

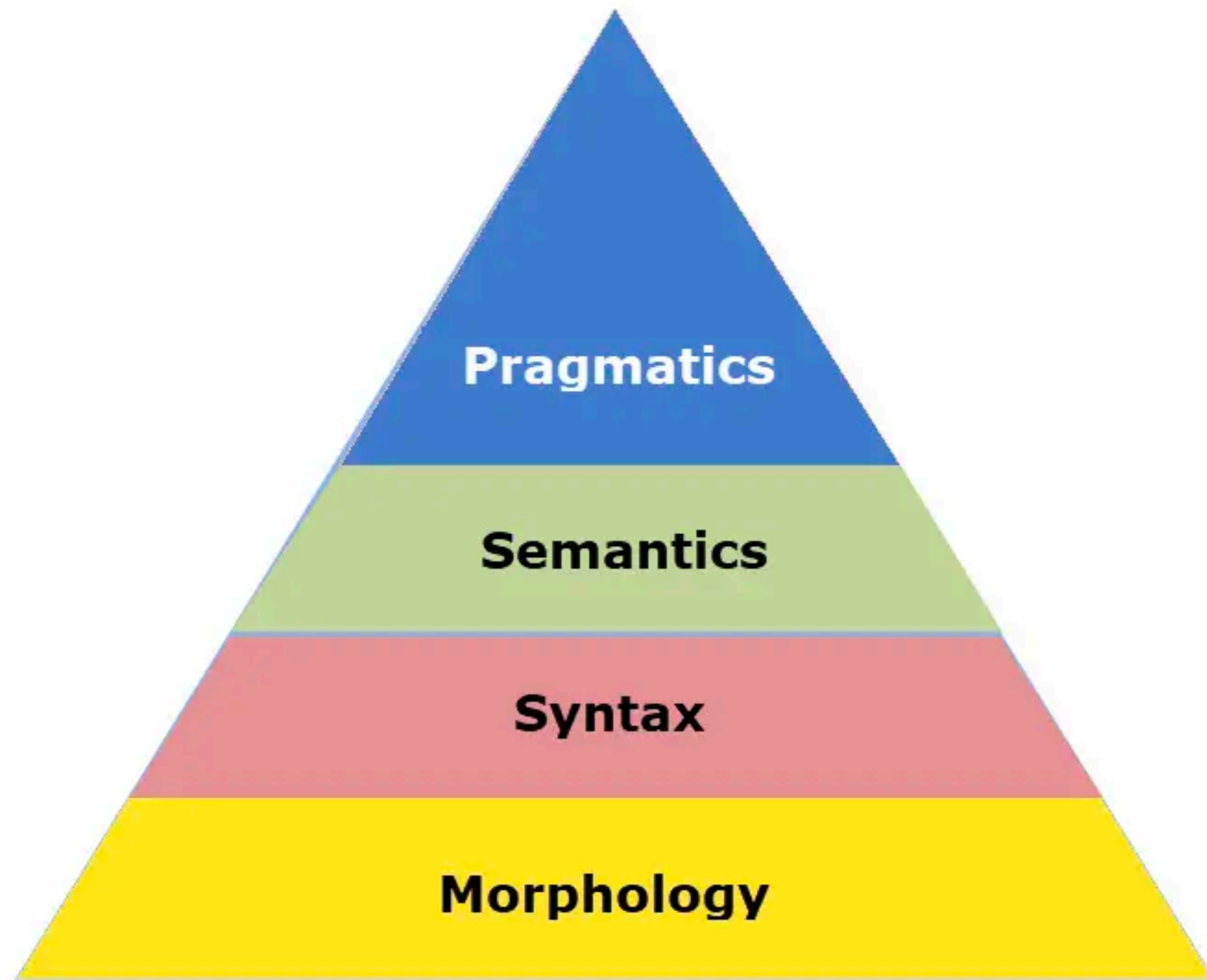
Lecture 8

Syntax - Structure of sentences

Zhizheng Wu

Agenda

- ▶ Recap
- ▶ Concept of syntax and constituency
- ▶ Context-free grammar
- ▶ Cocke-Kasami-Younger (CKY) algorithm



