

Programmierkurs Python II

Michaela Regneri

FR 4.7 Allgemeine Linguistik (Computerlinguistik)

Universität des Saarlandes

Sommersemester 2013



Übersicht

- Wiederholung: Elementare Parsing-Strategien
- CYK-Parser:
 - Parsing als dynamische Programmierung
 - Charts als kompakte Repräsentation von Teilergebnissen
 - Der Algorithmus: Erkennen & Parser
- Earley-Parser
 - ...kennt Ihr schon. :)

Probleme

- Die elementaren Parsing-Strategien (top-down, bottom-up) sind nicht auf allgemeine Grammatiken anwendbar
 - keine Tilgungsregeln, keine zyklischen Kettenregeln (BU)
 - keine (links-) rekursiven Regeln (TD)
- Lokale Ambiguität \Rightarrow Suche & Backtracking
 - Identische Teilergebnisse werden u.U. mehrfach berechnet
 - \Rightarrow Laufzeit exponentiell in der Eingabelänge (worst case)
- Lokale Ambiguität = Welche Regel bzw. Operation muss angewendet werden?

3

Beispielgrammatik

$S \rightarrow NP VP$	$DET \rightarrow an$
$NP \rightarrow DET N$	$DET \rightarrow the$
$NP \rightarrow POSS N$	$POSS \rightarrow his$
$NP \rightarrow NP PP$	$N \rightarrow elephant$
$PP \rightarrow P NP$	$N \rightarrow pajamas$
$VP \rightarrow V NP$	$N \rightarrow boy$
$VP \rightarrow VP PP$	$V \rightarrow shot$
	$P \rightarrow in$

4

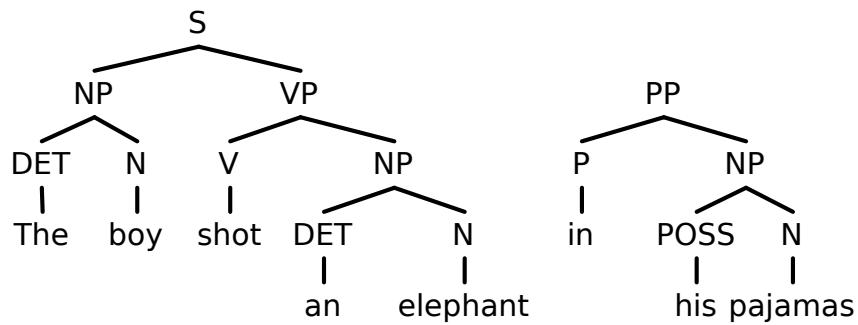
The boy shot an elephant in ...

⟨[], [the boy shot an elephant in his pajamas]⟩

⇒* ⟨[NP VP], [in his pajamas]⟩

⇒ ⟨[S], [in his pajamas]⟩

⇒* ⟨[S PP], []⟩ ⇒ **Backtracking**



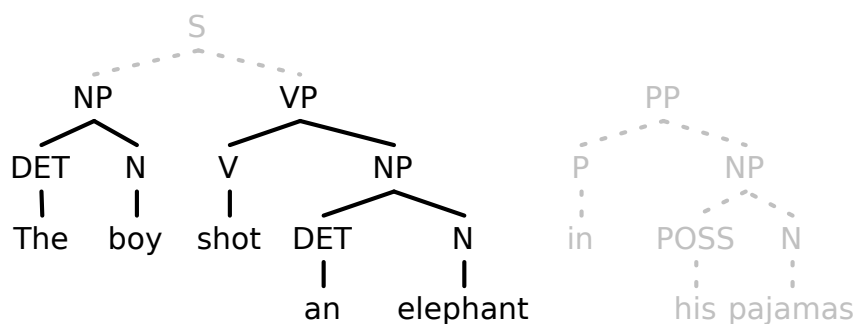
(Shift-Reduce)

5

The boy shot an elephant in ...

