

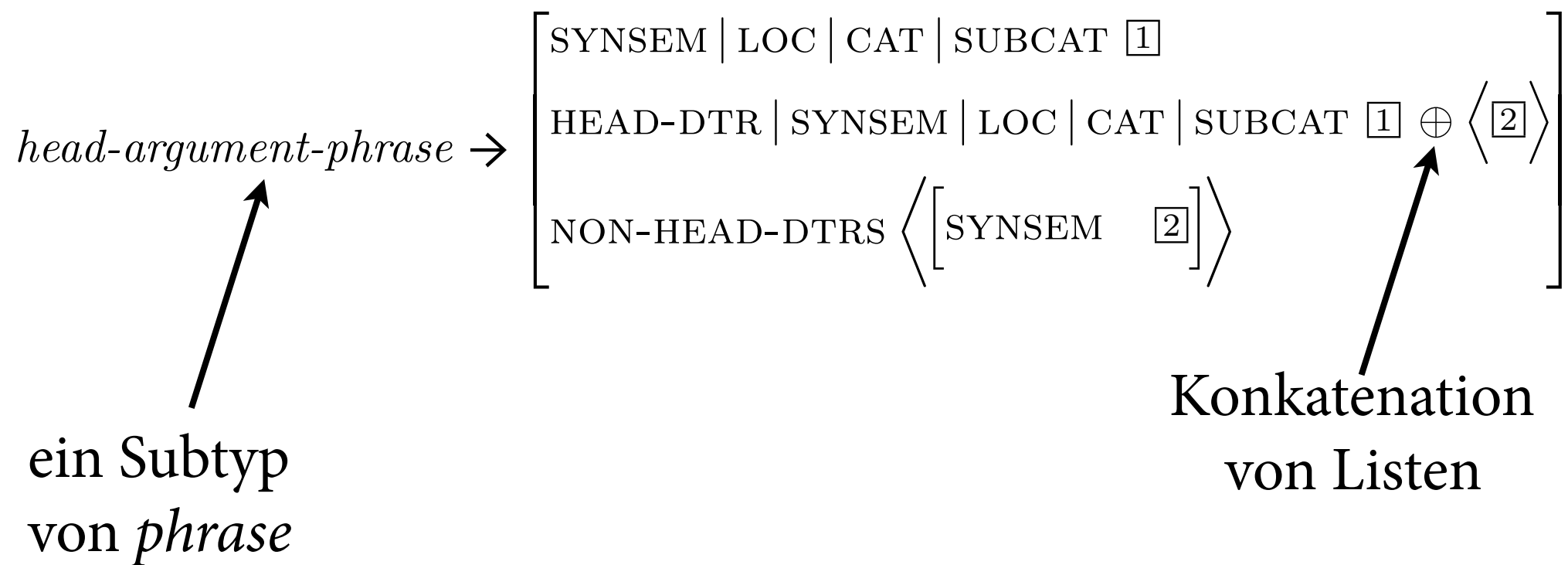
# Parsing von unifikationsbasierten Grammatikformalismen

Vorlesung “Grammatikformalismen”  
Alexander Koller

18. Juli 2017

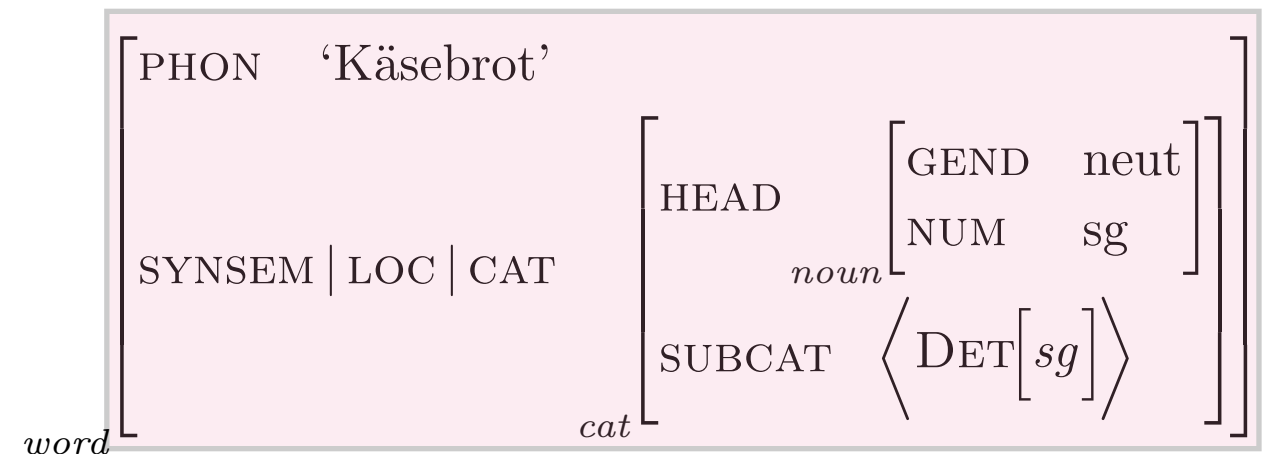
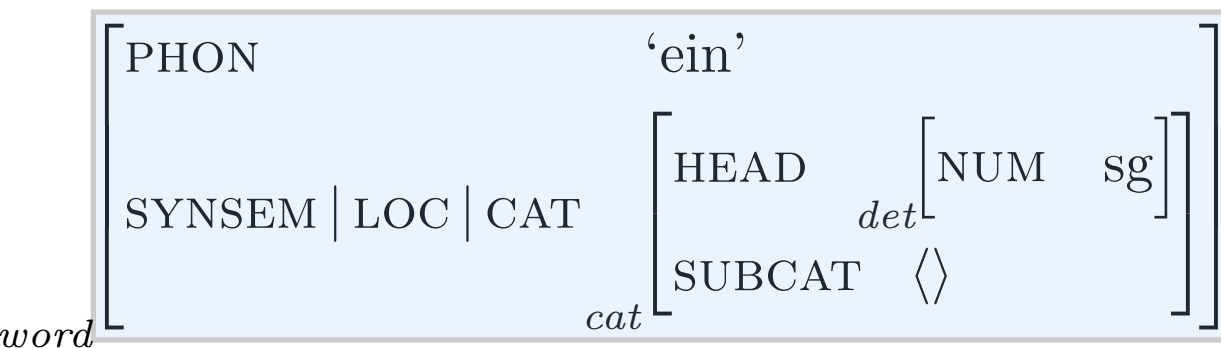
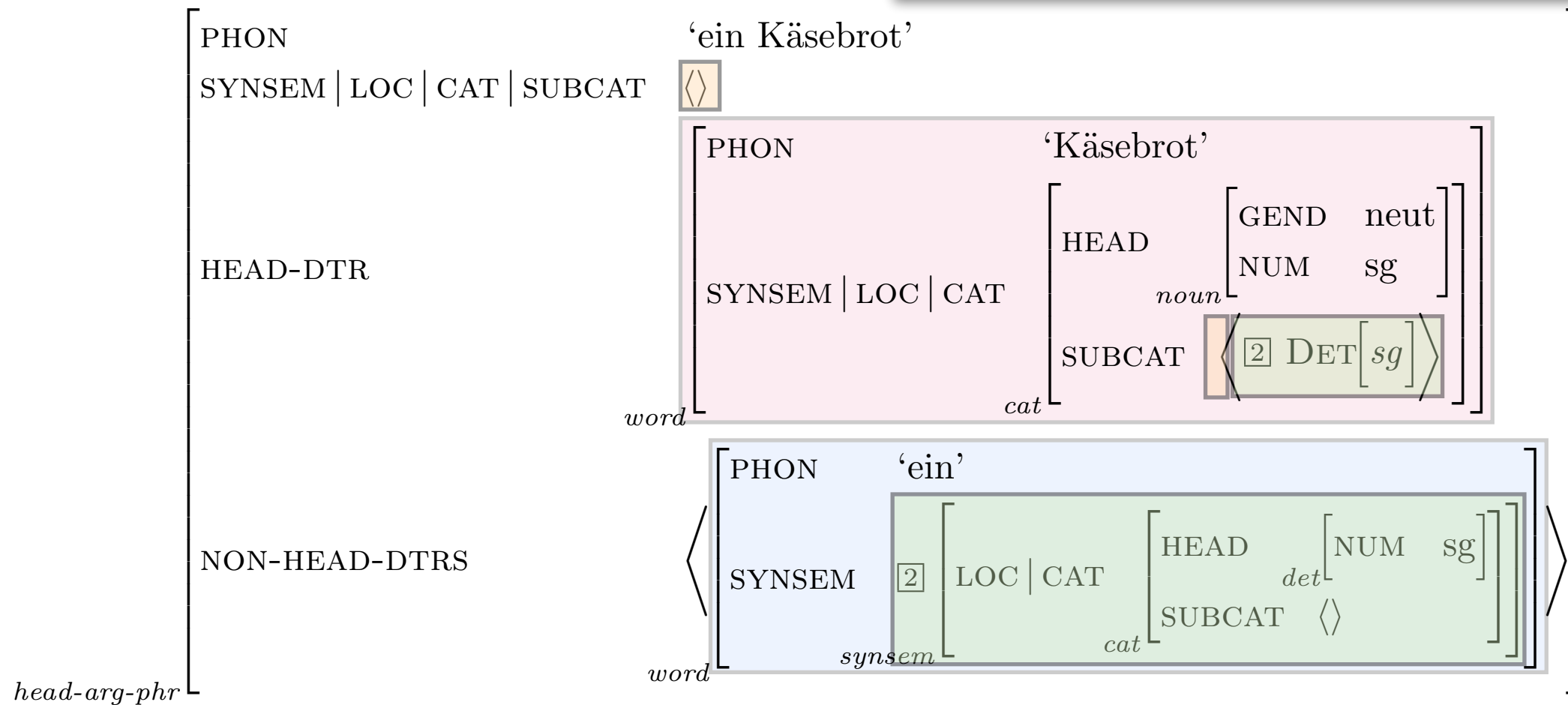
# Das Kopf-Argument-Schema

- Wichtigstes Schema: Kombination von Kopf mit Argument mit dem *Kopf-Argument-Schema*.



# Kopf + Argument

$$head-arg-phr \rightarrow \left[ \begin{array}{l} SYNSEM | LOC | CAT | SUBCAT \boxed{1} \\ HEAD-DTR | SYNSEM | LOC | CAT | SUBCAT \boxed{1} \oplus \langle \boxed{2} \rangle \\ NON-HEAD-DTRS \langle [ SYNSEM \boxed{2} ] \rangle \end{array} \right]$$



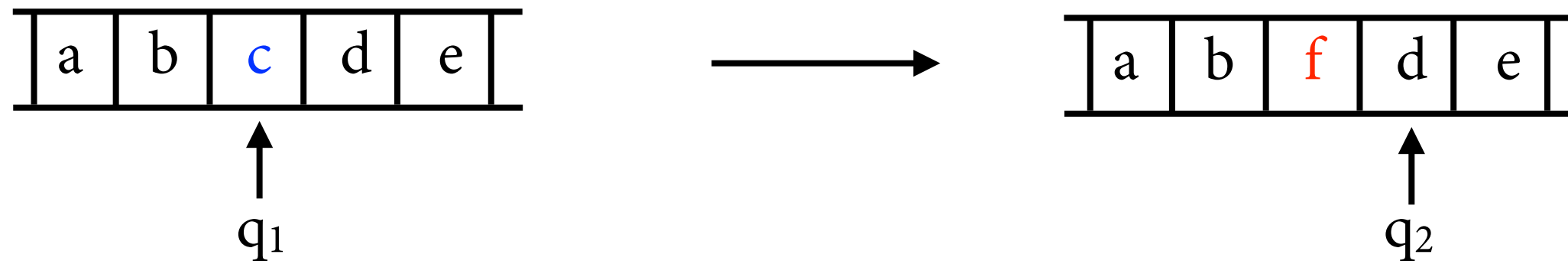
# Expressivität

- Wir haben uns überzeugt, dass HPSG linguistische Phänomene sehr elegant beschreiben kann.
- Welchen Preis in formaler Expressivität bezahlen wir dafür?

# Chomsky-Hierarchie

Typ	Grammatiken	Automaten
0	Typ-0-Grammatiken	Turingmaschinen
1	kontextsensitive	linear platzbeschränkte Turingmaschinen
2	kontextfreie	Kellerautomaten
3	reguläre	endliche Automaten

# Codierung von Turing-Maschinen


$$\begin{bmatrix} \text{LEFT} & \langle b, a \rangle \\ \text{RIGHT} & \langle c, d, e \rangle \\ \text{STATE} & q_1 \end{bmatrix}$$

$$\begin{bmatrix} \text{LEFT} & \langle f, b, a \rangle \\ \text{RIGHT} & \langle d, e \rangle \\ \text{STATE} & q_2 \end{bmatrix}$$

# Codierung von Turing-Maschinen

$$\begin{bmatrix} \text{LEFT} & \langle b, a \rangle \\ \text{RIGHT} & \langle c, d, e \rangle \\ \text{STATE} & q_1 \end{bmatrix}$$


$$\begin{bmatrix} \text{LEFT} & \langle f, b, a \rangle \\ \text{RIGHT} & \langle d, e \rangle \\ \text{STATE} & q_2 \end{bmatrix}$$

Regelschema  
(vereinfacht):

*rstep* →

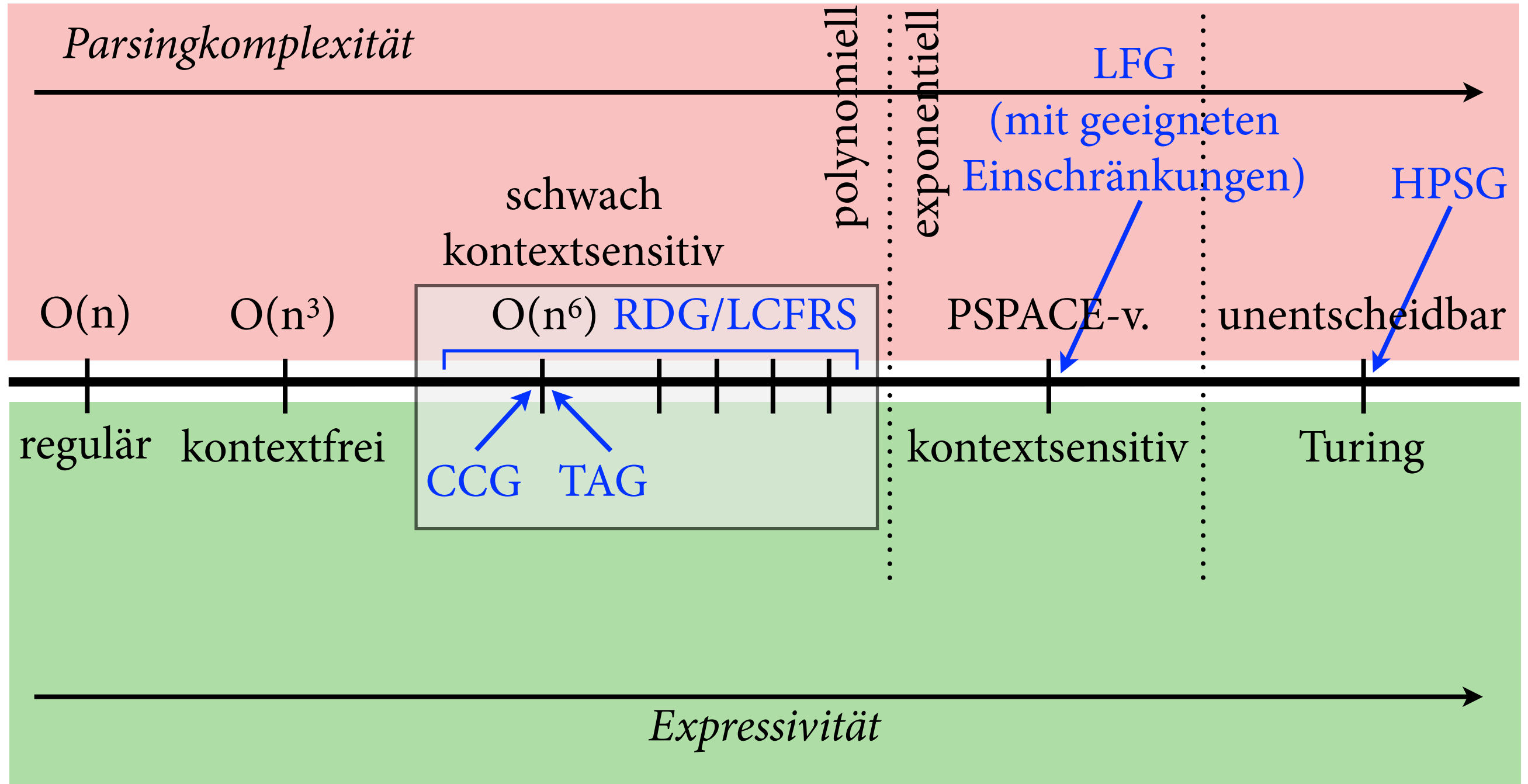
$$\begin{bmatrix} \text{LEFT} & \begin{bmatrix} \text{FIRST} & f \\ \text{REST} & \boxed{1} \end{bmatrix} \\ \text{RIGHT} & \boxed{2} \\ \text{STATE} & q_2 \\ \text{DTR} & \begin{bmatrix} \text{LEFT} & \boxed{1} \\ \text{RIGHT} & \begin{bmatrix} \text{FIRST} & c \\ \text{REST} & \boxed{2} \end{bmatrix} \\ \text{STATE} & q_1 \end{bmatrix} \end{bmatrix}$$

# Codierung: Grundidee

- Kann Berechnung einer Turingmaschine mit zwei Stacks simulieren.
- In HPSG-Featurestrukturen kann man zwei Stacks mit zwei listenwertigen Features darstellen.
  - ▶ geht in TAG, CCG nicht:  
dort endliche Wertebereiche für Features
- Daher ist Wortproblem von HPSG unentscheidbar.



