

Syntaktische Analyse

- Satzbildung durch Wortkomposition
- Ausgangspunkt für Bedeutungsanalyse
 - (1) *Der Hund fraß den Knochen*
 - (2) *Der Knochen wurde vom Hund gefressen*
- Struktur aus Syntaxregeln hilft der Bedeutungsfindung im Gegensatz zu:
 - Es ist immer das 2. Substantiv, das gefressen wird.
- Zerlegung der Bedeutungsfindung auf Basis von syntaktischen Zerlegungen (etwa von Phrasen)
 - (3) *sp[Der Hase mit den langen Ohren] erfreute sich an sp[einem großen grünen Salatblatt]*
- Bedeutungsfindung durch Komposition der Teilbedeutungen
- Komposition abhängig von gewählter Clustering
 - (4) *Tim sah Maria mit dem Fernglas*
 - (5) *Ich sah den Kölner Dom auf dem Flug nach Frankfurt*

51

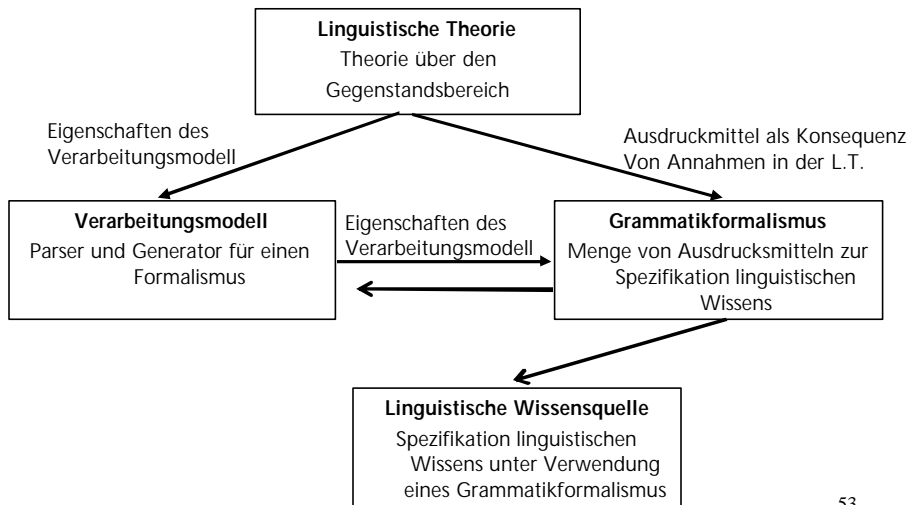
Beschreibungsformalismus

Definition: Ein Grammatikformalismus ist ein Formalismus zur Repräsentation sprachlichen Wissens auf den verschiedenen linguistischen Beschreibungsebenen.

- Kriterien zur Ausgestaltung eines Grammatikformalismus unter Berücksichtigung der Beziehung von Grammatikformalismus zu linguistischer Theorie:
 - **Ausdrucksstärke:** Welche Phänomenbereiche können beschrieben werden?
 - **Algorithmisierbarkeit:** Welche Verarbeitungsmodelle werden unterstützt? Welche Algorithmisierbarkeit ist möglich?
 - **Linguistische Motiviertheit:** Umfang direkt zu verwendender linguistisch abgesicherter Ausdrucksmittel? Inwieweit sind Ausdrucksmittel direkt an spezifische linguistische Theorie gebunden?

