

Syntactic Parsing



Background

2

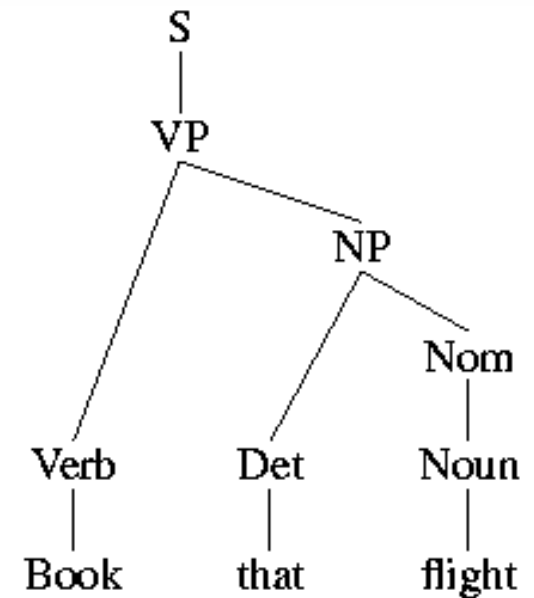
- Syntactic parsing
 - The task of recognizing a sentence and assigning a syntactic structure to it
- Since CFGs are a declarative formalism, they do not specify how the parse tree for a given sentence should be computed.
- Parse trees are useful in applications such as
 - Grammar checking
 - Semantic analysis
 - Machine translation
 - Question answering
 - Information extraction

Parsing as Search

3

- The parser can be viewed as searching through the space of all possible parse trees to find the correct parse tree for the sentence.
- ***How can we use the grammar to produce the parse tree?***

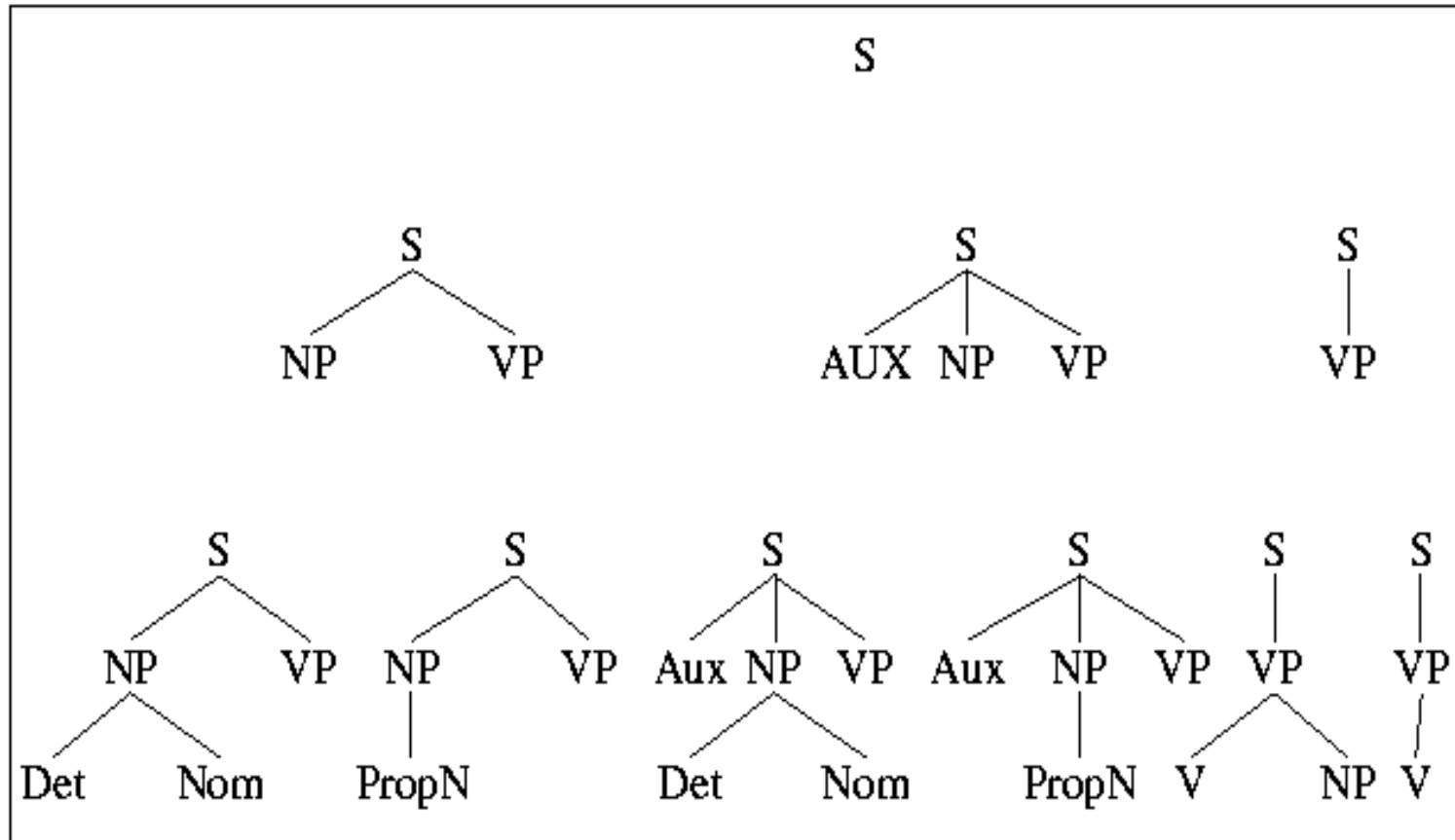
$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid money$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid prefer$
$NP \rightarrow Det Nominal$	$Aux \rightarrow does$
$Nominal \rightarrow Noun$	
$Nominal \rightarrow Noun Nominal$	$Prep \rightarrow from \mid to \mid on$
$NP \rightarrow Proper-Noun$	$Proper-Noun \rightarrow Houston \mid TWA$
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	$Nominal \rightarrow Nominal PP$