# Fazit



# Parsing mit Featurestrukturen

- Grammatikalität von HPSG-Analysen bisher über Constraints definiert:
  - ▶ FS mit Subtyp von *sign* und korrektem Wert von PHON
  - ▶ alle Knoten sort-resolved und wohlgeformt (d.h. jeder Knoten durch Lexikoneintrag oder spezifischen Typ von Phrase motiviert)
- Heute: CKY-artiger Parser für HPSG.
- Grundideen übertragen sich auf andere unifikationsbasierte Formalismen, z.B. LFG.

# Grundidee

- Parsingschema für Bottom-Up-Parser:

$$\frac{[i,j, F_1] \quad [j,k, F_2]}{[i,k, C(\tau) \sqcup [\text{hdtr } F_1] \sqcup [\text{nhdtr } F_2]]} \tau$$

- Einbettende Features `hdtr` und `nhdtr` hängen von Phrasentyp  $\tau$  ab.
- FS  $C(\tau)$  kommt aus Typconstraint für  $\tau$ . Unifikation kann fehlschlagen; dann Regel nicht anwendbar.

# Grundidee

- Parsingschema für Bottom-Up-Parser:

$$\frac{[i,j, F_1] \quad [j,k, F_2]}{[i,k, C(\tau) \sqcup [\text{hdtr } F_1] \sqcup [\text{nhdtr } F_2]]} \tau$$

FS  $F_2$  ist Analyse von Substring  $j$ - $k$

- Einbettende Features  $\text{hdtr}$  und  $\text{nhdtr}$  hängen von Phrasentyp  $\tau$  ab.
- FS  $C(\tau)$  kommt aus Typconstraint für  $\tau$ . Unifikation kann fehlschlagen; dann Regel nicht anwendbar.

# Grundidee

- Parsingschema für Bottom-Up-Parser:

$$\frac{[i,j, F_1] \quad [j,k, F_2]}{[i,k, C(\tau) \sqcup [\text{hdtr } F_1] \sqcup [\text{nhdtr } F_2]]} \tau$$

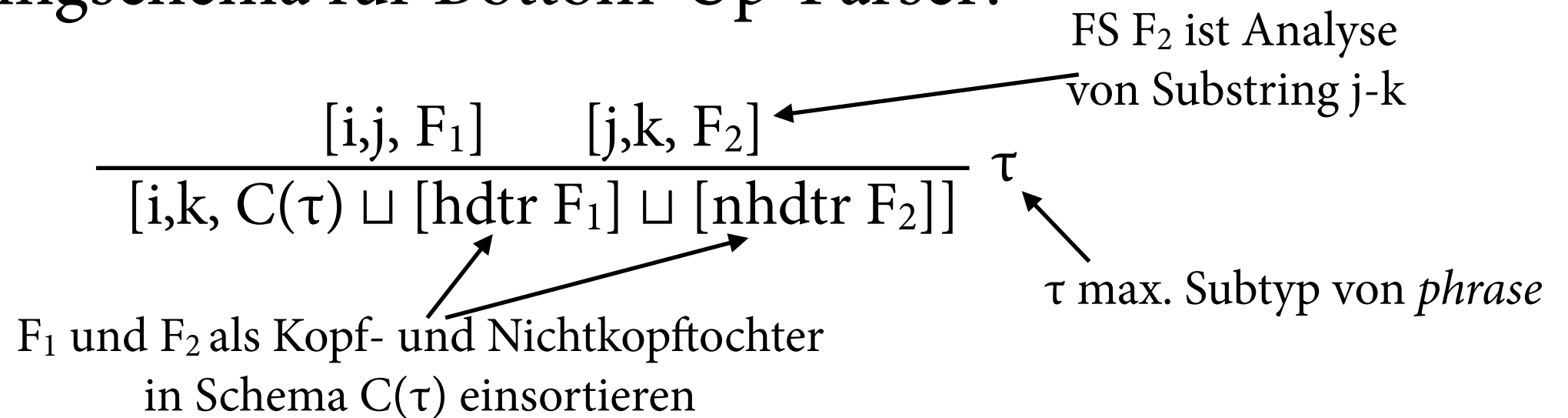
FS  $F_2$  ist Analyse von Substring  $j$ - $k$

$\tau$  max. Subtyp von *phrase*

- Einbettende Features  $\text{hdtr}$  und  $\text{nhdtr}$  hängen von Phrasentyp  $\tau$  ab.
- FS  $C(\tau)$  kommt aus Typconstraint für  $\tau$ . Unifikation kann fehlschlagen; dann Regel nicht anwendbar.

# Grundidee

- Parsingschema für Bottom-Up-Parser:



- Einbettende Features  $\text{hdtr}$  und  $\text{nhdtr}$  hängen von Phrasentyp  $\tau$  ab.
- FS  $C(\tau)$  kommt aus Typconstraint für  $\tau$ . Unifikation kann fehlschlagen; dann Regel nicht anwendbar.

# Beispiel


ein Käsebrot

isst

isst

ein Käsebrot

# Beispiel

ein Käsebrot

isst

<p><i>word</i> [ PHON 'isst' SYNSEM   LOC   CAT <i>cat</i> [ HEAD <i>verb</i> [ NUM sg ] SUBCAT &lt; NP[<i>nom</i>], NP[<i>acc</i>] &gt; ] ]</p>	

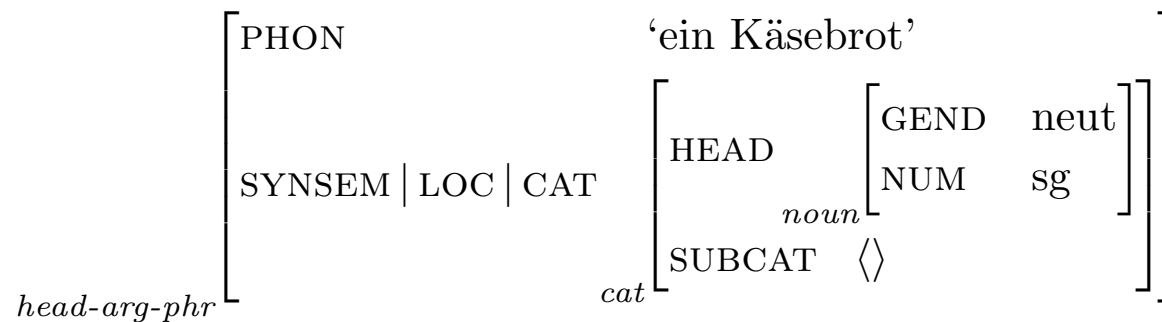
isst

ein Käsebrot

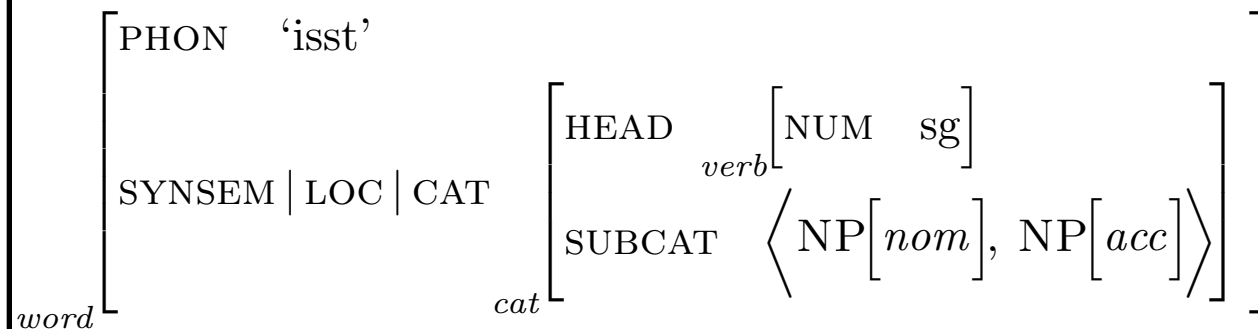


# Beispiel

ein Käsebro



isst

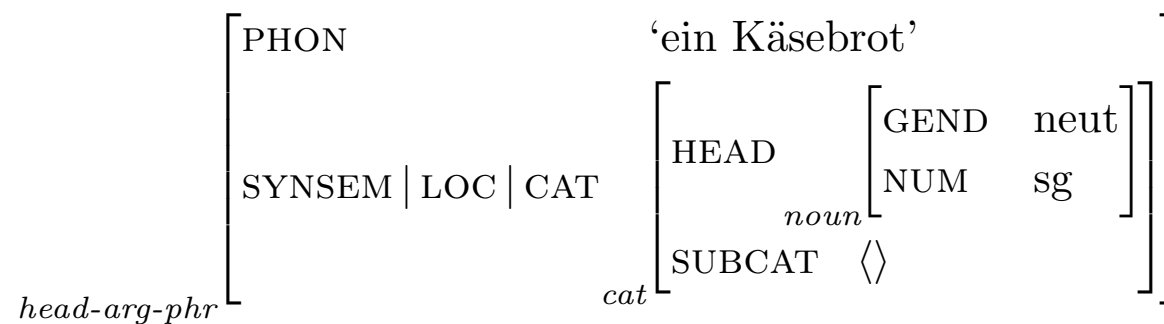
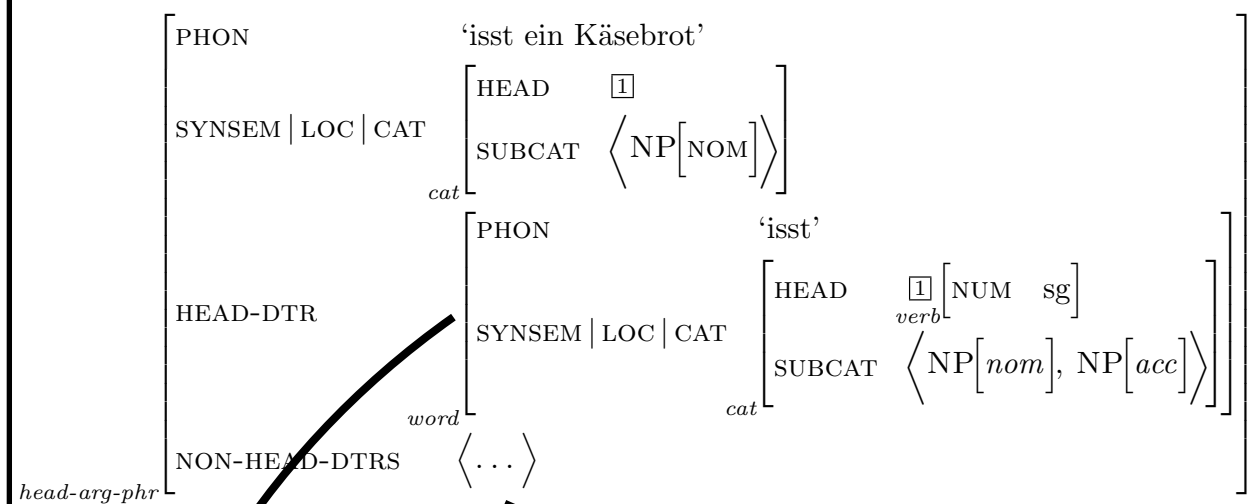


isst

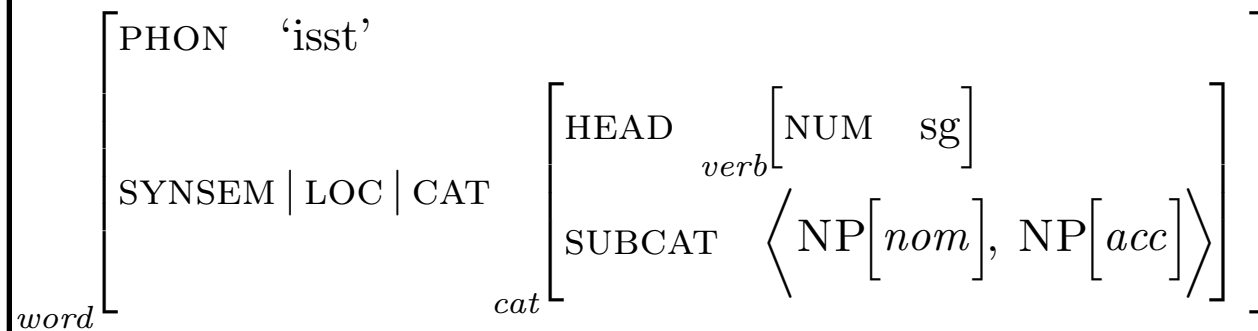
ein Käsebro

# Beispiel

ein Käsebrot



isst



isst

ein Käsebrot

# Algorithmische Fragen

- *Gleichheit testen*: Wenn ich  $[i,k,F]$  und  $[i,k,G]$  ableiten kann, muss ich effizient überprüfen, ob  $F$  und  $G$  die gleiche Featurestruktur sind.
  - ▶ Wenn nein, muss neues Item auf Agenda eingetragen werden.
- *Unifikation*: Muss  $C(\tau) \sqcup [\text{hdtr } F_1] \sqcup [\text{nhdtr } F_2]$  effizient ausrechnen.
- Mindestens 90% der Laufzeit eines unifikationsbasierten Parsers gehen in diese zwei Funktionen.

# Subsumption von TFS

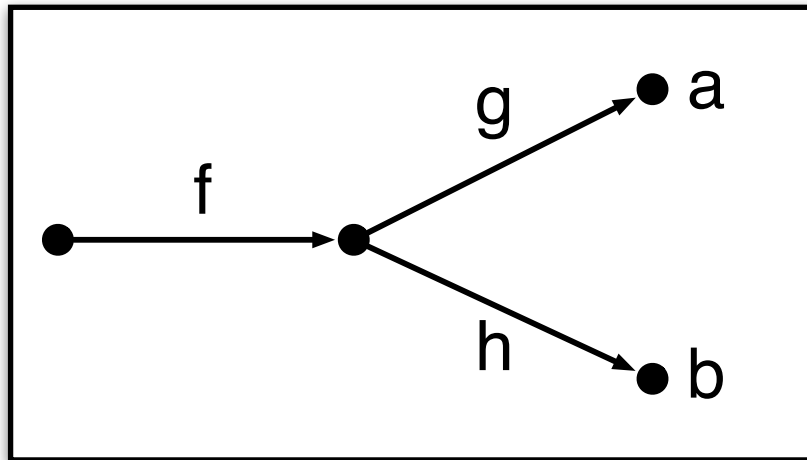
- TFS  $F$  *subsumiert* TFS  $G$  ( $F \sqsubseteq G$ ), wenn alle Informationen in  $F$  auch in  $G$  vorhanden sind:
  - ▶ Pfade:  $F(\pi)$  definiert  $\Rightarrow G(\pi)$  definiert
  - ▶ Reentrancy:  $F(\pi) = F(\pi') \Rightarrow G(\pi) = G(\pi')$
  - ▶ Subtypen:  $\theta_F(F(\pi)) \sqsubseteq \theta_G(G(\pi))$
- *Unifikation*:  $F = F_1 \sqcup F_2$  ist FS (falls existiert), so dass
  - ▶ subsumiert beide:  $F_1 \sqsubseteq F, F_2 \sqsubseteq F$
  - ▶ erfindet nichts dazu: für alle  $F'$  mit  $F_1 \sqsubseteq F', F_2 \sqsubseteq F'$  gilt  $F \sqsubseteq F'$ .

# Subsumption testen

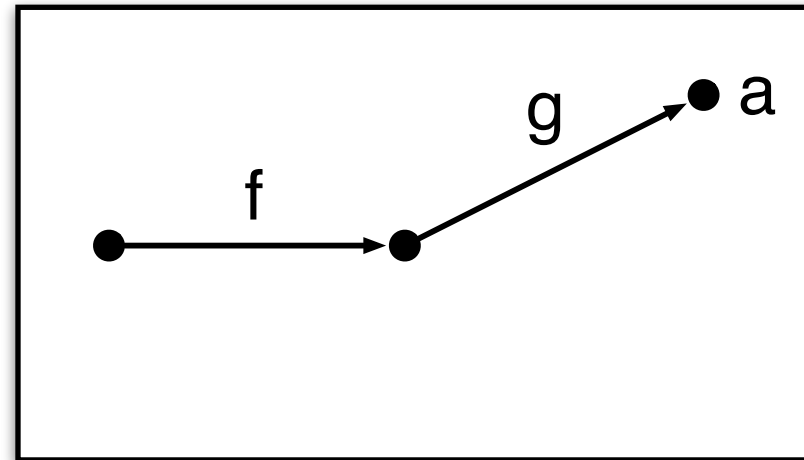
- Algorithmus von Malouf et al. 2000: Gegeben  $F$  und  $G$ , teste gleichzeitig, ob  $F \sqsubseteq G$  und ob  $G \sqsubseteq F$ .
  - ▶ Spezialfall:  $F = G$  gdw.  $F \sqsubseteq G$  und  $G \sqsubseteq F$ .
- Traversiere beide FS rekursiv.
  - ▶ Speichere für jeden Knoten  $u$  in  $F$  einen *Pointer* zu einem Knoten  $v$  in  $G$ , und umgekehrt.
  - ▶ Erkenne  $F \not\sqsubseteq G$ , falls:  $u$  mehr Features hat als  $v$ ; Typen von  $u$ ,  $v$  nicht zusammenpassen;  $u$  beim zweiten Besuch schon einen Pointer zu einem Knoten  $\neq v$  hatte (und umgekehrt).