52

Beschreibungsformalismus



53

Grammatik

- Grammatik als formales Beschreibungssystem für Syntaxregeln zur Definition der Wortorganisation
- Menge von Ersetzungsregeln für eine Kette von Symbolen (=Satz)
- Symbole eingeordnet in syntaktische Kategorien
 - Nomen, Verb, Adjektiv, ...
 - Nominalphrase, Verbalphrase
- Satz ist regulär, wenn ausgehend von einem Startsymbol der Satz mit Hilfe der Regeln erzeugt werden kann (Parsen)
- Was soll die Grammatik tun?
 - Unterscheidung zwischen „korrekten“ und „inkorrekten“ Sätzen
 - Der Junge wirft den Ball. (+)
 - Der wirft Junge Ball den. (-)
 - Zuweisung einer „sinnvollen“ Struktur
 - (Der Junge) (wirft den Ball)
 - (Der) (Junge wirft) (den Ball)
- Eigenschaften der Grammatik: Kompakt, modular
- Syntaxanalytiker (Parser) überprüft Satz auf Grammatikkonformität

54

Grammatik

- Linguistische Strukturen und Prozesse
 - Kompositionalität
 - Interpretation einer Struktur lässt sich aus Teilen und ihrer Zusammensetzung ableiten.
 - Interpretation jedes Teils kann lokal erfolgen.
 - Bedingt analytische (paradigmatische Betrachtung).
 - Deklarativität
 - Trennung zwischen Wissen und Verarbeitung.
 - Verarbeitung kann nach rein strukturellen Gesichtspunkten erfolgen
 - Unterstützung von analytischen und generativen Prozessen.
- Konkatenation
 - Lineare Aneinanderreihung einzelner Elemente.
 - Häufigste Operation über linguistischen Daten.
 - Bei natürlichen Sprachen finden sich auch Tilgung und Ersetzung.

Chomskys Hierarchie der Grammatiken

Regeln zum Ausdruck der korrekten grammatikalischen Struktur einer Sprache für Analyse und Synthese

Format der hierarchischen Chomsky Grammatiken

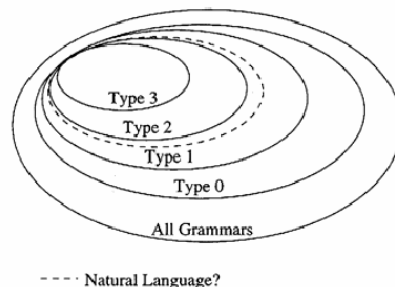
$A B C \dots \rightarrow E F G \dots$

Typ 0: Transformation Grammar
 Regeln der Form: *beliebig* \rightarrow *beliebig*
 Mächtigkeit: allgemeine Turing Maschine

Typ 1: kontextsensitive Grammatik
 Regeln der Form: $A B C \rightarrow A D C$
 Mächtigkeit: Linear Bound Automata

Typ 2: kontextfreie Grammatik
 Regeln der Form: $A \rightarrow B C D \dots$
 Mächtigkeit: Push Down Stack Automata

Typ 3: Regulär or Rechtslinear
 Regeln der Form: $A \rightarrow x B$ und $A \rightarrow x$ wobei x terminal
 Mächtigkeit: Finite State Automata



Grammatikformalismen

Phrasenstrukturgrammatiken Beispiel:

(7) NP \rightarrow Det Adj N

- (7) beschreibt einfache Sätze wie
 - „der alte Mann“
 - „ein junger Mann“
 - „dem großen Haus“
- Menge von Regeln wie (7) bilden eine Phrasenstrukturgrammatik (PSG).
- PSG sind kontextfreie Grammatiken.

PSG:

| | |
|--------------|--|
| Grammatik | $G = (\Sigma, \Phi, P, S)$ |
| Terminale | $\Sigma = \{alte, dem, der, ein, großen, Haus, junger, Mann, Vater, \dots\}$ |
| Nonterminale | $\Phi = \{Adj, Det, N, NP, \dots\}$ |
| Regeln | $P = \left\{ \begin{array}{l} S \rightarrow NP VP \mid \dots, NP \rightarrow Det N \mid Det Adj N \mid \dots \\ Det \rightarrow dem \mid der \mid ein \mid \dots, Adj \rightarrow alte \mid junger \mid \dots, \\ N \rightarrow Haus \mid Mann \mid Vater \mid \dots \end{array} \right\}$ |
| Startsymbol | S |

