nicht).
- $A \sqcup B = C \iff A \sqsubseteq C \text{ und } B \sqsubseteq C \text{ und } \forall D \in \mathbf{Type}. A \sqsubseteq D \wedge B \sqsubseteq D \implies C \sqsubseteq D$
(Vgl. Mengenvereinigung und Untermengenbeziehung)

Beispiel: Typunifikation



Noch ein Beispiel: Typunifikation



nicht-Genitiv \sqcup Nominativ = Nominativ

Nom-Akk \sqcup Dativ = *undefiniert*

Unifikation (\sqcup) für Merkmalstrukturen

- Idee: Unifikation ebenfalls kleinste obere Schranke bzgl. der Subsumptionsrelation \sqsubseteq auf Merkmalstrukturen
- Algorithmus in zwei Schritten:
 - 1 Identifiziere korrespondierende (äquivalente) Knoten
 - 2 Unifiziere deren Typen

Formale Definition: Identifikation (Schritt 1)

Für Merkmalstrukturen $F = (Q, \bar{q}, \theta, \delta)$, $F' = (Q', \bar{q}', \theta', \delta')$ mit $Q \cap Q' = \emptyset$ sei die *Äquivalenzrelation* \equiv wie folgt definiert:

- $\bar{q} \equiv \bar{q}'$
- $\delta(f, q) \equiv \delta'(f, q')$ wenn beide Seiten definiert und $q \equiv q'$

Formale Definition: Typunifikation (Schritt 2)

Die Unifikation von F und F' ist dann wie folgt definiert:

$$F \sqcup F' = ((Q \cup Q')/\equiv, [\bar{q}]_{\equiv}, \theta^{\equiv}, \delta^{\equiv})$$

mit

$$\theta^{\equiv}([q]_{\equiv}) = \bigsqcup \{(\theta \cup \theta')(q') \mid q' \equiv q\}$$

und

$$\delta^{\equiv}(f, [q]_{\equiv}) = \begin{cases} [(\delta \cup \delta')(f, q)]_{\equiv} & \text{falls } (\delta \cup \delta')(f, q) \text{ definiert} \\ \text{undefiniert} & \text{sonst} \end{cases}$$

Notation (für \equiv Äquivalenzrelation über X)

- $[x]_{\equiv} = \{y \in X \mid y \equiv x\}$
- $X/\equiv = \{[x]_{\equiv} \mid x \in X\}$

Beispiel: (Formale) Unifikation

$$q_1 \left[\begin{array}{cc} q_2 \text{ CAT} & N \\ q_3 \text{ AGR} & \left[\begin{array}{cc} q_4 \text{ NUM} & Sg \\ q_5 \text{ CAS} & \textit{nicht-Gen} \end{array} \right] \end{array} \right]$$

$$q_6 \left[\begin{array}{cc} q_7 \text{ ORTH} & \textit{Hund} \\ q_8 \text{ AGR} & \left[\begin{array}{cc} q_9 \text{ NUM} & Sg \\ q_{10} \text{ CAS} & \textit{Nom} \end{array} \right] \end{array} \right]$$

1 Identifikation korrespondierender Knoten

- $q_1 \equiv q_6$ (Initialisierung)
- Nach 1 Schritt mit δ :
 - $q_3 \equiv q_8$
- Nach 2 Schritten mit δ :
 - $q_4 \equiv q_9$
 - $q_5 \equiv q_{10}$

Beispiel: (Formale) Unifikation

$$q_1 \left[\begin{array}{cc} q_2 \text{ CAT} & N \\ q_3 \text{ AGR} & \left[\begin{array}{cc} q_4 \text{ NUM} & Sg \\ q_5 \text{ CAS} & \text{nicht-Gen} \end{array} \right] \end{array} \right] \quad q_6 \left[\begin{array}{cc} q_7 \text{ ORTH} & \text{Hund} \\ q_8 \text{ AGR} & \left[\begin{array}{cc} q_9 \text{ NUM} & Sg \\ q_{10} \text{ CAS} & \text{Nom} \end{array} \right] \end{array} \right]$$