# Beispiel

$F_1$



$F_2$

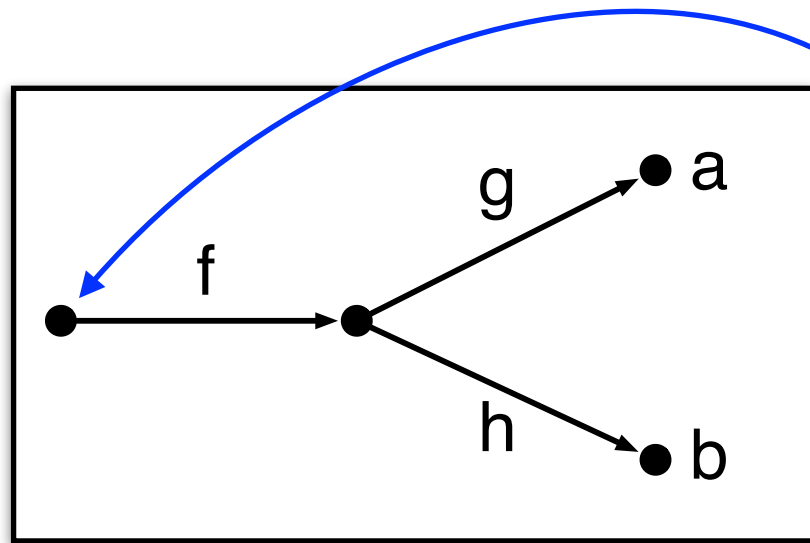


$F_1 \sqsubseteq F_2 ?$

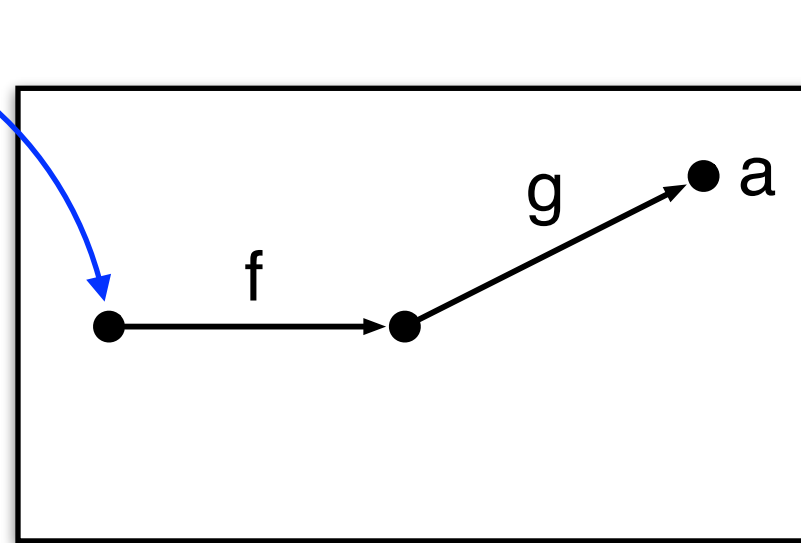
$F_2 \sqsubseteq F_1 ?$

# Beispiel

$F_1$



$F_2$

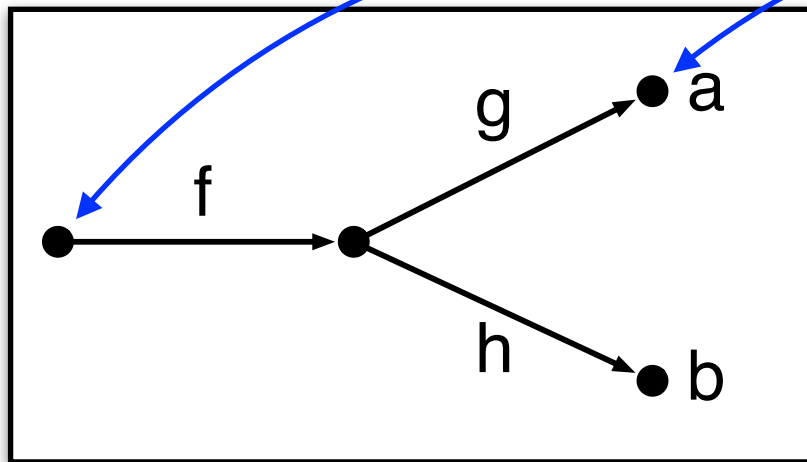


$F_1 \sqsubseteq F_2 ?$

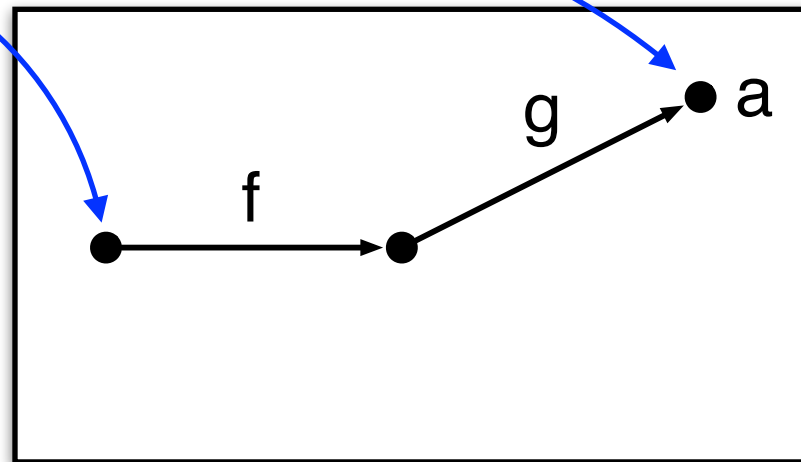
$F_2 \sqsubseteq F_1 ?$

# Beispiel

$F_1$



$F_2$



$F_1 \sqsubseteq F_2 ?$

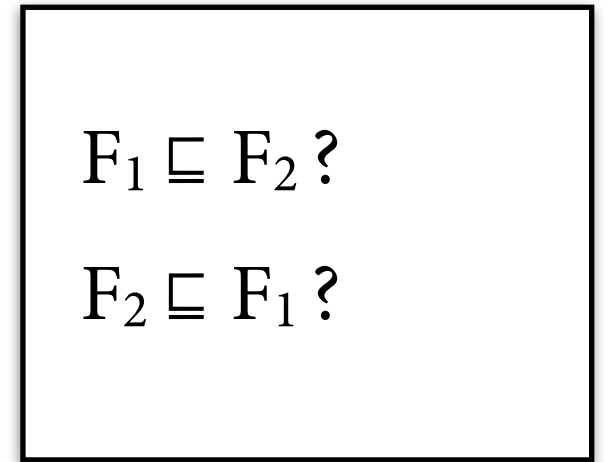
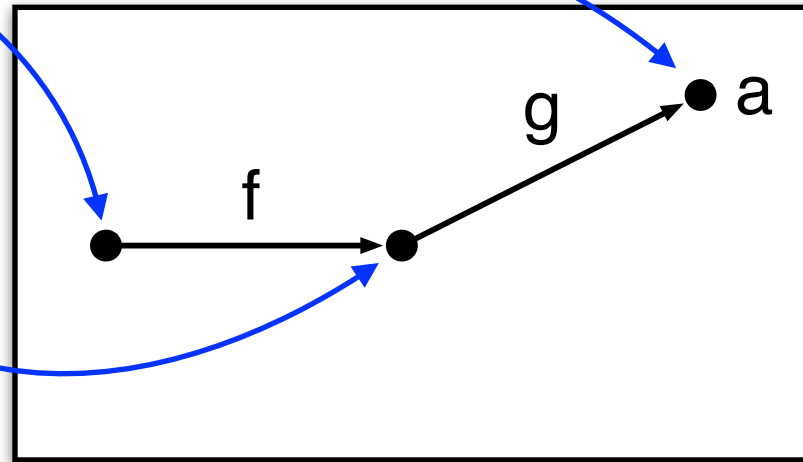
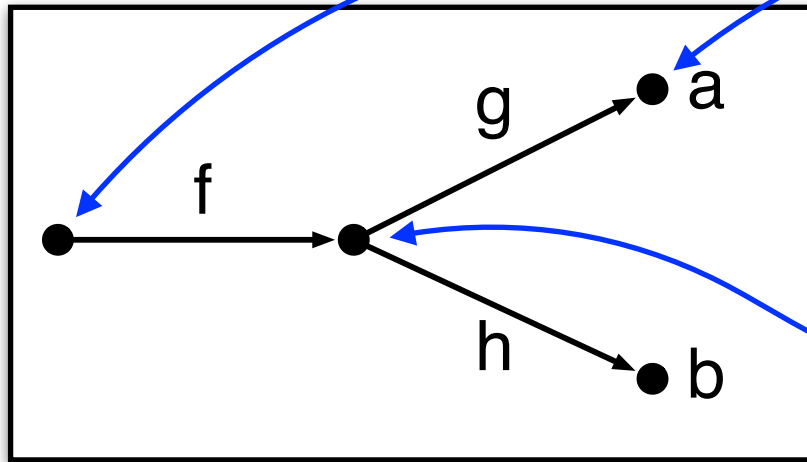
$F_2 \sqsubseteq F_1 ?$



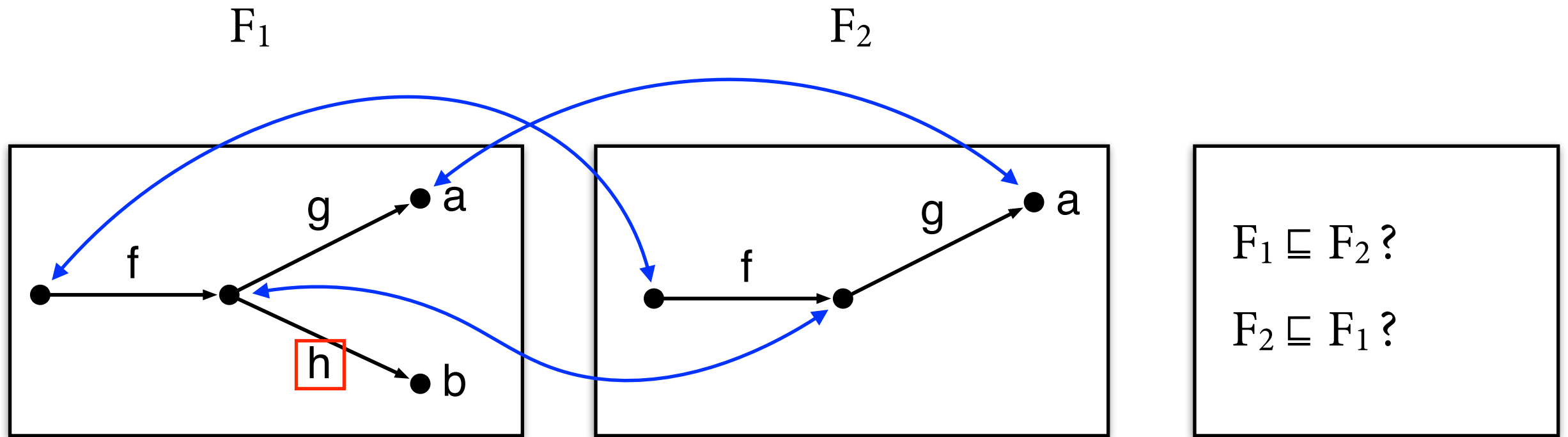
# Beispiel

$F_1$

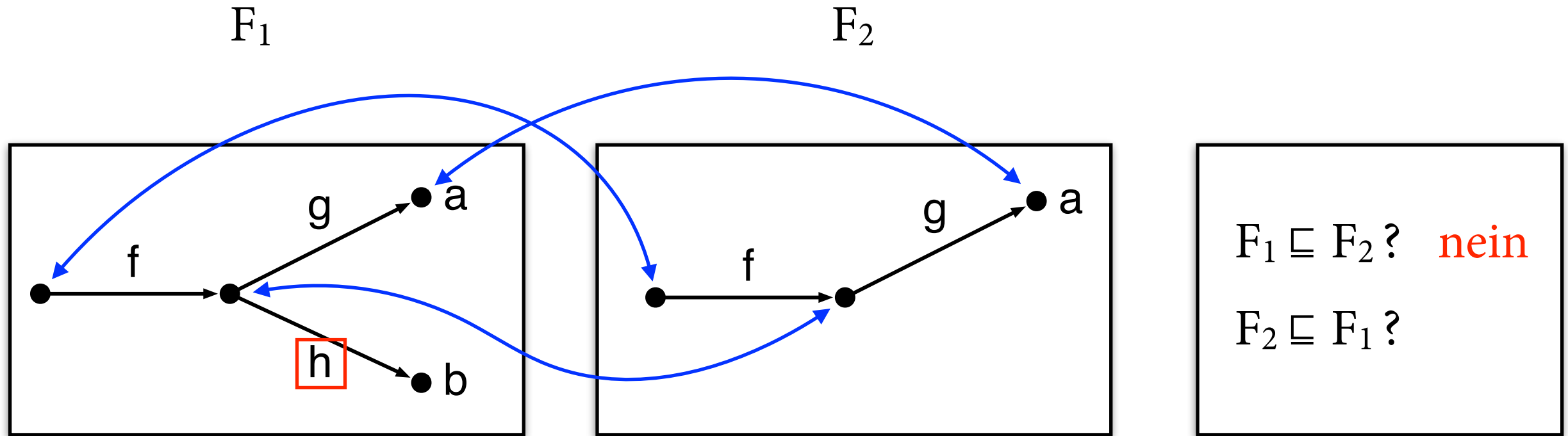
$F_2$



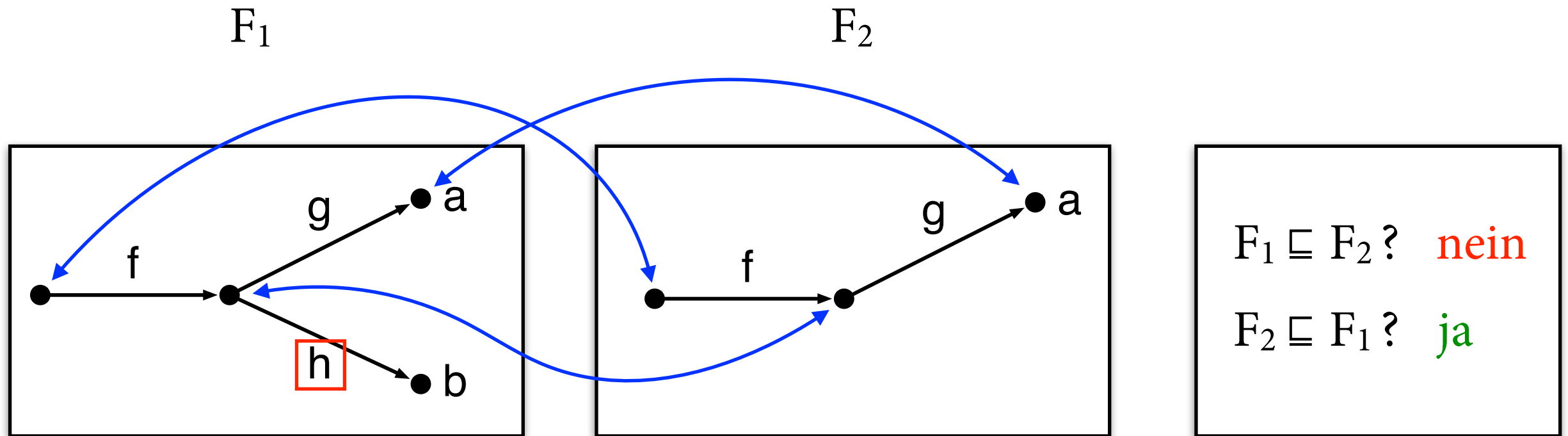
# Beispiel



# Beispiel



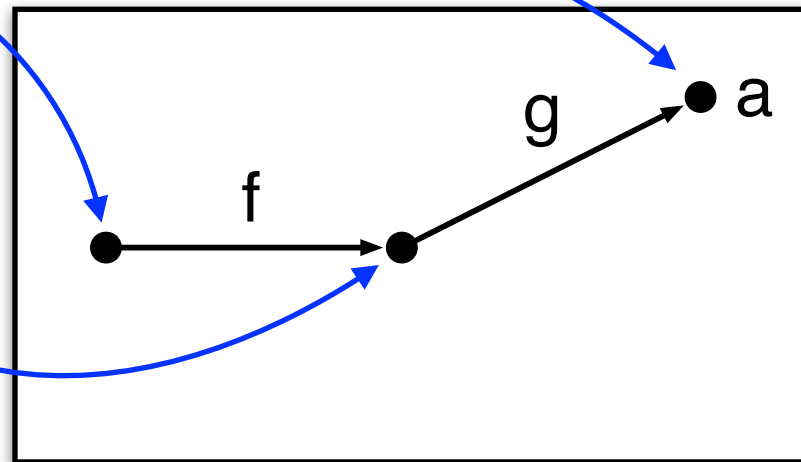
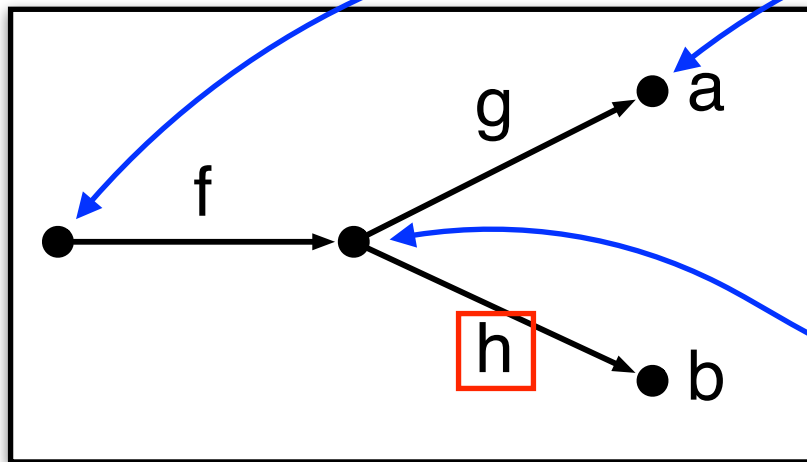
# Beispiel



