

# Parsing Grundlagen und Chart Parsing

Katja Markert

Institut für Computerlinguistik  
Uni Heidelberg  
markert@cl.uni-heidelberg.de  
mit einigen Folien von Yannick Versley

December 12, 2019

- 1 Was ist Parsing?
- 2 Top-Down und Bottom-Up Parsing
- 3 Chart Parsing: CKY

- 1 Was ist Parsing?
- 2 Top-Down und Bottom-Up Parsing
- 3 Chart Parsing: CKY

## Parsing

Gegeben eine CFG. Wie kann für einen Satz  $s$  automatisch

- 1 das Erkennungsproblem gelöst werden?
- 2 die syntaktische Analyse durchgeführt werden?

Eingabe: Ein Satz  $s$  und eine CFG

Ausgabe: Ein oder mehrere Phrasenstrukturbäume, falls  $s$  grammatisch. *False* sonst.

- Grammar Checking
- Zwischenstufe für tiefe semantische Verarbeitung  
(Informationsextraktion, Question-Answering)
  - [S [NP Who ] [VP won [NP the Nobel Prize for Physics ] in 2007 ] ? ]
  - [S [NP France's Albert Fert and Germany's Peter Gruenberg ] [VP win [NP the 2007 Nobel Prize for Physics ] ] it was announced Tuesday Oct 9 2007 . ]

- 1 Was ist Parsing?
- 2 Top-Down und Bottom-Up Parsing**
- 3 Chart Parsing: CKY

Zwei Informationsquellen beschränken die Suche:

- 1 Eingabesatz: Jeder Parsebaum muss die Wörter des Eingabesatzes als Blätter (Terminalsymbole) enthalten
- 2 Startsymbol der Grammatik: Jeder Parsebaum muss Startsymbol als Wurzelknoten besitzen

Daraus ergeben sich zwei Grundlegende Parsingstrategien:  
Bottom-Up (datengetrieben) und Top-Down (zielgetrieben)

Expandiert Grammatikregeln ausgehend vom Startsymbol. Die Suche ist erfolgreich, wenn die Blätter des Parsebaums genau mit den Wörtern des Eingabestrings übereinstimmen.

## Recursive Descent Algorithmus

- top-down
- left-right
- erschöpfende depth-first Suche mit Backtracking
- Regeln in Grammatikreihenfolge

$S \rightarrow NP VP$

$S \rightarrow Aux NP VP$

$S \rightarrow VP$

$NP \rightarrow Pronoun$

$NP \rightarrow Proper-Noun$

$NP \rightarrow Det Nominal$

$Nominal \rightarrow Noun$

$Nominal \rightarrow Nominal Noun$

$Nominal \rightarrow Nominal PP$

$VP \rightarrow Verb$

$VP \rightarrow Verb NP$

$VP \rightarrow Verb NP PP$

$VP \rightarrow Verb PP$

$VP \rightarrow VP PP$

$PP \rightarrow Preposition NP$

$Det \rightarrow that | this | a$

$Noun \rightarrow book | flight | meal | money$

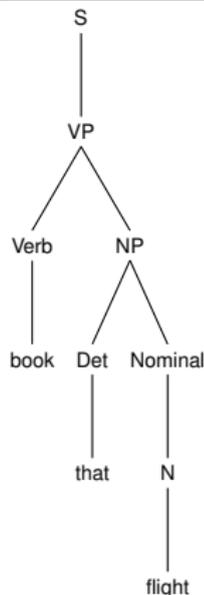
$Verb \rightarrow book | include | prefer$

$Pronoun \rightarrow I | she | me$

$Proper-Noun \rightarrow Houston | TWA$

$Aux \rightarrow does$

$Preposition \rightarrow from | to | on | near | through$



Beispielsatz  $S_1$ : book that flight

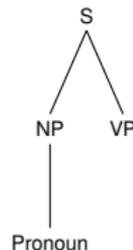
# Beispiel Recursive Descent für *book that flight*