⟨[], [the boy shot an elephant in his pajamas]⟩

⇒* ⟨[NP VP], [in his pajamas]⟩



(Shift-Reduce)

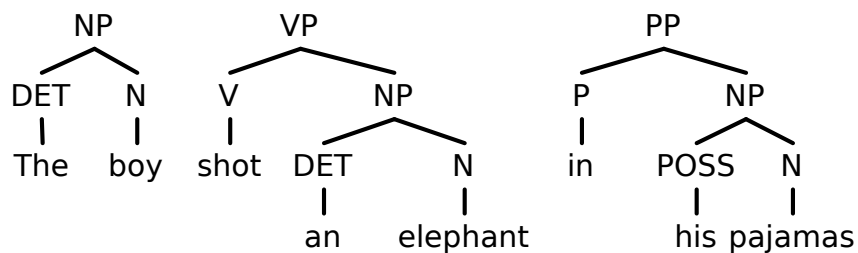
6

The boy shot an elephant in ...

⟨[], [the boy shot an elephant in his pajamas]⟩

⇒* ⟨[NP VP], [in his pajamas]⟩

⇒* ⟨[NP VP PP], []⟩



(Shift-Reduce)

7

The boy shot an elephant in ...

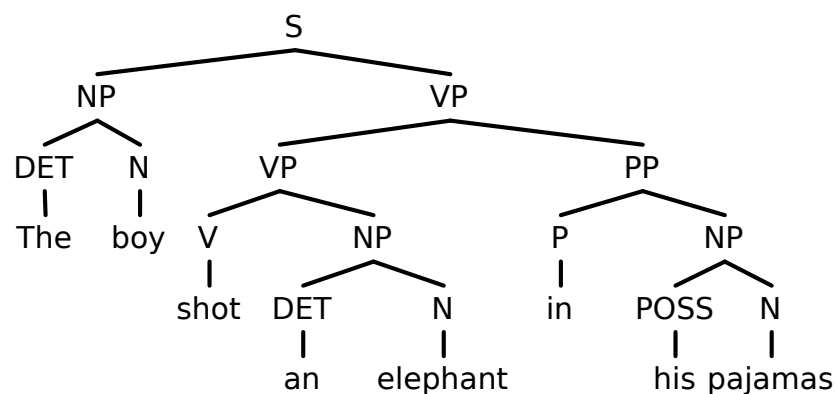
⟨[], [the boy shot an elephant in his pajamas]⟩

⇒* ⟨[NP VP], [in his pajamas]⟩

⇒* ⟨[NP VP PP], []⟩

⇒* ⟨[NP VP], []⟩

⇒* ⟨[S], []⟩

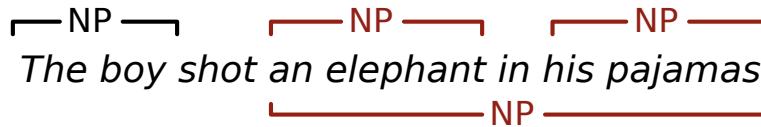


(Shift-Reduce)

8

Dynamische Programmierung

- **Kontextfreie Grammatiken:** die Anwendbarkeit einer Regel in einer Ableitung ist unabhängig vom Kontext.



- **Chart-Parsing:** Speichere bereits analysierte Teilergebnisse (Konstituenten) in einer „Chart.“
- **Charts** sind kompakte Repräsentation aller (lokal) möglichen Teilkonstituenten der Eingabekette.

9

Chart-Parsing

- **Chart-Parsing:** speichere bereits analysierte Teilergebnisse (Konstituenten) in einer „Chart.“
 - Charts aka. „well-formed substring table“
- **Charts können enthalten:**
 - Konstituenten, die bereits gefunden wurden
 - Hypothesen darüber, welche Konstituenten gefunden werden können (z.B. Earley-Algorithmus).
- **Verschiedene Chart-Parser:**
 - CYK, Earley, ...