2 Typunifikation

- $Q_U = \{\{q_1, q_6\}, \{q_2\}, \{q_7\}, \{q_3, q_8\}, \{q_4, q_9\}, \{q_5, q_{10}\}\}$
- $\bar{q}_U = \{q_1, q_6\}$
- $\theta^{\equiv}(\{q_2\}) = N, \theta^{\equiv}(\{q_7\}) = \text{Hund}, \theta^{\equiv}(\{q_3, q_8\}) = \mathbf{fs},$
 $\theta^{\equiv}(\{q_4, q_9\}) = Sg, \theta^{\equiv}(\{q_5, q_{10}\}) = \text{Nom}, \theta^{\equiv}(\{q_1, q_6\}) = \mathbf{fs}$
- $\delta(\text{CAT}, \{q_1, q_6\}) = \{q_2\}, \delta(\text{ORTH}, \{q_1, q_6\}) = \{q_7\}, \dots$

Lemma

Wenn $F \sqcup F'$ definiert ist, dann ist $F \sqcup F' \in \mathcal{F}$ eine Merkmalstruktur.

Theorem

$F \sqcup F'$ ist die *kleinste obere Schranke* von F und F' in $(\mathcal{F}, \sqsubseteq)$, falls F und F' eine obere Schranke haben.

Für Beweise siehe (Carpenter:Log-TyFeat).

1.4. Bedingungen

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - **Bedingungen**
- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

Pfade

- Sequenzen von Merkmalen werden *Pfade* genannt.
- **Path** = **Feat**^{*} sei die Menge aller Pfade.
- Für $p \in \mathbf{Path}$, $F \in \mathcal{F}$ sei $F@p$ der Knoten in F , den man am Ende von Pfad p erhält.

Beispiele

- AGR-NUM
- SYN-SBJ-AGR-NUM
- ORTH
- ε (der leere Pfad)

Definition (Beschreibung Desc)

Die Menge der Beschreibungen über **Type** und **Feat** sei die kleinste Menge, die folgende Bedingungen erfüllt:

- $A \in \mathbf{Desc}$, für alle $A \in \mathbf{Type}$
- $p : d \in \mathbf{Desc}$, für $p \in \mathbf{Path}$, $d \in \mathbf{Desc}$
- $x \in \mathbf{Desc}$, für alle $x \in \mathbf{Var}$
- $d \wedge e \in \mathbf{Desc}$, für $d, e \in \mathbf{Desc}$

Beispiel

- $\text{AGR-NUM} : Sg$
- $\text{SYN-SBJ} : \boxed{1} \wedge \text{SEM-AGT} : \boxed{1}$

Erfülltheit

Die Erfülltheitsrelation \models^α zwischen Merkmalstrukturen und Beschreibungen ist gegeben durch:

- Für $A \in \mathbf{Type}$, $F \models^\alpha A \iff A \sqsubseteq \theta(\bar{q})$
- $F \models^\alpha p : d \iff F@p \models^\alpha d$
- Für $x \in \mathbf{Var}$, $F \models^\alpha x \iff \alpha(x) = F$
- $F \models^\alpha d \wedge e \iff F \models^\alpha d \text{ und } F \models^\alpha e$

Sei F eine Merkmalstruktur.

$$F = \left[\begin{array}{ll} \text{CAT} & \boxed{1}N \\ \text{POS} & \boxed{2} \\ \text{AGR} & \left[\begin{array}{ll} \text{NUM} & Sg \\ \text{CAS} & \textit{Nominativ} \end{array} \right] \end{array} \right]$$

$$\alpha(\boxed{1}) = \alpha(\boxed{2})$$

Welche Beschreibungen aus Desc erfüllt F ?