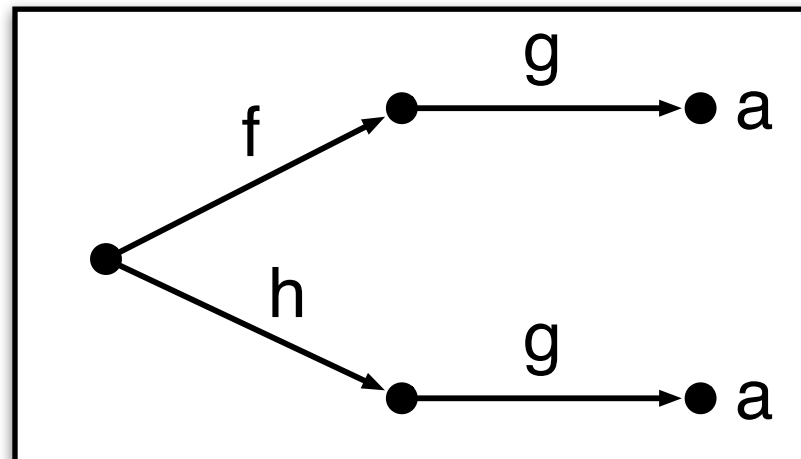
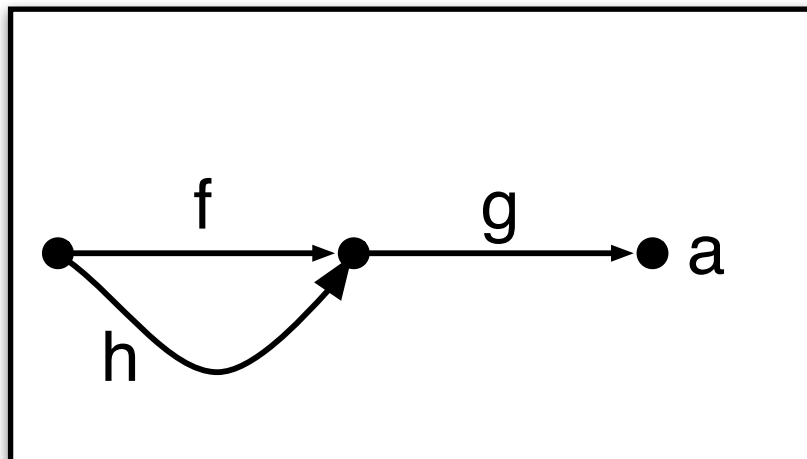
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**

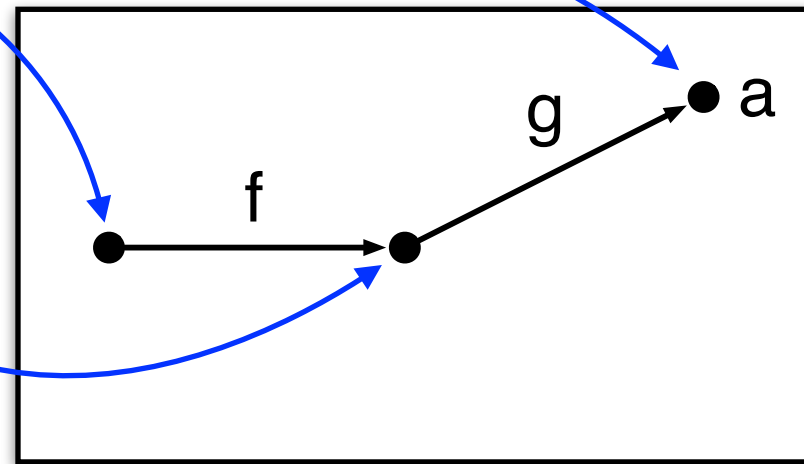
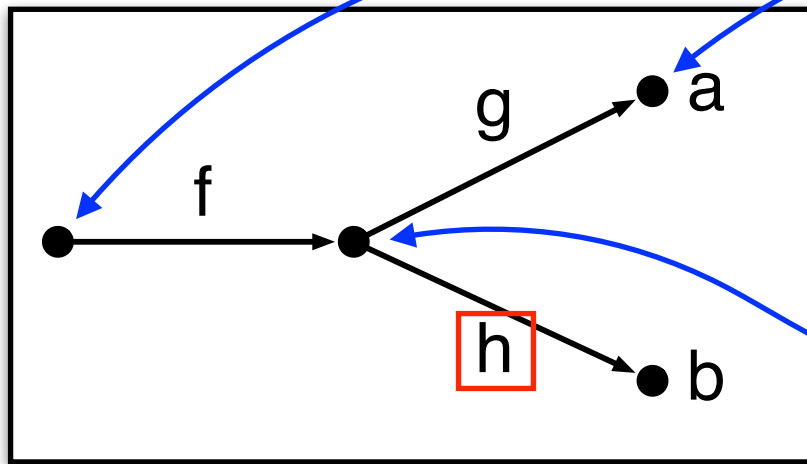


$F_1 \sqsubseteq F_2$ ?  
 $F_2 \sqsubseteq F_1$ ?

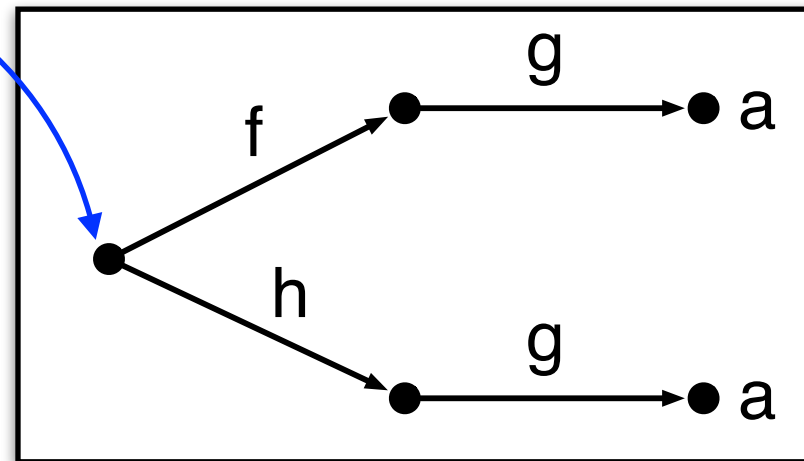
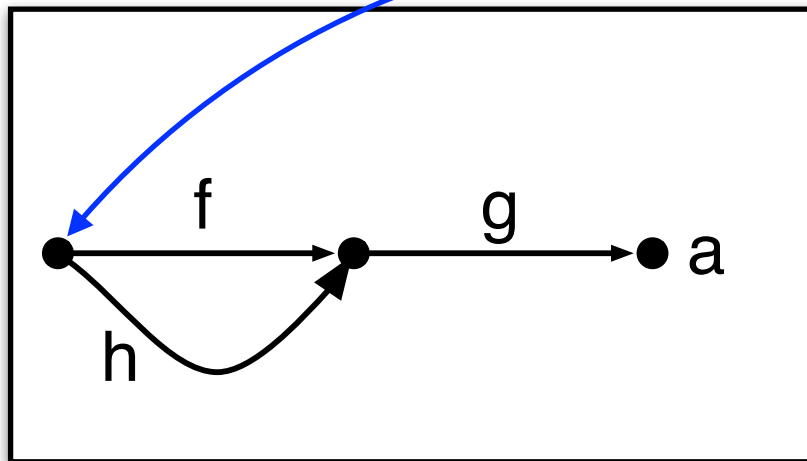
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**

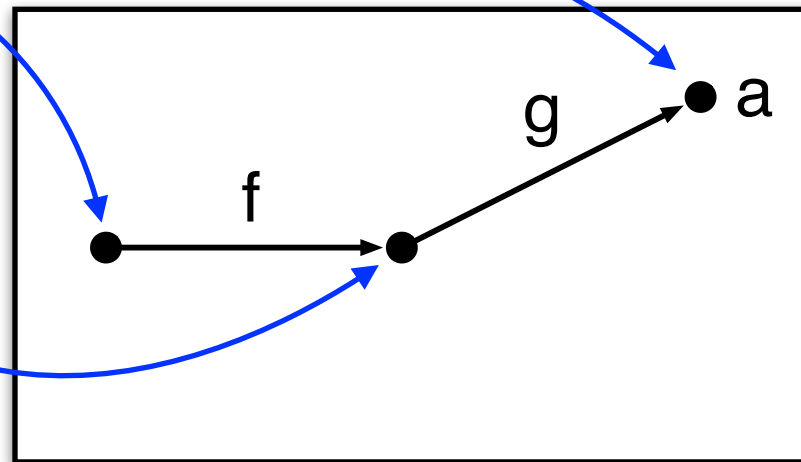
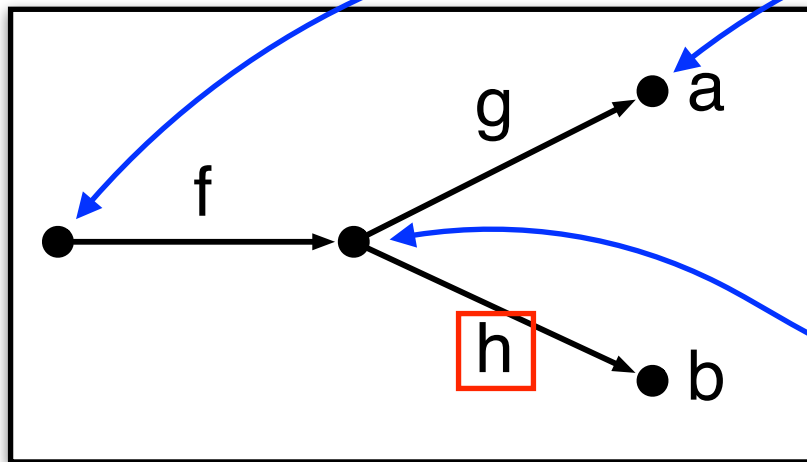


$F_1 \sqsubseteq F_2$ ?  
 $F_2 \sqsubseteq F_1$ ?

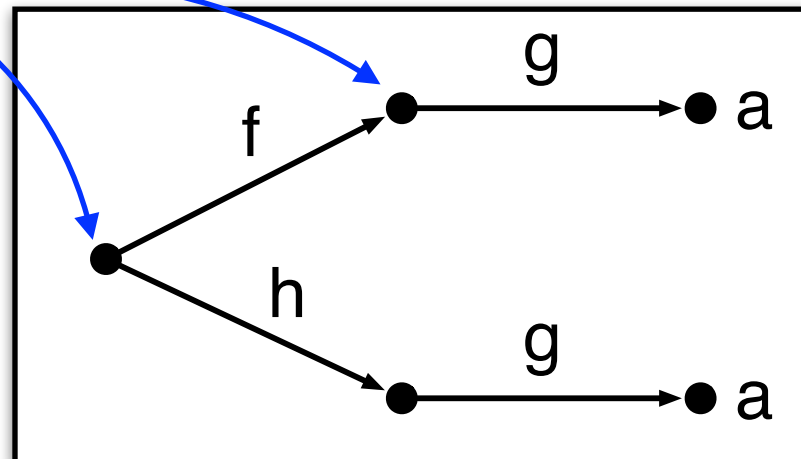
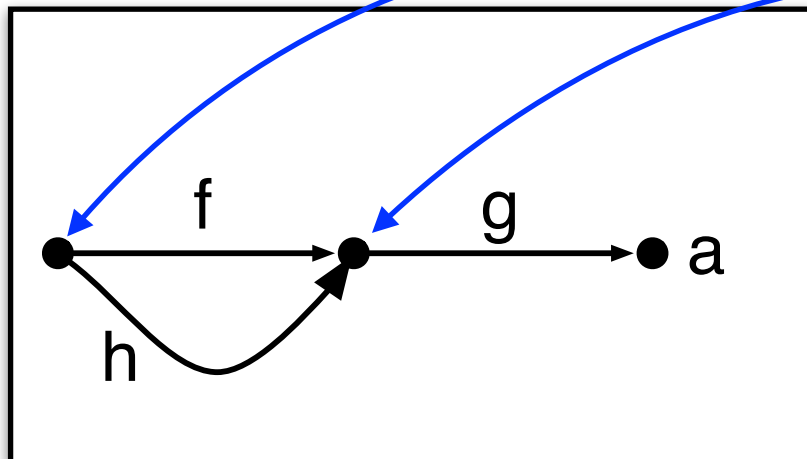
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**

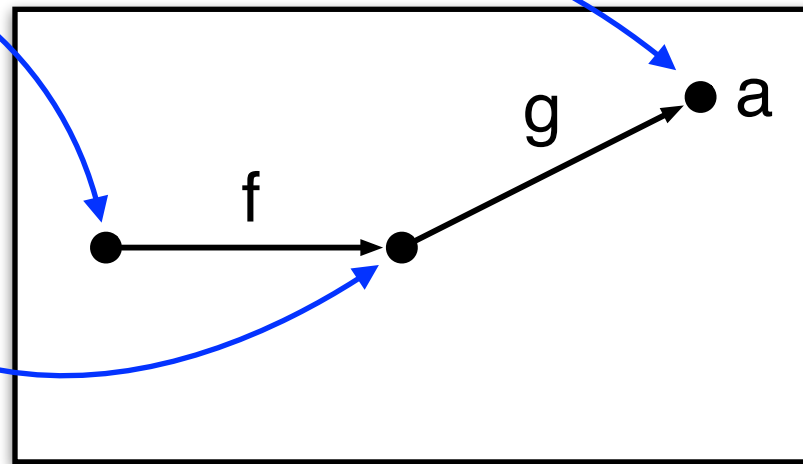
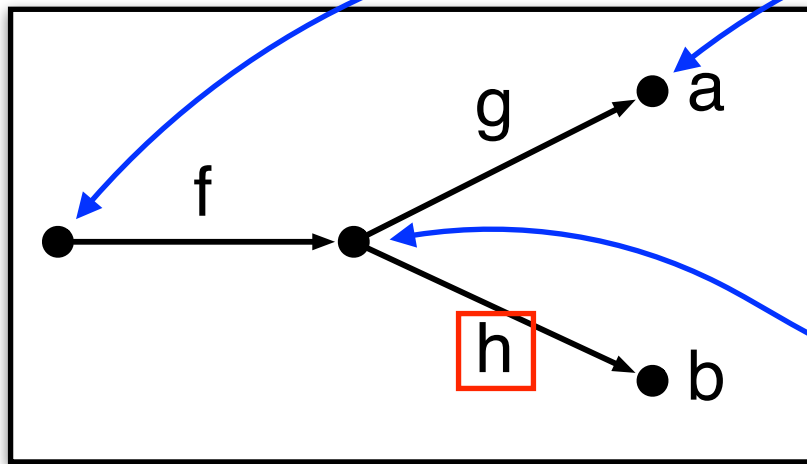


$F_1 \sqsubseteq F_2$ ?  
 $F_2 \sqsubseteq F_1$ ?

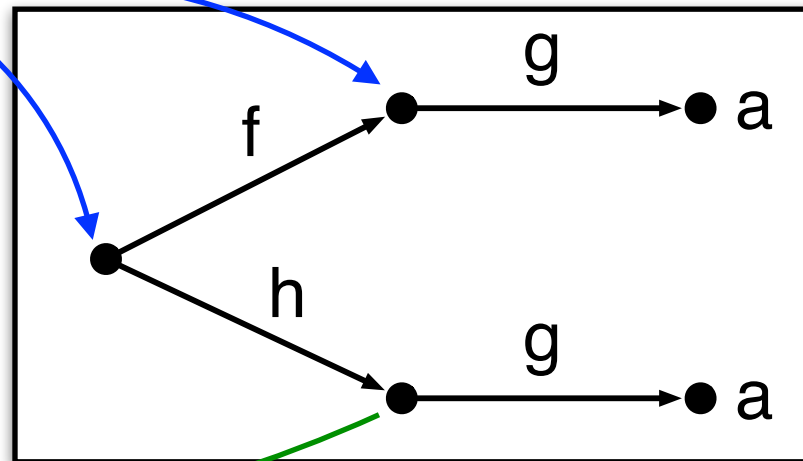
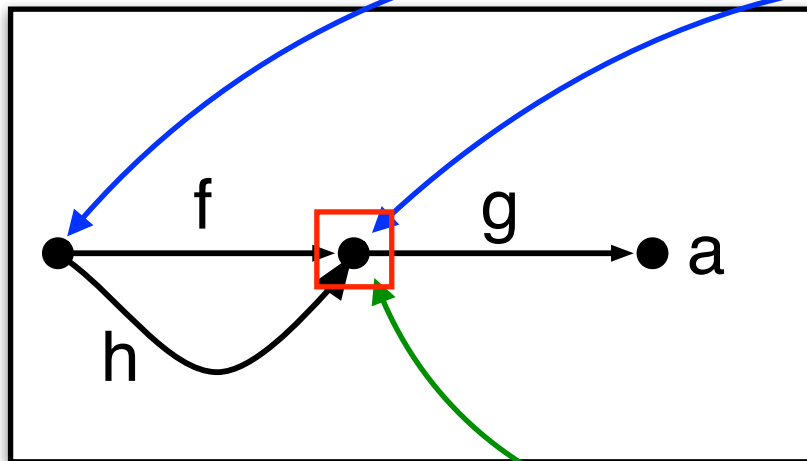
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**



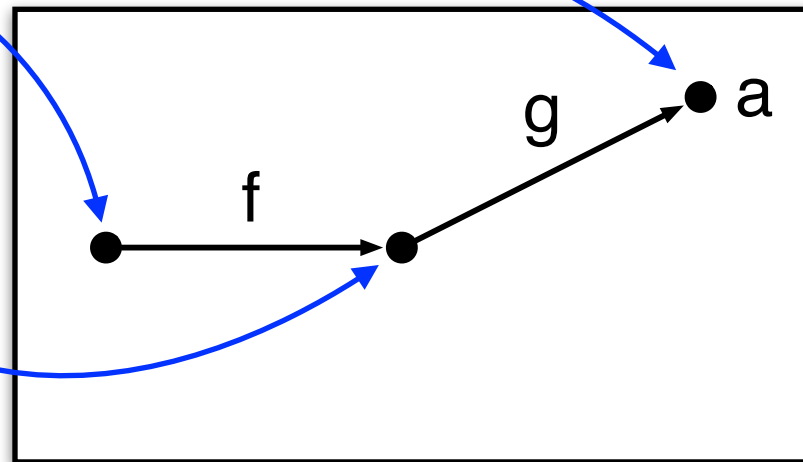
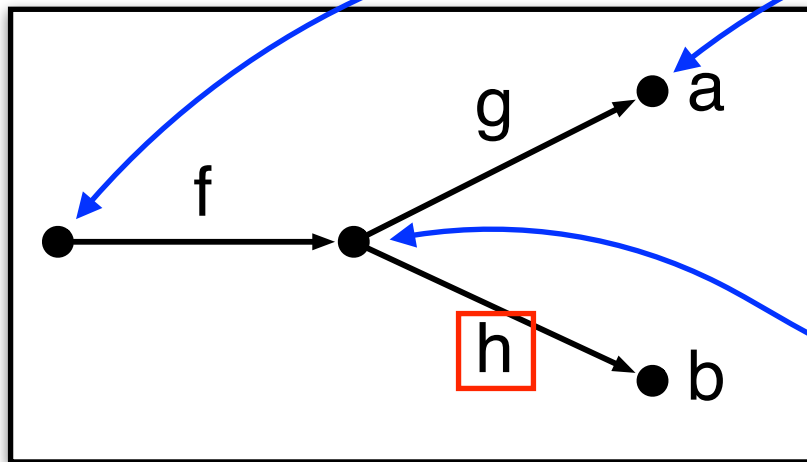
$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ?