10

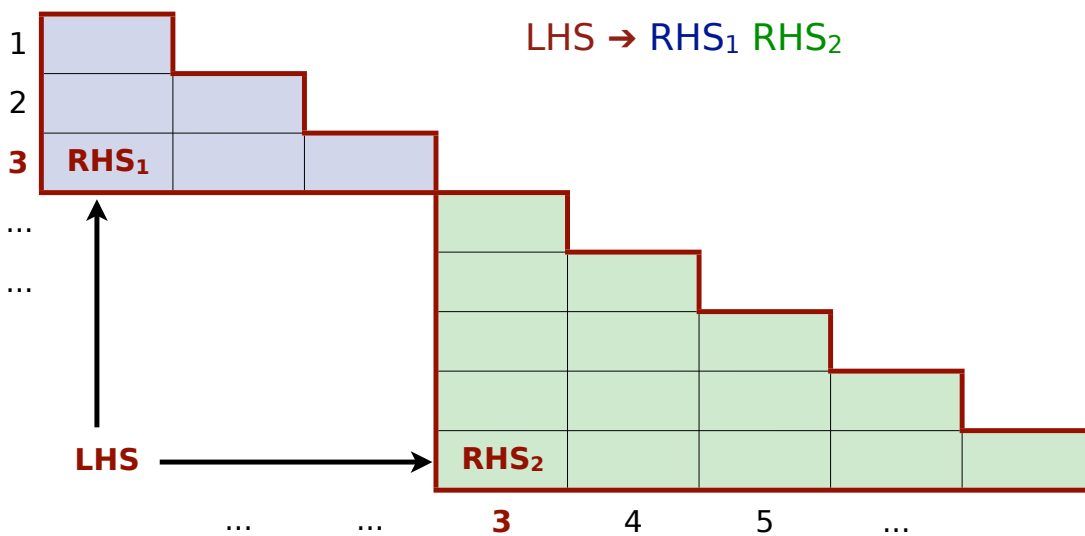
Charts

1	DET							
2	NP	N						
3	∅	∅	V					
4	∅	∅	∅	DET				
5	S	∅	VP	NP	N			
6	∅	∅	∅	∅	∅	P		
7	∅	∅	∅	∅	∅	∅	POSS	
8	S	∅	VP	NP	∅	PP	NP	N

$A \in T[i, j]$ gdw. $A \Rightarrow^* w_{i+1} \dots w_j$

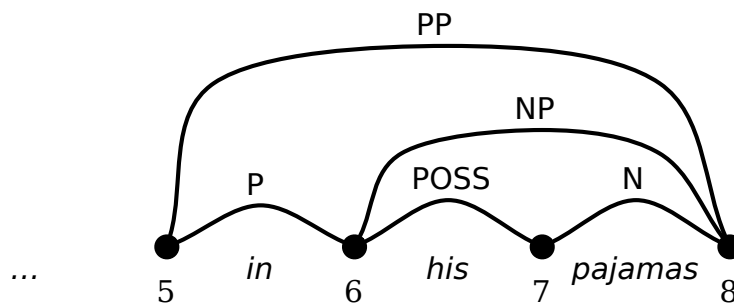
0 *The* 1 *boy* 2 *shot* 3 *an* 4 *elephant* 5 *in* 6 *his* 7 *pajamas* 8

Charts



Charts als Graph

- **Knoten** repräsentieren Positionen in der Eingabekette
- **Kanten** zwischen zwei Knoten repräsentieren Teilketten der Eingabe
 - Kante $n_i \rightarrow n_j$ gdw. $A \Rightarrow^* w_{i+1} \dots w_j$



13

CYK Algorithmus

- CYK (Cocke, Younger, Kasami) ist ein einfacher, chart-basierter bottom-up Parser.
- Die Grammatik muss in Chomsky-Normalform vorliegen:
 - $A \rightarrow w$ (w Terminalsymbol)
 - $A \rightarrow B C$ (B und C Nichtterminalsymbole)
 - $S \rightarrow \epsilon$ (S Startsymbol, nur wenn $\epsilon \in L$)
- Anmerkung: hier nehmen wir an, dass $\epsilon \notin L$, die Grammatik enthält also keine Regel $S \rightarrow \epsilon$

14

Grundlegende Idee

- Wenn
 - $B \Rightarrow^* w_i \dots w_{j-1}$
 - $C \Rightarrow^* w_j \dots w_{k-1}$
 - $A \rightarrow B C$
- Dann
 - $A \Rightarrow^* w_i \dots w_{k-1}$

15

CYK (Erkenner, Pseudo-code)

```
CYK(G, w1 ... wn):  
  for i in 1 ... n:  
    T[i-1, i] = { A | A → wi ∈ R }  
    for j in i - 2 ... 0:  
      T[j, i] = ∅  
      for k in j + 1 ... i - 1:  
        T[j, i] = T[j, i] ∪  
          { A | A → B C, B ∈ T[j,k], C ∈ T[k, i] }  
  if S ∈ T[0, n] then return True else return False
```