Parsing as Search

4

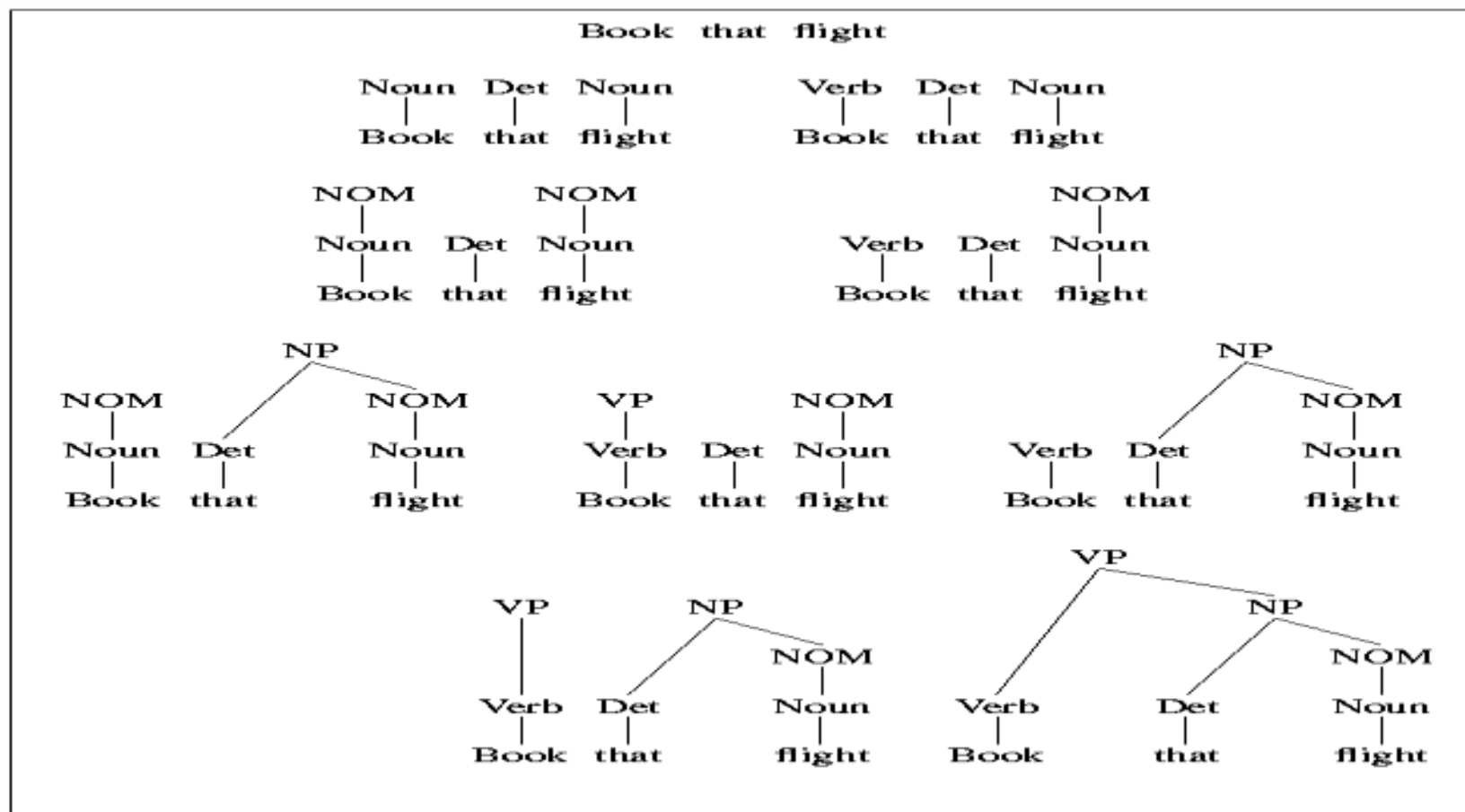
- Top-down parsing



Parsing as Search

5

- Bottom-up parsing



Parsing as Search

6

● Comparisons

- The top-down strategy never wastes time exploring trees that cannot result in an S .
- The bottom-up strategy, by contrast, trees that have no hope to leading to an S , or fitting in with any of their neighbors, are generated with wild abandon.
- Spend considerable effort on S trees that are not consistent with the input.