- $F \models^\alpha N$ Nein!
 - $F \models^\alpha \text{CAT} : N$ Ja!
 - $F \models^\alpha \text{AGR-CAS} : \textit{nicht-Genitiv}$ Ja!
 - $F \models^\alpha \text{POS} : N$ Ja!
- Denn: *nicht-Genitiv* \sqsubseteq *Nominativ*

MGSat (allgemeinster Erfüller)

Zu jeder *konsistenten* (widerspruchsfreien) Beschreibung $d \in \mathbf{Desc}$ gibt es eine Merkmalstruktur $MGSat(d) \in \mathcal{F}$ mit der Eigenschaft

$$\forall F \in \mathcal{F}. F \models d \iff MGSat(d) \sqsubseteq F$$

Konstruktion

- Für $A \in \mathbf{Type}$: $MGSat(A) = [A]$
- $MGSat(f_1 f_2 \dots f_n : d) = \left[f_1 \quad \left[f_2 \quad \dots \left[f_n \quad MGSat(d) \right] \right] \right]$
- Wenn $\mathbf{Var} = \mathbb{N}$, dann $MGSat(1) = [\underline{1}]$
- $MGSat(d \wedge e) = MGSat(d) \sqcup MGSat(e)$

Grammatikregel mit Constraint

$NP[CAS=?y] \rightarrow DET[GEN=?x, CAS=?y] \ N[GEN=?x]$

Bedingungen als Beschreibungen

- $type : NP \wedge CAS : \boxed{2}$
- $type : DET \wedge GEN : \boxed{1} \wedge CAS : \boxed{2}$
- $type : N \wedge GEN : \boxed{1}$

Bedingungen als Merkmalstrukturen

$$\begin{bmatrix} type & NP \\ CAS & \boxed{2} \end{bmatrix} \rightarrow \begin{bmatrix} type & DET \\ GEN & \boxed{1} \\ CAS & \boxed{2} \end{bmatrix} \begin{bmatrix} type & N \\ GEN & \boxed{1} \end{bmatrix}$$

2. *Implementierung

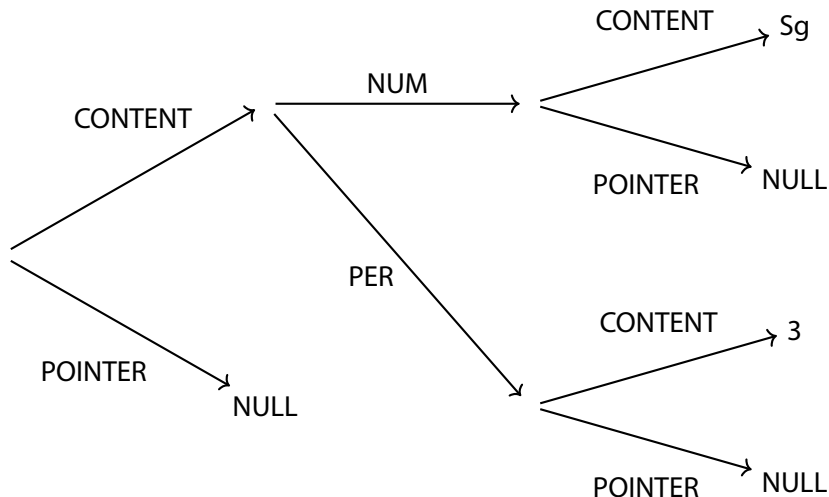
- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - Bedingungen

