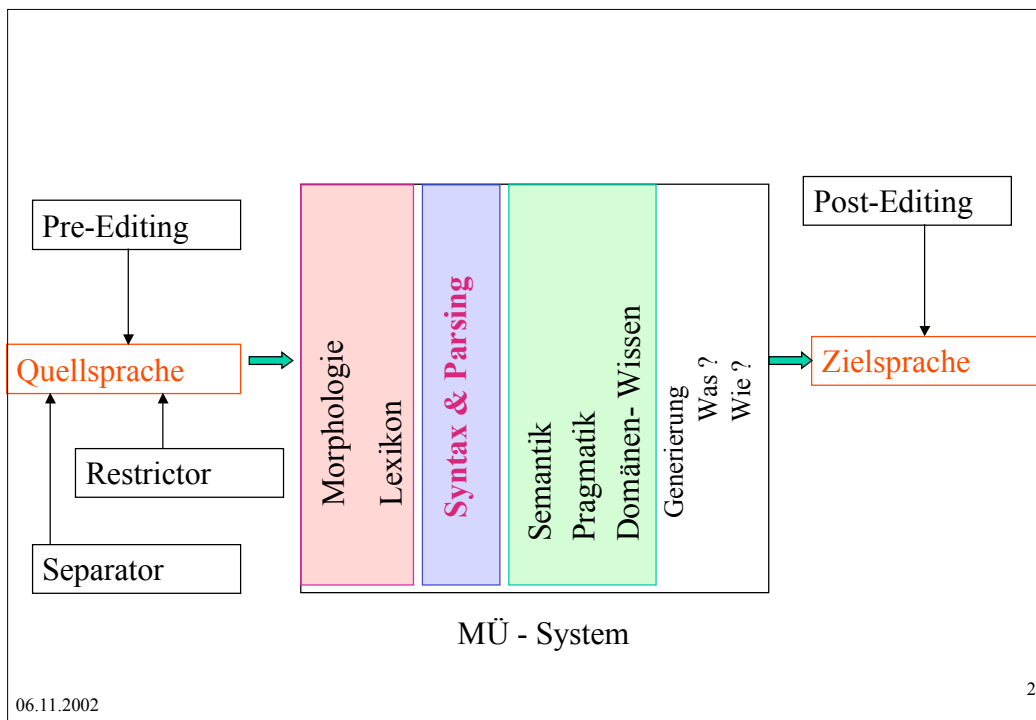


# Syntax und Parsing

06.11.2002

1



06.11.2002

2

## Definitionen -1-

- **Syntax:** Beschreibung der Regeln zur korrekten Kombination der lexikalischen Elemente einer Sprache.
- **Parsing:** Maschinelle (formale) syntaktische Analyse

06.11.2002

3


## Definitionen -2- Konstituenten

- NP - used to refer to things, places, concepts, events, qualities
  - single pronouns
  - proper nouns
  - common nouns = head
  - head + specifiers and (or) qualifiers
    - » specifiers : ordinals, cardinals, determiners, possessives, quantifiers
    - » qualifiers : adjectives or nouns used as modifiers
- VP
  - verb = head
  - head and complements
- AdjP
  - single adjectives
  - degree modifier
- PP
  - Preposition + other phrase
- AdvP : indicators of degree, locations, frequency etc.

06.11.2002

4

## Gliederung

- Syntax 
- Parsing
- Syntax und Parsing in MÜ Systeme

06.11.2002

5

## Zwei Sichten auf die Syntax

1. Im Zentrum steht die syntaktische Struktur als abstrakte Beziehung zwischen präterminalen Symbolen. Diese Beziehung kann dargestellt werden als

- Ersetzungsregel
- Netzübergänge oder
- Erkennungsprozeduren

Ein Lexikon stellt die Verbindung her zwischen den "Satzmustern" und den Wörtern

Regelgrammatik

2. Im Zentrum stehen Wortklassen oder Lexikoneinträge, die die Verbindbarkeit von Wörtern oder Wortklassen beschreiben. Durch Verbinden der Strukturinformation zu immer größeren Einheiten wird eine Satzbeschreibung erreicht

Lexikalische Grammatik

In beiden Grammatiktypen können Abhängigkeits- und Konstituenzstrukturen dargestellt werden.

06.11.2002

6

## Grammatische Beschreibungstypen

- Patterns
- Rekursive Netze (RTN)
- Kontextfreie Grammatik
- Dependenzgrammatik
- Konstituentenstrukturgrammatik
- Funktionale Grammatik
- Unifikationsgrammatik
- Kasusgrammatik

06.11.2002

7

## Patterns -1-

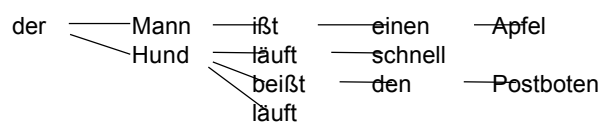
### Patterns

#### ● Direkte

- wörtliche "die linke hintere Maschine jetzt anfahren"
- offene "(&c (das) (ist) \$\$)" matcht alle Permutationen und akzeptiert ein beliebiges Folgewort: *das ist ein Fehler* und *ist das eine Maschine*
- lexikalische "Det, Adj, Adv, Adj, N, V" *Der rote sehr helle Knopf leuchtet*

#### ● mit Variablen "die 1 2 3 jetzt anfahren"

### Wortübergangsgraph (Strukturierte Pattern)



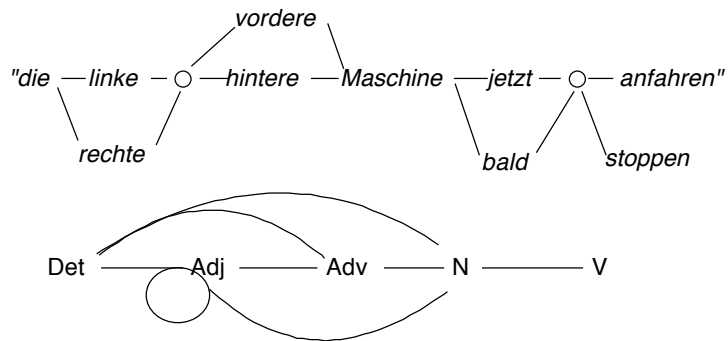
06.11.2002

8

## Patterns -2-

### Strukturierte Patterns

Übergangsnetz z (Endlicher Automat, Zustandsgraph)



Regulärer Ausdruck

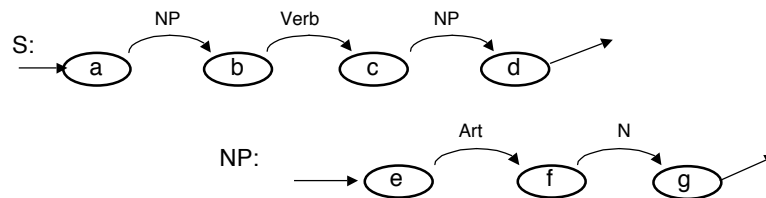
Det (Adj\*) (Adv) N V

06.11.2002

9

## Rekursive Netze

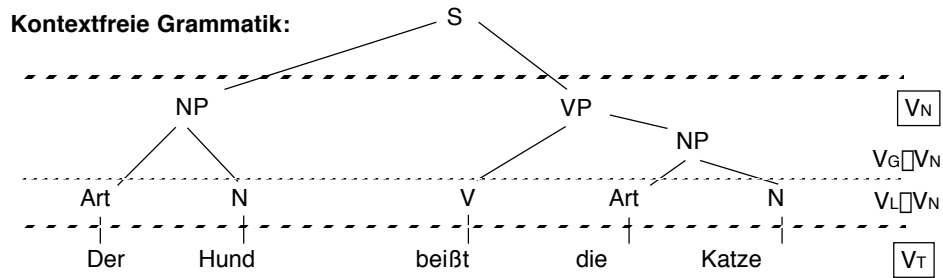
### Rekursives Netz (RTN)



06.11.2002

10

## Kontextfreie Grammatik



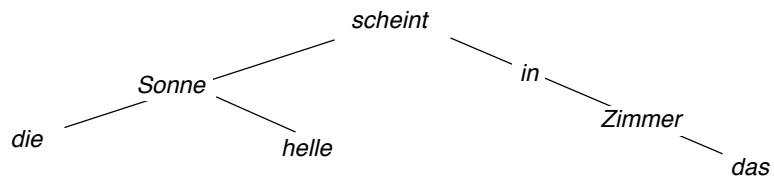
Verarbeitungsfähige systematische Methode

06.11.2002

11

## Abhängigkeitsgrammatik -1-

Abhängigkeit (Beispiel: Head und Modifizier-Grammatik):

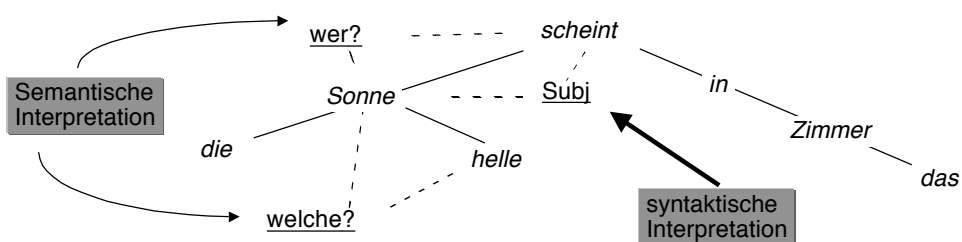


06.11.2002

12

## Dependenzgrammatik-2-

Dependenz (Beispiel: Head und Modifier-Grammatik):



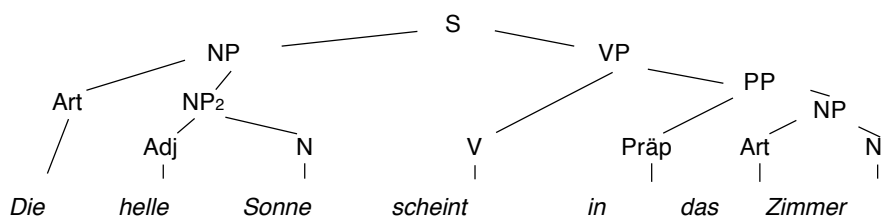
Eigenschaften: Stellt die syntaktische Abhängigkeit vom Verb dar, keine Projektion auf Oberflächenkette  
basiert auf der Verbalenz (mit obligatorischen und freien Ergänzungen)

06.11.2002

13

## Konstituentenstrukturgrammatik

Konstituenz (Beispiel IC-Analyse):



Eigenschaften: Unmittelbare Bestandteile des Satzes werden beschrieben, Oberflächenstruktur bleibt sichtbar, Meist binäre Teilung, über Teilbarkeit entscheidet Substitutionstest (Minimalpaare), Probleme bei Diskontinuität von Elementen.

06.11.2002

14

## Konstituenten -Findungsprozeduren

- **Substitutionstest:** Elemente, die man für einander ersetzen kann, ohne daß sich an der Grammatikalität des Satzes etwas ändert, sind Konstituenten.

