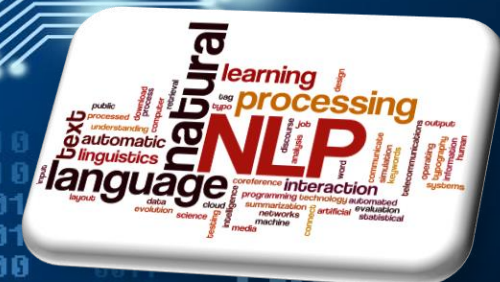




ΣΥΛΟΓΗ ΔΕΛΤΙΩΝ

Natural Language Processing

Dr. Sandra Wahid



Syntactic Parsing

- It is the task of assigning a **syntactic structure** to a sentence.
- It is useful in applications such as:
 - Grammar checking: sentence that cannot be parsed may have grammatical errors (or at least be hard to read).
 - Semantic analysis
 - Machine translation
 - Question answering: for example to answer the question: “Which flights to Denver depart before the Seattle flight?”
 - we’ll need to know that the questioner wants a list of flights going to Denver, not flights going to Seattle.

Two views of syntactic structures

Constituency Parsing

Dependency Parsing

Constituency

- Syntactic constituency is the idea that **groups of words** can behave as single units, or constituents.
- Example: **noun phrase** “a sequence of words surrounding at least one noun” form constituents. Why??

Harry the Horse
the Broadway coppers
they

a high-class spot such as Mindy's
the reason he comes into the Hot Box
three parties from Brooklyn

All these are examples of noun phrases

they can all appear in **similar syntactic environments**, for example, before a verb.

three parties from Brooklyn *arrive...*
a high-class spot such as Mindy's *attracts...*
the Broadway coppers *love...*
they *sit*

Context-Free Grammars

- The most widely used formal system for modeling constituent structure in English and other natural languages is the Context-Free Grammar (CFG).
- Also called Phrase-Structure Grammars.
- A context-free grammar consists of a set of **rules or productions**, each of which expresses the ways that symbols of the language can be grouped and ordered together, and a **lexicon** of words and symbols.
- Example: the following productions express that an NP (or noun phrase) can be composed of either a ProperNoun or a determiner (Det) followed by a Nominal, a Nominal in turn can consist of one or more Nouns.

$$\begin{aligned} NP &\rightarrow Det\ Nominal \\ NP &\rightarrow ProperNoun \\ Nominal &\rightarrow Noun \mid Nominal\ Noun \end{aligned}$$

- Context-free rules can be hierarchically embedded, so we can combine the previous rules with others that express facts about the lexicon:

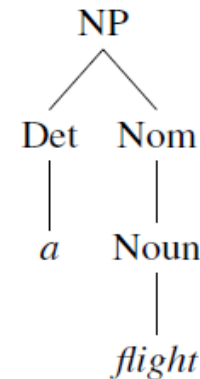
$$\begin{aligned} Det &\rightarrow a \\ Det &\rightarrow the \\ Noun &\rightarrow flight \end{aligned}$$

Context-Free Grammars

- The symbols that are used in a CFG are divided into two classes:
 - **Terminals:** the symbols that correspond to words in the language, the lexicon is the set of rules that introduce these terminal symbols.
 - **Non-terminals:** the symbols that express abstractions over these terminals.
- Format of context-free rule: $\langle \text{a single non-terminal symbol} \rangle \rightarrow \langle \text{ordered list of one or more terminals and non-terminals} \rangle$
- CFG can be used to generate a set of strings. This sequence of rule expansions is called a **derivation** of the string of words. It is common to represent a derivation by a **parse tree**.

$NP \rightarrow Det\ Nominal$
 $NP \rightarrow ProperNoun$
 $Nominal \rightarrow Noun \mid Nominal\ Noun$

$Det \rightarrow a$
 $Det \rightarrow the$
 $Noun \rightarrow flight$



A parse tree for “a flight”.