S



<i>S</i> → <i>NP VP</i>	<i>Det</i> → <i>that   this   a</i>
<i>S</i> → <i>Aux NP VP</i>	<i>Noun</i> → <i>book   flight   meal   money</i>
<i>S</i> → <i>VP</i>	<i>Verb</i> → <i>book   include   prefer</i>
<i>NP</i> → <i>Pronoun</i>	<i>Pronoun</i> → <i>I   she   me</i>
<i>NP</i> → <i>Proper-Noun</i>	<i>Proper-Noun</i> → <i>Houston   TWA</i>
<i>NP</i> → <i>Det Nominal</i>	<i>Aux</i> → <i>does</i>
<i>Nominal</i> → <i>Noun</i>	<i>Preposition</i> → <i>from   to   on   near   through</i>
<i>Nominal</i> → <i>Nominal Noun</i>	
<i>Nominal</i> → <i>Nominal PP</i>	
<i>VP</i> → <i>Verb</i>	
<i>VP</i> → <i>Verb NP</i>	
<i>VP</i> → <i>Verb NP PP</i>	
<i>VP</i> → <i>Verb PP</i>	
<i>VP</i> → <i>VP PP</i>	
<i>PP</i> → <i>Preposition NP</i>	

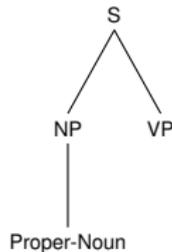
# Beispiel Recursive Descent für *book that flight*



Backtracking nötig, da POS Mismatch

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid money$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston \mid TWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from \mid to \mid on \mid near \mid through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	

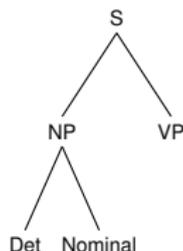
# Beispiel Recursive Descent für *book that flight*



Backtracking nötig, da POS Mismatch

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid money$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston \mid TWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from \mid to \mid on \mid near \mid through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	

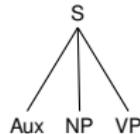
# Beispiel Recursive Descent für *book that flight*



Backtracking nötig, da POS Mismatch

$S \rightarrow NP VP$	$Det \rightarrow that   this   a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book   flight   meal   money$
$S \rightarrow VP$	$Verb \rightarrow book   include   prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I   she   me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston   TWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from   to   on   near   through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	

# Beispiel Recursive Descent für *book that flight*



Backtracking nötig, da POS Mismatch

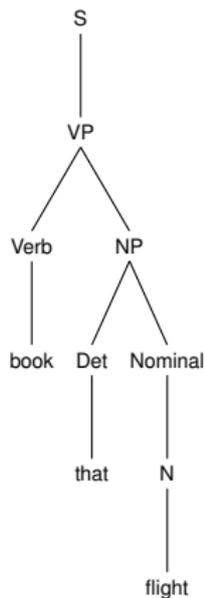
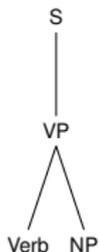
$S \rightarrow NP VP$	$Det \rightarrow that   this   a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book   flight   meal   money$
$S \rightarrow VP$	$Verb \rightarrow book   include   prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I   she   me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston   TWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from   to   on   near   through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	

# Beispiel Recursive Descent für *book that flight*



Backtracking nötig, das zuviel Input übrig

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid money$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Pronoun$	$Pronoun \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston \mid TWA$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	$Preposition \rightarrow from \mid to \mid on \mid near \mid through$
$Nominal \rightarrow Nominal Noun$	
$Nominal \rightarrow Nominal PP$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	
$VP \rightarrow Verb NP PP$	
$VP \rightarrow Verb PP$	
$VP \rightarrow VP PP$	
$PP \rightarrow Preposition NP$	



Endlich (nach zwei weiteren Backtrackings an NP)

- Die Blätter des generierten Baumes werden mit den Wörtern des Eingabesatzes abgeglichen. Ist das aktuelle Eingabewort kompatibel mit dem POS, so ist die verwendete Regel erfolgreich
- Können die Blätter des konstruierten Parsebaums mit den Wörtern des Eingabestrings vollständig abgeglichen werden, so ist der konstruierte Baum ein valider Parsebaum.

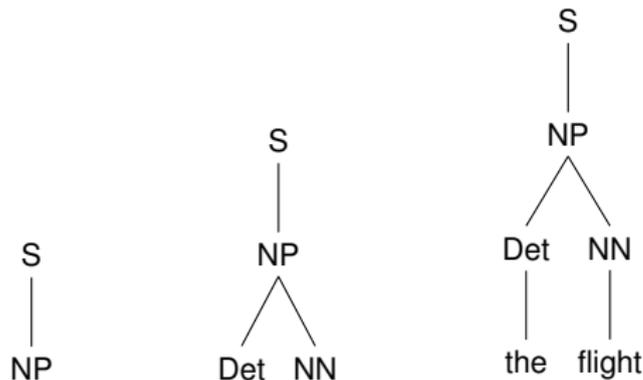
- Fehlschläge erst bei Abgleich mit Blättern sichtbar
- Probleme mit POS-Tagging Mismatch können durch eine Kombination mit Bottom-Up (left-corner parsing) vermieden werden
- Beim systematischem Rücksetzen müssen identische Teilstrukturen immer wieder berechnet werden (s. nächstes Beispiel)
- Diese Ineffizienz kann durch Bottom-up Kombination nicht vermieden werden