*Die langen Winterabende / Die Winterabende / Winterabende versetzen mich in Melancholie*

- **Pronominalisierungstest:** Was sich zusammen pronominalisieren läßt, ist eine Konstituente

*Ede und Caro wohnen in der Maxstraße. Sie haben ein Kind.*

- " -

*.Dort wohnen nur reiche Leute*

- **Fragetest:** Wonach man fragen kann, ist eine Konstituente

*Heute kommt meine Mutter zu Besuch. Wer kommt heute zu Besuch?*

- **Koordinationstest:** Was sich koordinieren läßt, ist eine Konstituente

*Ede gab eine sehr kluge und keineswegs triviale Antwort*

06.11.2002

15

## Funktionale Grammatiken - 1-

- Alle grammatischen Fakten können als constraints beschrieben werden:
  - $S \square NP VP$  kann man lesen als die Bedingung, daß ein korrekter Satz aus NP und VP bestehen muß (top-down gelesen),
  - $Person = Nom Sg$  kann man lesen als die Bedingung, daß das Wort *Person* nur eine syntaktische Stelle einnehmen kann, die die Eigenschaft *Nom Sg* hat (bottom-up gelesen).
  - für jede syntaktische Kategorie (klassisch: Nonterminale) gibt es eine Auswahl an Subkategorien und an relevanten Merkmalen.

06.11.2002

16



## Funktionale Grammatiken -2-

- Mit einem solchen constraint-Formalismus werden alle Stufen der Syntax beschrieben:
  - Satz,
  - Kategorien und
  - Lexeme.
- Ein constraint satisfaction system baut aus Teilen eine Gesamtinterpretation zusammen oder generiert aus Teilen eine Gesamtstruktur

06.11.2002

17

## Unifikationsgrammatiken

- Beschreibung von syntaktischen Merkmalen durch Attribut-Wert-Matrix (attribute terms / feature terms)
- Zentrale Operation: Unifikation

- **Grammatik-Formalismen:**

|                                      |         |                             |
|--------------------------------------|---------|-----------------------------|
|                                      | DCG     | Pereira/Shieber             |
|                                      | PATR II | Shieber, Karttunen          |
| Grammatik-Theorien:                  |         |                             |
| Functional Unification Grammar       | FUG     | Kay                         |
| Lexical-Functional Grammar           | LFG     | Kaplan/Bresnan              |
| Generalized Phrase Structure Grammar | GPSG    | Gazdar/Klein/<br>Pullum/Sag |
| Head-Driven Phrase Structure Grammar | HPSG    | Pollard/Sag                 |

06.11.2002

18

## Vorteile unifikationsorientierter Grammatiken

**einheitliche Operation** für Merkmale aller Ebenen (Phonologie, Morphologie, Syntax, Semantik)

**Deklarative Form** daher lesbar ohne Kenntnis der Prozesse

**Robustheit** weil korrekte Teilstrukturen definiert und durch Vorhersagbarkeit von Wirkungen allein aus den im Grammatikformalismus notierten Regeln

**Modularität** durch hierarchische Strukturierung der Information

**Homometrie** d.h. kleine Änderungen in den Regeln haben geringe Wirkungen auf die Resultate, große Änderungen, große Wirkungen

## Eigenschaften unifikationsbasierter Grammatiken

**oberflächenorientiert** d.h. es ist eine direkte Charakterisierung der faktischen Ordnung der Oberflächenstruktur - der aktuellen Anordnung der Zeichen in einem Satz - möglich

**interpretativ** in dem Sinne, daß mit den Zeichenketten wohldefinierte Beschreibungen assoziiert werden

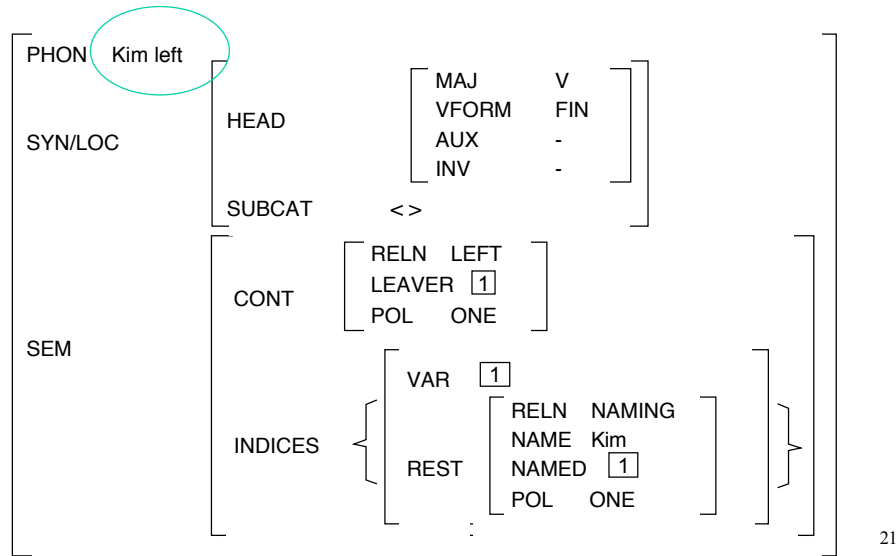
**induktiv** durch rekursive Definition der Assoziation von Beschreibungen mit Zeichenketten, so daß neue Beschreibungen aus Beschreibungen von Teilketten auf wohldefinierte Weise erzeugt werden, wobei die Kombinationsoperationen von Zeichenketten vorgegeben sind

**deklarativ** die Assoziation zwischen Zeichenketten und Beschreibungen ist in einer Weise definiert, die nur angibt, welche Assoziationen zulässig sind, nicht, wie sie berechnet werden.

**merkmalsbasiert** durch Mengen von Paaren aus Merkmalsnamen und (ggf. strukturierten) Merkmalswerten (f-Strukturen)

# HPSG

Carl Pollard  
Ivan A. Sag:  
Information-Based Syntax and Semantics  
Stanford 1987. S. 93



06.11.2002

21

## Syntaktische Mehrdeutigkeit

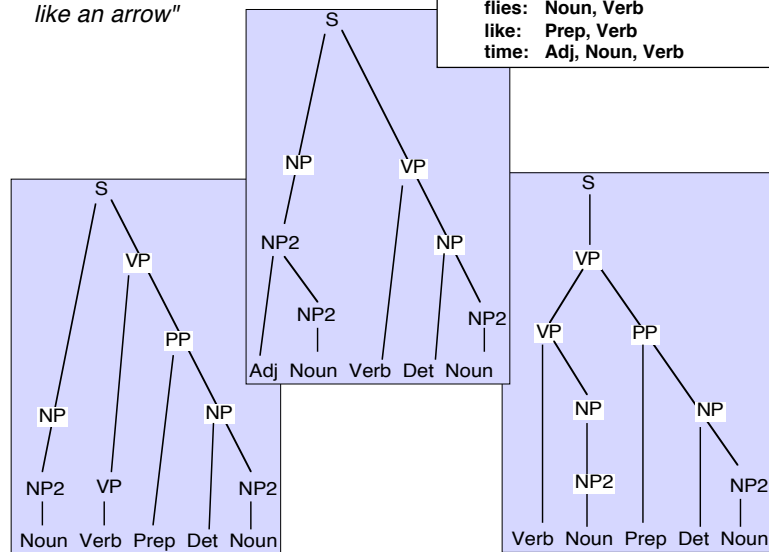
Grammatik

- S  $\emptyset$  NP VP
- (Aussagen)
- S  $\emptyset$  VP (Befehle)
- NP  $\emptyset$  Det NP2
- NP  $\emptyset$  NP2
- NP  $\emptyset$  NP PP
- NP2  $\emptyset$  Noun
- NP2  $\emptyset$  Adj NP2
- PP  $\emptyset$  Prep NP
- VP  $\emptyset$  Verb
- VP  $\emptyset$  Verb NP
- VP  $\emptyset$  VP PP

"time flies  
like an arrow"

Lexikon

- an: Det
- arrow: Noun
- flies: Noun, Verb
- like: Prep, Verb
- time: Adj, Noun, Verb



Nach Winograd  
06.11.2002

## Schwierigkeiten für deutsche Grammatiken

### 1. Flexions- und Kongruenzregeln

- Person-Numerus-Kongruenz bei Verb und Subjekt: *"ich gebe"*
- Kasus-Genus-Numerus-Kongruenz in NPs: *"einem großen Bären"*
- Syntaktische Restriktionen des Verbs für Ergänzungen: *"jemanden einer Sache verdächtigen"*
- artikelbestimmte Flexion des Adjektivs: *"ein großer / der große .."*

### 2. Variable Satzgliedstellung für einige Teilstrukturen (NPs)

### 3. Diskontinuierliche Komponenten

- Verbalklammer: *"er wird nächste Woche ankommen"*
- Fernstellung des Nebensatzes: *"Die Meinung ist falsch, daß das einfach ist"*


### 4. Koordination

- Präpositional-Attribute: *"Der Dicke in meinem Abteil mit dem riesigen Koffer auf dem Sitz neben sich"*
- Adjektiv-Attribute: *"die phantastische, überwältigende, überzeugende Lösung"*
- Nominalgruppen: *"Vorlesungen, Seminare und andere Veranstaltungen"*

06.11.2002

23

## Gliederung

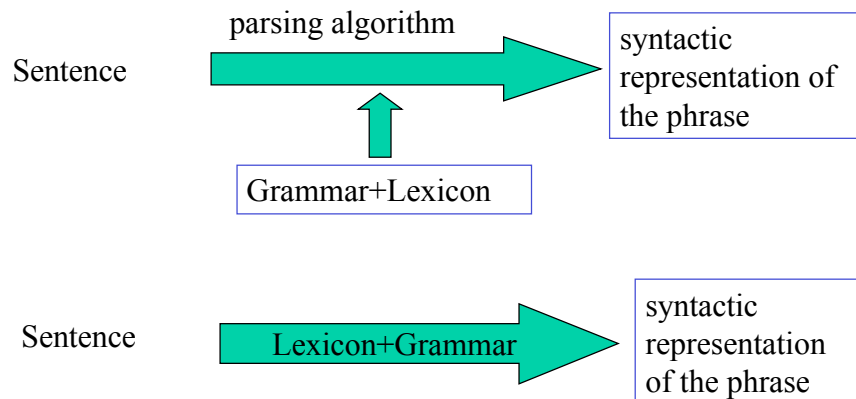
- Syntax
- Parsing 
- Syntax und Parsing in MÜ Systeme

06.11.2002

24

## Parser

Two types:



06.11.2002

25

## Grammatik und Parser

nach Hellwig

- Trennung von Grammatik und Parser
  - Gründe:
    - Grammatik existiert unabhängig von Programm und Daten
    - Grammatik ist eine Beschreibung der Objektsprache in deklarativer Form (unabhängig vom Erkennungsprozeß)
    - Parser fußt auf der Syntax des Grammatikformalismus, nicht auf den Inhalten der Grammatik
    - Interpretierender Parser
  - Vorteile: Getrennte Erweiterbarkeit, Lesbarkeit der Grammatik durch Linguisten, Parserbau während des Grammatikaufbaus, Anwendung auf mehrere Sprachen, Verwendung der Grammatik beim Generieren.
- Integration von Grammatik und Parser
  - Gründe:
    - Grammatik ist in das Programm eingebettet
    - Formulierung der Sprachbeschreibung mit Blick auf die Verarbeitung
    - Grammatik ist eine Sprachbeschreibung in prozeduraler Form (Erkennungsanweisungen)
    - Prozeduraler Parser
  - Vorteil: Effizienz, Formulierung ohne Grammatikformalismus möglich