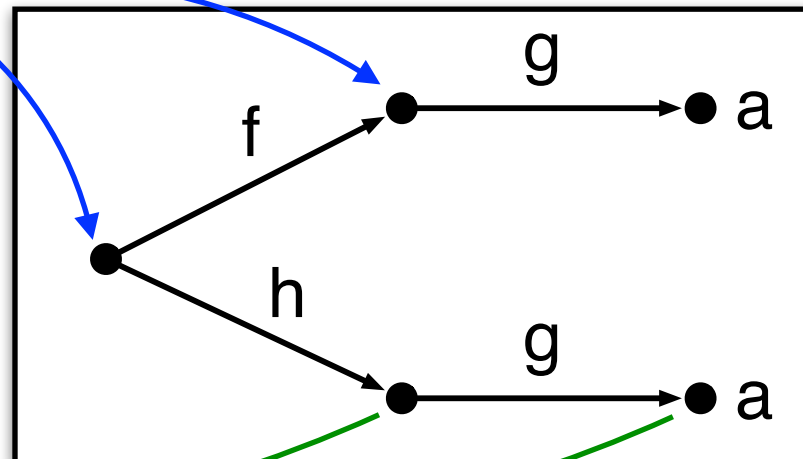
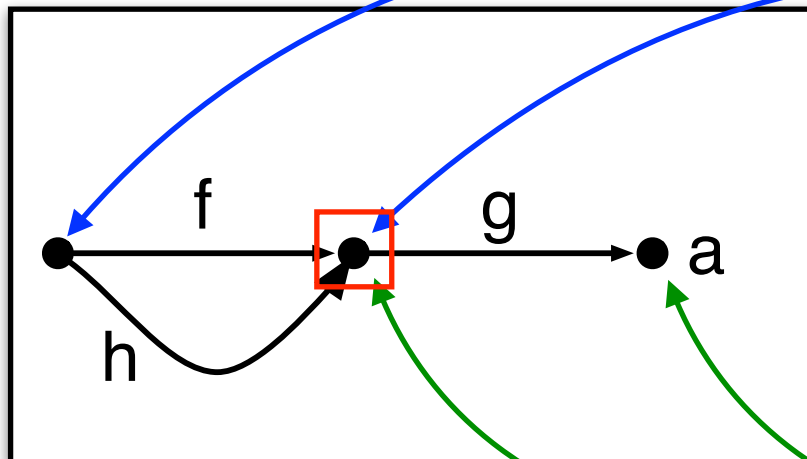
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**

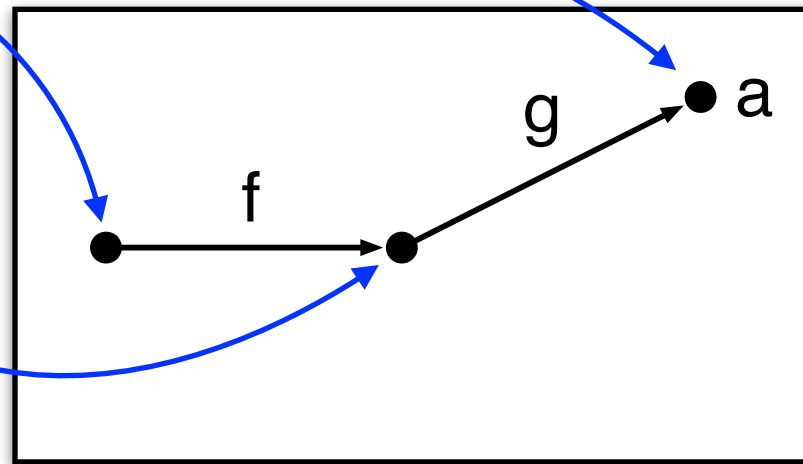
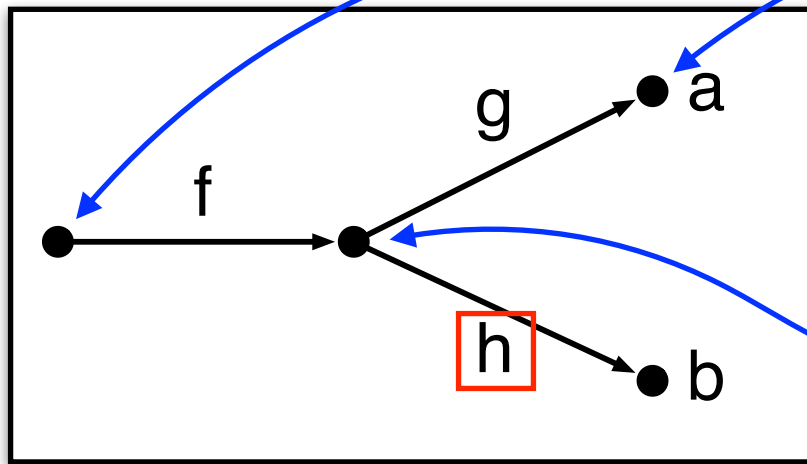


$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ?

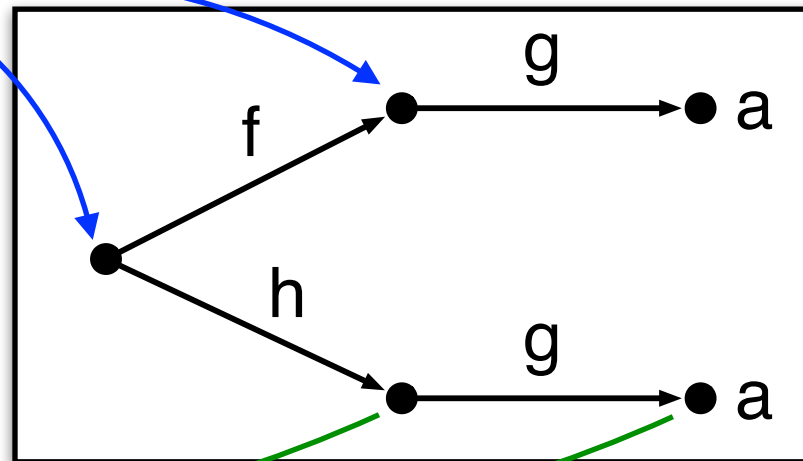
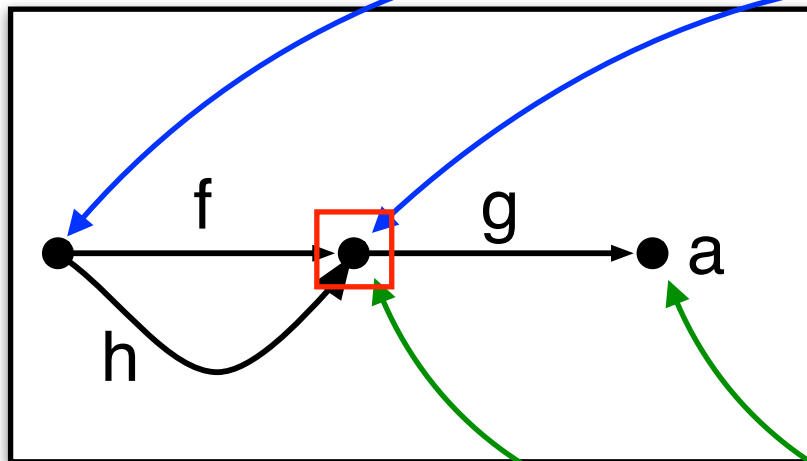
# Beispiel

$F_1$

$F_2$



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**



$F_1 \sqsubseteq F_2$ ? **nein**  
 $F_2 \sqsubseteq F_1$ ? **ja**

# Unifikation

- Grundidee: Rekursiv durch Featurestrukturen gehen und Pointer auf korrespondierende Knoten setzen.
  - ▶ Dabei wird Inhalt von  $F_2$  nach  $F_1$  kopiert.
  - ▶ In ähnlicher Weise wie bei Subsumption auf Konsistenz der Information testen.
- Rolle des Kopierens:
  - ▶ Im worst case kopiert Algorithmus fast die ganze FS  $F_2$ , bevor er erkennt, dass die FS nicht unifiziert werden können.
  - ▶ Schnellere Algorithmen verzögern und vermeiden Kopieren so gut wie möglich, z.B. Tomabechi (1991).

# Unifikationsalgorithmus

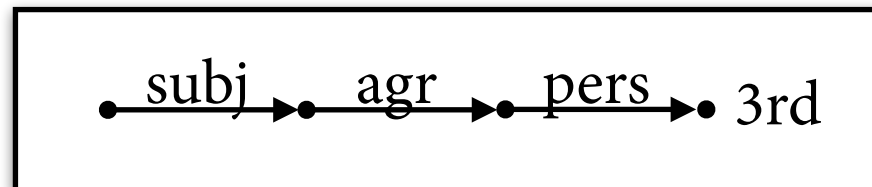
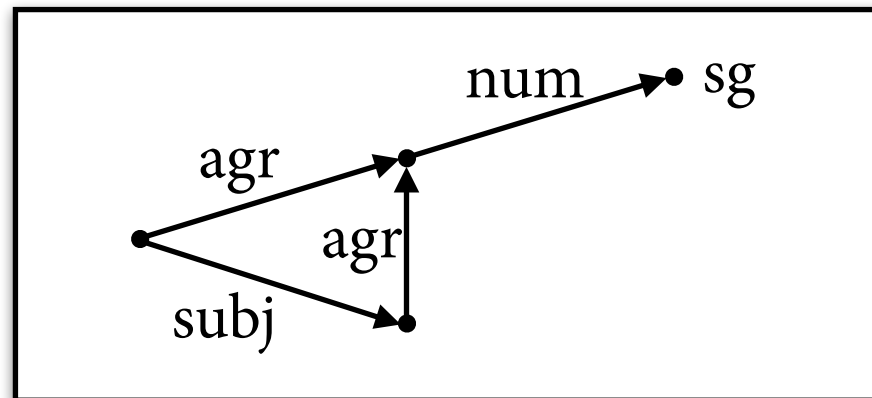
```
function UNIFY(f1, f2) returns fstructure or failure

f1-real ← Real contents of f1
f2-real ← Real contents of f2

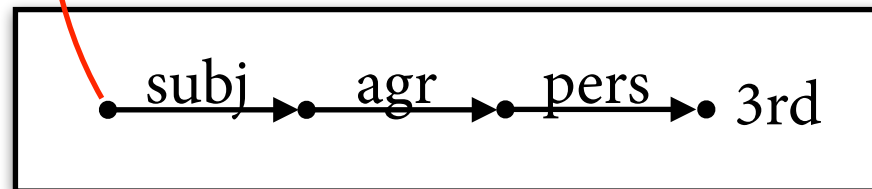
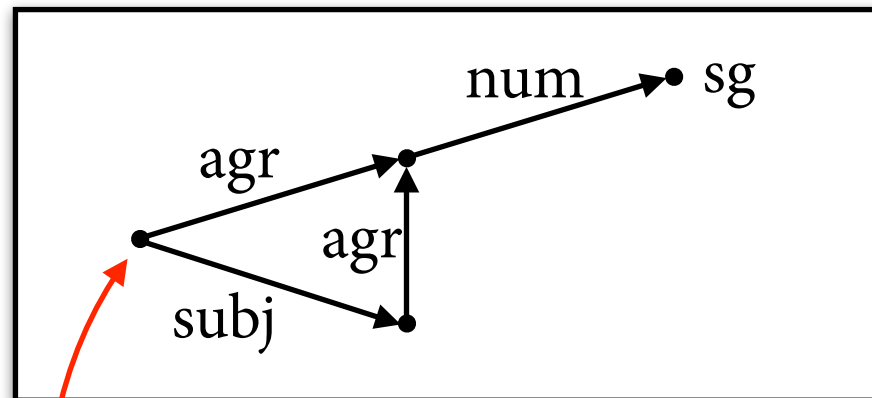
if f1-real is null then
  f1.pointer ← f2
  return f2
else if f2-real is null then
  f2.pointer ← f1
  return f1
else if f1-real and f2-real are identical then
  f1.pointer ← f2
  return f2
else if both f1-real and f2-real are complex feature structures then
  f2.pointer ← f1
  for each feature in f2-real do
    other-feature ← Find or create
                        a feature corresponding to feature in f1-real
    if UNIFY(feature.value, other-feature.value) returns failure then
      return failure
  return f1
else return failure
```

(aus, aber nicht von, Jurafsky & Martin)