Beispiel für Ineffizienz trotz gutem POS Matching

S  $\rightarrow$  NP  
NP  $\rightarrow$  Det NN  
NP  $\rightarrow$  NP PP  
NP  $\rightarrow$  NN  
PP  $\rightarrow$  IN NP  
NN  $\rightarrow$  {flight, London }  
IN  $\rightarrow$  to  
Det  $\rightarrow$  the

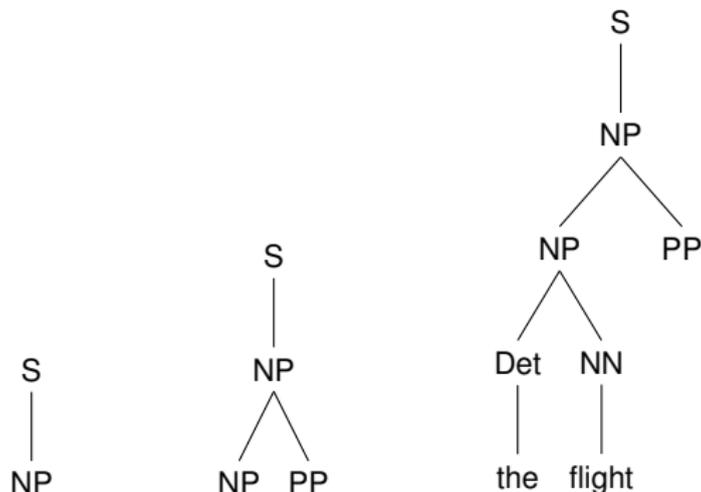
Satz: *the flight to London*

*the flight to London*



PP kann nicht inkorporiert werden → Backtrack

*the flight to London*



NP *the flight* wurde doppelt abgeleitet. Komplexität im schlimmsten Fall exponentiell zur Eingabelänge

$$\begin{aligned} S &\rightarrow NP \\ NP &\rightarrow NP PP \\ NP &\rightarrow NNS \\ PP &\rightarrow IN NP \\ NNS &\rightarrow \{\text{men, hats}\} \\ IN &\rightarrow \text{with} \end{aligned}$$

Parsen Sie damit Top-Down *men with hats*

Linksrekursion:  $A \rightarrow A X$ . Unendliche Schleife. Kann vermieden werden, wenn man Grammatik anders schreibt.

$$\begin{aligned} NP &\rightarrow NP PP \\ NP &\rightarrow NNS \end{aligned}$$
$$\begin{aligned} NP &\rightarrow NNS T \\ T &\rightarrow PP T \\ T &\rightarrow \varepsilon \end{aligned}$$

# Recursive Descent Probleme: Rekursion

$$\begin{aligned} S &\rightarrow NP \\ NP &\rightarrow NP PP \\ NP &\rightarrow NNS \\ PP &\rightarrow IN NP \\ NNS &\rightarrow \{\text{men, hats}\} \\ IN &\rightarrow \text{with} \end{aligned}$$