57

Beispiel für eine einfache deutsche Grammatik

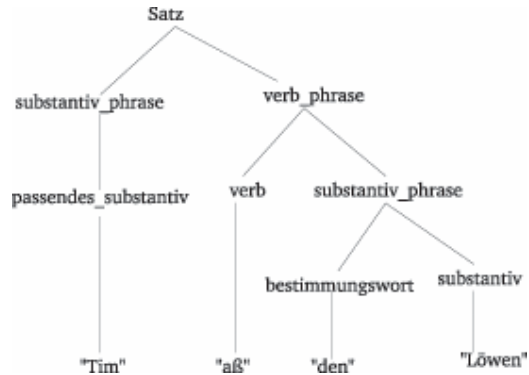
Grammatik in DCG (definite clause grammar oder definiten Klauselgrammatik):

| | | |
|----------------------|-----|--------------------------------|
| satz | --> | substantiv_phrase, verb_phrase |
| substantiv_phrase | --> | passendes_substantiv |
| substantiv_phrase | --> | bestimmungswort, substantiv |
| verb_phrase | --> | verb, substantiv_phrase |
| passendes_substantiv | --> | [Maria] |
| passendes_substantiv | --> | [Tim] |
| substantiv | --> | [Löwen] |
| substantiv | --> | [Keks] |
| verb | --> | [aß] |
| verb | --> | [küsste] |
| bestimmungswort | --> | [den] |

58

Ableitungsbaum

Ableitungsbaum für „Tim aß den Löwen“



Grammatikformalismen

Eine einfache Grammatik mit zahlenmäßiger Übereinstimmung zwischen Subjekt und Verb:

```

satz          --> substantiv_phrase(Num),
                verb_phrase(Num)
substantiv_phrase(Num) --> passendes_substantiv(Num)
substantiv_phrase(Num) --> bestimmungswort (Num),
                substantiv(Num)
verb_phrase(Num)    --> verb(Num), substantiv_phrase ( )
passendes_substantiv(s) --> [maria]
substantiv(s)      --> [löwe]
substantiv(p)      --> [löwen]
best(s)            --> [den]
best(p)            --> [die]
verb(s)            --> [isst]
verb(p)            --> [essen]
  
```

Grammatikformalismen

- Natürliche Sprache hat großes Vokabular.
- Syntaktisch sich gleich verhaltende Wörter werden zu prälexikalischen Kategorien (Präterminale) zusammengefasst.
- Zuordnung von Wörtern zu Kategorien durch ein Lexikon.
 - (8) NP \rightarrow N | Det N | AP N | Det AP N
 - AP \rightarrow Adj | Adj AP
 - (9) NP \rightarrow (Det) Adj* N
- Ausdrucksfähigkeit wird nicht erhöht (schwache Äquivalenz zur PSG ohne Lexikon, Rekursion (8) wird durch Iteration (9) ersetzt.)

PSG mit Lexikon:

| | |
|--------------|---|
| Grammatik | $G' = (\Sigma', \Phi', P', S')$ |
| Präterminale | $\Sigma' = \{Det, Adj, N, \dots\}$ |
| Nonterminale | $\Phi' = \{S, VP, NP, \dots\}$ |
| Regeln | $P' = \{S \rightarrow NP VP \dots, NP \rightarrow Det N Det Adj N \dots\}$ |
| Startsymbol | S |
| Lexikon | $L = \left\{ \begin{array}{l} dem[Det], der[Det], ein[Det], alte[Adj], großen[Adj], \\ junger[Adj], Haus[N], Mann[N], Vater[N], \dots \end{array} \right\}$ |