Problems with the Basic Top-Down Parser

7

- Problems with the top-down parser
 - Left-recursion
 - Ambiguity
 - Inefficiency reparsing of subtrees
- Introducing the Earley and CYK algorithm

Ambiguity

8

- Common structural ambiguity

- Attachment ambiguity

- ✦ A sentence has an attachment ambiguity if a particular constituent can be attached to the parse tree at **more than one place**.
 - ✦ Various kinds of **adverbial phrases** are also subject to this kind of ambiguity

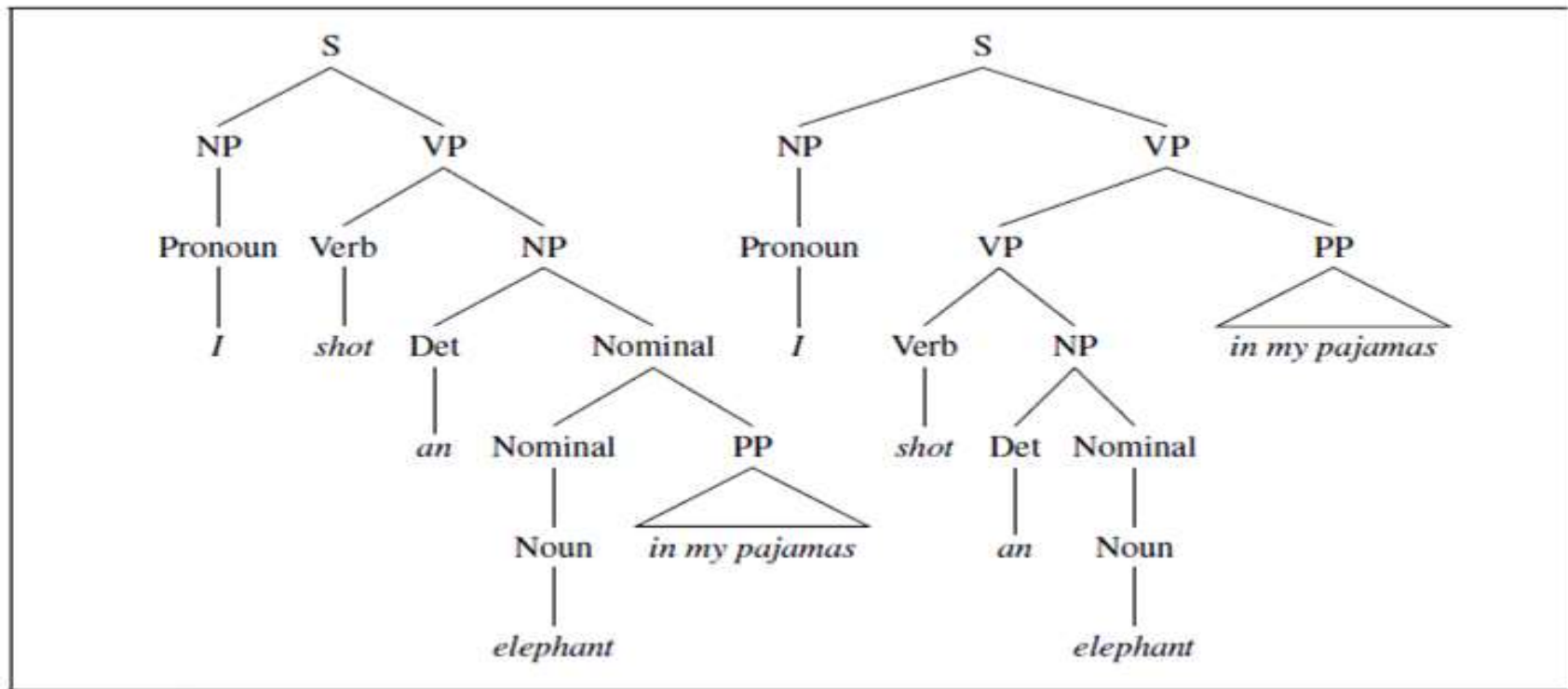
- Coordination ambiguity

- ✦ Different sets of phrases can be conjoined by a conjunction like **and**.
 - ✦ For example, the phrase old men and women can be bracketed as [old [men and women]], referring to old men and old women, or as [old men] and [women], in which case it is only the men who are old.