Natural Language Processing Pyramid

Open class ("content") words

Nouns

Proper

Janet
Italy

Common

cat, cats
mango

Verbs

Main

eat
went

Auxiliary

can
had

Adjectives

old green tasty

Adverbs

slowly yesterday

Numbers

122,312
one

Interjections *Ow hello*

... more

Closed class ("function")

Determiners *the some*

Conjunctions *and or*

Pronouns *they its*

Prepositions *to with*

Particles *off up*

... more

Part-of-speech tagging is a disambiguation process

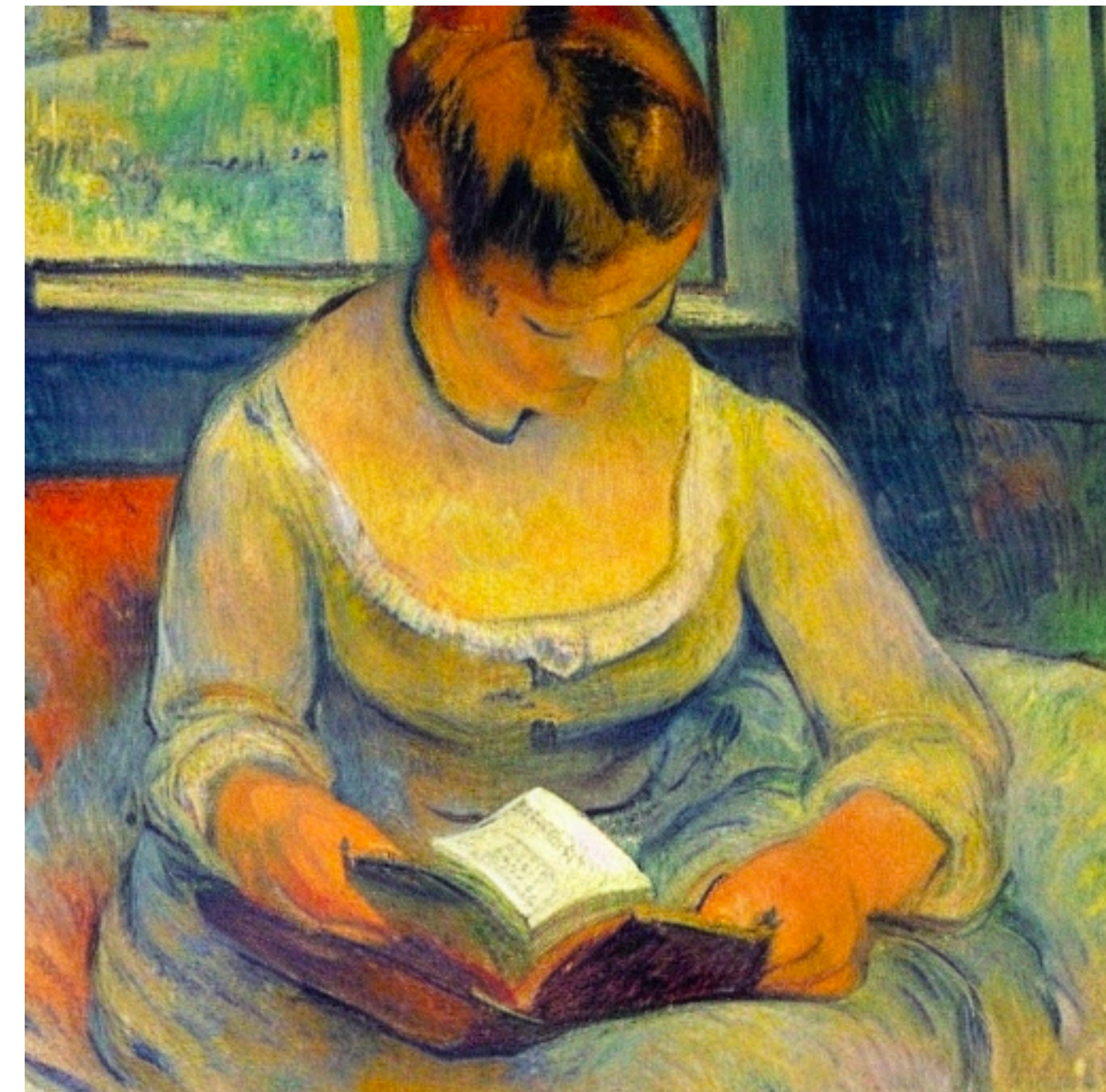
Verb or Noun?