Context-Free Grammars

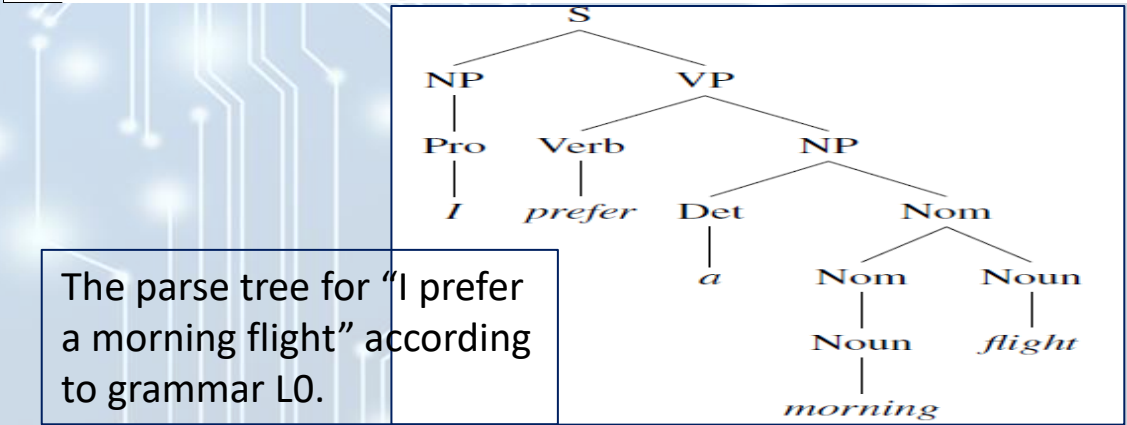
- The formal language defined by a CFG is the set of strings that are derivable from the designated **start symbol**.
- Each grammar must have **one** designated start symbol, which is often called S where S is usually interpreted as the **“sentence” node**.

Grammar Rules	Examples
$S \rightarrow NP VP$ <div>verb phrase</div>	I + want a morning flight
$NP \rightarrow$ Pronoun Proper-Noun Det Nominal Nominal \rightarrow Nominal Noun Noun	I Los Angeles a + flight morning + flight flights
$VP \rightarrow$ Verb Verb NP Verb NP PP Verb PP <div>prepositional phrase</div>	do want + a flight leave + Boston + in the morning leaving + on Thursday
$PP \rightarrow$ Preposition NP	from + Los Angeles

The grammar for \mathcal{L}_0 , with example phrases for each rule.

$Noun \rightarrow$	flights flight breeze trip morning
$Verb \rightarrow$	is prefer like need want fly do
$Adjective \rightarrow$	cheapest non-stop first latest other direct
$Pronoun \rightarrow$	me I you it
$Proper-Noun \rightarrow$	Alaska Baltimore Los Angeles Chicago United American
$Determiner \rightarrow$	the a an this these that
$Preposition \rightarrow$	from to on near in
$Conjunction \rightarrow$	and or but

The lexicon for \mathcal{L}_0 .



Formal Definition of Context-Free Grammar

- A context-free grammar G is defined by four parameters: N , Σ , R , S (technically this is a “4-tuple”).

N a set of **non-terminal symbols** (or **variables**)
 Σ a set of **terminal symbols** (disjoint from N)
 R a set of **rules** or productions, each of the form $A \rightarrow \beta$,
where A is a non-terminal,
 β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$
 S a designated **start symbol** and a member of N

- The following conventions are followed:

Capital letters like A , B , and S

S

Lower-case Greek letters like α , β , and γ

Lower-case Roman letters like u , v , and w

Non-terminals

The start symbol

Strings drawn from $(\Sigma \cup N)^*$

Strings of terminals

Grammar Equivalence and Normal Form

- A formal language is defined as a (possibly infinite) set of strings of words.
- This suggests that we could ask if two grammars are equivalent by asking if they generate the same set of strings. In fact, it is possible to have two distinct context-free grammars generate the same language.
- We usually distinguish **two kinds of grammar equivalence**:
 - **strong equivalence**: two grammars are strongly equivalent if they generate the same set of strings and if they assign the same phrase structure to each sentence (allowing merely for renaming of the non-terminal symbols).
 - **weak equivalence**: two grammars are weakly equivalent if they generate the same set of strings but do not assign the same phrase structure to each sentence.
- It is sometimes useful to have a normal form for grammars, in which each of the productions takes a particular form.

Chomsky Normal Form (CNF)

- A context-free grammar is in CNF if:
 - ϵ -free (no empty rules).
 - each production is either of the form $A \rightarrow BC$ or $A \rightarrow a$.
(each rule either has two non-terminal symbols or one terminal symbol.)
- Any context-free grammar can be converted into a weakly equivalent Chomsky normal form grammar.
 - For example, a rule of the form: $A \rightarrow BCD$ can be converted into the following two CNF rules: $A \rightarrow BX$ and $X \rightarrow CD$
 - This conversion is done to perform efficient parsing.
- Conversion to CNF steps:
 1. Get rid of all ϵ productions.
 2. Get rid of all productions where RHS is one variable.
 3. Replace every production that is too long by shorter productions.
 4. Move all terminals to productions where RHS is one terminal.

Chomsky Normal Form (CNF)

1) Eliminate ϵ Productions:

- Determine the nullable variables (those that generate ϵ)
- Go through all productions, and for each, omit every possible subset of nullable variables. For example, if $P \rightarrow AxB$ with both A and B nullable, add productions $P \rightarrow xB \mid Ax \mid x$.
- After this, delete all productions with empty RHS.