Ambiguity

9

PP attachment ambiguity



Ambiguous because the phrase *in my pajamas* can be part of the NP headed by *elephant* or a part of the verb phrase headed by *shot*.

Ambiguity

10

We saw the Eiffel Tower flying to Paris.

- The gerundive-VP *flying to Paris* can be
 - part of a gerundive sentence, or
 - an adjunct modifying the VP

The CYK Algorithm

11

● *The membership problem*

○ Problem:

✦ Given a context-free grammar \mathbf{G} and a string \mathbf{w}

□ $\mathbf{G} = (V, \Sigma, P, S)$ where

- V finite set of variables
- Σ (the alphabet) finite set of terminal symbols
- P finite set of rules
- S start symbol (distinguished element of V)
- V and Σ are assumed to be disjoint

□ \mathbf{G} is used to generate the string of a language

○ Question:

✦ Is \mathbf{w} in $\mathbf{L(G)}$?

The CYK Algorithm

12

- J. Cocke
- D. Younger,
- T. Kasami
 - Independently developed an algorithm to answer this question.

Dynamic Programming

13

- DP search methods fill tables with partial results and thereby
 - Avoid doing avoidable repeated work
 - Solve in polynomial time
 - Efficiently store ambiguous structures with shared sub-parts.
- We'll cover two approaches that roughly correspond to bottom-up and top-down approaches.
 - CKY
 - Earley

The CYK Algorithm Basics

14

- The Structure of the rules in a Chomsky Normal Form grammar
- Uses a “dynamic programming” or “table-filling algorithm”
- Based on bottom-up parsing and requires first normalizing the grammar.
- **Earley parser** is based on top-down parsing and does not require normalizing grammar but is more complex.
- More generally, **chart parsers** retain completed phrases in a chart and can combine top-down and bottom-up search.

Chomsky Normal Form

15

- *Normal Form* is described by a set of conditions that each rule in the grammar must satisfy
 - Context-free grammar is in CNF if each rule has one of the following forms:
 - $A \rightarrow BC$ at most 2 symbols on right side
 - $A \rightarrow a$, or terminal symbol
 - $S \rightarrow \lambda$ null string
- where $B, C \in V - \{S\}$

Construct a Triangular Table

16

- Each row corresponds to one length of substrings
 - Bottom Row – Strings of length 1
 - Second from Bottom Row – Strings of length 2
 - .
 - .
 - Top Row – string 'w'

Construct a Triangular Table

17

- $X_{i,i}$ is the set of variables A such that $A \sqsubseteq w_i$ is a production of G
- Compare at most n pairs of previously computed sets:
 $(X_{i,i}, X_{i+1,j}), (X_{i,i+1}, X_{i+2,j}) \dots (X_{i,j-1}, X_{j,j})$

Construct a Triangular Table

18

$X_{1,5}$				
$X_{1,4}$	$X_{2,5}$			
$X_{1,3}$	$X_{2,4}$	$X_{3,5}$		
$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,5}$	
$X_{1,1}$	$X_{2,2}$	$X_{3,3}$	$X_{4,4}$	$X_{5,5}$
w1	w2	w3	w4	w5

Table for string 'w' that has length 5

Construct a Triangular Table

19

$X_{1,5}$				
$X_{1,4}$	$X_{2,5}$			
$X_{1,3}$	$X_{2,4}$	$X_{3,5}$		
$X_{1,2}$	$X_{2,3}$	$X_{3,4}$	$X_{4,5}$	
$X_{1,1}$	$X_{2,2}$	$X_{3,3}$	$X_{4,4}$	$X_{5,5}$
w1	w2	w3	w4	w5

Looking for pairs to compare

Example CYK Algorithm

20

- Show the CYK Algorithm with the following example:
 - CNF grammar G
 - ✦ $S \rightarrow AB \mid BC$
 - ✦ $A \rightarrow BA \mid a$
 - ✦ $B \rightarrow CC \mid b$
 - ✦ $C \rightarrow AB \mid a$
 - w is baaba
 - Question Is **baaba** in $L(G)$?