Verb or Noun?



She is **reading** a book about **Reading**





One morning I shot an elephant in my pajamas

https://www.youtube.com/watch?v=NfN_gcjGoJo

One morning I shot an elephant in my pajamas

How he got into my pajamas I don't know

Syntax is not Morphology

- ▶ Morphology deals with the internal structure of words
- ▶ Syntax deals with combinations of words

- ▶ Morphology is usually irregular
- ▶ Syntax has its irregularities, but it is usually regular
 - Syntax is mostly made up of general rules that apply across-the-board

Constituency

- ▶ One way of viewing the structure of a sentence is as a collection of nested constituents
- ▶ Constituent: a group of neighboring words relate more closely to one another than to other words in the sentence
- ▶ Constituents larger than a word are called phrases
 - Noun phrases
 - Prepositional phrases
 - Verb phrases
- ▶ Phrases can contain other phrases

Noun phrase (NP)

- ▶ a phrase that has a noun or pronoun as its head or performs the same grammatical function as a noun
 - The elephant arrived
 - It arrived.
 - Elephants arrived.
 - The big pretty elephant arrived.
 - The elephant she loves arrived.

Prepositional phrase (PP)

- ▶ I arrived on Tuesday.
 - ▶ I arrived in March.
 - ▶ I arrived under the leaking roof.
-
- ▶ Every prepositional phrase contains a noun phrase

Verb phrase

- ▶ A verb phrase in English consists of a verb followed by assorted other things
 - VP → Verb NP
 - I prefer an afternoon lecture
 - VP → Verb NP PP
 - have a lecture in the afternoon
 - VP → Verb PP
 - Teaching on Tuesday

Is a string constituent?

- ▶ Substitution test
 - Can the string be replaced by a single word?
- ▶ Movement test
 - Can the string be moved around in the sentence?
- ▶ Answer test
 - Can the string be the answer to a question?

He talks [in class]