2) Eliminate Variable Unit Productions:

- A unit production is where RHS has only one variable.
- Consider production $A \rightarrow B$. Then for every production $B \rightarrow \alpha$, add the production $A \rightarrow \alpha$.
- Repeat until done (but don't re-create a unit production already deleted).

3) Replace Long Productions by Shorter Ones:

- For example, if have production $A \rightarrow BCD$, then replace it with $A \rightarrow BE$ and $E \rightarrow CD$.

4) Move Terminals to Unit Productions:

- For every terminal on the right of a non-unit production, add a substitute variable.
- For example, replace production $A \rightarrow bC$ with productions $A \rightarrow BC$ and $B \rightarrow b$.

Chomsky Normal Form (CNF)

- Example1:

Consider the CFG:

$$S \rightarrow aXbX$$

$$X \rightarrow aY \mid bY \mid \varepsilon$$

$$Y \rightarrow X \mid c$$

- 1]** The variable X is nullable; and so therefore is Y .
After elimination of ε , we obtain:

$$S \rightarrow aXbX \mid abX \mid aXb \mid ab$$

$$X \rightarrow aY \mid bY \mid a \mid b$$

$$Y \rightarrow X \mid c$$

- 2]** After elimination of the unit production $Y \rightarrow X$,
we obtain:

$$S \rightarrow aXbX \mid abX \mid aXb \mid ab$$

$$X \rightarrow aY \mid bY \mid a \mid b$$

$$Y \rightarrow aY \mid bY \mid a \mid b \mid c$$

Chomsky Normal Form (CNF)

- Example1:

3]&4] Now, break up the RHSs of S ; and replace a by A , b by B and c by C wherever not units:

$$S \rightarrow EF \mid AF \mid EB \mid AB$$

$$X \rightarrow AY \mid BY \mid a \mid b$$

$$Y \rightarrow AY \mid BY \mid a \mid b \mid c$$

$$E \rightarrow AX$$

$$F \rightarrow BX$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$C \rightarrow c$$

Chomsky Normal Form (CNF)

- Example2:

Convert the following CFG into Chomsky Normal Form:

$$\begin{aligned} S &\rightarrow AbA \\ A &\rightarrow Aa \mid \varepsilon \end{aligned}$$

1]
$$\begin{aligned} S &\rightarrow AbA \mid bA \mid Ab \mid b \\ A &\rightarrow Aa \mid a \end{aligned}$$

3] The second step does not apply. After the third step, one has:

$$\begin{aligned} S &\rightarrow TA \mid bA \mid Ab \mid b \\ A &\rightarrow Aa \mid a \\ T &\rightarrow Ab \end{aligned}$$

4]
$$\begin{aligned} S &\rightarrow TA \mid BA \mid AB \mid b \\ A &\rightarrow AC \mid a \\ T &\rightarrow AB \\ B &\rightarrow b \\ C &\rightarrow a \end{aligned}$$

Chomsky Normal Form (CNF)

• Example3:

$S \rightarrow NP VP$	$N \rightarrow people$
$VP \rightarrow V NP$	$N \rightarrow fish$
$VP \rightarrow V NP PP$	$N \rightarrow tanks$
$NP \rightarrow NP NP$	$N \rightarrow rods$
$NP \rightarrow NP PP$	$V \rightarrow people$
$NP \rightarrow N$	$V \rightarrow fish$
$NP \rightarrow e$	$V \rightarrow tanks$
$PP \rightarrow P NP$	$P \rightarrow with$

$S \rightarrow NP VP$	2. Rewrite as:	$N \rightarrow people$
$VP \rightarrow V NP$	$S \rightarrow NP VP$	$N \rightarrow fish$
$VP \rightarrow V NP PP$	$S \rightarrow VP$	$N \rightarrow tanks$
$NP \rightarrow NP NP$	3. Repeat	$N \rightarrow rods$
$NP \rightarrow NP PP$		$V \rightarrow people$
$NP \rightarrow N$		$V \rightarrow fish$
$NP \rightarrow e$	1. Remove empty rules	$V \rightarrow tanks$
$PP \rightarrow P NP$		$P \rightarrow with$

$S \rightarrow NP VP$		$N \rightarrow people$
$S \rightarrow VP$	1. Remove unary rule	$N \rightarrow fish$
$VP \rightarrow V NP$	}	$N \rightarrow tanks$
$VP \rightarrow V$		$S \rightarrow V NP$
$VP \rightarrow V NP PP$		$S \rightarrow V$
$VP \rightarrow V PP$		$S \rightarrow V NP PP$
$NP \rightarrow NP NP$		$N \rightarrow rods$
$NP \rightarrow NP$	2. Rewrite for rules where VP appears on the left	$V \rightarrow people$
$NP \rightarrow NP PP$		$V \rightarrow fish$
$NP \rightarrow PP$		$V \rightarrow tanks$
$NP \rightarrow N$		$P \rightarrow with$
$PP \rightarrow P NP$		
$PP \rightarrow P$		