Constructing The Triangular Table

21

					<div>S □ AB BC</div> <div>A □ BA a</div> <div>B □ CC b</div> <div>C □ AB a</div>
{B}	{A, C}	{A, C}	{B}	{A, C}	
b	a	a	b	a	

Calculating the Bottom ROW

Constructing The Triangular Table

22

- $X_{1,2} = (X_{i,i}, X_{i+1,j}) = (X_{1,1}, X_{2,2})$
- $\square \quad \{B\}\{A,C\} = \{BA, BC\}$
- Steps:
 - Look for production rules to generate BA or BC
 - There are two: S and A
 - $X_{1,2} = \{S, A\}$

S	\square	AB		BC
A	\square	BA		a
B	\square	CC		b
C	\square	AB		a

Constructing The Triangular Table

23

{S, A}				
{B}	{A, C}	{A, C}	{B}	{A, C}
b	a	a	b	a

Constructing The Triangular Table

24

- $X_{2,3} = (X_{i,i}, X_{i+1,j}) = (X_{2,2}, X_{3,3})$
- $\square \quad \{A, C\}\{A, C\} = \{AA, AC, CA, CC\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There is one: B
 - $X_{2,3} = \{B\}$

S	\square	AB		BC
A	\square	BA		a
B	\square	CC		b
C	\square	AB		a

Constructing The Triangular Table

25

{S, A}	{B}			
{B}	{A, C}	{A, C}	{B}	{A, C}
b	a	a	b	a

Constructing The Triangular Table

26

- $X_{3,4} = (X_{i,i}, X_{i+1,j}) = (X_{3,3}, X_{4,4})$
- $\square \quad \{A, C\}\{B\} = \{AB, CB\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There are two: S and C
 - $X_{3,4} = \{S, C\}$

S	\square	AB		BC
A	\square	BA		a
B	\square	CC		b
C	\square	AB		a

Constructing The Triangular Table

27

{S, A}	{B}	{S, C}		
{B}	{A, C}	{A, C}	{B}	{A, C}
b	a	a	b	a

Constructing The Triangular Table

28

- $X_{4,5} = (X_{i,i}, X_{i+1,j}) = (X_{4,4}, X_{5,5})$
- $\square \quad \{B\}\{A, C\} = \{BA, BC\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There are two: S and A
 - $X_{4,5} = \{S, A\}$

S \square	AB BC
A \square	BA a
B \square	CC b
C \square	AB a

Constructing The Triangular Table

29

{S, A}	{B}	{S, C}	{S, A}	
{B}	{A, C}	{A, C}	{B}	{A, C}
b	a	a	b	a

Constructing The Triangular Table

30

- $X_{1,3} = (X_{i,i}, X_{i+1,j}) (X_{i,i+1}, X_{i+2,j})$
 $= (X_{1,1}, X_{2,3}), (X_{1,2}, X_{3,3})$
- $\square \quad \{B\}\{B\} \cup \{S, A\}\{A, C\} = \{BB, SA, SC, AA, AC\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There are NONE: S and A
 - $X_{1,3} = \emptyset$
 - **no elements in this set (empty set)**

S	\square	AB		BC
A	\square	BA		a
B	\square	CC		b
C	\square	AB		a

Constructing The Triangular Table

31

\emptyset				
$\{S, A\}$	$\{B\}$	$\{S, C\}$	$\{S, A\}$	
$\{B\}$	$\{A, C\}$	$\{A, C\}$	$\{B\}$	$\{A, C\}$
b	a	a	b	a

Constructing The Triangular Table

32

- $X_{2,4} = (X_{i,i}, X_{i+1,j}) (X_{i,i+1}, X_{i+2,j})$
 $= (X_{2,2}, X_{3,4}), (X_{2,3}, X_{4,4})$
- $\square \quad \{A, C\}\{S, C\} \cup \{B\}\{B\} = \{AS, AC, CS, CC, BB\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There is one: B
 - $X_{2,4} = \{B\}$

S	\square	AB BC
A	\square	BA a
B	\square	CC b
C	\square	AB a

Constructing The Triangular Table

33

\emptyset	$\{B\}$			
$\{S, A\}$	$\{B\}$	$\{S, C\}$	$\{S, A\}$	
$\{B\}$	$\{A, C\}$	$\{A, C\}$	$\{B\}$	$\{A, C\}$
b	a	a	b	a

Constructing The Triangular Table

34

- $X_{3,5} = (X_{i,i}, X_{i+1,j}) (X_{i,i+1}, X_{i+2,j})$
 $= (X_{3,3}, X_{4,5}), (X_{3,4}, X_{5,5})$
- $\square \{A,C\}\{S,A\} \cup \{S,C\}\{A,C\}$
 $= \{AS, AA, CS, CA, SA, SC, CA, CC\} = Y$
- Steps:
 - Look for production rules to generate Y
 - There is one: B
 - $X_{3,5} = \{B\}$

S	\square	AB		BC
A	\square	BA		a
B	\square	CC		b
C	\square	AB		a

Constructing The Triangular Table

35

\emptyset	$\{B\}$	$\{B\}$		
$\{S, A\}$	$\{B\}$	$\{S, C\}$	$\{S, A\}$	
$\{B\}$	$\{A, C\}$	$\{A, C\}$	$\{B\}$	$\{A, C\}$
b	a	a	b	a

Final Triangular Table

36

$\{S, A, C\}$	<div> <div>□</div> <div>$X_1,$</div> <div>5</div> </div>			
\emptyset	$\{S, A, C\}$			
\emptyset	$\{B\}$	$\{B\}$		
$\{S, A\}$	$\{B\}$	$\{S, C\}$	$\{S, A\}$	
$\{B\}$	$\{A, C\}$	$\{A, C\}$	$\{B\}$	$\{A, C\}$
b	a	a	b	a

- Table for string 'w' that has length 5
- The algorithm populates the triangular table

Example (Result)

37

● Is baaba in $L(G)$?