He talks there

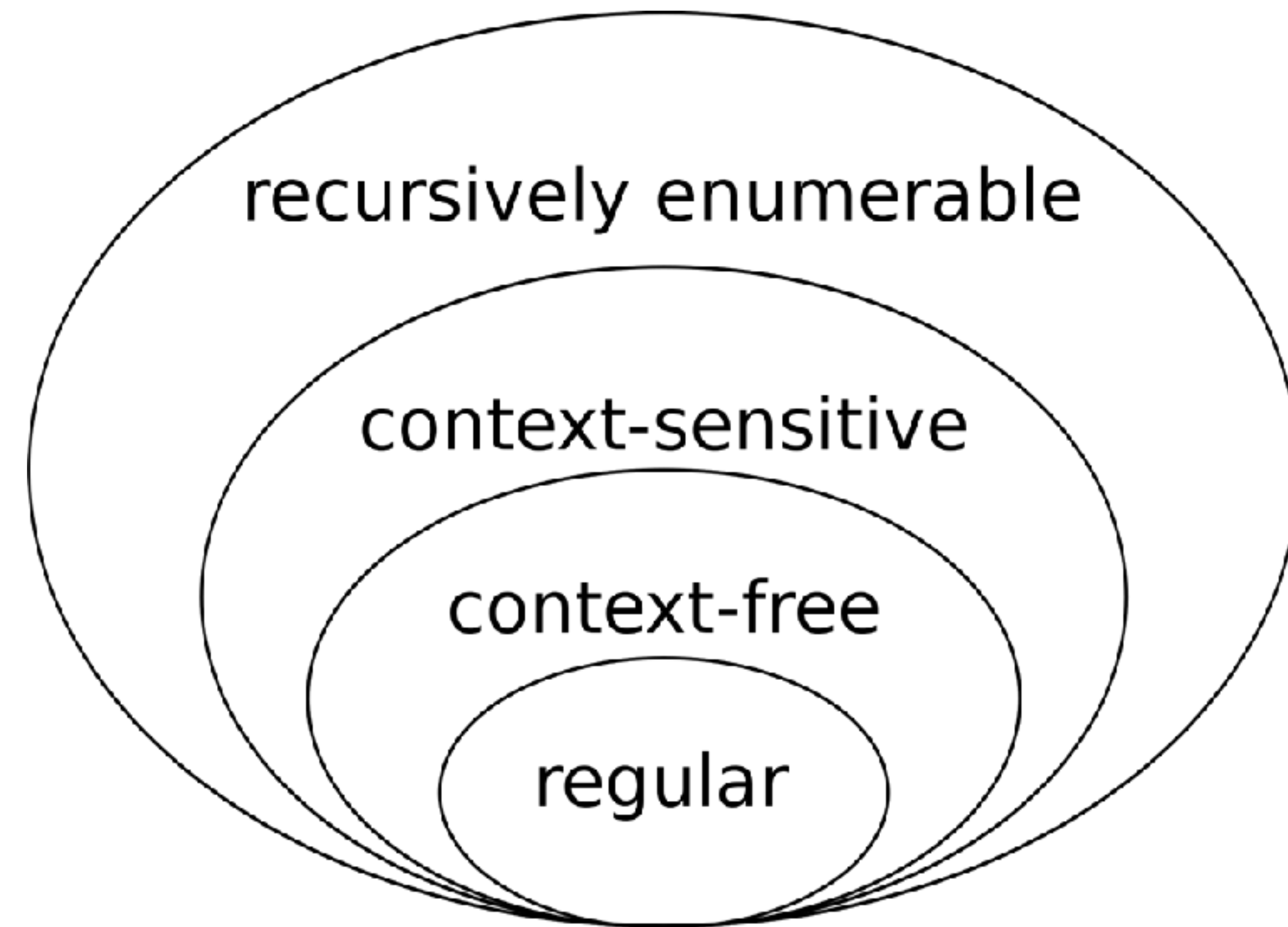
[In class], he talks

Where does he talks?

[In class]

Chomsky hierarchy

- ▶ Type-0 grammars include all formal grammars
- ▶ Type-1 grammars generate context-sensitive languages
- ▶ Type-2 grammars generate the context-free languages
- ▶ Type-3 grammars generate the regular languages, which can be described using regular expressions



Context-free grammar

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

Rules or productions

- ▶ Context-free
 - production rules are independent of the context
 - There is no context in the left hand side (LHS) of rules