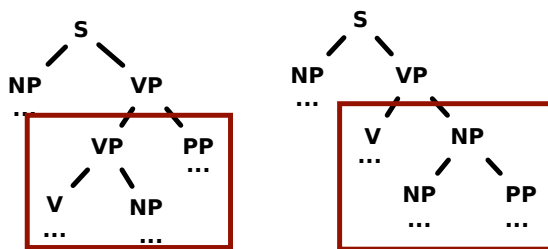
16

Erkenner → Parser

- Speichere zu jedem Nichtterminal NT in der Chart die "spans", die verwendet werden können, um NT abzuleiten
- Liste von Listen nötig, da Teilketten mehrere Ableitungsbäume haben können (Ambiguität)

CYK (Parser)

...	...						
3	...	V					
4	...	∅	DET				
5	...	VP [(2,3-4)]	NP [(3,4)]	N			
6	...	∅	∅	∅	P		
7	...	∅	∅	∅	∅	POSS	
8	...	VP [(2-4,5-7) (2, 3-7)]	NP [(3-4,5-7)]	∅	PP [(5,6-7)]	NP [(6,7)]	N
	...	2	3	4	5	6	7
		boy	shot	an	elephant	in	his pajamas



Chomsky Normalform

- Zu jeder kontextfreien Grammatik G gibt es eine äquivalente kontextfreie Grammatik G' in Chomsky Normalform.
- Algorithmus:
 - Eliminierung von ε -Regeln
 - Eliminierung "zu kurzer" Regeln ("Kettenregeln")
 - Eliminierung zu langer Regeln (mehr als 2 NTe)

19

Eliminierung "zu langer" Regeln

Linksbinarisierung(G):

```
while  $G$  enthält Regel  $A \rightarrow A_1 A_2 A_3 \dots A_k, k \geq 3$   
  entferne die Regel aus  $G$   
  neue Regel:  $\langle A_1, \dots, A_{k-1} \rangle \rightarrow A_1 \dots A_{k-1}$   
  neue Regel:  $A \rightarrow \langle A_1, \dots, A_{k-1} \rangle A_k$ 
```

Rechtsbinarisierung(G):

```
while  $G$  enthält Regel  $A \rightarrow A_1 A_2 A_3 \dots A_k, k \geq 3$   
  entferne die Regel aus  $G$   
  neue Regel:  $\langle A_2, \dots, A_k \rangle \rightarrow A_2 \dots A_k$   
  neue Regel:  $A \rightarrow A_1 \langle A_2, \dots, A_k \rangle$ 
```

20

Implementierungsvarianten

- $T[i,j] = T[i,j] \cup \{ A \mid A \rightarrow B C, B \in T[i,k], C \in T[k,j] \}$
 - \Rightarrow kann verschieden implementiert werden
- **Variante 1**
 - Iteriere über alle Regeln $A \rightarrow B C$
 - Prüfe, ob $B \in T[i,k]$ und $C \in T[k,j]$
- **Variante 2**
 - Iteriere über alle $B \in T[i,k]$
 - Iteriere über alle Regeln $A \rightarrow B C$
 - Prüfe, ob $C \in T[k, j]$

21

Implementierungsvarianten

- $T[i,j] = T[i,j] \cup \{ A \mid A \rightarrow B C, B \in T[i,k], C \in T[k,j] \}$
 - \Rightarrow kann verschieden implementiert werden
- **Variante 3**
 - Iteriere über alle $C \in T[k,j]$
 - Iteriere über alle Regeln $A \rightarrow B C$
 - Prüfe, ob $A \in T[i,k]$
- **Variante 4**
 - Iteriere über alle $B \in T[i,k]$ und $C \in T[k,j]$
 - Prüfe, ob es Regel $A \rightarrow B C$ gibt

22

Song &al. (2008)

- Ein paar konkrete Zahlen zur Grammatikgröße

	# of Symbols	# of Rules
Original	72	14,971
Right	10,654	25,553
Left	12,944	27,843
Head	11,798	26,697
Compact	3,644	18,543
Ours	8,407	23,306

Table 3: Grammar size of different binarizations

23

Song &al. (2008)

- Ein paar konkrete Laufzeitresultate:

Binarization	Constituents	Time (s)
Right	241,924,229	5,747
Left	193,238,759	3,474
Head	166,425,179	3,837
Compact	94,257,478	2,302
Ours	52,206,466	2,182

Table 4: Performance on test set

- Das Test-Set umfasst 90% der Sätze im Wall Street Journal mit maximal 40 Worten.