06.11.2002

26

## Parsingstrategien - 1 -

|                         |                   |    |                        |
|-------------------------|-------------------|----|------------------------|
| • <b>Architektur</b>    | parallel          | vs | sequentiell            |
| • <b>Startpunkt</b>     | top-down          | vs | bottom-up              |
| • <b>Regelanwendung</b> | deterministisch   | vs | nicht-deterministisch  |
| • <b>Auswahl</b>        | first guess       | vs | informierte Auswahl    |
| • <b>Vorgehen</b>       | Breitensuche      | vs | Tiefensuche            |
| • <b>Richtung</b>       | rechts ---> links | vs | links ---> rechts      |
|                         |                   | vs | Insel-Parsing          |
| • <b>Steuerung</b>      | syntaktisch       | vs | semantisch             |
|                         |                   | vs | semantisch-pragmatisch |

06.11.2002

27

## Parsingstrategien - 2 -

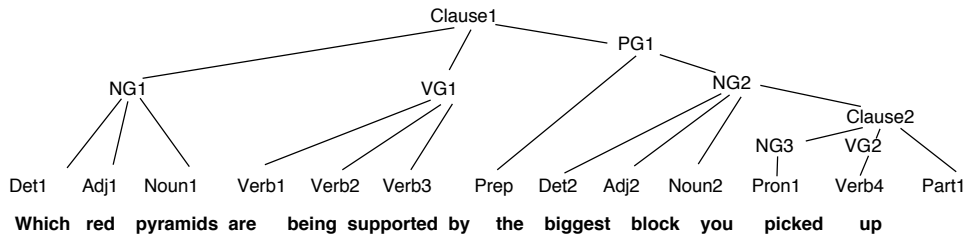
|  |                   |    |                     |
|--|-------------------|----|---------------------|
| • <b>Zielstruktur</b>                      | syntaktisch       | vs | semantisch          |
| • <b>Formalisiertheit</b>                  | ad hoc            | vs | uniform             |
| • <b>Ebenenspezifität</b>                  | ebenengebunden    | vs | systemorientiert    |
| • <b>Genauigkeit</b>                       | partiell genau    | vs | formal beschreibbar |
| • <b>Informiertheit</b>                    | blind             | vs | vorherschauend      |
| • <b>Betrachtungsraum</b>                  | wortweise         | vs | gruppenweise        |
| • <b>Ambiguitäts-</b><br><b>behandlung</b> | erfolgsorientiert | vs | exhaustiv           |

06.11.2002

28

## Parsing-Ergebnis: SHRDLU

## Syntakt. Relationen



Clause1 (*Which red pyramids are being supported by the biggest block you picked up?*) Clause  
Major Question Wh- Transitive Passive (Present in Present) Subject-Question Agent

NG1 (*Which red pyramids*) NG Determined Indefinite Question Plural Subject

- Det1 (*Which*) Determiner Question-Determiner Singular Plural
- Adj1 (*red*) Adjective
- Noun1 (*pyramids*) Noun Plural

VG1 (*are being supported*) VG Finite Passive (Present in Present)

- Verb1 (*are*) Verb Auxiliary Intensive Be Plural Present
- Verb2 (*being*) Verb Auxiliary Intensive Be Ing
- Verb3 (*supported*) Verb Transitive Past-Participle

06.11.2002

29

## Parsing-Ergebnis: GUS Semantische Slot Filler

"I want to go to San Diego on May 28"

(Client Declare

(Case for *want* / e (Tense Present)

Agent = Dialog.Client.Person

Event = (Case for *go* (Tense Present)

Agent = Dialog.Client.Person

To-Place = (Case for City

Name = *San Diego*)

Date = (Case for Date

Month = *May*

Day = *28* )))

Nach Winograd

06.11.2002

30

## Parsing-Ergebnis: DB-Interface Suchausdruck in einer DB-Sprache

"Liste die Namen aller Lieferanten, die mindestens alle die Teile liefern, die auch vom Lieferanten S2 geliefert werden"

```

SELECT          UNIQUE S#
FROM            SP SP X
WHERE           NOT EXISTS
                (SELECT *
                 FROM SP SP Y
                 WHERE S# = 'S2'
                 AND      NOT EXISTS
                         (SELECT *
                          FROM SP
                          WHERE S# = SP X. S#
                          AND P# = SP Y. P#))
    
```



06.11.2002

nach Date

31

## Parsing-Ergebnis: Tabelle

| REF        | SCALE    | PHASE     | YIELD                   | TEMP      |
|------------|----------|-----------|-------------------------|-----------|
| para.1     | small    | solid     | 77%                     | -78 to 20 |
| TIME       | ENERGY   | APPARATUS | FEATURES                |           |
|            | cooling  |           | IR. N MR. MS            |           |
| REG. NO    | FUNCTION | AMT.      | AUTHOR ID               |           |
| 78624-62-1 | product  | 2.70 g    | 7a                      |           |
| 78624-61-0 | reactant |           | 6a                      |           |
| 13274-48-6 | reactant | 1.24 g    | N-methyltriazolinedione |           |
|            | solvent  | 80 ml     | pentane                 |           |
|            | solvent  | 40 ml     | ethyl acetate           |           |

Formular-Repräsentation des Systems SIE für eine chemische Reaktionsbeschreibung

06.11.2002

32



## Parsing-Ergebnis: Semantische Repräsentation

```
[ request (referent(_5747))
  presuppose (exists (_4340)) ]
  some (_4340)
[ unique (_4407)
  single (_4407)
  instance (_4407, person)
  propval (person,_4407,sex,male)
  [ some (_4725)
    [ unique (_5033)
      single (_5033)
      instance (_5033,project)
      propval (project,_5033,name,str (LOKI)) ]
```

"who is the man that leads the LOKI project?"

```
instance (_4725,leading)
  propval (leading,_4725,theta,_5033)
  propval (leading,_4725,alpha,_4407)
  topic (_4407) ] ]
instance (_4340,identity)
  propval (identity,_4340,alpha,_4407)
  propval (identity,_4340,theta,_5747)
  topic (_4407)
[ some (_5747)
  single (_5747)
  instance (_5747,person) ]
```

06.11.2002

33

## Pattern matching -1-

- Pattern = ein syntaktischer Rahmen für lexikalisch-semantische Äquivalenzklassen.
- Patterns beschreiben häufige Ausdrücke einer Sprache.

Welche Übungen gibt es im Montag studium ?

06.11.2002

34

## Pattern matching -2-

- Sammlung von Patterns spezifizieren
  - eine feste Wortreihenfolge
  - Plätze für variable Wörter
- Ein Satz kann verarbeitet werden, wenn es (mindestens) ein gültiges Pattern für diesen Satz gibt.
- Für einen Satz kann man ein oder mehrere Patterns benutzen
- Probleme: Syntaktisch unkorrekte Sätze werden auch verarbeitet
  - z. B: Welche am Montag gibt es ?

06.11.2002

35

## Sämtliche im System SIR (Raphael 1967) verwendeten Patterns

- Is \_ \_ ?
- How many \_ does \_ have?
- How many \_ does \_ own?
- What is the \_ of \_ ?
- How many \_ are parts of \_ ?
- How many \_ are there on \_ ?
- \_ has as a part one \_ .
- There is one \_ on \_ .
- \_ is just to the left of \_ .
- Is \_ just to the left of \_ ?
- \_ is just to the right of \_ .
- Is \_ just to the right of \_ ?
- \_ is \_ .
- \_ has \_ .
- \_ owns \_ .
- Where is \_ ?
- Is \_ part of \_ ?
- Does \_ own \_ ?
- \_ is \_ part of \_ .
- There are \_ on \_ .
- \_ is to the left of \_ .
- Is \_ to the left of \_ ?
- \_ is to the right of \_ .
- Is \_ to the right of \_ ?

06.11.2002

36

## Pattern Operationen in DYPAR

- Optionales Element ?
- Disjunktion |
- Wertaufnehmende Variable :=
- Wertabgebende Variable =
- Beliebiges Einzelwort \$
- Beliebige Zahl \$ n
- 0-n Wiederholungen \*
- 1-n Wiederholungen +
- n Wiederholungen ≠
- Bis hin zu &u
- Rest der Eingabe \$ r
- Alle Permutationen &c
- konsumierende Negation ~
- nicht konsumierende Negation &n
- Ersetzung von Ausdruck durch Funktionswert &i

06.11.2002

nach Wahlster

37

## DYPAR Patterns

### Beispiele:

(&u show) konsumiert die Eingabe bis zum ersten "show" ("*Could you please show*" fi show)

(&c (Das)(ist)(so)) matcht z.B. "*So ist das*" oder "*Ist das so*".

(!varname := (&i true <positive Antwort>)) fi alle positiven Antworten werden durch true dargestellt.

(!word):= ~<article> bindet alles außer der Negation an word

06.11.2002

38

## BNF-Syntaxdefinition der in DYPAR möglichen Patterns

|                                 |  |
|---------------------------------|--|
| <pattern> ::= nil               | <term> ::= &c <pattern-list>                   |
| <pattern> ::= (<pattern1>)      | <term> ::= &i VALUE <pattern>                  |
| <pattern> ::= <term>            | <term> ::= (%f LISP-FUNCTION (VARIABLE-NAME)*) |
| <pattern> ::= <term> <pattern1> | <term> ::= (= VARIABLE-NAME <pattern>)         |
| <term> ::= ATOMIC CONSTANT      | <term> ::= (= VARIABLE-NAME)                   |
| <term> ::= \$                   | <term> ::= (&m <list>)                         |
| <term> ::= NON-TERMINAL         | <n> ::= 1 2 3                                  |
| <term> ::= FUNCTION             | <list> ::= ATOMIC CONSTANT                     |
| <term> ::= (≠ <n ><pattern>)    | <list> ::= ATOMIC CONSTANT <list>              |
| <term> ::= ? <pattern>          | <pattern-list> ::= (<pattern-list>)            |
| <term> ::= * <pattern>          | <pattern-list> ::= (<pattern>)                 |
| <term> ::= + <pattern>          | <pattern-list> ::= (<pattern> <pattern-list>)  |
| <term> ::= &u <pattern>         |  |
| <term> ::= ! <pattern-list>     |  |

06.11.2002

39

## Vorteile und Nachteile von Pattern - Matchern

- |   |   |
|---|---|
| <ul style="list-style-type: none"> <li>● Vorteile:</li> <li>● schnell (für zeitkritische Anwendungen geeignet)</li> <li>● einfach (intuitiv formulierbar, von BS oder anderen Anwendungen bekannt)</li> <li>● bei kleinem Umfang gut wartbar</li> <li>● lokale Effekte</li> </ul> | <ul style="list-style-type: none"> <li>● Nachteile:</li> <li>● liefern keine Strukturbeschreibung</li> <li>● ungeeignet für einige syntaktische Phänomene (offene Wortfolge, Abhängigkeiten, ...)</li> <li>● umständlich für gleiche Strukturen an verschiedenen Stellen</li> <li>● bei größerem Umfang unübersichtlich und schwer wartbar</li> </ul> |
|---|---|