Grammar Rules

$S \rightarrow NP VP$

$NP \rightarrow Pronoun$

| $Proper-Noun$

| $Det Nominal$

$Nominal \rightarrow Nominal Noun$

| $Noun$

$VP \rightarrow Verb$

| $Verb NP$

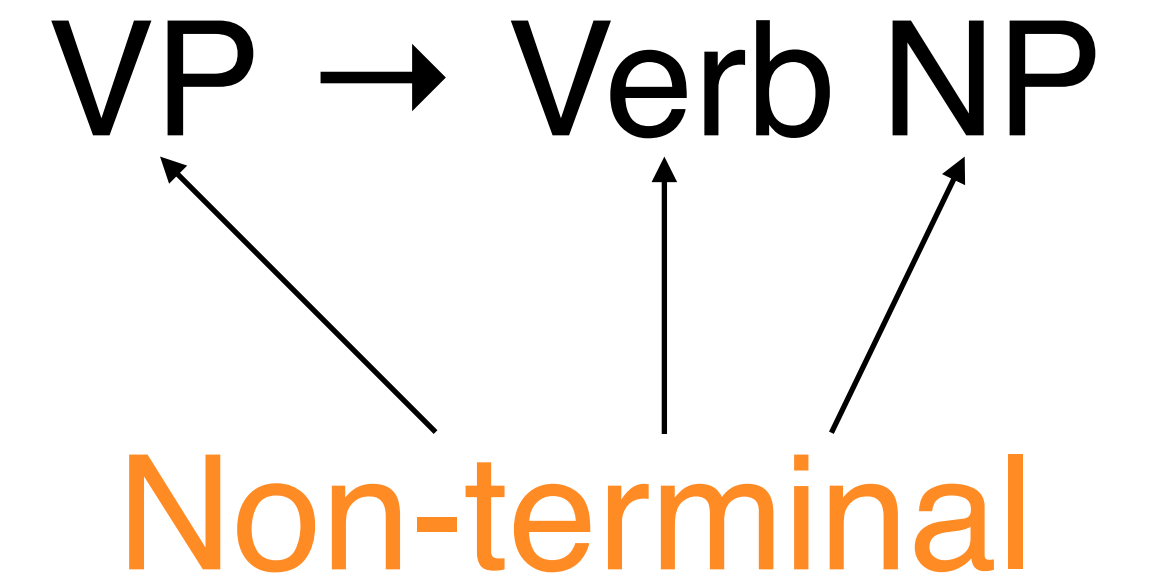
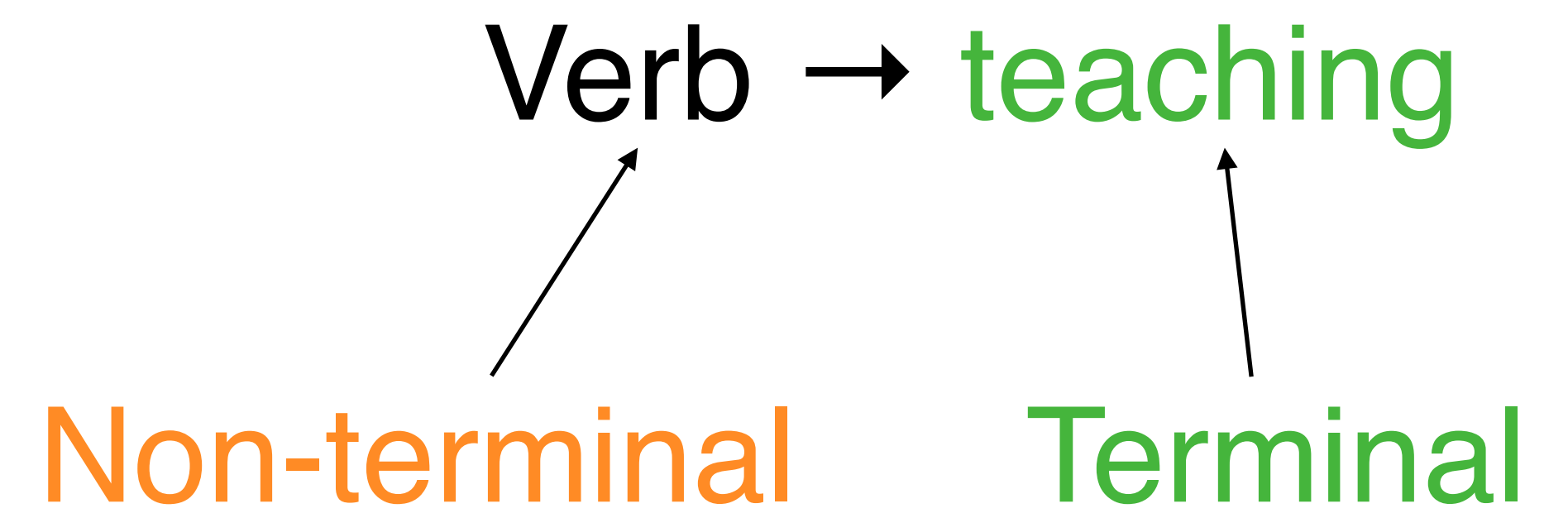
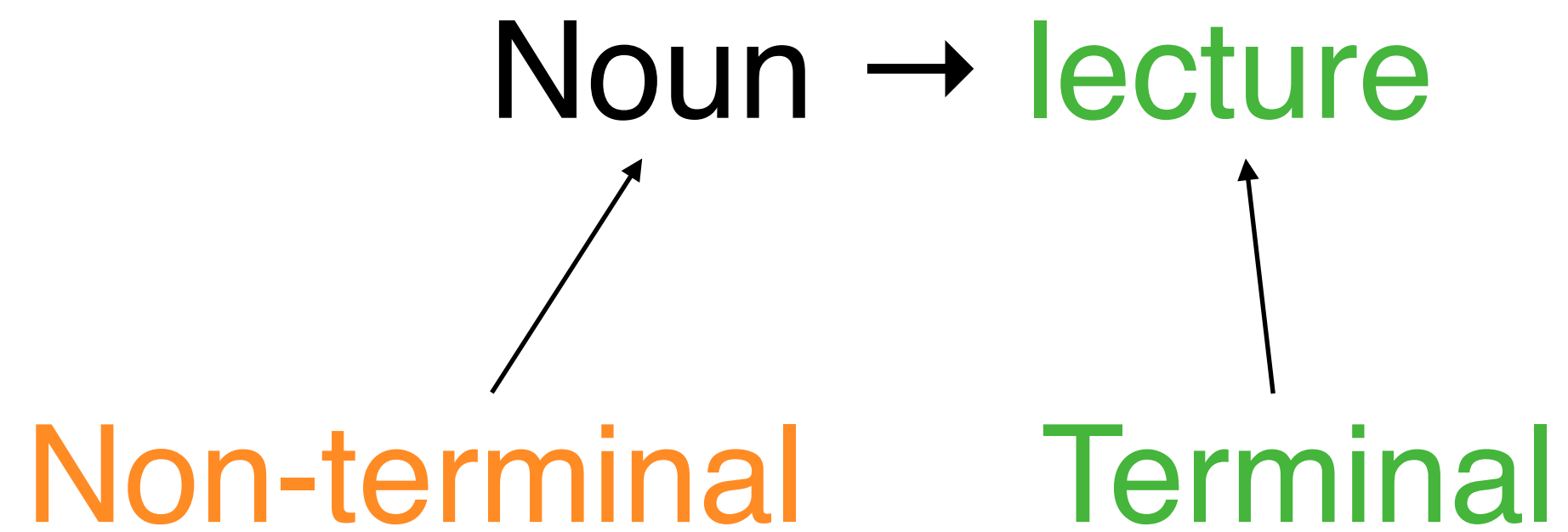
| $Verb NP PP$