Chomsky Normal Form (CNF)

• Example3:

$S \rightarrow NP VP$	$N \rightarrow people$	
$VP \rightarrow V NP$	$N \rightarrow fish$	
$S \rightarrow V NP$	$N \rightarrow tanks$	
$VP \rightarrow V$	$N \rightarrow rods$	
$VP \rightarrow V NP PP$	$V \rightarrow people$	$VP \rightarrow people$
$S \rightarrow V NP PP$	$S \rightarrow people$	
$VP \rightarrow V PP$	$V \rightarrow fish$	$VP \rightarrow fish$
$S \rightarrow V PP$	$S \rightarrow fish$	
$NP \rightarrow NP NP$	$V \rightarrow tanks$	$VP \rightarrow tanks$
$NP \rightarrow NP$	$S \rightarrow tanks$	
$NP \rightarrow NP PP$	$P \rightarrow with$	
$NP \rightarrow PP$		
$NP \rightarrow N$		
$PP \rightarrow P NP$		
$PP \rightarrow P$		

Keep removing unaries

$S \rightarrow NP VP$	$N \rightarrow people$	$NP \rightarrow people$
$VP \rightarrow V NP$	$N \rightarrow fish$	$NP \rightarrow fish$
$S \rightarrow V NP$	$N \rightarrow tanks$	$NP \rightarrow tanks$
$VP \rightarrow V NP PP$	$N \rightarrow rods$	$NP \rightarrow rods$
$S \rightarrow V NP PP$	$V \rightarrow people$	
$VP \rightarrow V PP$	$S \rightarrow people$	
$S \rightarrow V PP$	$VP \rightarrow people$	
$NP \rightarrow NP NP$	$V \rightarrow fish$	
$NP \rightarrow NP$	$S \rightarrow fish$	
$NP \rightarrow NP PP$	$VP \rightarrow fish$	
$NP \rightarrow PP$	$V \rightarrow tanks$	
$NP \rightarrow N$	$S \rightarrow tanks$	
$PP \rightarrow P NP$	$VP \rightarrow tanks$	
$PP \rightarrow P$	$P \rightarrow with$	

Keep removing unaries

Chomsky Normal Form (CNF)

- Example3:

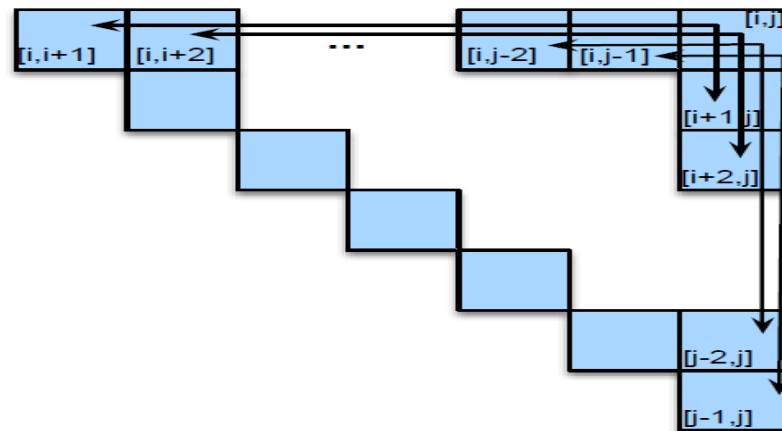
$S \rightarrow NP VP$	Done with unary rules	$NP \rightarrow people$
$VP \rightarrow V NP$		$NP \rightarrow fish$
$S \rightarrow V NP$		$NP \rightarrow tanks$
$VP \rightarrow V NP PP$	$VP \rightarrow V @VP_P$	$NP \rightarrow rods$
$S \rightarrow V NP PP$	$@VP_P \rightarrow NP PP$	$V \rightarrow people$
$VP \rightarrow V PP$		$S \rightarrow people$
$S \rightarrow V PP$	Replace ternary rule with two binary rules by adding a new non-terminal symbol	$VP \rightarrow people$
$NP \rightarrow NP NP$		$V \rightarrow fish$
$NP \rightarrow NP PP$		$S \rightarrow fish$
$NP \rightarrow P NP$		$VP \rightarrow fish$
$PP \rightarrow P NP$		$V \rightarrow tanks$
		$S \rightarrow tanks$
		$VP \rightarrow tanks$
		$P \rightarrow with$
		$PP \rightarrow with$