Yes

We can see the S in the set X_{1n} where 'n' = 5

We can see the table

the cell $X_{15} = (S, A, C)$ then

if $S \in X_{15}$ then baaba $\in L(G)$

CYK Algorithm



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

for $j \leftarrow$ **from** 1 **to** LENGTH(*words*) **do**

$table[j - 1, j] \leftarrow \{A \mid A \rightarrow words[j] \in grammar\}$

for $i \leftarrow$ **from** $j - 2$ **downto** 0 **do**

for $k \leftarrow i + 1$ **to** $j - 1$ **do**

$table[i, j] \leftarrow table[i, j] \cup$

$\{A \mid A \rightarrow BC \in grammar,$

$B \in table[i, k],$

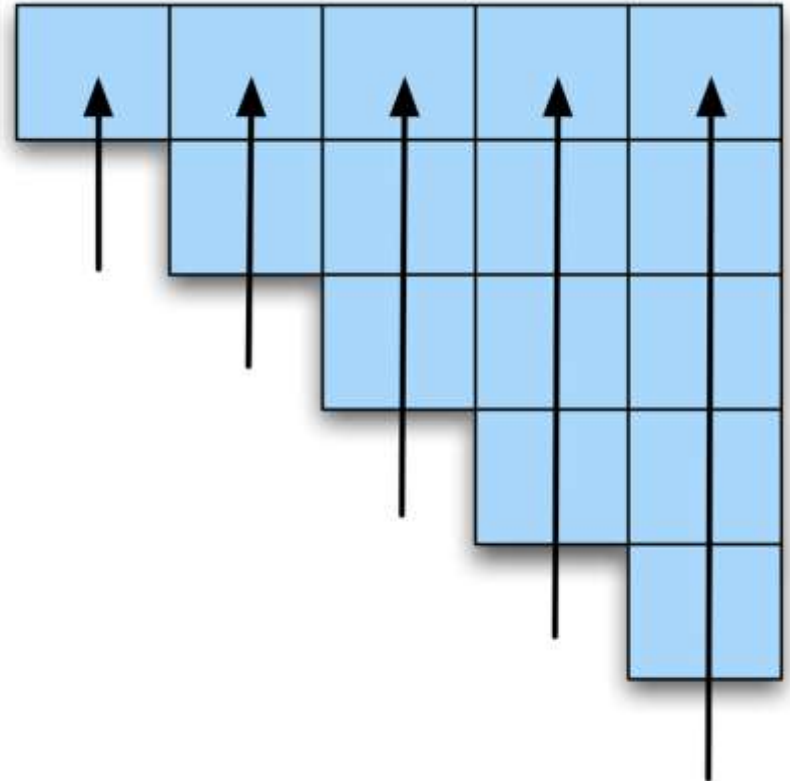
$C \in table[k, j]\}$

Example



Book the flight through Houston

S, VP, Verb Nominal, Noun [0,1]		S,VP,X2 [0,3]		S,VP,X2 [0,5]
	Det [1,2]	NP [1,3]		NP [1,5]
		Nominal, Noun [2,3]		Nominal [2,5]
			Prep [3,4]	PP [3,5]
				NP, Proper- Noun [4,5]



The Earley Algorithm



- Chart entries represent three type of constituents
 - predicted constituents (top-down predictions)
 - Scan in-progress constituents (we're in the midst of ...)
 - completed constituents (we've found ...)
- Progress in parse represented by **Dotted Rules**
 - Position of • indicates type of constituent
 - $_0$ Book $_1$ that $_2$ flight $_3$
 - S --> • VP, [0,0] (predicting VP)
 - NP --> Det • Nom, [1,2] (finding NP)
 - VP --> V NP •, [0,3] (found VP)
 - [x,y] tells us where the state begins (x) and where the dot lies (y) wrt the input

The Earley Algorithm



o Book₁ that₂ flight₃

S --> • VP, [0,0]

- o First 0 means S constituent begins at the start of the input
- o Second 0 means the dot here too

NP --> Det • Nom, [1,2]