61

Grammatikformalismen

- Regel (7) NP \rightarrow Det Adj N enthält zwei Arten von Informationen:
 - Dominanzinformation: NP besteht aus Det, Adj und N.
 - Präzedenzinformation: Det kommt vor Adj, Adj vor N.
- Dadurch erhält man zwei disjunkte Mengen Beschreibungsregeln:
 - Dominanzregeln: Aus welchen Subphrasen bestehen Phrasen.
 - Präzedenzregeln: Partielle Ordnung der Phrasen; Reihenfolge der Phrasen, die innerhalb der Dominanzregeln nicht verletzt werden dürfen.
- Diese Art der Grammatik nennt sich ID/LP-Grammatik [Gazdar und Pullum, 1981; Gazdar et al., 1985] (immediate dominance/linear-precedence)

ID/LP-Grammatik:

| | |
|-----------------|---|
| Grammatik | $G = (\Sigma', \Phi', D, P, S)$ |
| Präterminale | $\Sigma' = \{Det, Adj, N, \dots\}$ |
| Nonterminale | $\Phi' = \{S, VP, NP, \dots\}$ |
| Dominanzregeln | $D = \{S \rightarrow NP, VP \dots; NP \rightarrow Det, N Det, Adj, N; \dots\}$ |
| Präzedenzregeln | $P = \{NP < VP, Det < Adj, Det < N, Adj < N, \dots\}$ |
| Startsymbol | S |
| Lexikon | $L = \left\{ \begin{array}{l} dem[Det], der[Det], ein[Det], alte[Adj], großen[Adj], \\ junger[Adj], Haus[N], Mann[N], Vater[N], \dots \end{array} \right\}$ |

62

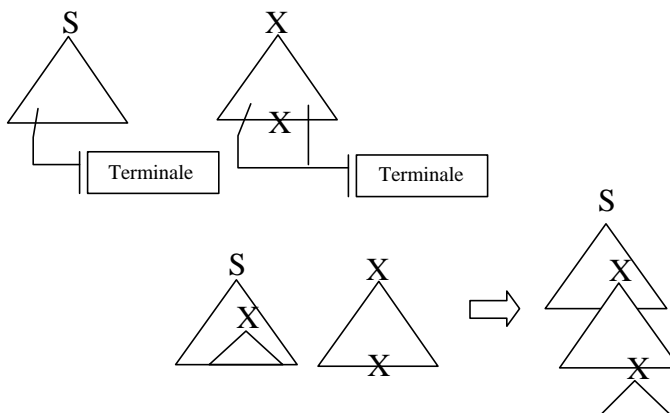
Grammatikformalismen

- Baumadjunktionsgrammatiken (Tree Adjoining Grammar – TAG) [Joshi, 1985]
 - Konkatenation als zweidimensionale Baumstruktur.
 - Statt Produktionsregeln gibt es zwei Baumarten:
 - Initialbäume: Wurzel Startsymbol S und Blätter ausschließlich Terminalsymbole
 - Auxiliarbäume: Blätter sind beliebige Anzahl von Terminalsymbolen und genau ein Nichtterminal, dieses ist das gleiche wie die Wurzel des AB.
 - Adjunktion als Standardoperation: AB kann einen beliebigen Teilbaum ersetzen wenn er die Wurzeln matchen
 - Jeder Initialbaum repräsentiert einen gültigen Satz der Sprache.
 - Adjunktion erzeugt ebenfalls gültige Sätze.

63

Grammatikformalismen

- Baumadjunktionsgrammatiken (TAG)



64

Grammatikformalismen

- **Baumadjunktionsgrammatiken (TAG):**

$G = (I, A)$

$I =$ Menge der Initialbäume

$A =$ Menge der Auxiliarbäume

Initialbäume sind Bäume, deren Wurzel das Startsymbol und deren Blätter alle Terminalsymbole sind.

Die Blätter der **Auxiliarbäume** sind Terminalsymbole bis auf genau ein Nichtterminal, das mit dem Wurzelknoten gleich ist.

Mit Hilfe der **Adjunktion** können Bäume zu komplexeren Bäumen zusammengesetzt werden.