$S \rightarrow NP VP$	$NP \rightarrow people$
$VP \rightarrow V NP$	$NP \rightarrow fish$
$S \rightarrow V NP$	$NP \rightarrow tanks$
$VP \rightarrow V @VP_V$	$NP \rightarrow rods$
$@VP_V \rightarrow NP PP$	$V \rightarrow people$
$S \rightarrow V @S_V$	$S \rightarrow people$
$@S_V \rightarrow NP PP$	$VP \rightarrow people$
$VP \rightarrow V PP$	$V \rightarrow fish$
$S \rightarrow V PP$	$S \rightarrow fish$
$NP \rightarrow NP NP$	$VP \rightarrow fish$
$NP \rightarrow NP PP$	$V \rightarrow tanks$
$NP \rightarrow P NP$	$S \rightarrow tanks$
$PP \rightarrow P NP$	$VP \rightarrow tanks$
	$P \rightarrow with$
	$PP \rightarrow with$

Constituency Parsing

CKY (Cocke-Kasami-Younger) Parsing Algorithm: A dynamic programming approach to represent all possible parses of the sentence if any

- First need to transform grammar into CNF: produces binary parse trees.
- Bottom-up parsing.
- Dynamic programming: save the results in a table/chart and re-use these results in finding larger constituent.
- A two-dimensional matrix can be used to encode the structure of an entire tree.
- For a sentence of length n , we will work with the upper-triangular portion of an $(n+1) \times (n+1)$ matrix.
- Each cell $[i,j]$ in this matrix contains the set of non-terminals that represent all the constituents that span positions i through j of the input.
- It follows then that the cell that represents the entire input resides in position $[0,n]$ in the matrix.
- We fill the upper-triangular matrix a column at a time working from left to right, with each column filled from bottom to top.
- Each cell $[i,j]$ is filled as follows:



\mathcal{L}_1 in CNF

$S \rightarrow NP VP$

$S \rightarrow XI VP$

$XI \rightarrow Aux NP$

$S \rightarrow book \mid include \mid prefer$

$S \rightarrow Verb NP$

$S \rightarrow X2 PP$

$S \rightarrow Verb PP$

$S \rightarrow VP PP$

$NP \rightarrow I \mid she \mid me$

$NP \rightarrow TWA \mid Houston$

$NP \rightarrow Det Nominal \mid$

$Nominal \rightarrow book \mid flight \mid meal \mid money$

$Nominal \rightarrow Nominal Noun$

$Nominal \rightarrow Nominal PP$

$VP \rightarrow book \mid include \mid prefer$

$VP \rightarrow Verb NP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

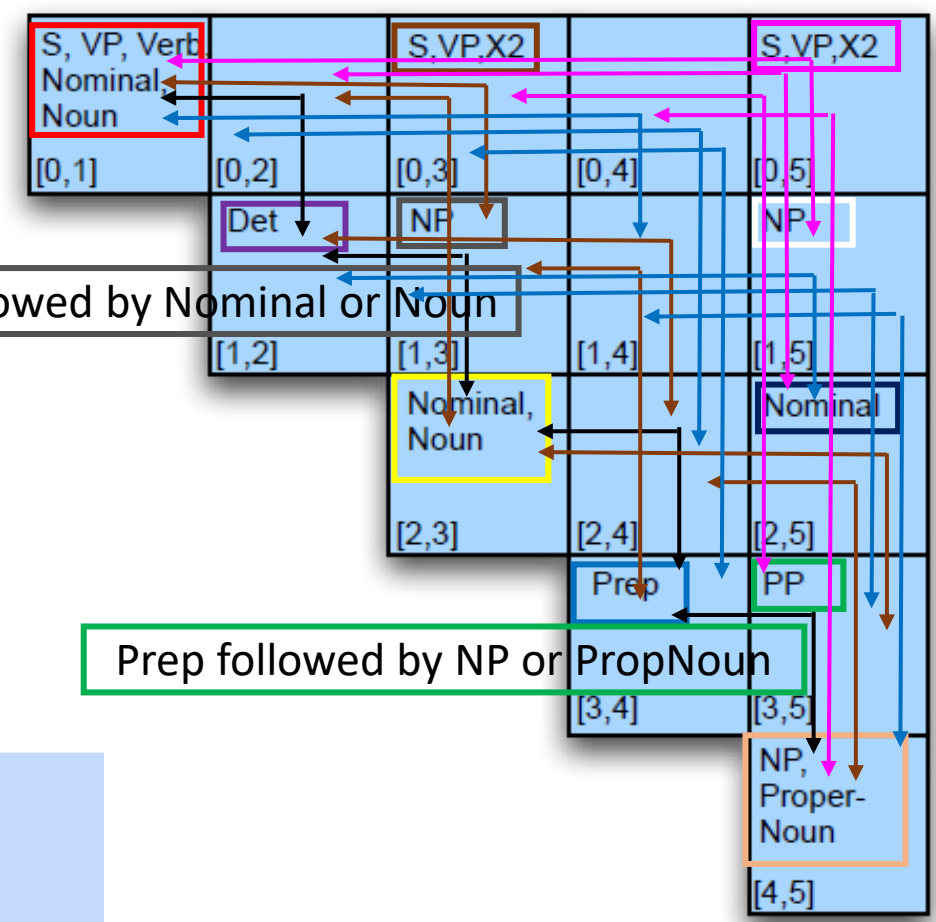
$VP \rightarrow Verb PP$

$VP \rightarrow VP PP$

$PP \rightarrow Preposition NP$

Cell [0,1] spans *Book*
 Cell [3,4] spans *through*
 Cell [2,5] spans *flight through Houston*
 Cell [1,4] spans *the flight through*

0 *Book* 1 *the* 2 *flight* 3 *through* 4 *Houston* 5



Det followed by Nominal or Noun

Prep followed by NP or PropNoun

Lexicon

$Det \rightarrow that \mid this \mid the \mid a$

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

$Verb \rightarrow book \mid include \mid prefer$

$Pronoun \rightarrow I \mid she \mid me$

$Proper-Noun \rightarrow Houston \mid NWA$

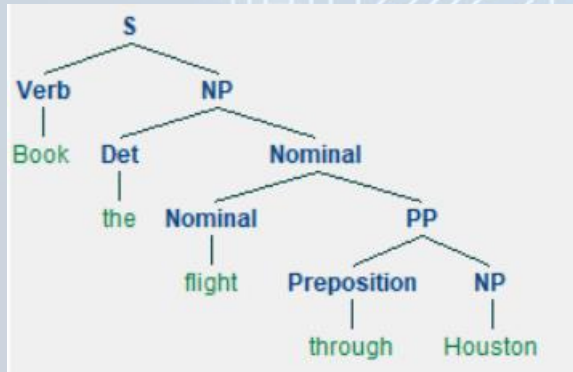
$Aux \rightarrow does$

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

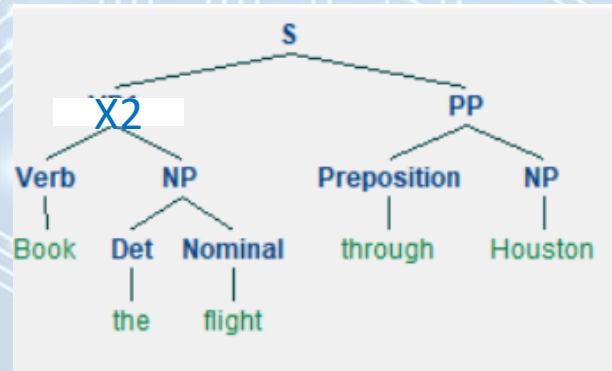
- If the cell $[0,n]$ contains S then a valid parse is available for the sentence.
- To get all possible parses for the sentence:
 - Choose an S from cell $[0,n]$ and then recursively retrieve its component constituents from the table.
 - Repeat for all S's in cell $[0,n]$.

• In the example:

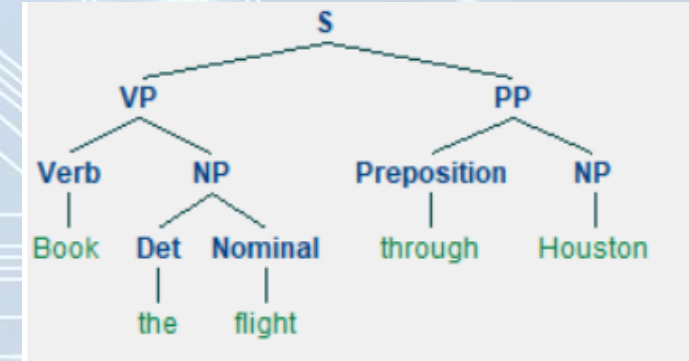
1. $S \rightarrow \text{Verb NP}$



2. $S \rightarrow X_2 \text{ PP}$



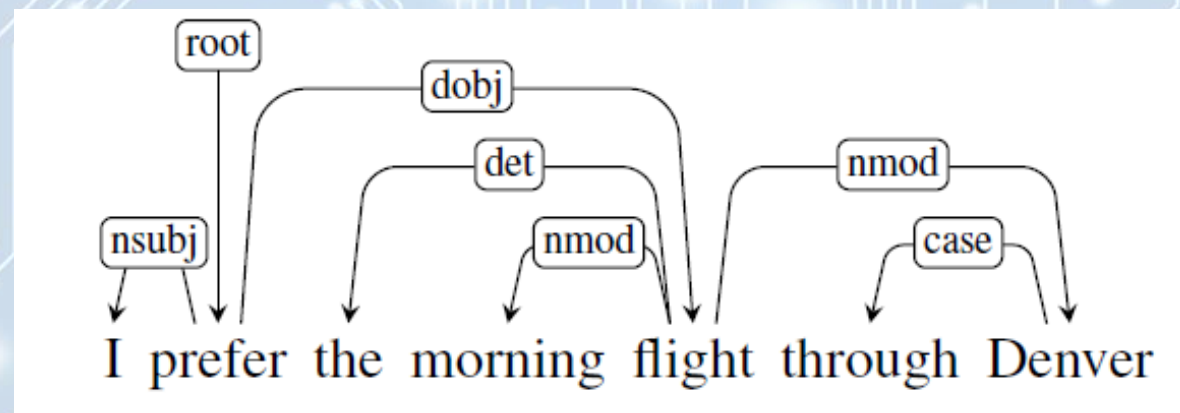
3. $S \rightarrow \text{VP PP}$



- There may be an exponential number of parses for a given sentence so there are augmentations to the method to retrieve only the **best parse**.

Dependency Parsing

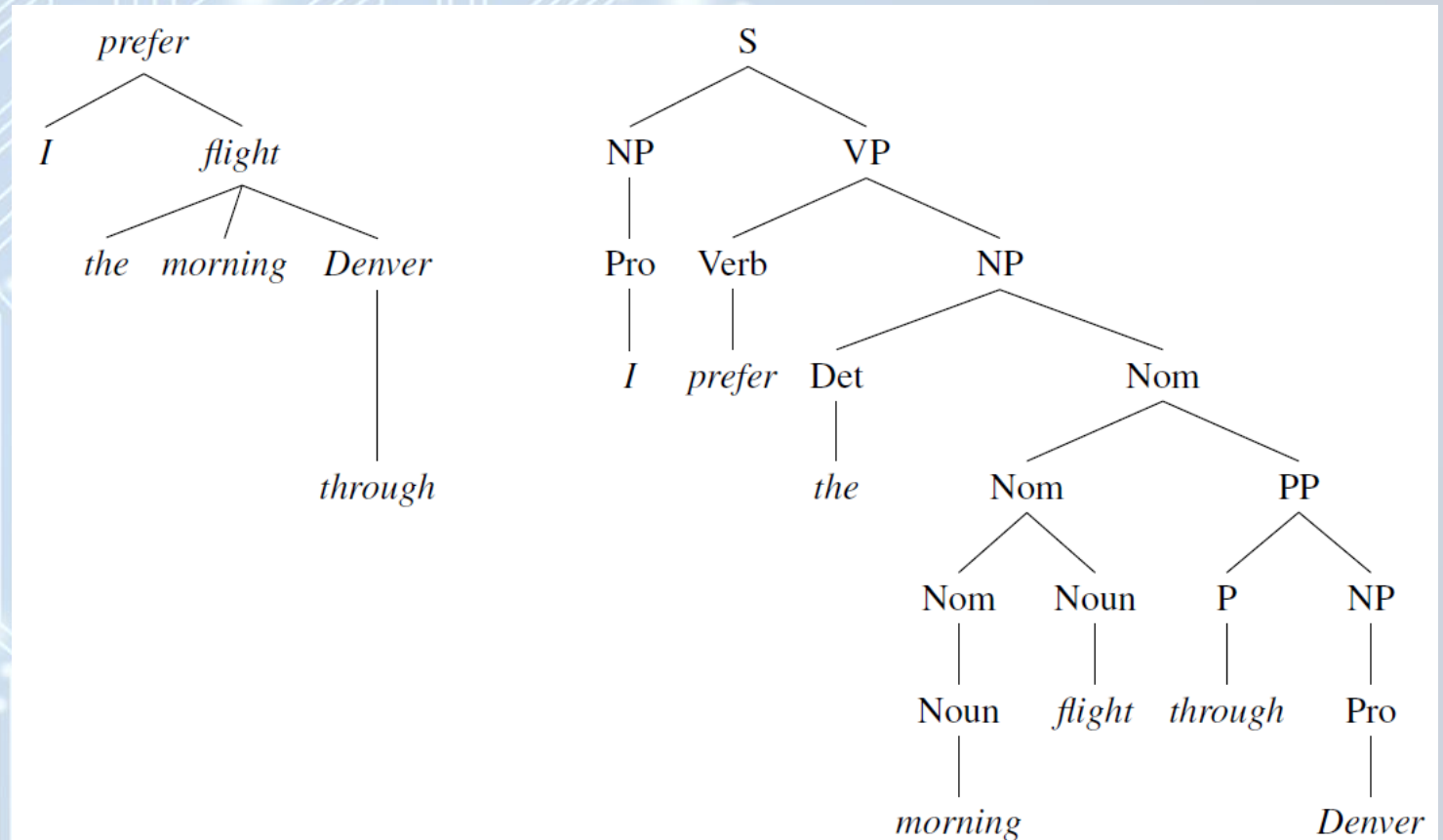
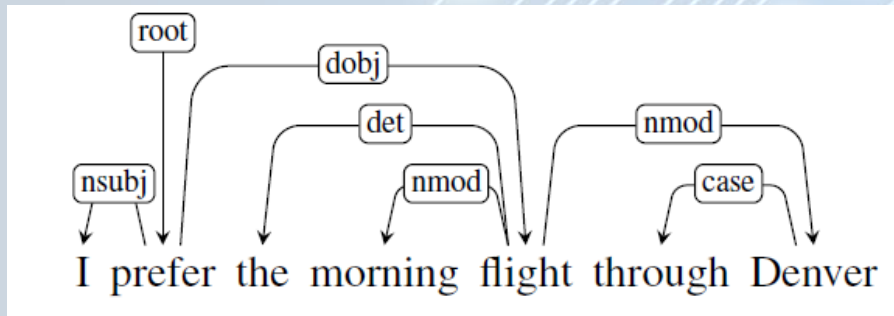
- So far, we have seen context-free grammars and constituent-based representations.
- Another important family of grammar formalisms called **dependency grammars**.
- The syntactic structure of a sentence is described solely in terms of directed **binary grammatical relations** between the words, as in the following **dependency parse**:



- The parse tree has directed, labeled arcs from *heads* to *dependents*.
- This is called **typed dependency structure** because the labels are drawn from a fixed inventory of grammatical relations.
- A root node explicitly marks the root of the tree, the head of the entire structure.

Dependency Parsing

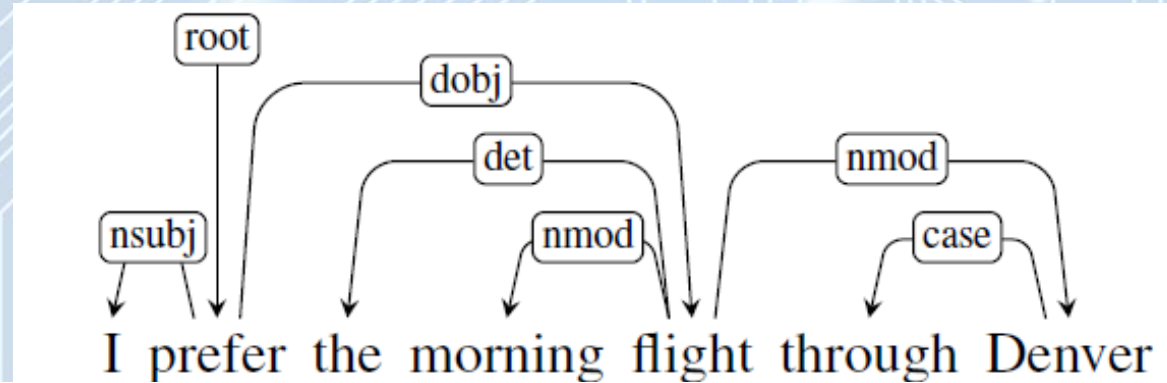
- Dependency and constituent analyses for *I prefer the morning flight through Denver*:



- The arguments to the verb *prefer* are directly linked to it in the dependency structure, while their connection to the main verb is more distant in the phrase-structure tree.
- Similarly, *morning* and *Denver*, modifiers of *flight*, are linked to it directly in the dependency structure.

Dependency Parsing

- The Universal Dependencies (UD) project provides an inventory of dependency relations that are **cross-linguistically applicable**.
- A subset of the UD relations from the following example:



Relation	Description
nsubj	Nominal subject
dobj	Direct object
det	Determiner
nmod	Nominal modifier
case	Prepositions, postpositions and other case markers

- Dependency parsing approaches include: **transition-based** parsing and **graph-based** parsers.



Thank You

0100

0100

0110

0110

1001

1001

0100

0100

1101

0100

0110

0100

0100

0101

0100

0110

0100

0100

0101

1010

0110

1001

0101

0100

1101

0101

1110

1010

0110

1001

0101

0100

0110

1001

0101

0100

1101

0101

1110

0110

1001