| $Verb PP$

$PP \rightarrow Preposition NP$

Terminal vs Non-terminal

- ▶ **Terminal:** The symbols that that correspond to words in the language
- ▶ **Non-terminal:** The symbols that express abstractions over these terminals



Lexicon: Terminal vs Non-terminal

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

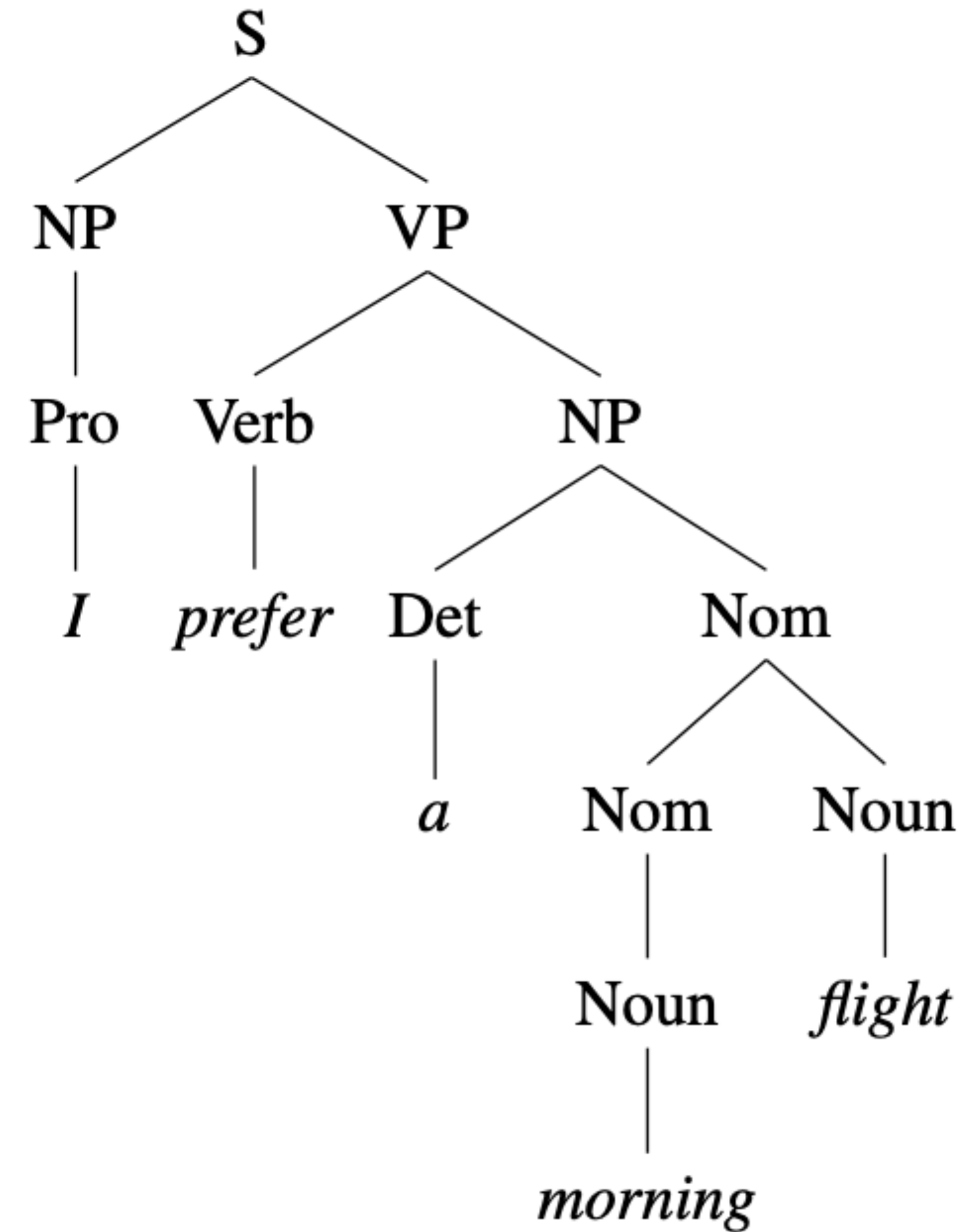
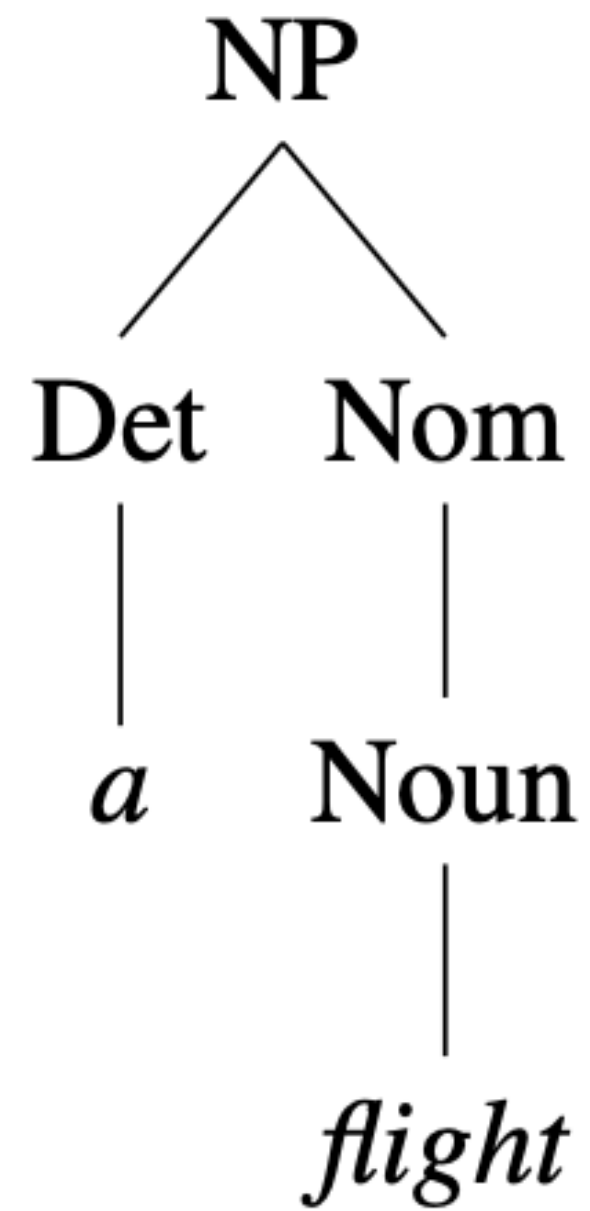
S: Start symbol