06.11.2002

40

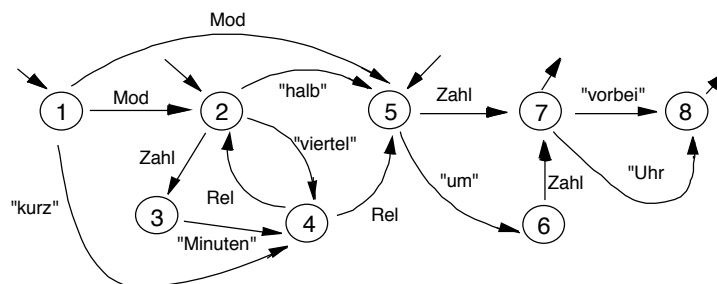
## Netzwerkgrammatiken

- Einfache Netzwerkgrammatiken (Endliche Automaten)  
Übergänge in **einem** Netz zwischen
  - Wörtern oder
  - lexikalischen Kategorien
- Rekursive Netzwerkgrammatiken (BTN / RTN)  
Übergänge in gestaffelten Netzen zwischen
  - Netzen
  - Wörtern oder
  - lexikalischen Kategorien
- Erweiterte Netzwerkgrammatiken (ATN)  
Übergänge in gestaffelten Netzen zwischen
  - Netzen
  - Wörtern oder
  - lexikalischen Kategorienunter :
  - Ausführung von Aktionen und
  - Aufbau einer Strukturbeschreibung

06.11.2002

41

## Endlicher Automat (BTN) für deutsche Zeitangaben



- Mod: eben, bald, gerade, gleich, genau
- Rel: vor nach
- Zahl: 1, 2, 3, ...59, 60.

06.11.2002

42

## CFG für deutsche Zeitangaben

S → Vor Zeit

Zeit → Zahl (Uhr)

Vor → (Abtön) (Mod)

Mod → Frac | Int Rel | Mod

Abtön → (Hecke) Mod

Frac: viertel, halb

Rel: vor, nach

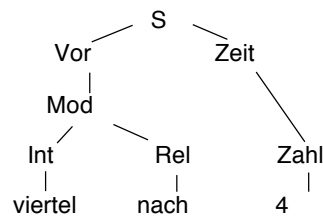
Int: viertel, kurz

Nach: vorbei

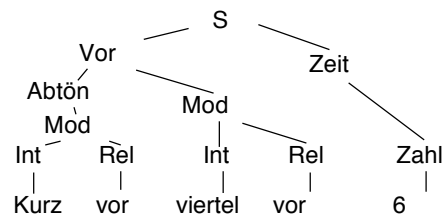
Uhr: Uhr

Zahl: 1,2,3,4,...

Hecke: fast, bald, gleich



06.11.2002



43

## Cocke / Kasami / Younger - Algorithmus - 1 -

Parser, die auf dem **Cocke**-Algorithmus beruhen, bearbeiten nur CFGs in Chomsky-Normalform und arbeiten

von links nach rechts

bottom-up

breitenorientiert (parallele Verfolgung aller Analysepfade)

in einem einzigen Durchgang (one-pass)

ohne Zurücksetzen (backtracking)

vor allem: binäre Regeln

Es wird für alle erreichten X mit allen zurückliegenden Y überprüft

Zunächst werden die einzelnen Wörter der Eingabekette von links nach rechts durchnummeriert.

nach Wahlster

44

06.11.2002

## Cocke / Kasami / Younger - Algorithmus - 2 -

Es wird eine Parsingtabelle (auch: Wegetabelle) aufgebaut:

- jedem Reduktionsschritt wird eindeutig eine Nummer zugeordnet,
- Das durch den Reduktionsschritt erzeugte nichtterminale Symbol und die angewandte Regel werden vermerkt,
- Mit LG (linke Grenze) und RG (rechte Grenze) wird relativ zur Nummerierung der Eingabekette jeweils der Teilstring der Eingabekette identifiziert, der von dem Reduktionsschritt erfaßt wird (nur bei lexikalischen Regeln ist  $RG = LG$ ),
- Mit LK (linke Komponente) und RK (rechte Komponente) wird spezifiziert, aus welchen Komponenten eine neue syntaktische Kategorie gewonnen wurde, indem die entsprechenden Nummern der Reduktionsschritte angegeben werden.

Das ist nötig, weil höhere Nonterminale tiefere Nonterminale zusammenfassen können, die keine direkte Verbindung zu Terminalen haben können

06.11.2002

## Cocke Algorithmus - 3 -

(1) Lies nächstes Wort  $w$  ein und notiere in der Parsingtabelle

- seine lexikalische Kategorie      - seine Position  $z$  ( $RG = LG$ )
- setze  $RK = LK = 0$

(2) Ab dem zweiten Eintrag der Parsingtabelle:

Solange es für die aktuelle Zeile in der Parsingtabelle mit der Kategorie  $X$  oder vorangehende Zeilen mit der Kategorie  $Y$  und  $RG(X) = Z$  solche Kategorien  $Y$  gibt, für die gilt

- $RG(Y) = LG(X) - 1$  (d.h.  $Y$  steht unmittelbar links von  $X$ )
- es gibt eine Regel  $A \rightarrow YX$

ergänze als neue aktuelle Zeile in der Parsingtabelle

- die Kategorie  $A$
- $LG(A) = LG(Y)$ ,  $RG(A) = RG(X)$
- $LK(A) = \langle \text{Nummer von } Y \text{ in Parsing-Tabelle} \rangle$
- $RK(A) = \langle \text{Nummer von } X \text{ in Parsing-Tabelle} \rangle$

Adjazenzregel

(3) Falls die neu eingetragene syntaktische Kategorie  $S$  ist und

$LG = 1$  sowie  $RG = \langle \text{Länge der Eingabekette} \rangle$

dann: Stop                      sonst: gehe nach (1)

06.11.2002

nach Wahlster 46

## Cocke-Algorithmus - Schritt 1 -

Beispiel nach Wahlster

Det eintragen mit :  
LG = I RG = I  
LK = 0 RK = 0

(L1) "Det  $\emptyset$  die, manche, einige"

1  
▼  
DET

*Die*      *Mädchen*      *sahen*      *die*      *Kinder*      *dort*  
I            II            III            IV            V            VI

06.11.2002

47

## Cocke-Algorithmus - Schritt 2 -

N eintragen mit :  
LG = II RG = II  
LK = 0 RK = 0

(L3) "N  $\emptyset$  Mädchen, Kinder"

1    2  
DET    N

*Die*      *Mädchen*      *sahen*      *die*      *Kinder*      *dort*  
I            II            III            IV            V            VI

06.11.2002

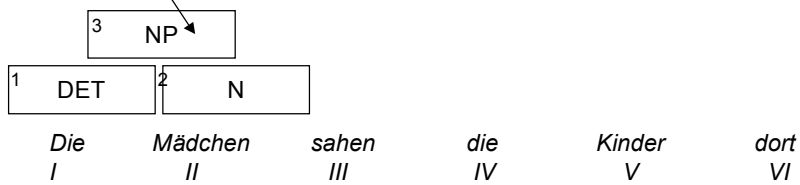
48



### Cocke-Algorithmus - Schritt 3 -

NP eintragen mit :  
 LG = I RG = II  
 LK = 1 RK = 2

(G10) "NP  $\emptyset$  DET N"



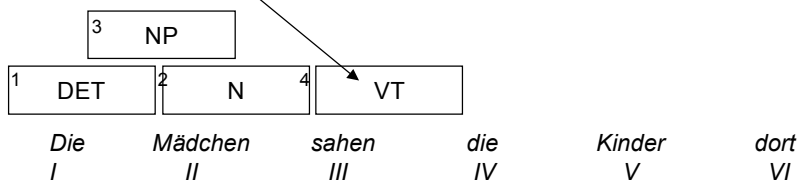
06.11.2002

49

### Cocke-Algorithmus - Schritt 4 -

VT eintragen mit :  
 LG = III RG = III  
 LK = 0 RK = 0

(L5) "VT  $\emptyset$  sahen, ehrten"

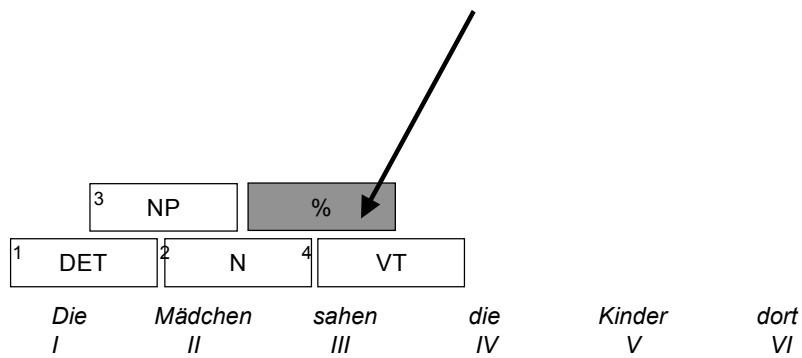


06.11.2002

50

### Cocke-Algorithmus - Schritt 4a -

Weiterer Schritt mit N und VT ist nicht möglich:  
 4a Kat = % LG = II RG = III LK = 2 RK = 4 R = %

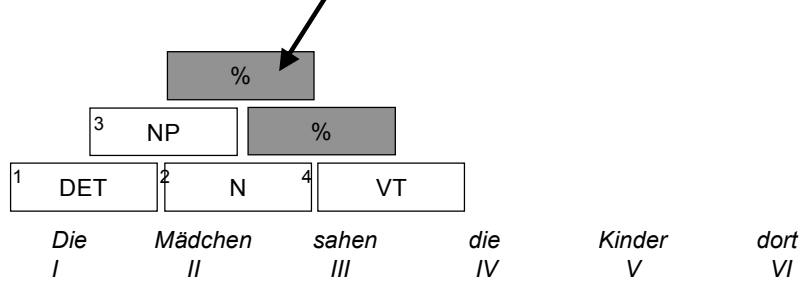


06.11.2002

51

### Cocke-Algorithmus - Schritt 4b -

Weiterer Schritt mit NP und VT nicht möglich:  
 4b Kat = % LG = I RG = III LK = 3 RK = 4 R = %



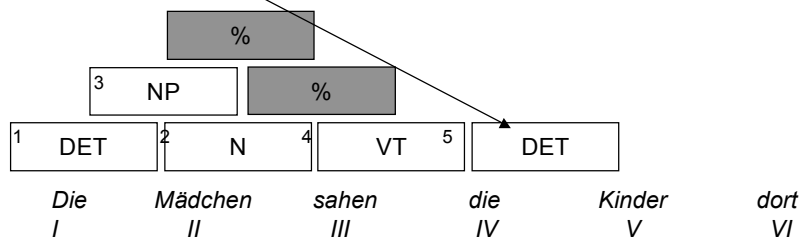
06.11.2002

52

## Cocke-Algorithmus - Schritt 5 -

DET eintragen mit :  
 LG = IV RG = IV  
 LK = 0 RK = 0

(L1) "Det Ø die, manche, einige"