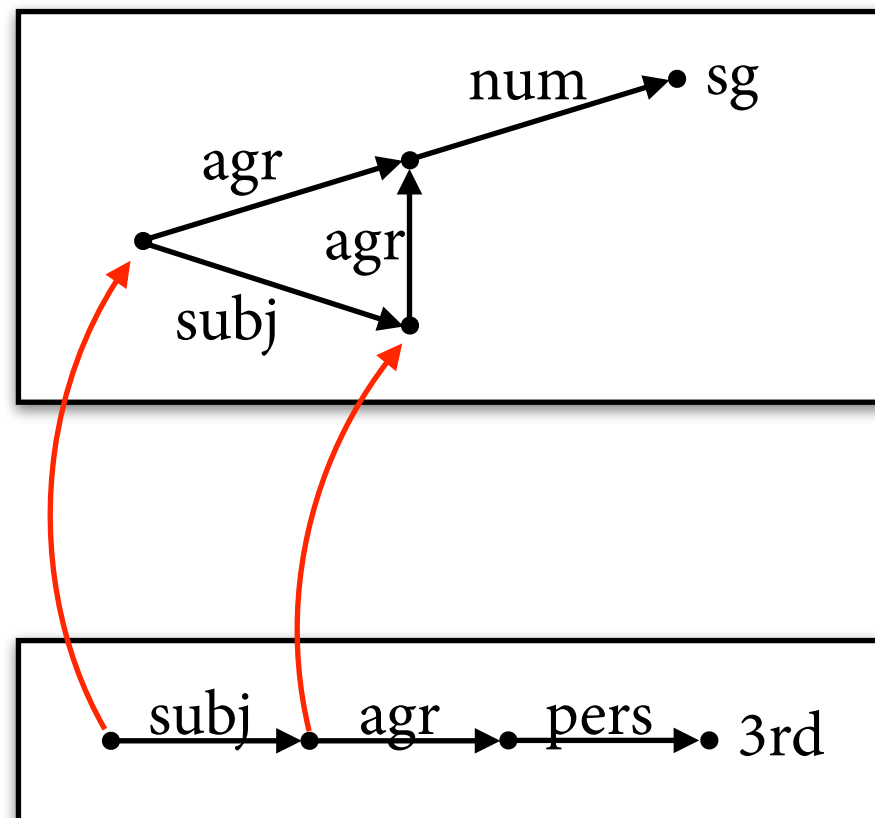
# Beispiel



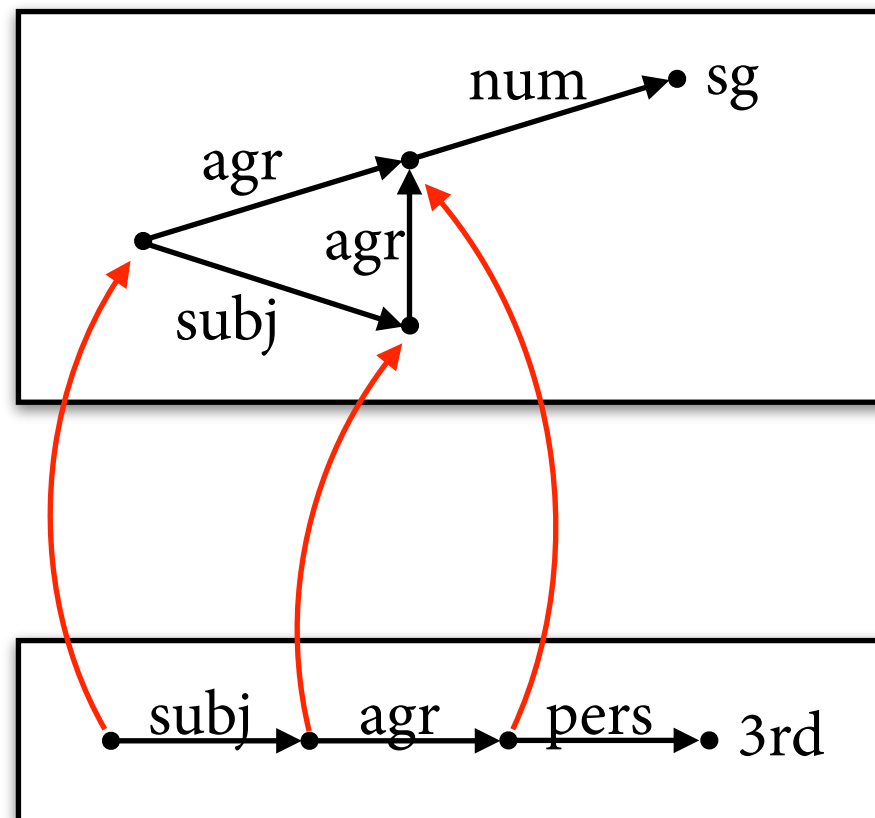
# Beispiel



# Beispiel

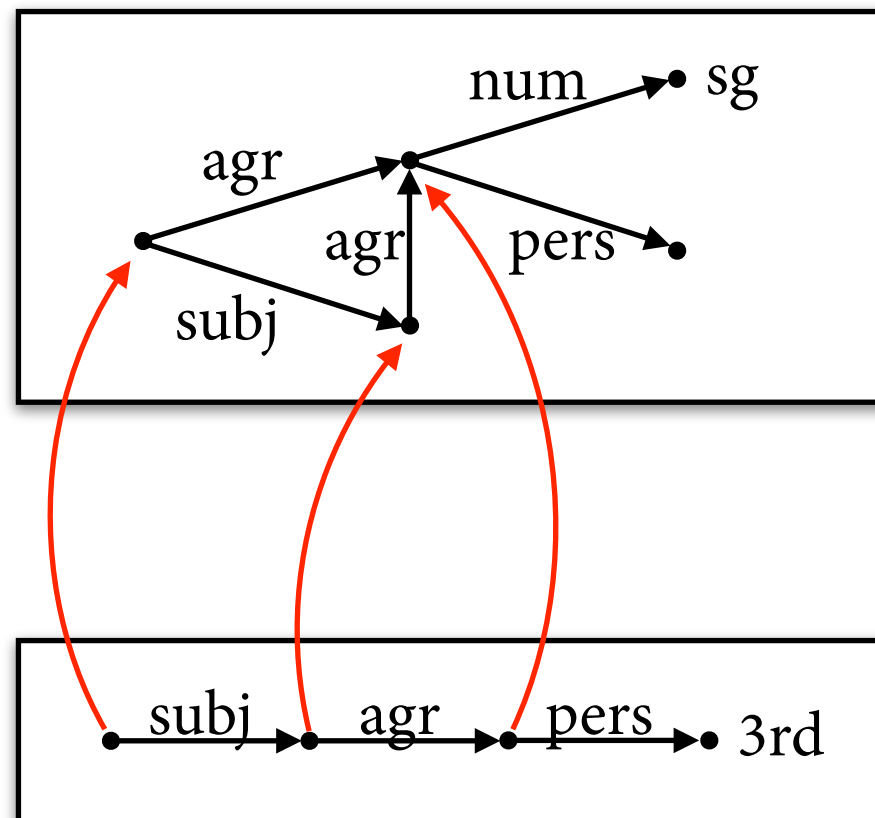


# Beispiel

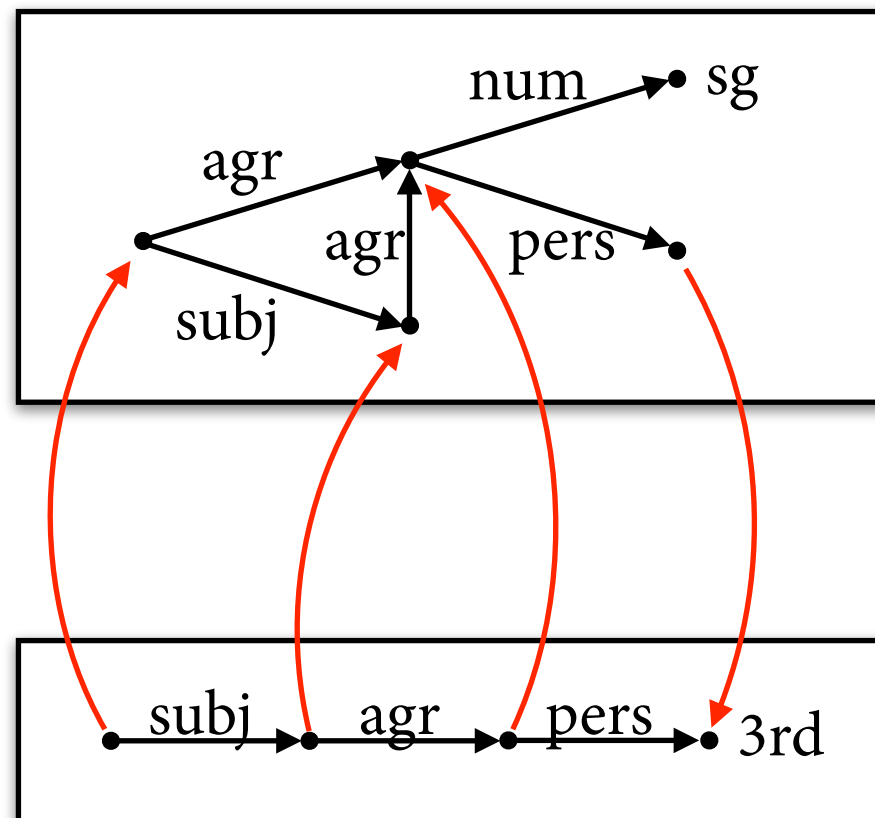




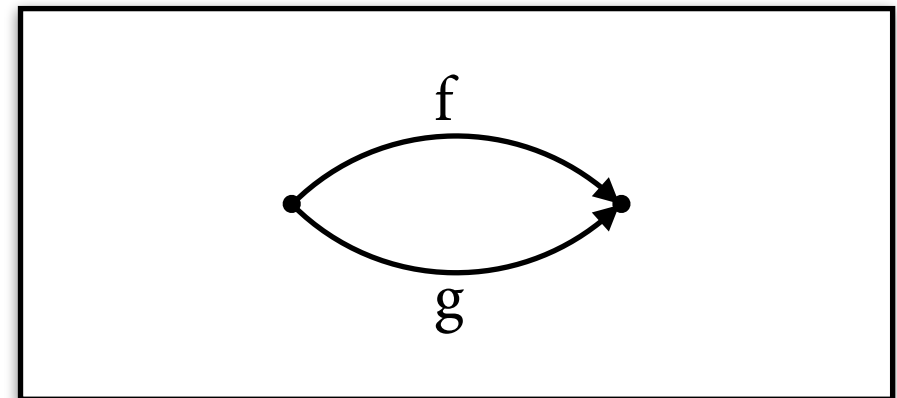
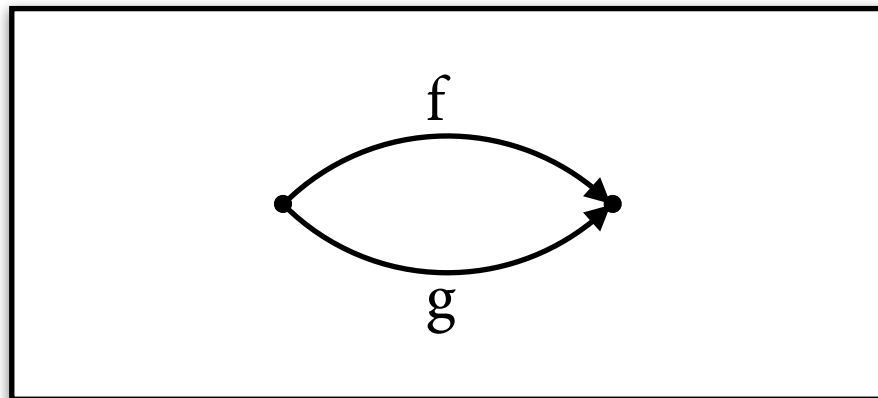
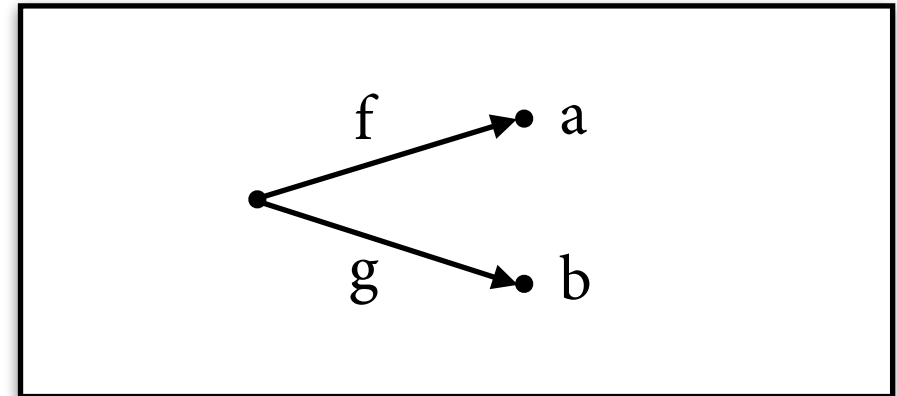
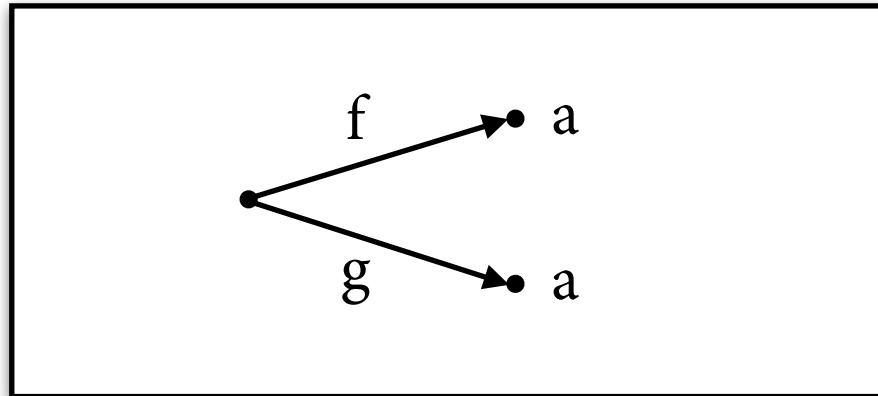
# Beispiel



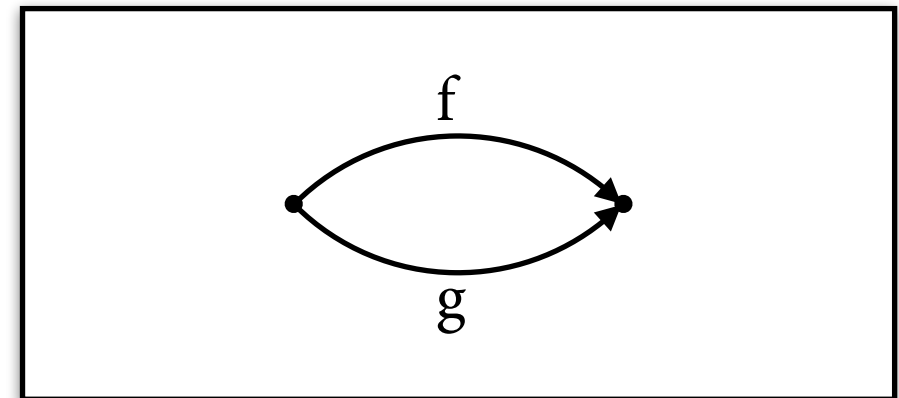
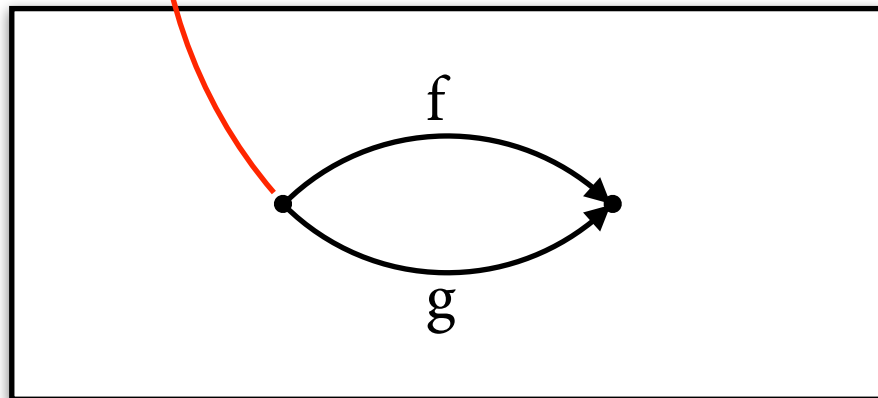
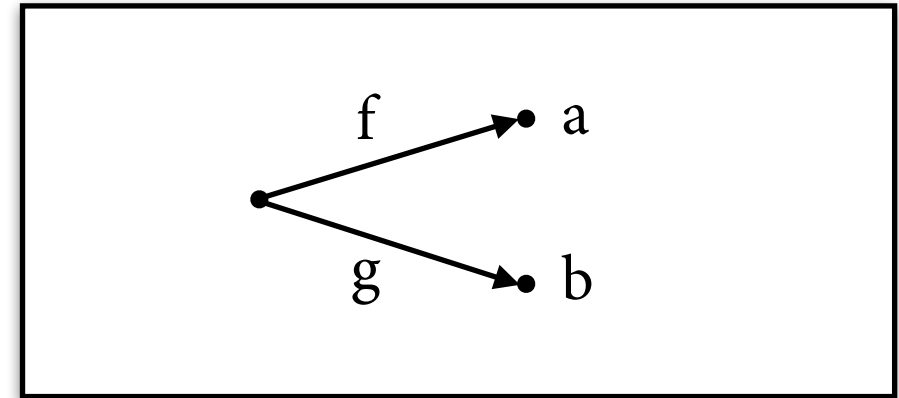
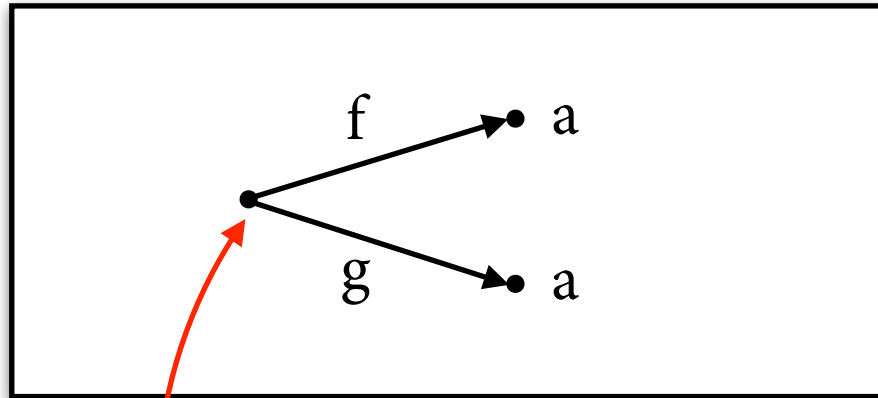
# Beispiel