- Mit TAGs können alle kontextfreien Sprachen dargestellt werden.
- TAGs können aber Bäume zuweisen, welche durch eine PSG nicht erzeugbar wären (schwache Äquivalenz).
- TAGs können aber so genannte mild kontextabhängige Sprachen akzeptieren.

65

Probleme mit PSG

- Unzureichende schwache generative Kapazität:
 - Natürliche Sprachen enthalten kontextsensitive Konstruktion. Zumindest scheinen einige Sprachen nicht kontextfreie Konstruktionen zu erlauben (-> TAGs).
 - Kontextfreie Grammatik wird immer übergenerieren.
- Unzureichende starke generative Kapazität:
 - Manche linguistische Daten sind mit PS-Regeln schwer darzustellen (Fernabhängigkeit, Übereinstimmung, ...)
 - Durch PS-Grammatik erzeugbare Strukturbeschreibungen sind linguistisch nicht adäquat (s. auch ID/LP und PSG).

66

Probleme mit PSG

Übereinstimmung (agreement) am Beispiel (7) NP -> Det Adj N:

- (7) beschreibt einfache Sätze wie
 - „der alte Mann“
 - „ein junger Mann“
 - „dem großen Haus“
- Aber auch
 - Der alter Vater
 - Ein junges Mann
 - Dem großen Häusern
- Regel (7) generiert über!
- Deutsch erfordert Übereinstimmung von Kasus, Genus und Numerus innerhalb einer NP. Adjektiva flektieren abhängig vom Determiner stark oder schwach (Paradigma).
- Dennoch können deutsche NPs kontextfrei beschrieben werden! Aufspaltung der Nichtterminalen (nach den agreement-Merkmalen) und Ersetzung von (7) durch Menge neuer Regeln.
 - Große Menge an Regeln.
 - Generalisierung geht verloren.

67

Probleme mit PSG

- Annotation als Einschränkungsmöglichkeit:

Regel: NP -> Det Adj N

Test: $Num_{Det} \cap Num_{Adj} \cap Num_N \neq \{ \}$

$Genus_{Det} \cap Genus_{Adj} \cap Genus_N \neq \{ \}$

$Casus_{Det} \cap Casus_{Adj} \cap Casus_N \neq \{ \}$

$Parad_{Det} \cap Parad_{Adj} \neq \{ \}$

Aktionen:

$Num_{NP} := Num_{Det} \cap Num_{Adj} \cap Num_N$

$Genus_{NP} := Genus_{Det} \cap Genus_{Adj} \cap Genus_N$

$Casus_{NP} := Casus_{Det} \cap Casus_{Adj} \cap Casus_N$

- Registerbelegung folgt aus Lexikoneinträgen:

Der: Num=Sing, Genus=mask, Casus=Nom, Parad.=stark

Num=Sing, Genus=fem, Casus=Gen, Parad.=stark

...

- Lexikonzugriff initialisiert die Register der prälexikalischen Kategorien.
- Aktionen propagieren die Werte zu den Nichtterminalen.

68

Spezielle Themen der KI

NLP Natural Language Processing Parsing

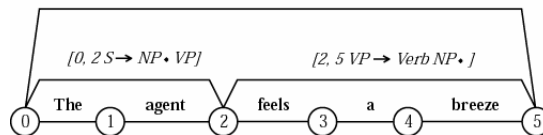
Parsing

Beispielllexikon und -grammatik für eine Sprache E0:

| | | |
|--|-----------------------------------|--|
| <i>Noun</i> → <i>stench</i> <i>breeze</i> <i>glitter</i> <i>nothing</i> <i>wumpus</i> <i>pit</i> <i>pits</i> <i>gold</i> <i>east</i> ... | | |
| <i>Verb</i> → <i>is</i> <i>see</i> <i>smell</i> <i>shoot</i> <i>feel</i> <i>stinks</i> <i>go</i> <i>grab</i> <i>carry</i> <i>kill</i> <i>turn</i> ... | | |
| <i>Adjective</i> → <i>right</i> <i>left</i> <i>east</i> <i>south</i> <i>back</i> <i>smelly</i> ... | | |
| <i>Adverb</i> → <i>here</i> <i>there</i> <i>nearby</i> <i>ahead</i> <i>right</i> <i>left</i> <i>east</i> <i>south</i> <i>back</i> ... | | |
| <i>Pronoun</i> → <i>me</i> <i>you</i> <i>I</i> <i>it</i> <i>S/HE</i> <i>Y'ALL</i> ... | <i>S</i> → <i>NP VP</i> | I + feel a breeze |
| <i>Name</i> → <i>John</i> <i>Mary</i> <i>Boston</i> <i>UCB</i> <i>PAJC</i> ... | <i>S Conjunction S</i> | I feel a breeze + and + I smell a wumpus |
| <i>Article</i> → <i>the</i> <i>a</i> <i>an</i> ... | <i>NP</i> → <i>Pronoun</i> | I |
| <i>Preposition</i> → <i>to</i> <i>in</i> <i>on</i> <i>near</i> ... | <i>Noun</i> | pits |
| <i>Conjunction</i> → <i>and</i> <i>or</i> <i>but</i> ... | <i>Article Noun</i> | the + wumpus |
| <i>Digit</i> → 0 1 2 3 4 5 6 7 8 9 | <i>Digit Digit</i> | 3 4 |
| | <i>NP PP</i> | the wumpus + to the east |
| | <i>NP RelClause</i> | the wumpus + that is smelly |
| | <i>VP</i> → <i>Verb</i> | stinks |
| | <i>VP NP</i> | feel + a breeze |
| | <i>VP Adjective</i> | is + smelly |
| | <i>VP PP</i> | turn + to the east |
| | <i>VP Adverb</i> | go + ahead |
| | <i>PP</i> → <i>Preposition NP</i> | to + the east |
| | <i>RelClause</i> → <i>that VP</i> | that + is smelly |

Parsing

- Parsing von links nach rechts kann backtracking erfordern:
 - „Have the students in section 4 of course 201 take the exam.“
 - „Have the students in section 4 of course 201 taken the exam?“
- Backtracking mit Reanalyse ist kostenintensiv.
- Resultate müssen zwischengespeichert werden
 - Im Beispiel kann die gelungene Analyse der NP „the students in section 4 of course 201“ in einen Zwischenspeicher - **chart** - abgelegt werden.
 - Parsing Algorithmen mit charts werden **Chartparser** genannt.
 - Alle Phrasen im Kontext eines Zweiges können auch in einem anderen Zweig des Suchraums verwendet werden.
 - Ein chart für einen n-Wort Satz besteht aus n+1 Knoten und Kanten



Chartparser

- Chartparser vereinen top-down und bottom-up:
 - Scanner (bottom-up) startet von den Worten aus und benutzt das Wort nur, um einen existierenden Kanteneintrag zu ergänzen.
 - Predictor (top-down) trägt neue Kanten in chart ein und bestimmt welche Symbole an welchen Positionen erfordert werden.
 - Completer (bottom-up) arbeitet wie Scanner und baut Konstituenten zur Ergänzung vorhandener Kanten auf.

```

function CHART-PARSE(string, grammar) returns chart
    chart[0] ← { [0, 0, S' → • S] }
    for v ← from 1 to LENGTH(string) do
        SCANNER(v, string[v])
    end
    return chart

procedure ADD-EDGE(edge)
    if edge in chart[END(edge)] then do nothing
    else
        push edge on chart[END(edge)]
        if COMPLETE?(edge) then COMPLETER(edge)
        else PREDICTOR(edge)

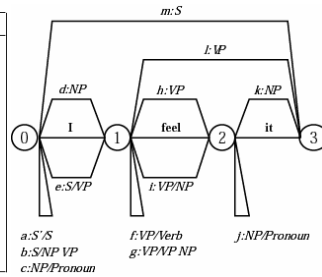
procedure SCANNER(j, word)
    for each [i, j, A → α • B β] in chart[j] do
        if word is of category B then
            ADD-EDGE([i, j+1, A → α B • β])
        end

procedure PREDICTOR([i, j, A → α • B β])
    for each (B → γ) in REWRITES-FOR(B, grammar) do
        ADD-EDGE([i, j, B → • γ])
    end

procedure COMPLETER([i, k, B → γ •])
    for each [i, j, A → α • B' β] in chart[j] do
        if B = B' then
            ADD-EDGE([i, k, A → α B' • β])
        end
    end
    
```