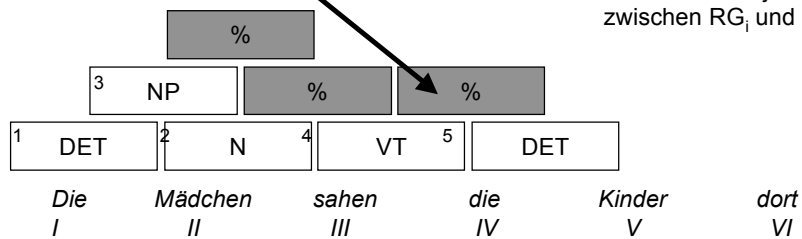
06.11.2002

53

## Cocke-Algorithmus - Schritt 5a -

Weiterer Schritt mit VT und DET nicht möglich:  
 5a Kat = % LG = IV RG = V LK = 0 RK = 0 R = %

DET mit N oder DET oder NP  
 wird nicht geprüft,  
 da verletzte Adjazenzregel  
 zwischen  $RG_1$  und  $LG_5$



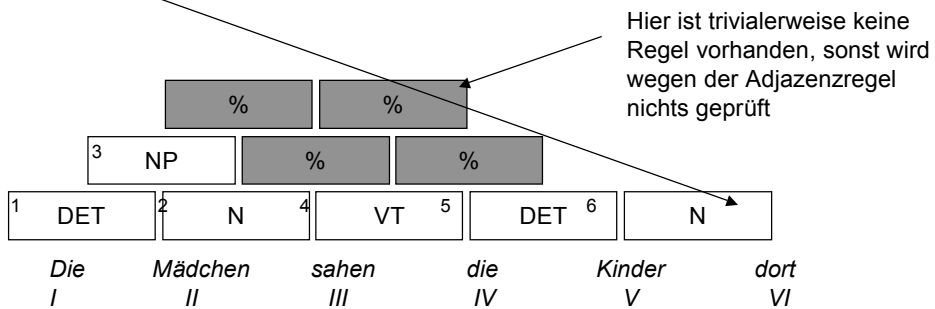
06.11.2002

54

## Cocke-Algorithmus - Schritt 6 -

N eintragen mit :  
 LG = V RG = V  
 LK = 0 RK = 0

(L3) "N  $\emptyset$  Mädchen, Kinder"



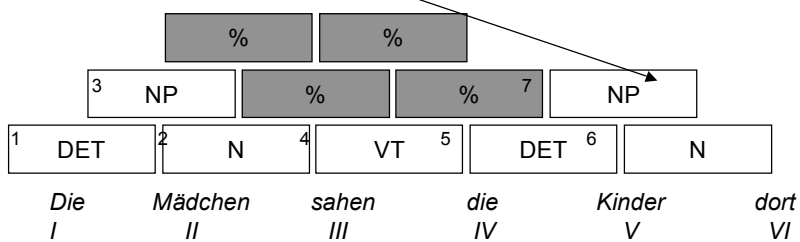
06.11.2002

55

## Cocke-Algorithmus - Schritt 7 -

NP eintragen mit :  
 LG = IV RG = V  
 LK = 5 RK = 6

(G10) "NP  $\emptyset$  DET N"



06.11.2002

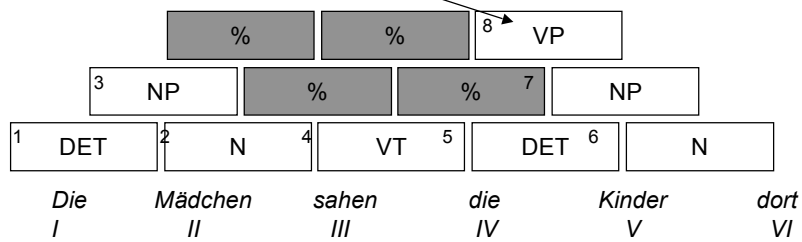
56

## Cocke-Algorithmus - Schritt 8 -

VP eintragen mit :  
LG = III RG = V  
LK = 4 RK = 7

(G9) "VP  $\emptyset$  VT NP"

Weitere Schritte verstoßen gegen die Adjazenzregel

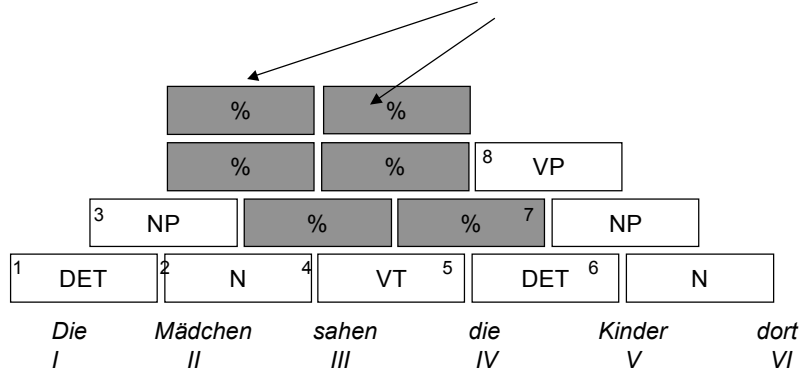


06.11.2002

57

## Cocke-Algorithmus - Schritt 8a und 8b -

VP mit N und DET entsprechen keinen Regeln



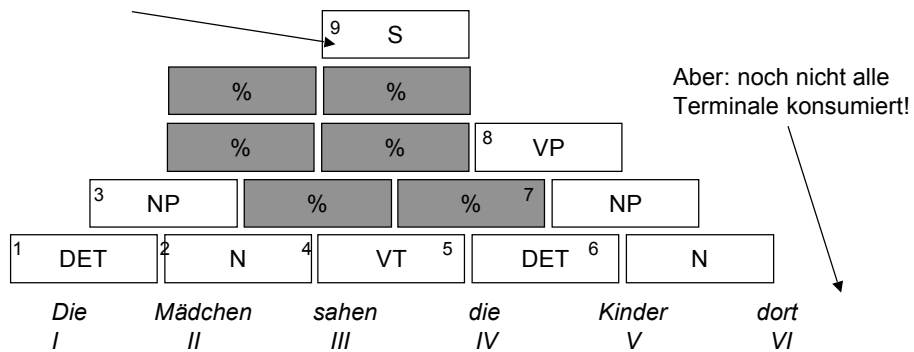
06.11.2002

58

## Cocke-Algorithmus - Schritt 9 -

S eintragen mit :  
 LG = I RG = V  
 LK = 3 RK = 8

(G3) "S  $\emptyset$  NP VP"



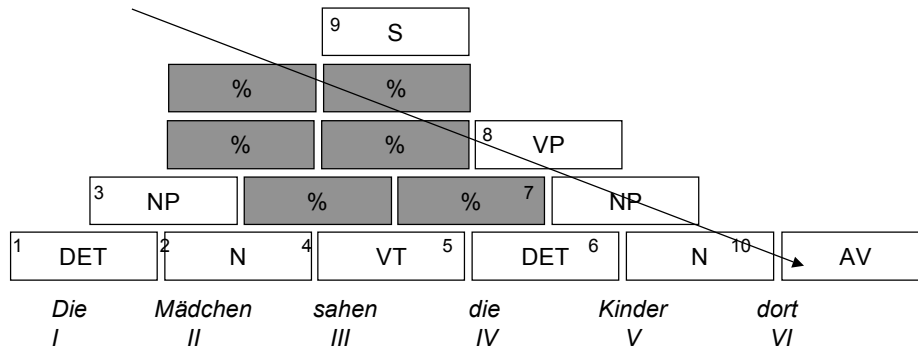
06.11.2002

59

## Cocke-Algorithmus - Schritt 10 -

AV eintragen mit :  
 LG = VI RG = VI  
 LK = 0 RK = 0

(L6) "AV  $\emptyset$  dort, hier"



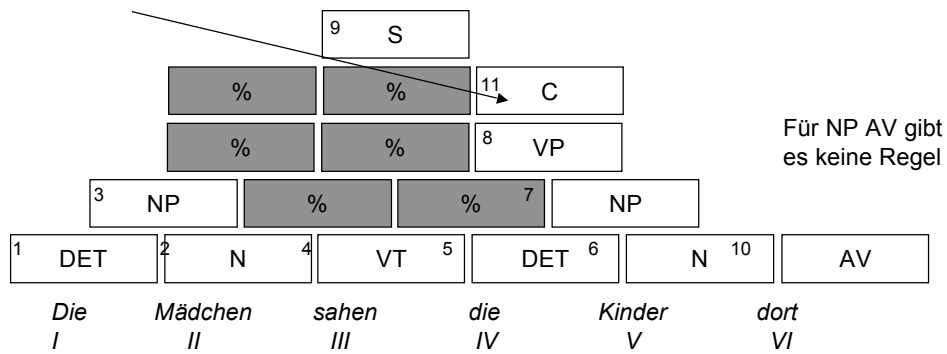
06.11.2002

60

## Cocke-Algorithmus - Schritt 11-

C eintragen mit :  
 LG = III RG = VI  
 LK = 8 RK = 10

(G8) "C Ø VP AV"



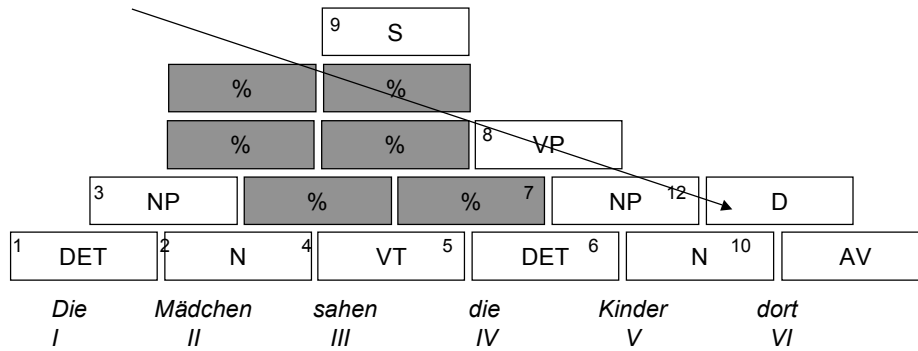
06.11.2002

61

## Cocke-Algorithmus - Schritt 12-

D eintragen mit :  
 LG = VI RG = VI  
 LK = 0 RK = 0

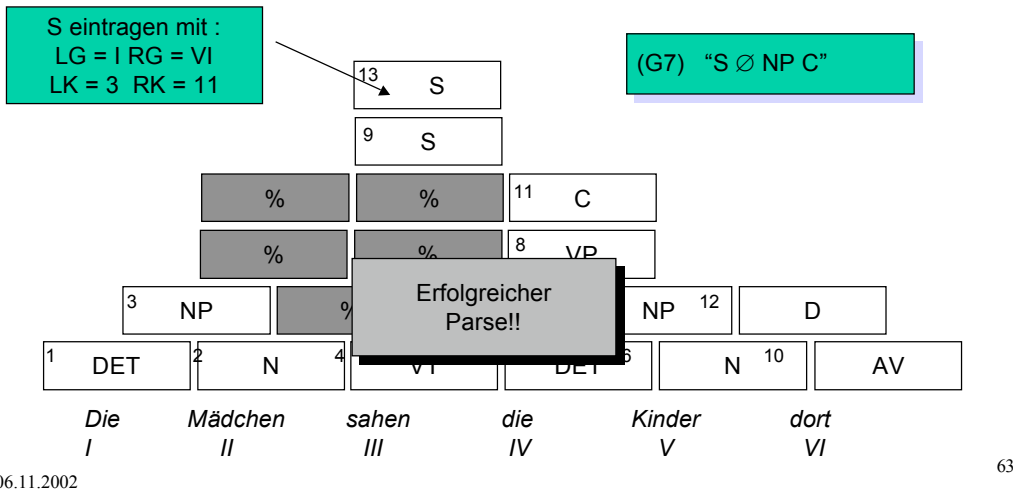
(G12) "D Ø N AV"