24

CYK

- CYK ist eingeschränkt auf Grammatiken in CNF
 - keine echte Einschränkung, Grammatiken können in CNF überführt werden
- CYK ist ein Bottom-Up-Parser und erzeugt unter Umständen viele nicht benötigte Konstituenten
- *Mary gave the man a book*
 - $N \rightarrow \text{man}$
 - $V \rightarrow \text{man} \Rightarrow$ CYK leitet für „man a book“ eine VP ab

25

Earley-Algorithmus

- Ein effizienter „aktiver“ Chart-Parser
 - Zeikomplexität $O(n^3)$ in der Eingabelänge n
- Für beliebige kontextfreie Grammatiken geeignet
 - Tilgungsregeln
 - Zyklische Kettenregeln
 - Links- und rechtsrekursive Regeln
- Mischung aus bottom-up und top-down
 - Analyserichtung links \rightarrow rechts

26

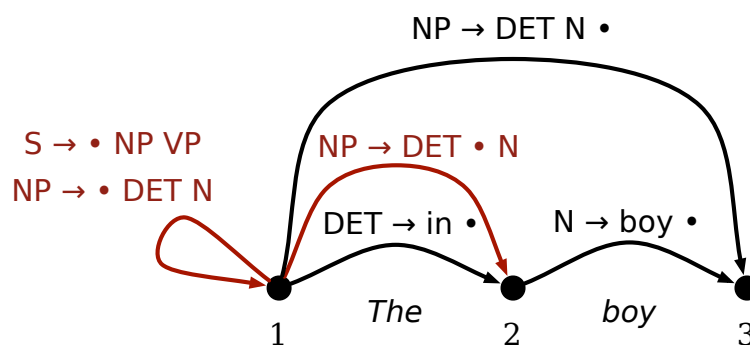
Earleys Algorithmus

- Passive Chart-Parser speichern nur vollständige Konstituenten als Teilergebnis in der Chart.
- Aktive Chart-Parser speichern auch Hypothesen.
- $A \rightarrow \mathbf{B_1 \dots B_n} \cdot \mathbf{C_{n+1} \dots C_k}$
 - $\mathbf{B_1 \dots B_2}$ wurden bereits gefunden
 - $\mathbf{C_{n+1} \dots C_k}$ müssen noch gefunden werden

27

Die Chart

- **Knoten** entsprechen Positionen in der Eingabe
- **Kanten** entsprechen gefundenen Teilkonstituenten. Auch unvollständige Konstituenten werden gespeichert.
 - Aktive Kanten $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle$
 - Passive Kanten $\langle A \rightarrow \alpha \cdot, i, j \rangle$

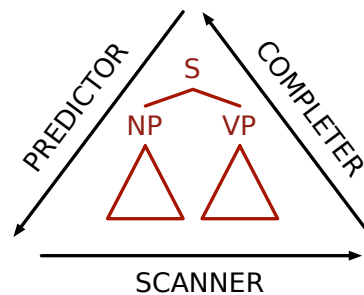


28

Earleys Algorithmus

- **Drei Prozeduren**

- Predictor
- Scanner
- Completer

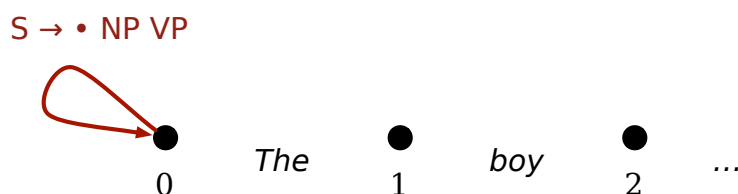


- Die drei Prozeduren werden für jeden Knoten in der Chart von links nach rechts aufgerufen.
- \Rightarrow jeweils solange, bis keine neuen Kanten zur Chart hinzugefügt werden