Chartparser

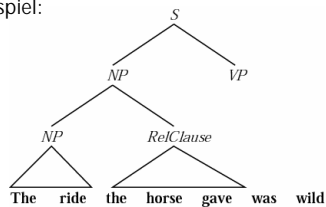
| Edge | Procedure | Derivation |
|------|----------------|---|
| a | INITIALIZER | [0, 0, $S' \rightarrow \bullet S$] |
| b | PREDICTOR(a) | [0, 0, $S \rightarrow \bullet NP VP$] |
| c | PREDICTOR(b) | [0, 0, $NP \rightarrow \bullet Pronoun$] |
| d | SCANNER(c) | [0, 1, $NP \rightarrow Pronoun \bullet$] |
| e | COMPLETER(b,d) | [0, 1, $S \rightarrow NP \bullet VP$] |
| f | PREDICTOR(e) | [1, 1, $VP \rightarrow \bullet Verb$] |
| g | PREDICTOR(e) | [1, 1, $VP \rightarrow \bullet VP NP$] |
| h | SCANNER(f) | [1, 2, $VP \rightarrow Verb \bullet$] |
| i | COMPLETER(g,h) | [1, 2, $VP \rightarrow VP \bullet NP$] |
| j | PREDICTOR(g) | [2, 2, $NP \rightarrow \bullet Pronoun$] |
| k | SCANNER(j) | [2, 3, $NP \rightarrow Pronoun \bullet$] |
| l | COMPLETER(i,k) | [1, 3, $VP \rightarrow VP NP \bullet$] |
| m | COMPLETER(e,l) | [0, 3, $S \rightarrow NP VP \bullet$] |



- Im Beispiel werden 13 Kanten in chart eingetragen, 8 sind unvollständig und 5 komplett.
- Predictor-Scanner-Completer Zyklen, etwa:
 - Kante (a) sucht nach S, um eine NP (b) mit folgenden Pronomen (c) vorherzusehen.
 - Scanner findet Pronomen (d) am richtigen Platz.
 - Completer kombiniert die unvollständige Kante (b) mit der vollständigen Kante (d) für eine neue Kante (e).

Chartparser

- Chartparser vermeiden die Erzeugung vieler Kanten bei einfachem bottom-up Vorgehen.
- Beispiel:



- bottom-up würde „ride the horse“ als VP labeln und später verwerfen.
- In der Beispielgrammatik EO ist aber keine Konkatenation von VP zu the; der Chartparser verhindert die unnötige VP Erzeugung.
- Algorithmen welche von links nach rechts arbeiten und die Erzeugung unnötiger Konstituenten erzeugen heißen **left-corner Parser**.
- Zeit- und Speicherbedarf bei Chartparsern ist **polynomial**.
 - $O(kn^2)$ Speicher für die Kanten, mit n: Anzahl der Worte und k: Grammatikkonstante.
 - Der Algorithmus terminiert, wenn keine Kanten mehr erzeugt werden können.
 - $O(n^3)$ Zeitaufwand im worst-case (Bestes Ergebnis für kontextfreie Sprachen).