Parsen Sie damit Top-Down *men with hats*

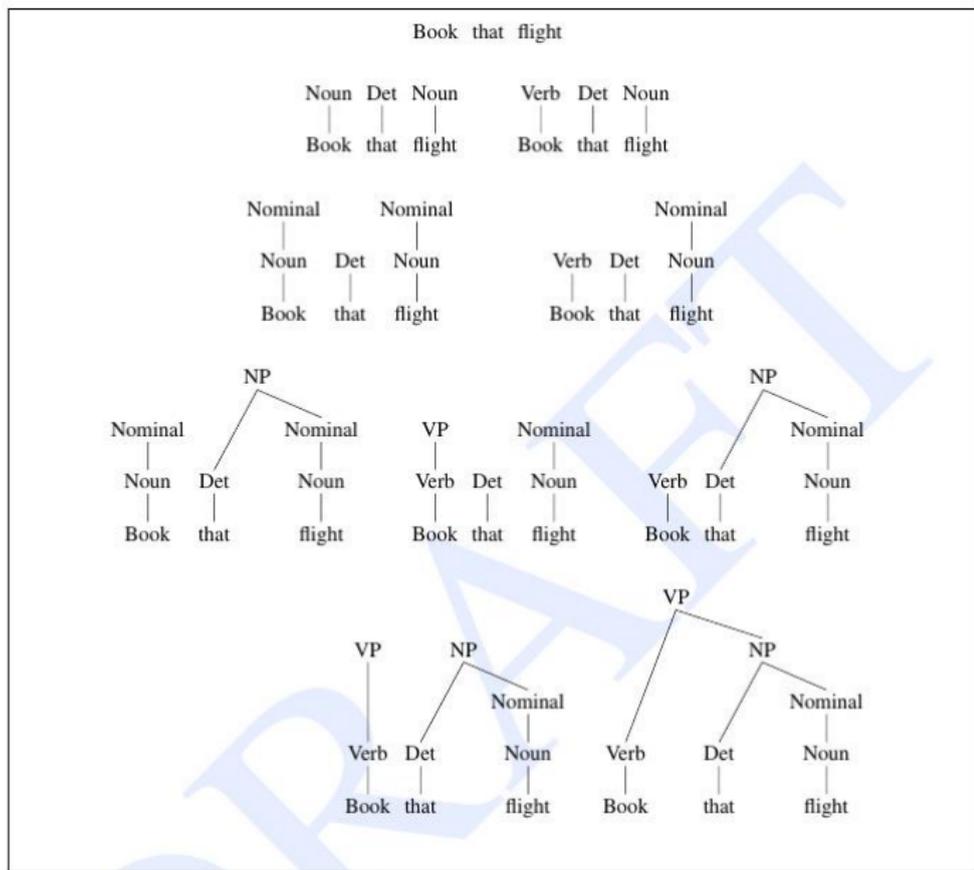
Linksrekursion:  $A \rightarrow A X$ . Unendliche Schleife. Kann vermieden werden, wenn man Grammatik anders schreibt.

$$\begin{aligned} NP &\rightarrow NP PP \\ NP &\rightarrow NNS \end{aligned}$$
$$\begin{aligned} NP &\rightarrow NNS T \\ T &\rightarrow PP T \\ T &\rightarrow \varepsilon \end{aligned}$$

Wie behandelt der Parser Ambiguität? Lösungen?

- ausgehend von der Eingabesequenz
- generiere Teilbäume auf Basis von Ersetzungsregeln (betrachte rechte Seite), bis das Startsymbol erreicht ist
- erfolgreich, wenn der Parsebaum die gesamte Eingabesequenz überspannt und einen einzigen Wurzelknoten enthält, der mit dem Startsymbol übereinstimmt

# Bottom-Up Parsing Beispiel



## Top-Down:

- 1 exploriert keine ungrammatikalischen Bäume
- 2 exploriert keine Bäume, die keinen Platz in einem Baum mit Wurzelknoten S finden können
- 3 generiert viele Bäume, die keinerlei Aussicht haben auf die Eingabewörter zu passen

## Bottom-Up:

- 1 exploriert nur solche Bäume, die mit der Eingabesequenz kompatibel sind
- 2 generiert viele Bäume, die keinerlei Aussicht haben, in das Startsymbol zu münden

## Beide:

- 1 Viele Teilbäume sind nicht Teil einer vollständigen validen Analyse
- 2 Viele Teilbäume werden wiederholt berechnet. → Dynamische Programmierung/Chart Parsing/CKY

- 1 Was ist Parsing?
- 2 Top-Down und Bottom-Up Parsing
- 3 Chart Parsing: CKY**

- Gesamtparse besteht aus validen Teilbäumen
- Dynamische Programmierung: Speichern von Lösungen für Teilprobleme in Chart zur Vermeidung wiederholter Berechnungen (polynomiale Laufzeit)

## Chart

Tabelle, in der Teilresultate der Analyse abgelegt werden. Einträge in Chart heißen Kanten. Eine Kante enthält mindestens:

- Satzabschnitt, auf den sich die Kante bezieht
- Angewandte Syntaxregel

- Eingabe:  $_0 \text{Book}_1 \text{ the}_2 \text{ flight}_3 \text{ through}_4 \text{ Houston}_5$
- Eine Kante der Chart:  $PP[3,5]$ . Überspannt die Phrase *through Houston*. Kann sich nur zusammensetzen aus  $[3,4]$  und  $[4, 5]$
- CYK Algorithmus: Erstellung von Chart-Kanten (Kasami 1965, Younger 1967, Cocke and Schwary 1970)
- CYK arbeitet nur mit Grammatiken in Chomsky-Normalform.