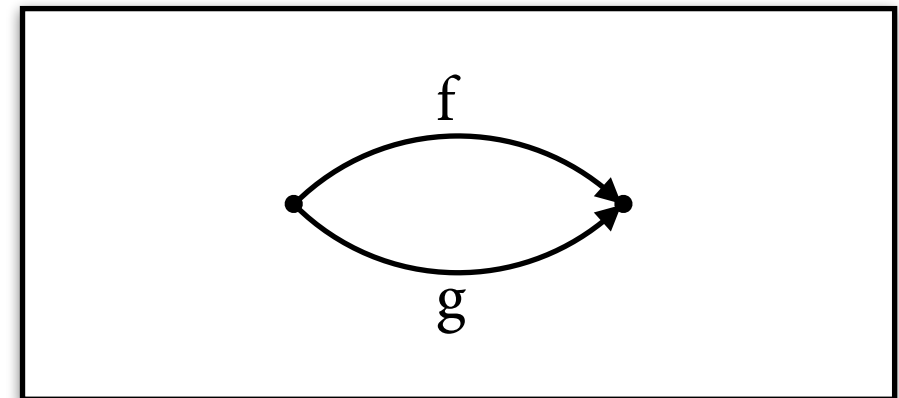
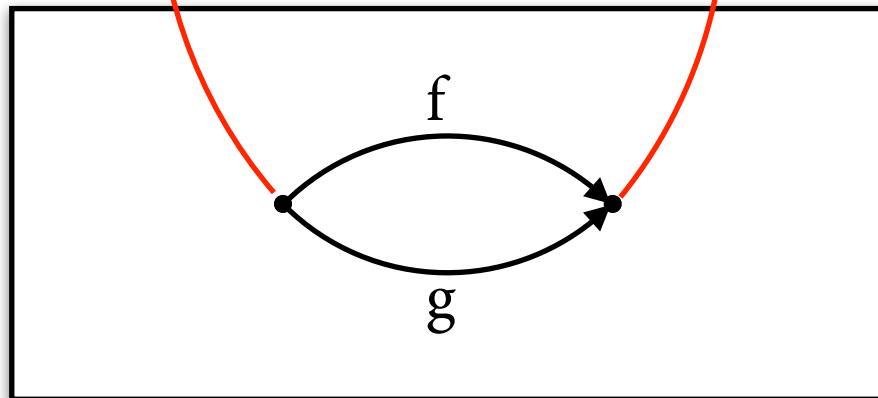
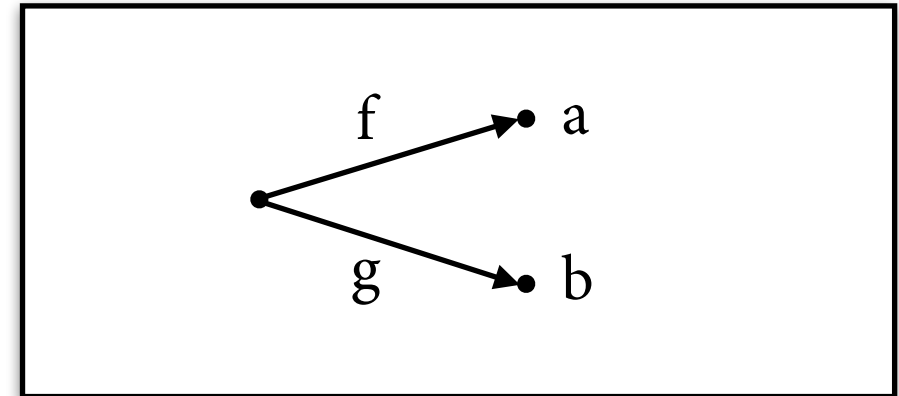
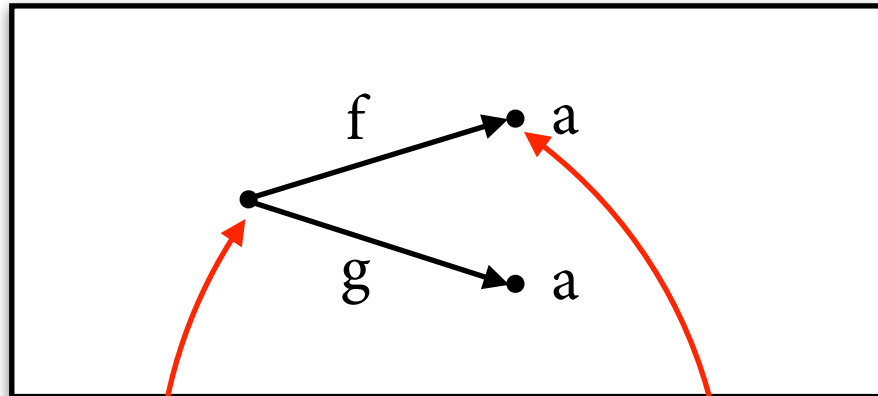
# Beispiel



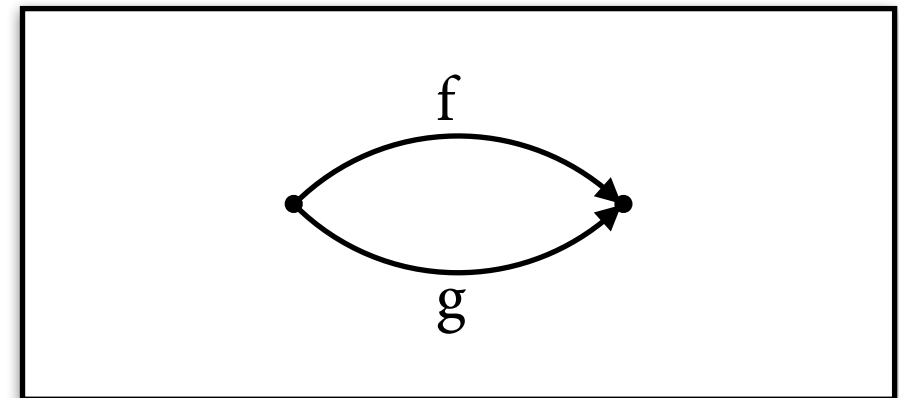
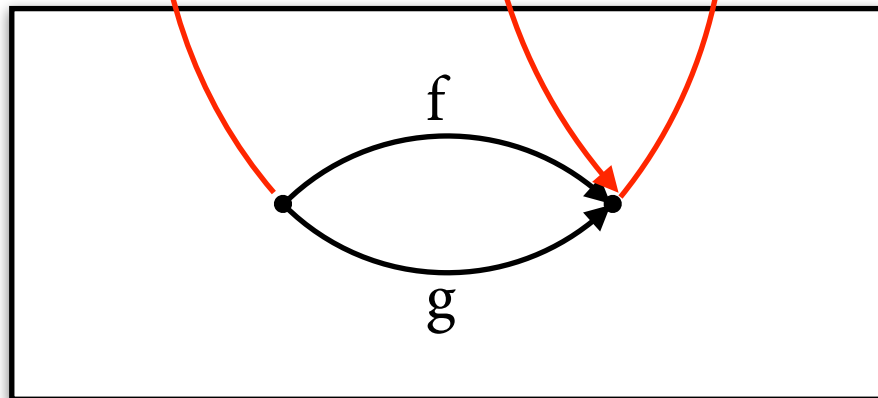
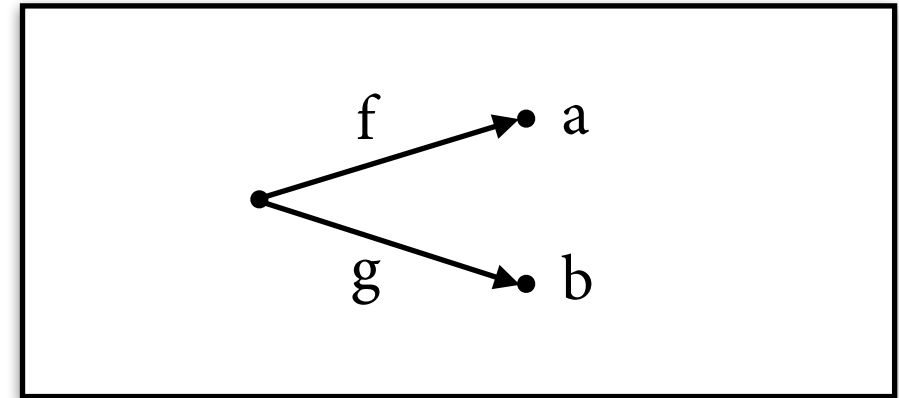
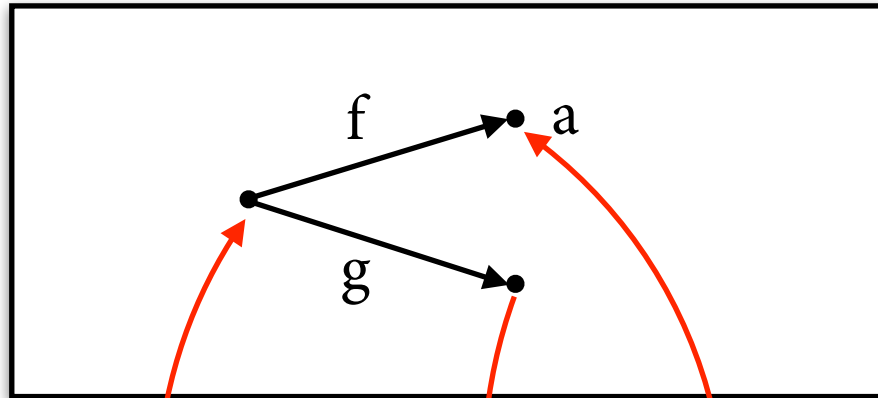
# Beispiel



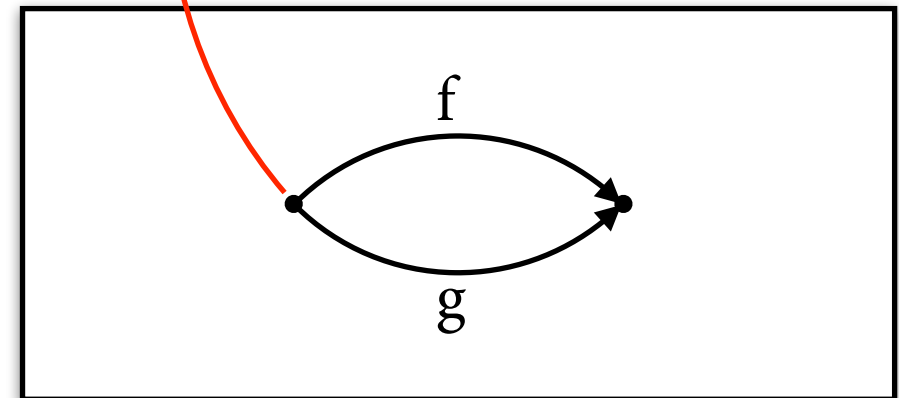
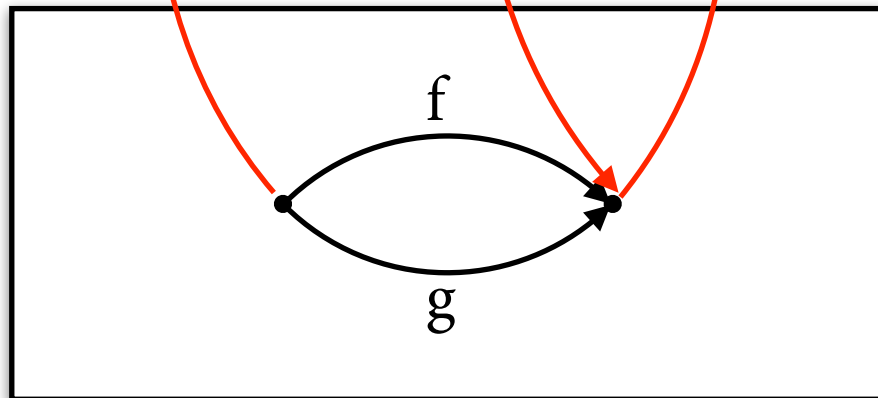
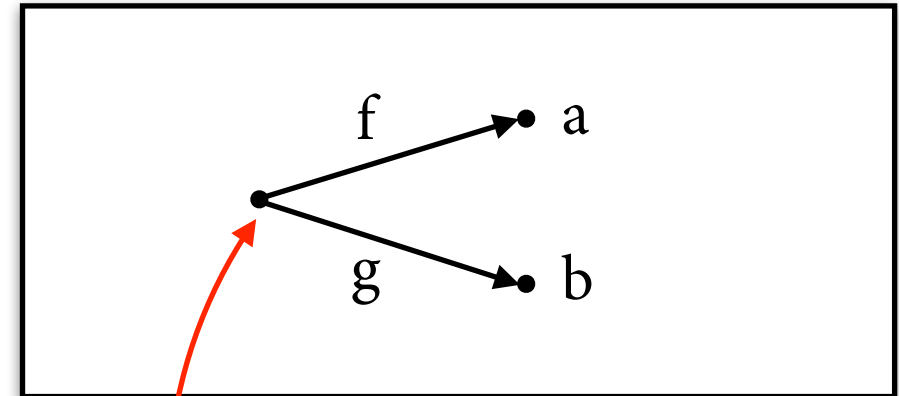
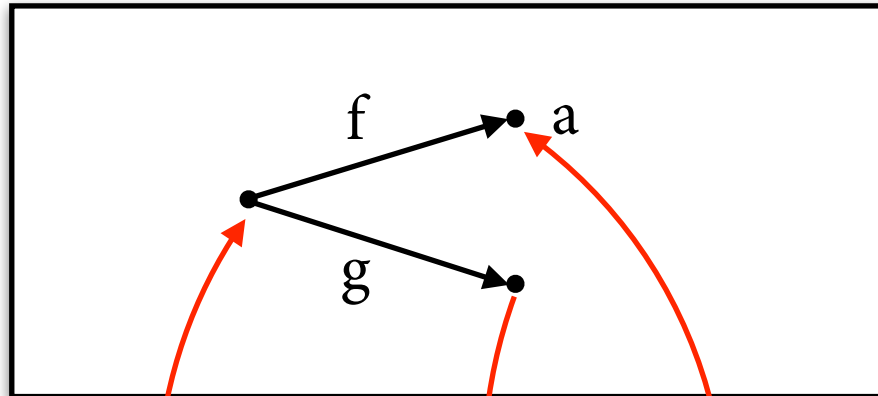
# Beispiel



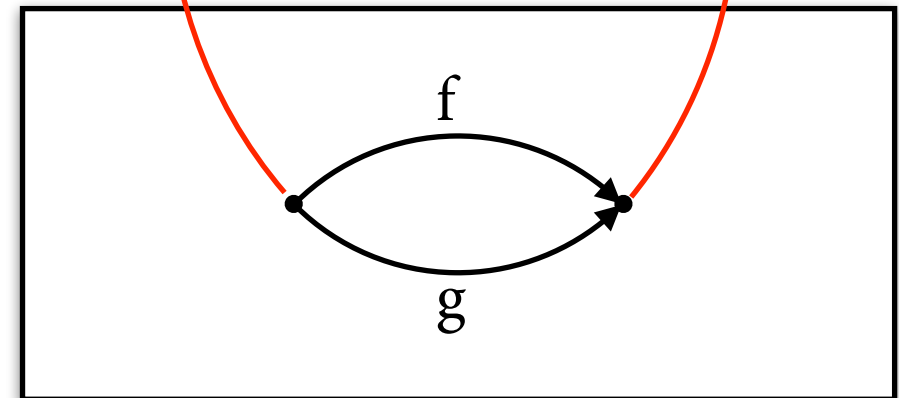
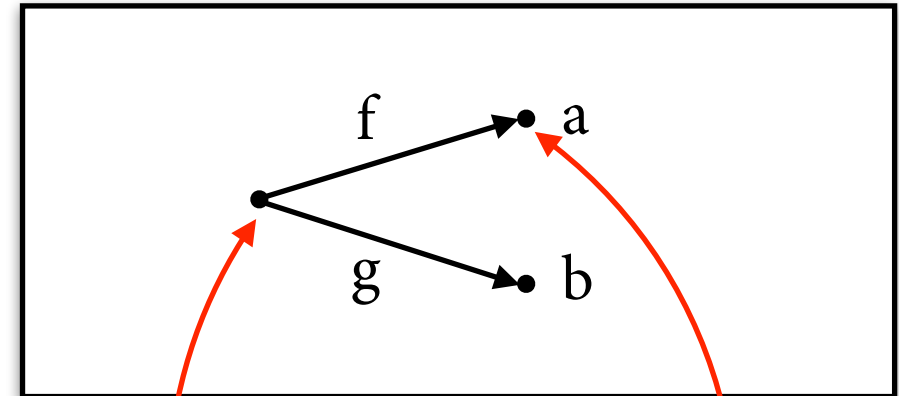
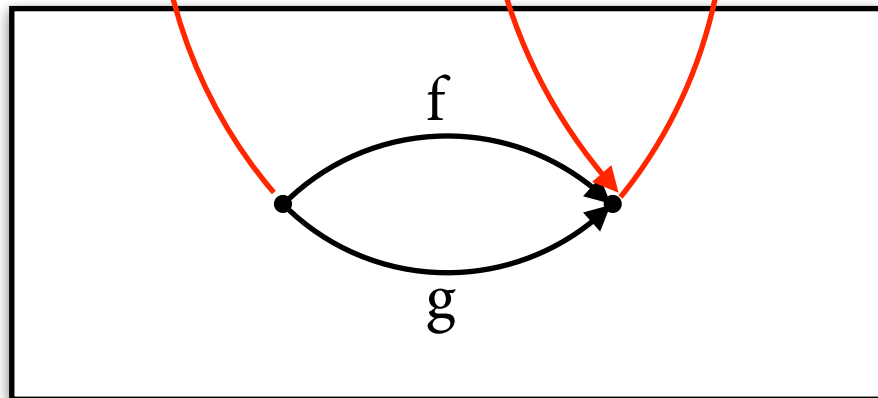
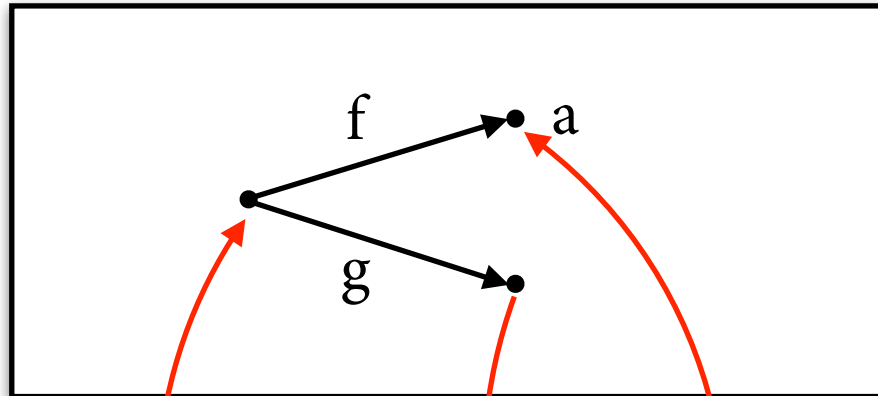
# Beispiel