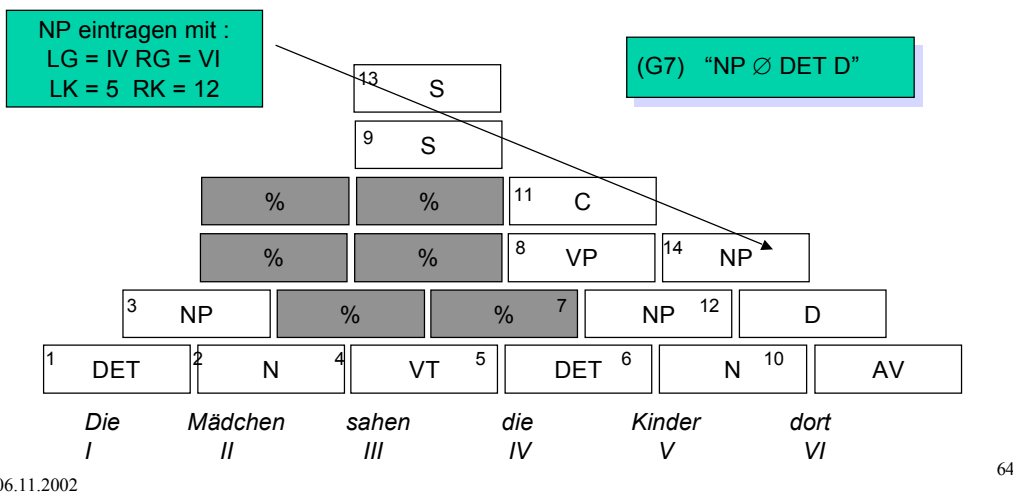
06.11.2002

62

### Cocke-Algorithmus - Schritt 13 -

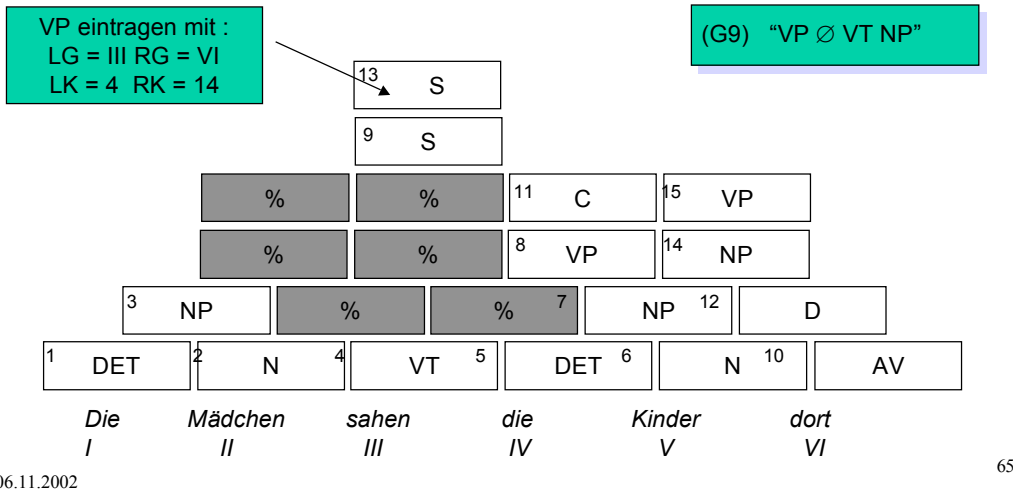


### Cocke-Algorithmus - Schritt 14 -

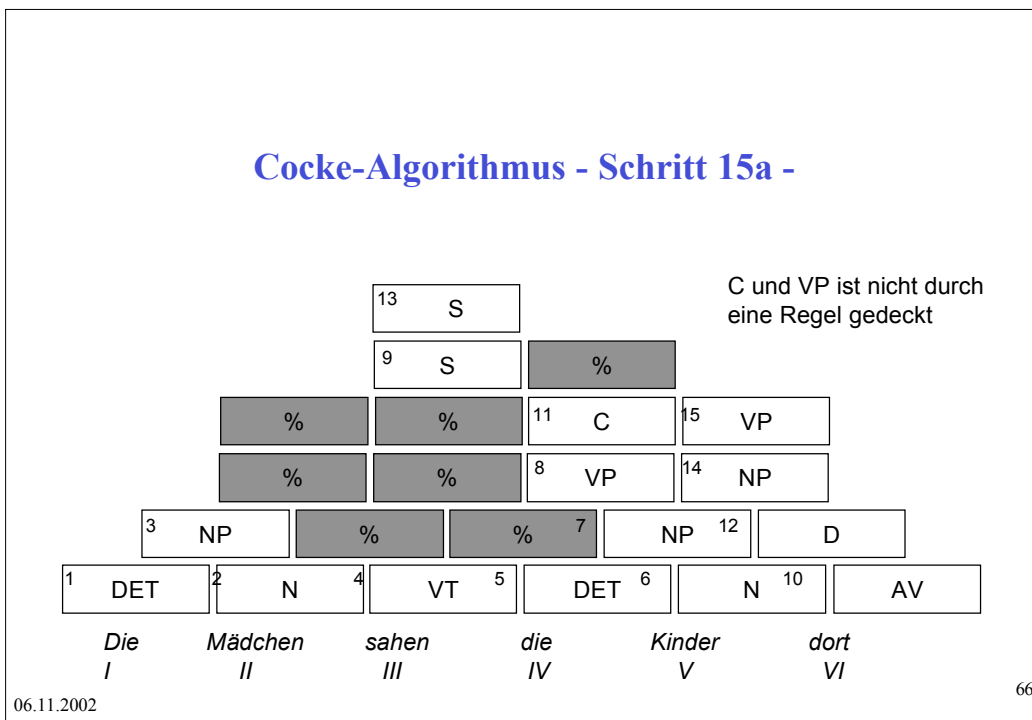




## Cocke-Algorithmus - Schritt 15 -



## Cocke-Algorithmus - Schritt 15a -





## Parsing Konstituentenstrukturgrammatik

- **top-down** : begins with the start symbol and searches through different paths to rewrite the symbols until the input sentence is generated or until all possibilities have been explored
- **bottom-up** : starts with the words in the sentence and uses the rewrite rules backward to reduce the sequence of symbols until it consists solely of S.

06.11.2002

69

## Chart-Parsing - 1 -

Von Martin Kay 1973 vorgeschlagen. Von Kaplan und Kay als General Syntactic Processor implementiert.

Eine Chart ist ein gerichteter bewerteter Graph, der im Verlauf des Parsing-Prozesses um immer weitere Kanten erweitert wird.

Zunächst wird die Chart initialisiert:

- Je einen Knoten an den Satzgrenzen und zwischen je zwei Wörtern der Eingabekette einsetzen, sowie
- Kanten einsetzen, die zwei Knoten links und rechts von einem terminalen Symbol verbinden und mit dem jeweils überspannten Wort markieren
- Kanten einsetzen, die zwei Knoten links und rechts von einem terminalen Symbol verbinden und mit der Wortklassenbezeichnung des durch sie überspannten Worts markieren. Können einem Wort mehrere alternative Kategorien zugeordnet werden, je eigene Kanten aufspannen.

06.11.2002

Nach Wahlster

70

## Chart-Parsing - 2 -

Die Aufgabe des Parsers besteht darin, für jede im Verlauf der Analyse der Eingabe gefundene Konstituente eine mindestens zwei Knoten überspannende Kante in die Chart einzuführen.

Die Analyse gilt als abgeschlossen, wenn eine die ganze Chart überspannende Kante eingerichtet ist.

Die Chart enthält dann alle möglichen syntaktischen Zerlegungen der Eingabe.

Bei unvollständigen Eingaben wird zwar keine den ganzen Satz überspannende Kante gefunden, aber die Chart enthält meist einige der Grammatik entsprechenden Konstituenten als Fragmente, die dann z.B. von einer Komponente zur Rekonstruktion von Ellipsen weiterverarbeitet werden können.

06.11.2002

71

## Chart-Parsing - 3 -

Man unterscheidet zwei Typen von Kanten in einem Chart-Parser:

- Aktive Kanten zur Darstellung noch unvollständiger Konstitute
- Inaktive Kanten zur Darstellung bereits vollständiger Konstitute

Wenn das Ende einer aktiven Kante A und der Anfang einer inaktiven Kante I an einem Knoten zusammenstoßen, wird folgendes Verfahren zur Generierung einer neuen Kante N ausgelöst:

- der Anfang von N ist der von A,
- das Ende von N ist das von I,
- die Kategorie von N ist die von A,
- der Inhalt von N hängt vom Inhalt von A und der Kategorie sowie dem Inhalt von I ab,
- N ist entweder Aktiv oder inaktiv, je nachdem, ob A damit vervollständigt wird oder nicht.

Nach Wahlster

06.11.2002

72

## Chart-Parsing - Beispiel -

|                                 |                              |
|---------------------------------|------------------------------|
| G1: S $\square$ NP VP           | G6: NP2 $\square$ NP2 PP     |
| G2: NP $\square$ Determiner NP2 | G7: PP $\square$ Preposition |
| G3: NP $\square$ NP2            | G8: VP $\square$ Verb        |
| G4: NP2 $\square$ Noun          | G9: VP $\square$ Verb NP     |
| G5: NP2 $\square$ Adjective NP2 | G10: VP $\square$ VP PP      |

"The rabbit with a saw nibbled on an orange"

Schritt: 1  
Eintrag:  $1S \emptyset 1NP VP$  (aus Initialisierung)  
Offene Kanten:  $1NP \emptyset 1Det NP2, 1NP \emptyset NP2$   
Komplette Kanten:  $1Det_2, \dots$

Schritt: 2  
Eintrag:  $1NP \emptyset 1Det NP2$  (aus Schritt 1)  
Neue offene Kante:  $1NP \emptyset Det 2NP2$   
Komplette Kanten:  $1Det_2, \dots$

06.11.2002

nach Winograd

73

## Chart-Parsing - Beispiel - 2-

Schritt: 3  
Eintrag:  $1NP \emptyset 1NP2$  (aus Schritt 1)  
Neue offene Kante:  $1NP2 \emptyset 1Noun, 1NP2 \emptyset 1ADJ NP2,$   
 $1NP2 \emptyset 1NP2 PP$   
Komplette Kanten:  $1Det_2, \dots$

Schritt: 4  
Eintrag:  $Det 2NP2$  (aus Schritt 2)  
Neue offene Kante:  $2NP2 \emptyset 2Noun, 2NP2 \emptyset 2ADJ NP2,$   
 $2NP2 \emptyset 2NP2 PP$   
Komplette Kanten:  $1Det_2, \dots$

Schritt: 5  
Eintrag:  $1NP2 \emptyset 1Noun$  (aus Schritt 3)  
Neue offene Kante: keine, da der remainder weder ein Nonterminal ist, noch mit der Eingabe matcht.  
Komplette Kanten:  $1Det_2, \dots$

06.11.2002

74

## Chart-Parsing - Beispiel - 3 -

Schritt: 6

Eintrag:  $1NP2 \emptyset 1Adj NP2$  (aus Schritt 3)

Neue offene Kante: keine, da der remainder weder ein Nonterminal ist, noch mit der Eingabe matcht.