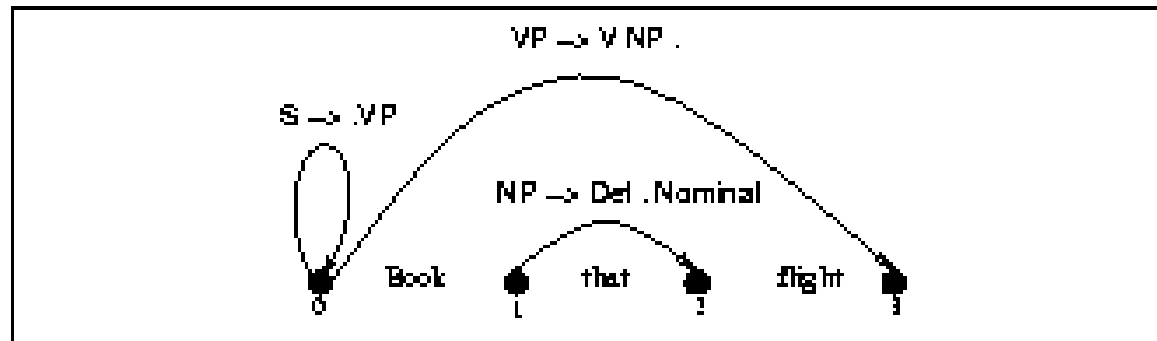
- o the NP begins at position 1
- o the dot is at position 2
- o Det has been successfully parsed
- o Nom predicted next

The Earley Algorithm



$VP \rightarrow V NP \bullet, [0,3]$

- Successful VP parse of entire input
- Graphical representation



Successful Parse



- Final answer is found by looking at last entry in chart
- If entry resembles $S \rightarrow \alpha \bullet [o, N]$ then input parsed successfully
- But ...
 - note that chart will also contain a record of all possible parses of input string, given the grammar
 - not just the successful one(s)

Parsing Procedure for the Earley Algorithm



- Move through each set of states in order, applying one of three operators to each state
 - **predictor**: add top-down predictions
 - **scanner**: read input and add corresponding state
 - **completer**: move dot to right when new constituent found
- No backtracking and no states removed: keep complete history of parse

Predictor



- New states represent top-down expectations
- Applied when **non part-of-speech non-terminals** are to the right of a dot

$S \rightarrow \bullet VP [0,0]$

- Adds new states to end of **current** chart
 - One new state for each expansion of the non-terminal in the grammar

$VP \rightarrow \bullet V [0,0]$

$VP \rightarrow \bullet V NP [0,0]$

Scanner



- New states for predicted part of speech.
- Applicable when part of speech is to the right of a dot
VP --> • V NP [0,0]
- Looks at current word in input
- If match, adds state(s) to **next** chart
VP --> V • NP [0,1]
 - i.e., we've **found** a piece of this constituent!

Completer



- We've found a constituent, so tell everyone waiting for this
- Applied when dot has reached right end of rule
NP --> Det Nom • [1,3]
- Find all states w/dot at 1 and expecting an NP
VP --> V • NP [0,1]
- Adds new (completed) state(s) to **current** chart
VP --> V NP • [0,3]

CFG for Fragment of English



S <input type="checkbox"/> NP VP	Det <input type="checkbox"/> that this
S <input type="checkbox"/> Aux NP VP	a <input type="checkbox"/> N <input type="checkbox"/> book flight meal
S <input type="checkbox"/> VP	money <input type="checkbox"/> V <input type="checkbox"/> book include
NP <input type="checkbox"/> Det Nom	prefer <input type="checkbox"/> Aux <input type="checkbox"/>
Nom <input type="checkbox"/>	does
Nom <input type="checkbox"/> N	Prep <input type="checkbox"/> from to on
Nom <input type="checkbox"/> NP	PropN <input type="checkbox"/> Houston TWA
<input type="checkbox"/> VP <input type="checkbox"/> PropN	Nom <input type="checkbox"/> Nom
<input type="checkbox"/> VP <input type="checkbox"/> V	PP <input type="checkbox"/> Prep
NP	NP

Book that flight (Chart [o])



- Seed chart with top-down predictions for S from [grammar](#)

$\gamma \rightarrow \bullet S$	[0,0]	Dummy start state
$S \rightarrow \bullet NP VP$	[0,0]	Predictor
$S \rightarrow \bullet Aux NP VP$	[0,0]	Predictor
$S \rightarrow \bullet VP$	[0,0]	Predictor
$NP \rightarrow \bullet Det Nom$	[0,0]	Predictor
$NP \rightarrow \bullet PropN$	[0,0]	Predictor
$VP \rightarrow \bullet V$	[0,0]	Predictor
$VP \rightarrow \bullet V NP$	[0,0]	Predictor

Parsing by Earley Algorithm



- When dummy start state is processed, it's passed to Predictor, which produces states representing every possible expansion of S, and adds these and every expansion of the left corners of these trees to bottom of Chart[0]
- When $VP \rightarrow \bullet V, [0,0]$ is reached, Scanner called, which consults first word of input, Book, and adds first state to Chart[1], $VP \rightarrow Book \bullet, [0,0]$