# Beispiel

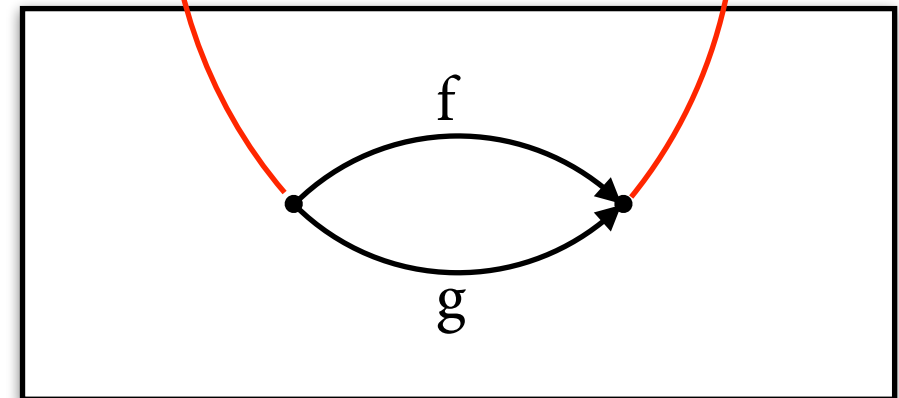
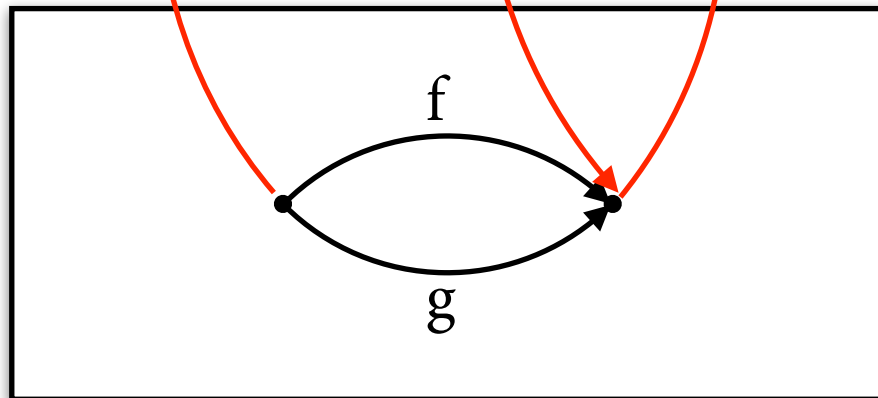
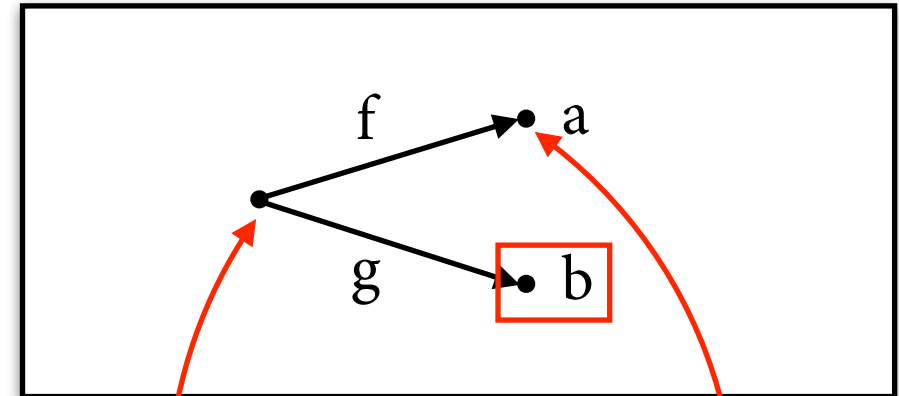
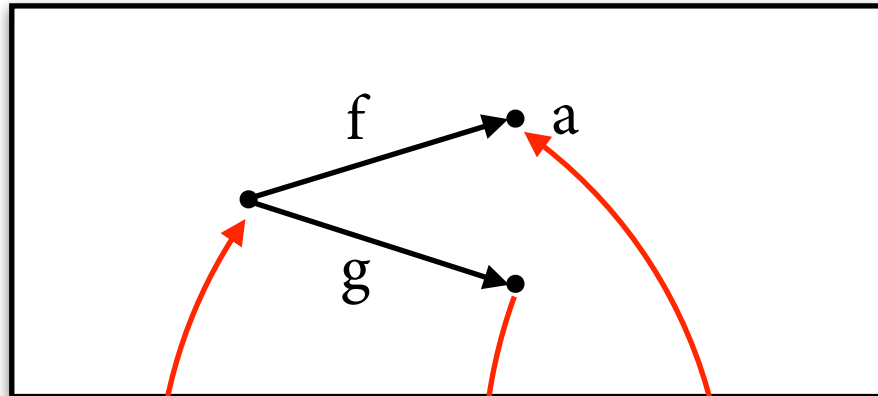


# Beispiel





# Beispiel



# Evaluation

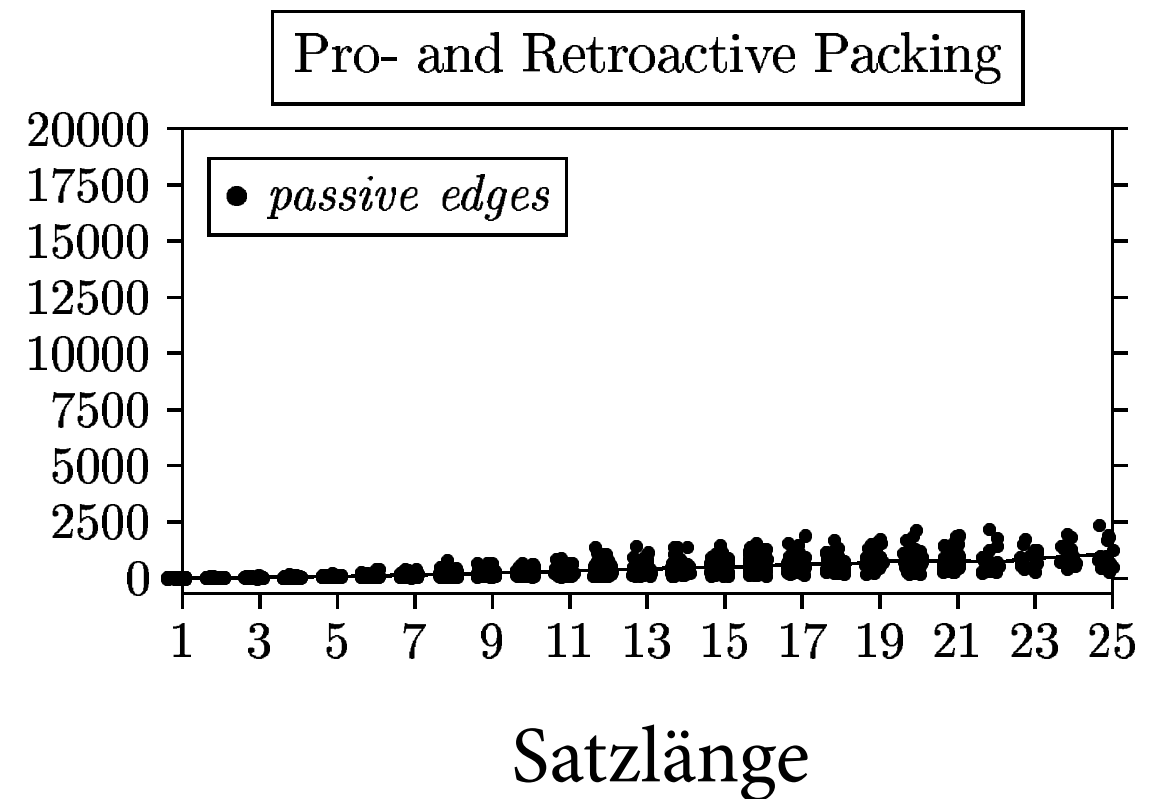
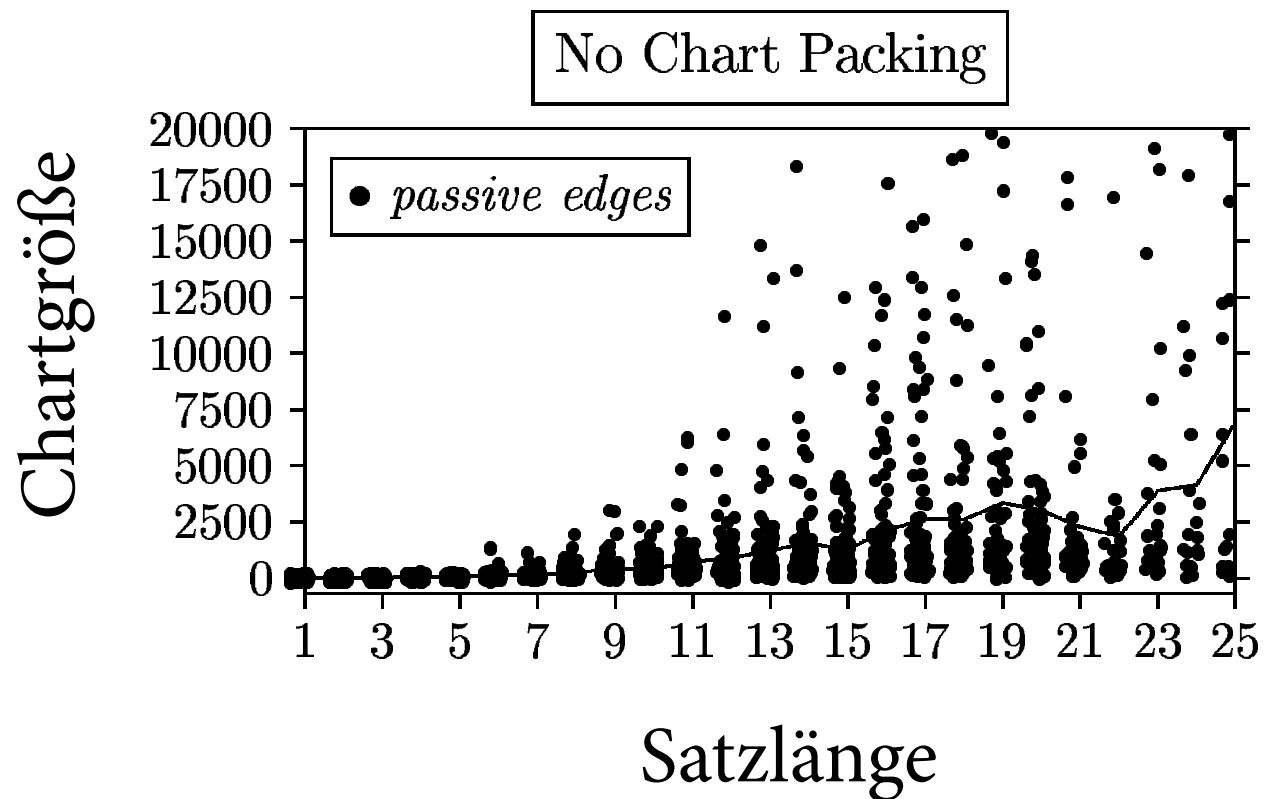
Unifier	tcpu		space (kb)	over (nodes)	copying		redundant (nodes)
	(s)	(s)			early (nodes)	(nodes)	
<i>quick check</i>	<i>on</i>	<i>off</i>	<i>on</i>	<i>on/off</i>	<i>on</i>	<i>off</i>	<i>on/off</i>
<i>unify1</i>	0.281	1.151	6,672	39.5	226.9	208.2	156.0
<i>unify2</i>	0.205	0.407	5,614	4.2	78.3	52.6	120.6
<i>unify3</i>	0.232	0.406	5,702	9.7	161.0	137.4	126.2
<i>tomabechi</i>	0.170	0.267	5,214	–	–	–	83.6
<i>tom-smart</i>	0.167	0.258	2,174	–	–	–	–

(Callmeier 2000; englisches Verbmobil-Testset mit LinGO-Grammatik; gemessen auf 500 MHz Pentium III)

# Ambiguitäten packen

- Unter bestimmten Umständen kann man zwei FSen als äquivalent gelten lassen.
  - ▶ wenn Parser  $(i,k,\sigma)$  kannte und  $(i,k,\tau)$  findet und  $\tau \sqsubseteq \sigma$  (d.h. ist allgemeiner), dann darf er  $(i,k,\sigma)$  löschen.
  - ▶ Chart bleibt kleiner  $\Rightarrow$  Parsing wird schneller.
- Löschen von Einträgen, die aus  $(i,k,\sigma)$  abgeleitet wurden, ist nicht trivial.

# Evaluation

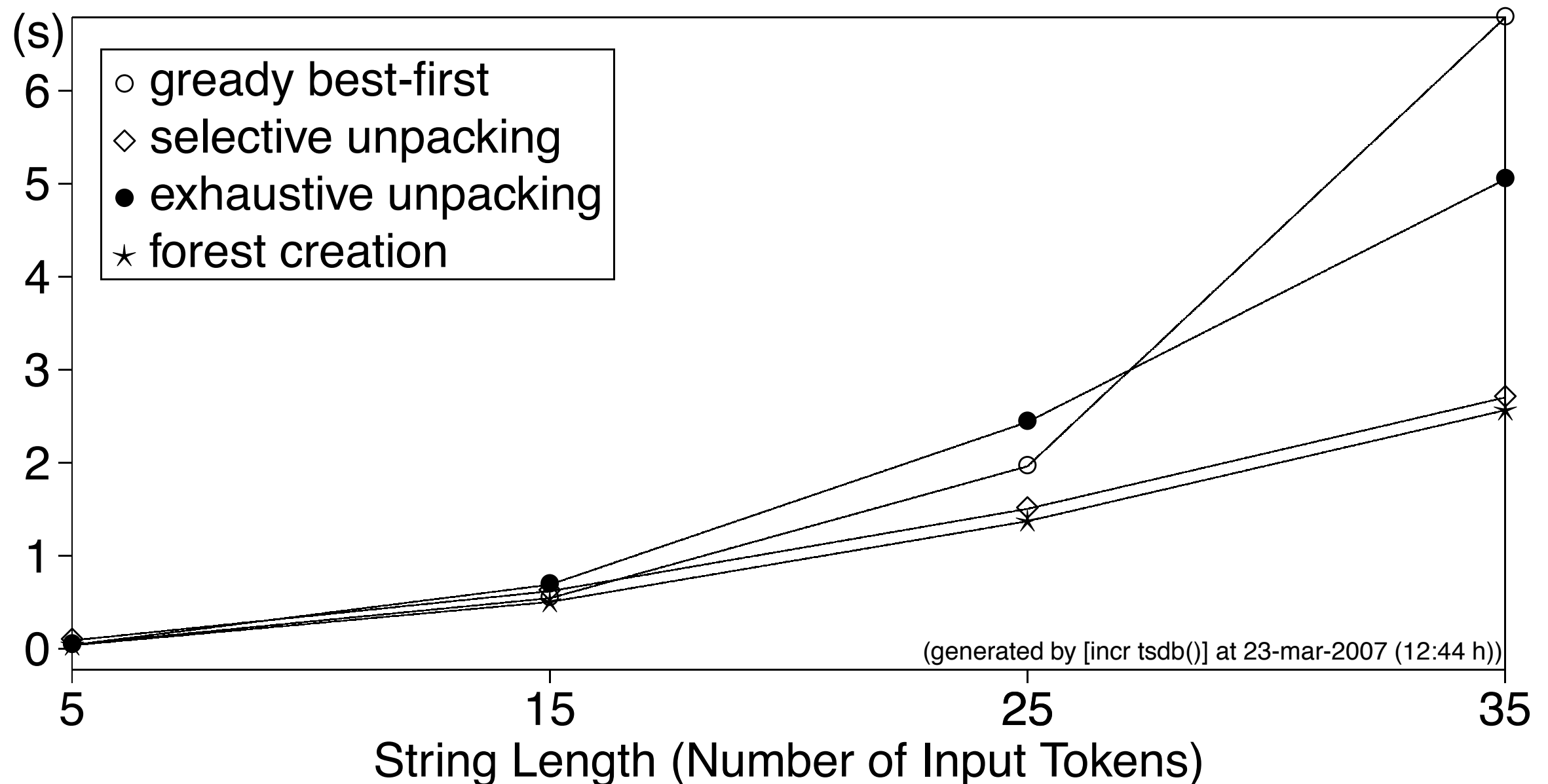


(Oepen & Carroll 2000; englisches Verbmobil-Testset mit LinGO-Grammatik)

# Approximation

- HPSG-Schemata definieren Phrasenstruktur. Kann man mit kfG approximieren.
  - ▶ [cat=V, subcat=<NPnom>] →  
[cat=V, subcat=<NPnom, NPacc>] [cat=N, subcat=<>]
- Mit (statistischer) kfG parsen; dann versuchen, beste Parsebäume in HPSG-Parses zu expandieren. (Kann fehlschlagen!)
- Unifikationen nur für Strukturen ausführen, die überhaupt eine Chance haben.

# Polynomial Time (Practical) Parsing



→ **Average Time for 20-Word Sentences around One Second**



# Stand der Kunst

- LKB: Grammatikentwicklungssystem.
  - ▶ verwendet Type Description Language (TDL), um Grammatiken aufzuschreiben
- TDL-Grammatiken kann man effizient mit dem PET-Parser verarbeiten.
  - ▶ LKB + PET + andere = DELPH-IN
  - ▶ DELPH-IN ist internationales Konsortium, das zueinander kompatible Parser, Grammatiken etc. entwickelt:  
<http://www.delph-in.net/>

# Zusammenfassung

- Expressivität von HPSG: turing-vollständig.
- Parsing von HPSG:
  - ▶ Wird in der Praxis dominiert von Zeit für Unifikation und Subsumptionstests.
  - ▶ Kann durch geschickte Algorithmen für diese Probleme in der Praxis effizient werden.
  - ▶ Algorithmen auch für andere Grammatikformalismen anwendbar, z.B. FTAG, LFG.



# Fazit