Komplette Kanten:  $1Det_2, \dots$

Schritt:7

Eintrag:  $1NP2 \emptyset 1NP2 PP$  (aus Schritt 3)

Neue offene Kante: keine, da bereits ein  $1NP2$  in der Chart ist. Zur Blockierung von Links-Rekursion

Komplette Kanten:  $1Det_2, \dots$

Schritt: 8

Eintrag:  $2NP2 \emptyset 2Noun$  (aus Schritt 4)

Neue offene Kante:  $2NP2 \emptyset Nouns$  (nicht in der Graphik)

Komplette Kanten:  $1Det_2, \dots$

06.11.2002

75

## Chart-Parsing - Beispiel - 4 -

Schritt: 9

Eintrag:  $2NP2 \emptyset 2Adj NP2$  (aus Schritt 4)

Neue offene Kante: keine, da der remainder weder ein Nonterminal ist, noch mit der Eingabe matcht.

Komplette Kanten:  $1Det_2, \dots$

Schritt: 10

Eintrag:  $2NP2 \emptyset 2NP2 PP$  (aus Schritt 4)

Neue offene Kante: keine, da bereits ein  $1NP2$  in der Chart ist. Zur Blockierung von Links-Rekursion

Komplette Kanten:  $1Det_2, \dots$

06.11.2002

76

## Chart-Parsing - Beispiel - 5 -

Schritt: 11

Eintrag:  ${}_2\text{NP}_2 \ \emptyset \ \text{Noun}_3$  (aus Schritt 8)

Neue offene Kante:  ${}_2\text{NP}_2 \ \emptyset \ \text{NP}_2 \ \text{PP}$ ,  ${}_1\text{NP} \ \emptyset \ \text{Det} \ \text{NP}_2$

Komplette Kanten:  ${}_1\text{Det}_2$ ,  ${}_2\text{NP}_3$ , ...

Schritt: 12

Eintrag:  ${}_2\text{NP}_2 \ \emptyset \ \text{NP}_2 \ \text{PP}$  (aus Schritt 11)

Neue offene Kante:  ${}_3\text{PP} \ \emptyset \ \text{Prep} \ \text{NP}$

Komplette Kanten:  ${}_1\text{Det}_2$ ,  ${}_2\text{NP}_3$ , ...

Schritt: 13

Eintrag:  ${}_1\text{NP} \ \emptyset \ \text{Det} \ \text{NP}_2$  (aus Schritt 11)

Neue offene Kante:  ${}_1\text{S} \ \emptyset \ \text{NP}_3 \ \text{VP}$

Komplette Kanten:  ${}_1\text{NP}_3$ ,  ${}_3\text{Prep}_4$ , ...

!

06.11.2002

77

## Chart-Parsing - Beispiel - 6 -

Schritt: 14

Eintrag:  ${}_3\text{PP} \ \emptyset \ \text{Prep} \ \text{NP}$  (aus Schritt 12)

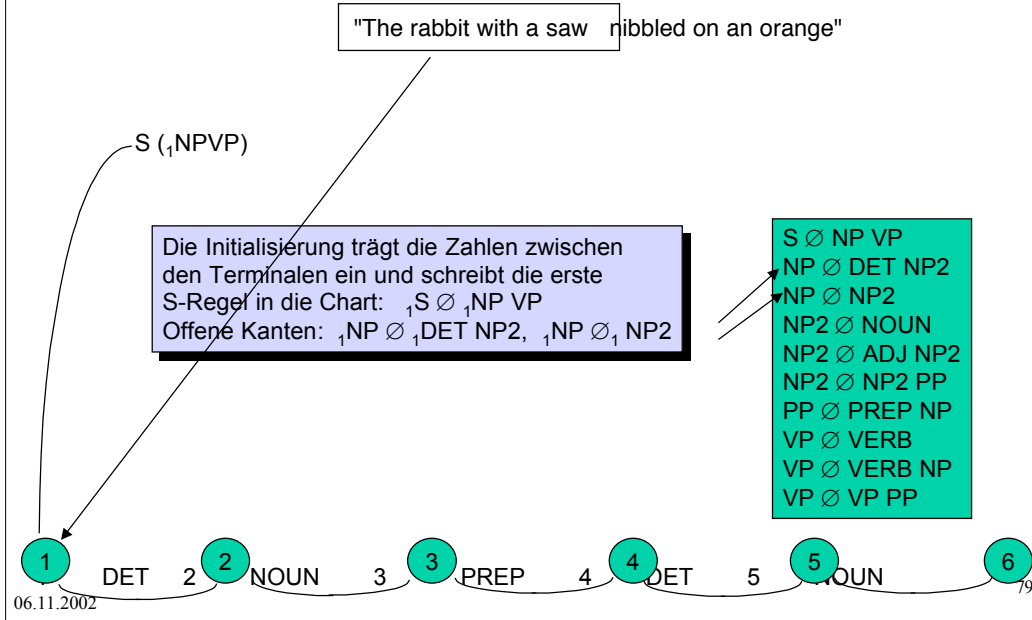
Neue offene Kante:  ${}_3\text{PP} \ \emptyset \ \text{Prep} \ \text{NP}$

Komplette Kanten:  ${}_1\text{NP}_3$ ,  ${}_3\text{Prep}_4$ , ...