- 2 **Implementierung*
 - **Merkmalstrukturen*
 - **Unifikation*

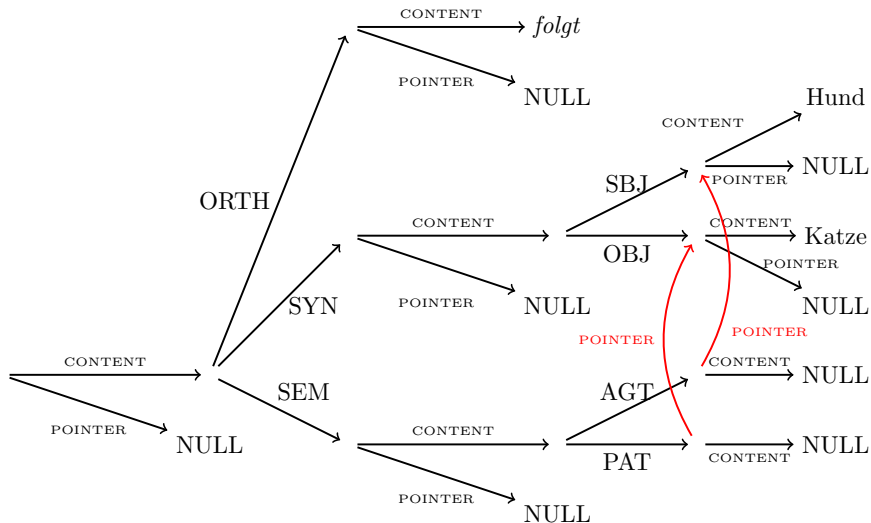
2.1. *Merkmalstrukturen

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - Bedingungen
- 2 *Implementierung
 - *Merkmalstrukturen
 - *Unifikation

Merkmalstrukturen: Implementierung mit Zeigern



Reentrance (Wiedereintritt)



2.2. *Unifikation

- 1 Formale Grundlagen
 - Merkmalstrukturen
 - Subsumption
 - Unifikation
 - Bedingungen
- 2 *Implementierung
 - *Merkmalstrukturen
 - *Unifikation

function unify(f1-orig, f2-orig)

 f1 \leftarrow deref(f1-orig)

▷ Verfolgen von .pointer

 f2 \leftarrow deref(f2-orig)

if f1, f2 \in **Type** **then**

return unifyValues(f1, f2)

▷ z. B. per Typhierarchie

if f1, f2 $\in \mathcal{F}$ **then**

for all feat2 \in f2 **do**

 feat1 \leftarrow erstelle oder finde entsprechendes Feature in f1

if unify(feats1, feat2) = failure **then**

return failure

return f1

...

function unify(f1-orig, f2-orig)

...

if f1 \in **ContPoint** \wedge f1.content is NULL **then**

 f1.pointer \leftarrow f2

return f2

if f2 \in **ContPoint** \wedge f2.content is NULL **then**

 f2.pointer \leftarrow f1

return f1

if f1, f2 \in **ContPoint** **then**

 f2.pointer \leftarrow f1

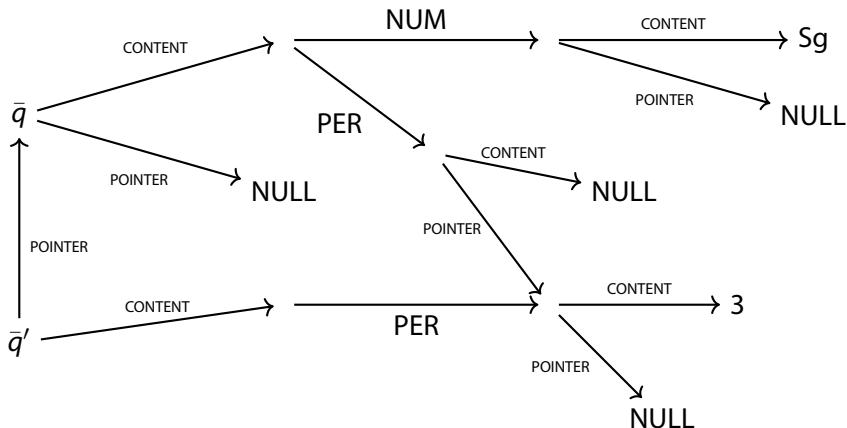
return unify(f1.content, f2.content)

return failure

▷ Kein Fall vorher hat zugetroffen

Unifikation: Implementierung mit Zeigern

$$\left[\text{NUM} \quad \text{Sg} \right] \sqcup \left[\text{PER} \quad 3 \right]$$



1 Formale Grundlagen

- Merkmalstrukturen
- Subsumption
- Unifikation
- Bedingungen

2 *Implementierung

- *Merkmalstrukturen
- *Unifikation