## Rechts in Chomsky-Normalform

$S \rightarrow NP VP$	$S \rightarrow NP VP$
$S \rightarrow Aux NP VP$	$S \rightarrow XI VP$
	$XI \rightarrow Aux NP$
$S \rightarrow VP$	$S \rightarrow book \mid include \mid prefer$
	$S \rightarrow Verb NP$
	$S \rightarrow X2 PP$
	$S \rightarrow Verb PP$
	$S \rightarrow VP PP$
$NP \rightarrow Pronoun$	$NP \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$NP \rightarrow TWA \mid Houston$
$NP \rightarrow Det Nominal$	$NP \rightarrow Det Nominal$
$Nominal \rightarrow Noun$	$Nominal \rightarrow book \mid flight \mid meal \mid money$
$Nominal \rightarrow Nominal Noun$	$Nominal \rightarrow Nominal Noun$
$Nominal \rightarrow Nominal PP$	$Nominal \rightarrow Nominal PP$
$VP \rightarrow Verb$	$VP \rightarrow book \mid include \mid prefer$
$VP \rightarrow Verb NP$	$VP \rightarrow Verb NP$
$VP \rightarrow Verb NP PP$	$VP \rightarrow X2 PP$
	$X2 \rightarrow Verb NP$
$VP \rightarrow Verb PP$	$VP \rightarrow Verb PP$
$VP \rightarrow VP PP$	$VP \rightarrow VP PP$
$PP \rightarrow Preposition NP$	$PP \rightarrow Preposition NP$

**Figure 13.8**  $\mathcal{L}_1$  Grammar and its conversion to CNF. Note that although they aren't shown here all the original lexical entries from  $\mathcal{L}_1$  carry over unchanged as well.

# Unsere Beispielgrammatik in CNF

S → NP VP

S → X1 VP

X1 → Aux NP

S → Verb NP

S → X2 PP

S → Verb PP

S → VP PP

NP → Det Nominal

Nominal → Nominal Noun

Nominal → Nominal PP

VP → Verb NP

VP → X2 PP

X2 → Verb NP

VP → VP PP

PP → Preposition NP

S → book |include |prefer

NP → I |she |me

NP → TWA|Houston

Nominal → book | flight | meal | money

VP → book|include|prefer

Verb → book |include|prefer

Noun → book |flight |meal |money

Proper Noun → TWA |Houston

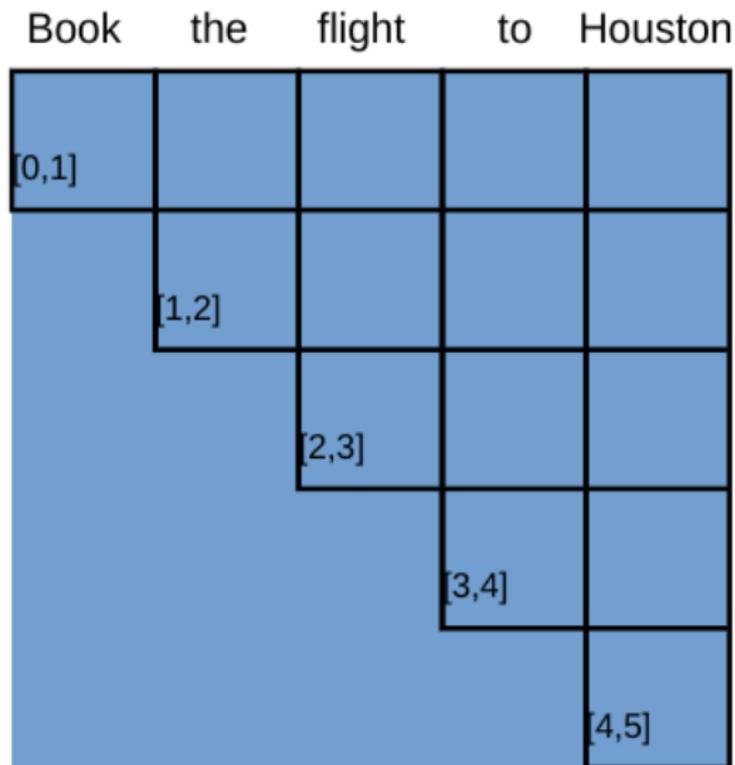
Pronoun → I | she | me

Aux → does

Det → that | this | a | the

Preposition → from | to | on |near |through

# CKY Algorithmus am Beispiel



# CKY Algorithmus am Beispiel

Book the flight to Houston

[0,1]	[0,2]	[0,3]	[0,4]	[0,5]
	[1,2]	[1,3]	[1,4]	[1,5]
		[2,3]	[2,4]	[2,5]
			[3,4]	[3,5]
				[4,5]