29

Initialisierung

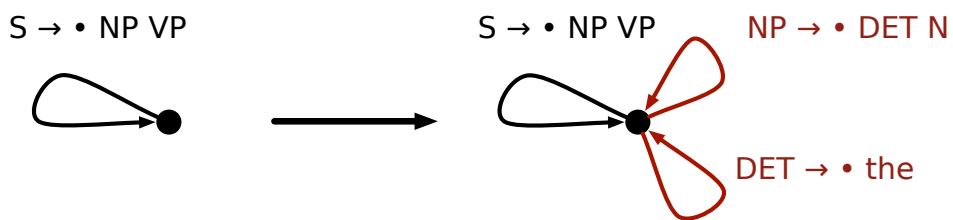
- Grammatik $G = \langle V, \Sigma, R, S \rangle$
- Eingabe $w_1 \dots w_n$
- Für jede Regel $S \rightarrow \alpha$
 - Füge $\langle S \rightarrow \cdot \alpha, 0, 0 \rangle$ zur Chart hinzu



30

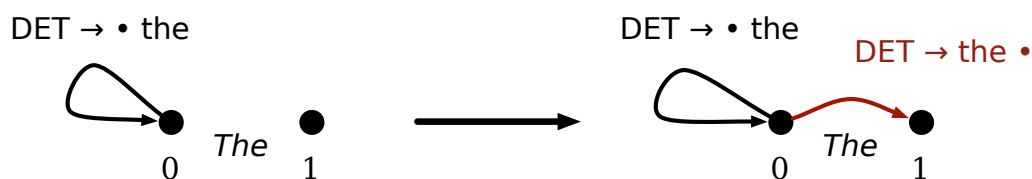
Predictor

- Der Predictor sagt voraus, welche Grammatikregeln zum Ziel führen könnten (top-down).
- Für neue Kanten $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle$
 - Für alle Regeln $B \rightarrow \gamma \in \text{Regeln}$
 - Neue Kante: $\langle B \rightarrow \cdot \gamma, j, j \rangle$



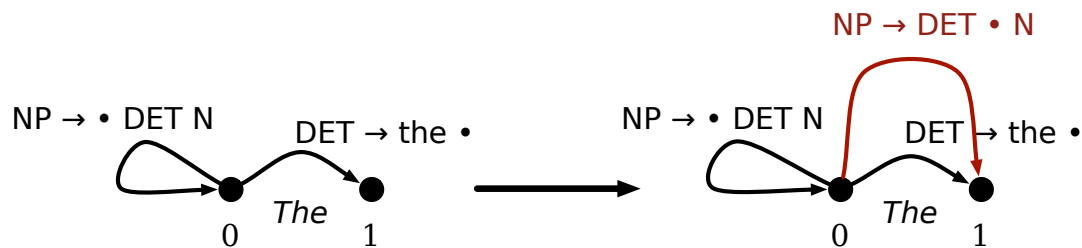
Scanner

- Der Scanner konsumiert Terminalsymbole aus der Eingabekette.
- Für neue Kanten $\langle A \rightarrow \alpha \cdot w_{j+1} \beta, i, j \rangle$
 - Neue Kante $\langle A \rightarrow \alpha w_{j+1} \cdot \beta, i, j + 1 \rangle$



Completer

- Der Completer vervollständigt aktive Kanten durch Kombination mit inaktiven Kanten (bottom-up).
- Für jede neue Kante $\langle A \rightarrow \alpha \cdot, j, k \rangle$
 - Für jede alte Kante $\langle B \rightarrow \gamma \cdot A \beta, i, j \rangle$
 - Neue Kante: $\langle B \rightarrow \gamma A \cdot \beta, i, k \rangle$



33

Übersicht

- **Scan**
 - $\langle A \rightarrow \alpha \cdot w_{j+1} \beta, i, j \rangle \Rightarrow \langle A \rightarrow \alpha w_{j+1} \cdot \beta, i, j + 1 \rangle$
- **Predict**
 - $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle, B \rightarrow \gamma \Rightarrow \langle B \rightarrow \cdot \gamma, j, j \rangle$
- **Complete**
 - $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle, \langle B \rightarrow \gamma \cdot, j, k \rangle \Rightarrow \langle A \rightarrow \alpha B \cdot \beta, i, k \rangle$