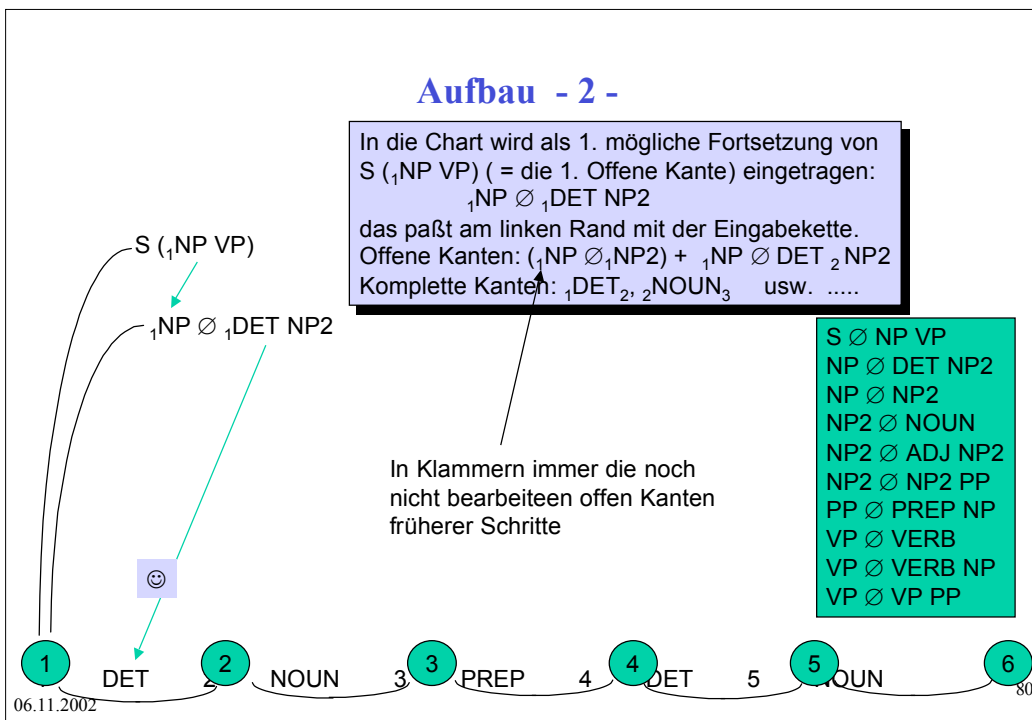
06.11.2002

78

## Aufbau der Chart - Beispiel -

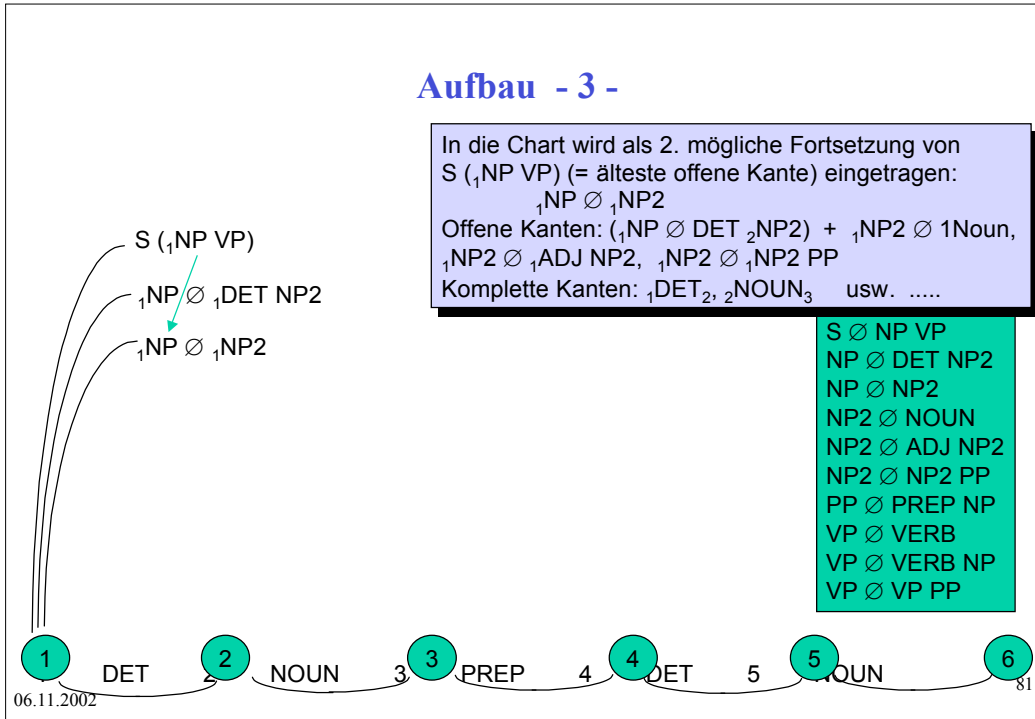


## Aufbau - 2 -

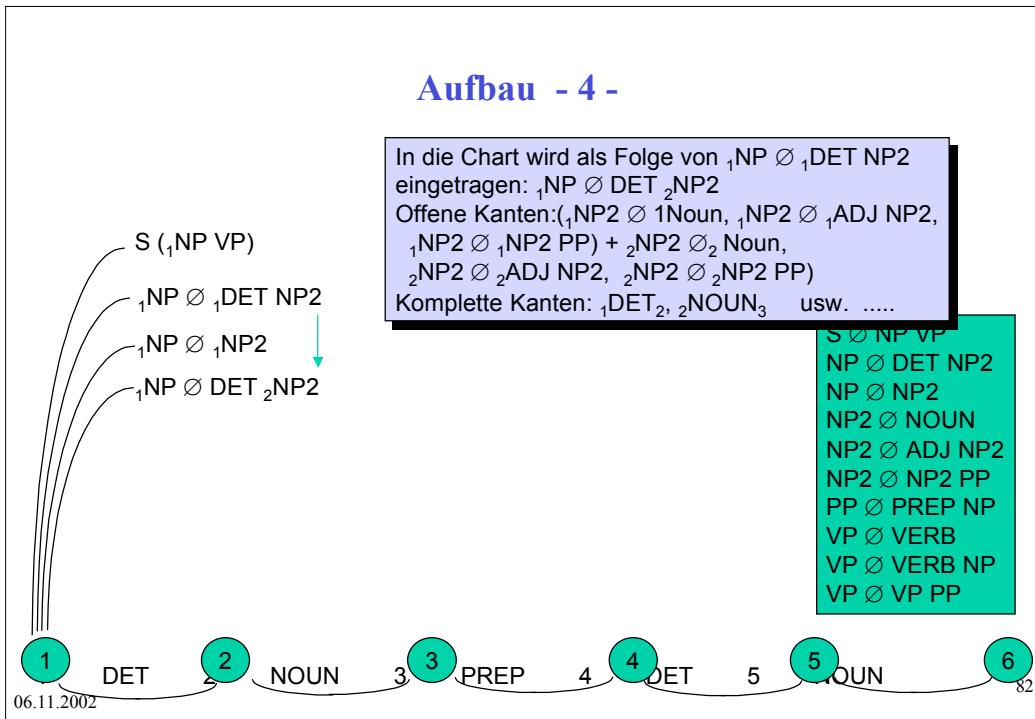




### Aufbau - 3 -



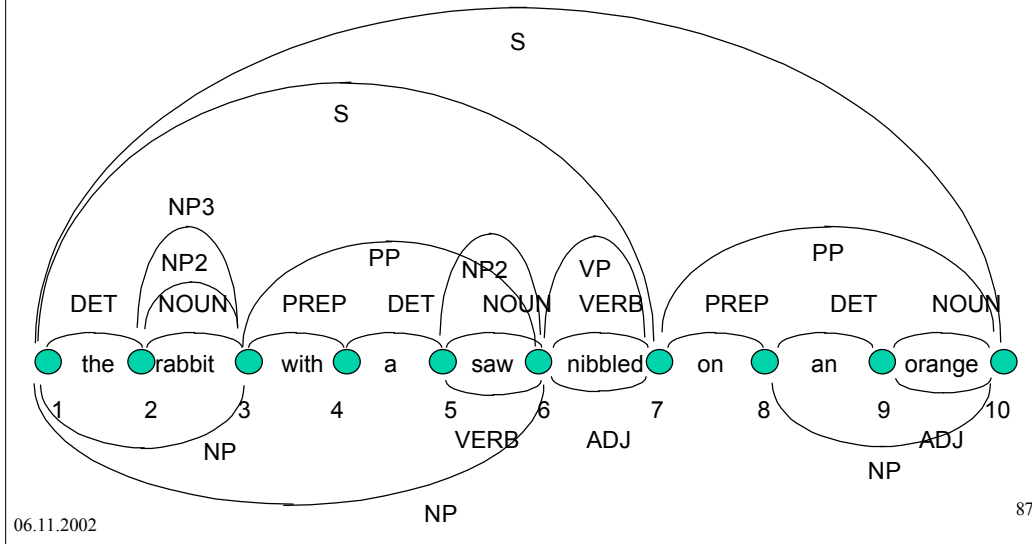
### Aufbau - 4 -







## Komplette Chart



## Chart-parsing vs. Traditionelle Methoden

- the same constituent is never used more than once
- in the worst case would build every possible constituent between every possible pair of positions
- more work in each step

*Example:* a sentence with 12 words :

brute force  $10^{12}$  operations

(1,000,000,000,000)

chart parser  $1000 \cdot 12^3$

(1,728,000)

on some examples, 500 000  
times faster

06.11.2002

88

## Parsing - problems

- partially incorrect phrases
  - Ex: An der Kupplung ist ein Kabelzug, **der ist beschädigt**.  
(looks like a relative clause, but wrong word order)
- ambiguities
  - Ex: Ich kann das rote Kabel am Vergaser nicht abziehen
- German separable verbs
  - Ex : Der Stecker, der an der Lampe ist, **hängt** lose **herunter**

06.11.2002

89

## Gliederung

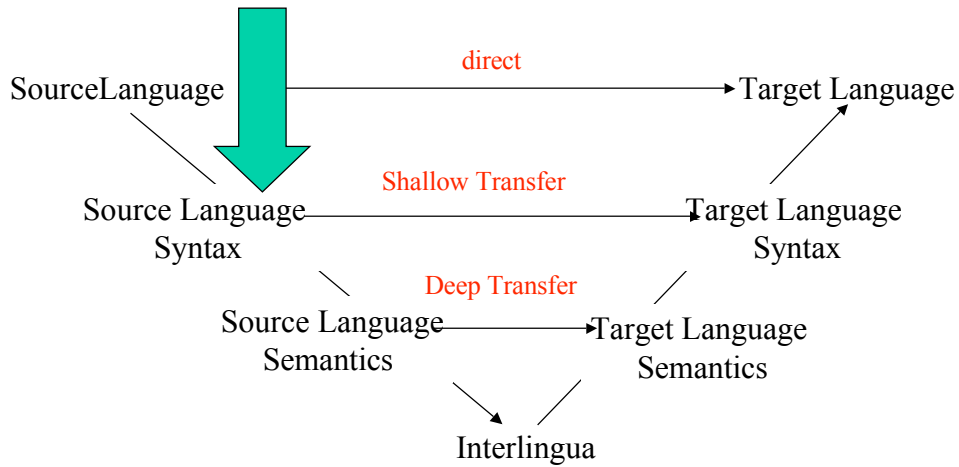
- Syntax
- Parsing
- Syntax und Parsing in MÜ Systeme



06.11.2002

90

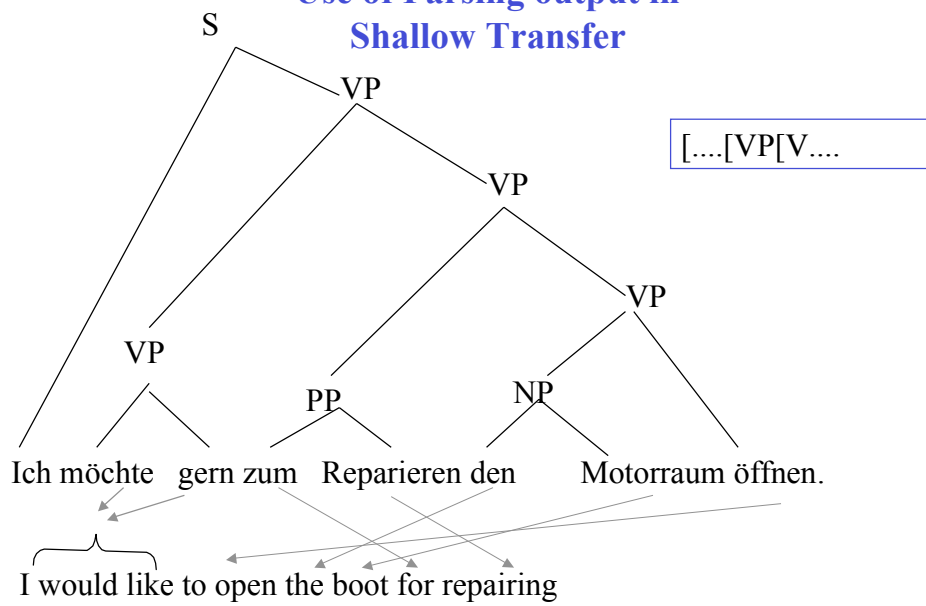
## Parsing for Machine Translation



06.11.2002

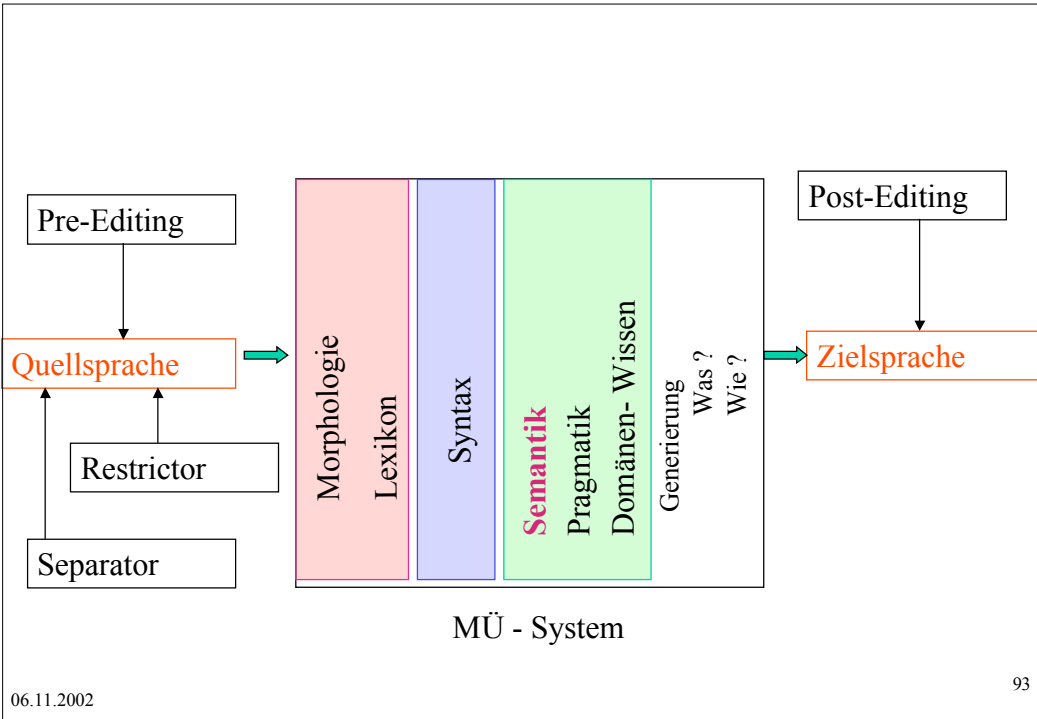
91

## Use of Parsing output in Shallow Transfer



06.11.2002

92



06.11.2002

93