Initialisierung: Auf Diagonale  $[j-1,j]$   
stehen mögliche Pos Tags  
(bottom-up)

Dann spaltenweise links nach  
Rechts, unten nach oben

Zelle  $[i,j]$  kann sich zusammensetzen  
aus  $[i,k]$  und  $[k,j]$  für alle  $k$   
zwischen  $i+1$  und  $j-1$

Binäre Zusammensetzung  $\rightarrow$  CNF

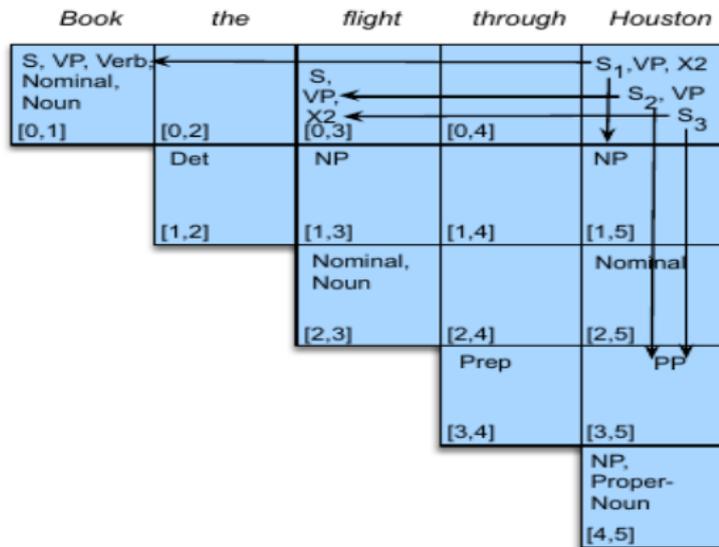
# CKY Algorithmus am Beispiel

Am Anfang der letzten Spalte:

<i>Book</i>	<i>the</i>	<i>flight</i>	<i>through</i>	<i>Houston</i>
S, VP, Verb, Nominal, Noun [0,1]	[0,2]	S,VP,X2 [0,3]	[0,4]	[0,5]
	Det [1,2]	NP [1,3]	[1,4]	NP [1,5]
		Nominal, Noun [2,3]	[2,4]	[2,5]
			Prep ← PP [3,4]	[3,5]
				NP, Proper- Noun [4,5]

]

Endergebnis (Bild aus Jurafsky und Martin):



- Obere rechte Ecke: im Buch wurde m.E. ein VP aus X<sub>2</sub> PP vergessen
- S<sub>2</sub> und S<sub>3</sub> sind äquivalent und beschreiben beide die Verb Attachment "book through Houston" Lesart.

```
function CKY-PARSE(words, grammar) returns table
```

```
for j ← from 1 to LENGTH(words) do
```

```
  table[j - 1, j] ← {A | A → words[j] ∈ grammar }
```

```
  for i ← from j - 2 downto 0 do
```

```
    for k ← i + 1 to j - 1 do
```

```
      table[i, j] ← table[i, j] ∪
```

```
        {A | A → BC ∈ grammar,
```

```
          B ∈ table[i, k],
```

```
          C ∈ table[k, j] }
```

**Figure 13.10** The CKY algorithm

- Die Chart berechnet alle möglichen Parsebäume
- Erlaubt das Packen von Ambiguitäten. Chart stellt lokale und globale Ambiguitäten dar.
- Phrasen werden angezeigt (nicht-leere Zellen)
- Laufzeit  $O(|R| \cdot n^3)$ , wenn  $n$  Satzlänge und  $R$  die Anzahl der Regeln. Grund: Teilbäume werden nachgeschaut und nicht neu geparkt
- Nachteil: braucht CNF

Jurafsky und Martin, 3rd Edition, Kapitel 13.1, 13.2

CKY-Demonstrator <http://lxmls.it.pt/2015/cky.html> (leider mit anderer Chartdarstellung)