34

Eine Chart

der Hund bellt

0		1		2		3	
$START \rightarrow \cdot S$	0,0	$D \rightarrow der \cdot$	0,1	$N \rightarrow Hund \cdot$	1,2	$V \rightarrow bellt \cdot$	2,3
$S \rightarrow \cdot NP VP$	0,0	$NP \rightarrow D \cdot N$	0,1	$NP \rightarrow D N \cdot$	0,2	$VP \rightarrow V \cdot$	2,3
$NP \rightarrow \cdot D N$	0,0			$S \rightarrow NP \cdot VP$	0,2	$S \rightarrow NP VP \cdot$	0,3
				$VP \rightarrow \cdot V$	2,2		

35

Der Algorithmus

- Neue Kante $\langle START \rightarrow \cdot S, 0, 0 \rangle$
- Für alle $j \in 1, \dots, n$
 - Für alle $\langle A \rightarrow \alpha \cdot \beta, i, j \rangle \in \text{Chart}$
 - Scan, wenn β mit einem Terminalsymbol beginnt
 - Predict, wenn β mit einem Nichtterminal beginnt
 - Complete, wenn $\beta = \epsilon$
- Akzeptiere, wenn $\langle START \rightarrow S \cdot, 0, n \rangle$

36

Implementierung

- Kanten $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle$ stellen wir dar als $(A, [B, \dots], i, j)$, d.h. α wird nicht explizit repräsentiert.
- Der Algorithmus (Skizze)
 - Auf einer Agenda (Liste) werden durch Anwendung von Scan, Predict, Complete erzeugte Kanten (zwischen-) gespeichert
 - Die Agenda enthält initial die Kante $(\text{START}, [\text{S}], 0, 0)$
 - In jedem Schritt entfernen wir eine Kante von der Agenda, wenden Scan, Predict oder Complete an und speichern die Kante abschließend in der Chart (Liste)
 - Durch Anwendung von Scan, Predict, Complete erzeugte Kanten werden auf der Agenda gespeichert

37

Implementierung

```
rules = [ ('S', ['NP', 'VP']),  
          ('NP', ['DET', 'N']),  
          ('VP', ['V', 'NP']),  
          ('VP', ['V']),  
          ('DET', ['der']),  
          ('DET', ['das']),  
          ('N', ['Student']),  
          ('N', ['Buch']),  
          ('V', ['liest']),  
          ('V', ['arbeitet'])  
        ]
```

38

Implementierung

```
def earley(start, rules, words):
    length = len(words)
    chart = []
    agenda = [('START', [start], 0, 0)]
    while agenda:
        (lhs, rhs, i, j) = agenda.pop()
        # scan, predict, complete (nächste Folie)
        chart.append((lhs, rhs, i, j))
    return ('START', [], 0, length) in chart
```

39

Implementierung

```
def earley(start, rules, words):
    ...
    while agenda:
        (lhs, rhs, i, j) = agenda.pop()
        if rhs:
            if j < length and rhs[0] == words[j]: # scan
                agenda.append((lhs, rhs[1:], i, j + 1))
            else: # predict
                for (_lhs, _rhs) in rules:
                    if rhs[0] == _lhs:
                        agenda.append((_lhs, _rhs, j, j))
        else:
            ...
```

40

Variante

- Im hier vorgestellten Verfahren wird nicht zwischen Lexikon und „eentlichen“ Regeln unterschieden.
- Regeln wie „DET → the“ werden vom Predictor behandelt
 - ⇒ es können Terminalsymbole vorhergesagt werden, die gar nicht in der Eingabekette vorkommen (ineffizient)
- **Sinnvolle Variante:**
 - Regeln wie „DET → the“ vom Scanner verarbeiten lassen:
 - Für neue Kanten $\langle A \rightarrow \alpha \cdot B \beta, i, j \rangle$
 - Neue Kante $\langle A \rightarrow \alpha B \cdot \beta, i, j + 1 \rangle$ wenn $B \rightarrow w_{i+1}$