Chart[1]



$V \rightarrow \text{book} \bullet$	[0,1]	Scanner
$VP \rightarrow V \bullet$	[0,1]	Completer
$VP \rightarrow V \bullet NP$	[0,1]	Completer
$S \rightarrow VP \bullet$	[0,1]	Completer
$NP \rightarrow \bullet \text{Det Nom}$	[1,1]	Predictor
$NP \rightarrow \bullet \text{PropN}$	[1,1]	Predictor

$V \rightarrow \text{book} \bullet$ passed to Completer, which finds 2 states in Chart[0] whose left corner is V and adds them to Chart[1], moving dots to right

Parsing by Earley Algorithm

Chart[0]		
$\gamma \rightarrow \bullet S$	[0,0]	Dummy start state
$S \rightarrow \bullet NP VP$	[0,0]	Predictor
$NP \rightarrow \bullet Det NOMINAL$	[0,0]	Predictor
$NP \rightarrow \bullet Proper-Noun$	[0,0]	Predictor
$S \rightarrow \bullet Aux NP VP$	[0,0]	Predictor
$S \rightarrow \bullet VP$	[0,0]	Predictor
$VP \rightarrow \bullet Verb$	[0,0]	Predictor
$VP \rightarrow \bullet Verb NP$	[0,0]	Predictor

Chart[1]		
$Verb \rightarrow book \bullet$	[0,1]	Scanner
$VP \rightarrow Verb \bullet$	[0,1]	Completer
$S \rightarrow VP \bullet$	[0,1]	Completer
$VP \rightarrow Verb \bullet NP$	[0,1]	Completer
$NP \rightarrow \bullet Det NOMINAL$	[1,1]	Predictor
$NP \rightarrow \bullet Proper-Noun$	[1,1]	Predictor

Chart[2]		
$Det \rightarrow that \bullet$	[1,2]	Scanner
$NP \rightarrow Det \bullet NOMINAL$	[1,2]	Completer
$NOMINAL \rightarrow \bullet Noun$	[2,2]	Predictor
$NOMINAL \rightarrow \bullet Noun NOMINAL$	[2,2]	Predictor

Chart[3]		
$Noun \rightarrow flight \bullet$	[2,3]	Scanner
$NOMINAL \rightarrow Noun \bullet$	[2,3]	Completer
$NOMINAL \rightarrow Noun \bullet NOMINAL$	[2,3]	Completer
$NP \rightarrow Det NOMINAL \bullet$	[1,3]	Completer
$VP \rightarrow Verb NP \bullet$	[0,3]	Completer
$S \rightarrow VP \bullet$	[0,3]	Completer
$NOMINAL \rightarrow \bullet Noun$	[3,3]	Predictor
$NOMINAL \rightarrow \bullet Noun NOMINAL$	[3,3]	Predictor

The Earley Algorithm

function EARLEY-PARSE(*words, grammar*) **returns** *chart*

ENQUEUE($(\gamma \rightarrow \bullet S, [0,0])$, *chart*[0])

for $i \leftarrow$ **from** 0 **to** LENGTH(*words*) **do**

for each *state* **in** *chart*[*i*] **do**

if INCOMPLETE?(*state*) **and**

 NEXT-CAT(*state*) is not a part of speech **then**

 PREDICTOR(*state*)

elseif INCOMPLETE?(*state*) **and**

 NEXT-CAT(*state*) is a part of speech **then**

 SCANNER(*state*)

else

 COMPLETER(*state*)

end

end

return(*chart*)

procedure PREDICTOR($(A \rightarrow \alpha \bullet B \beta, [i, j])$)

for each $(B \rightarrow \gamma)$ **in** GRAMMAR-RULES-FOR(*B, grammar*) **do**

 ENQUEUE($(B \rightarrow \bullet \gamma, [j, j])$, *chart*[*j*])

end

procedure SCANNER($(A \rightarrow \alpha \bullet B \beta, [i, j])$)

if $B \in$ PARTS-OF-SPEECH(*word*[*j*]) **then**

 ENQUEUE($(B \rightarrow \text{word}[j], [j, j + 1])$, *chart*[*j* + 1])

procedure COMPLETER($(B \rightarrow \gamma \bullet, [j, k])$)

for each $(A \rightarrow \alpha \bullet B \beta, [i, j])$ **in** *chart*[*j*] **do**

 ENQUEUE($(A \rightarrow \alpha B \bullet \beta, [i, k])$, *chart*[*k*])

end

procedure ENQUEUE(*state, chart-entry*)

if *state* is not already in *chart-entry* **then**

 PUSH(*state, chart-entry*)

end

THANK YOU