- ▶ 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
- ▶ S is usually interpreted as the “sentence” node

$S \rightarrow NP VP$

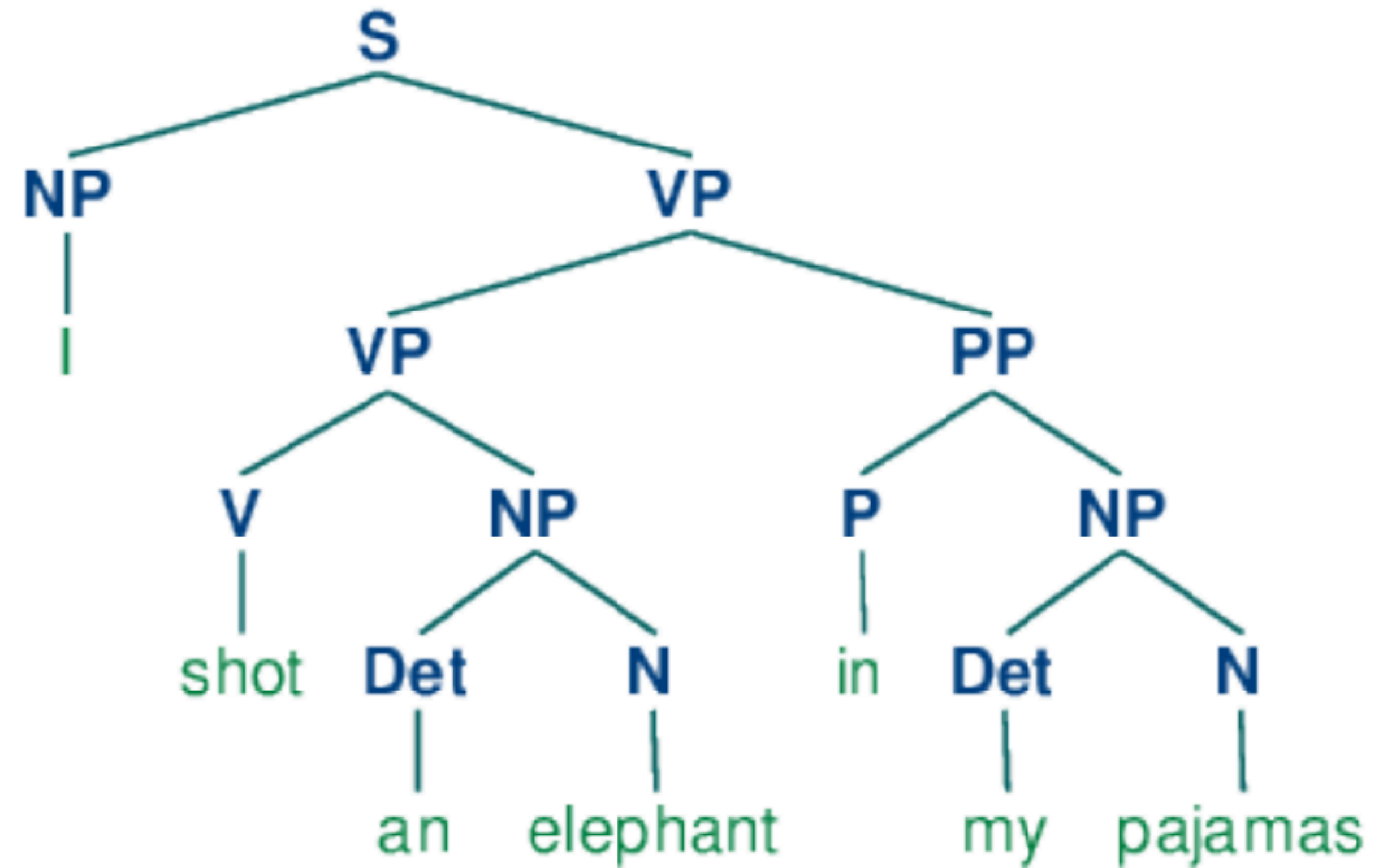
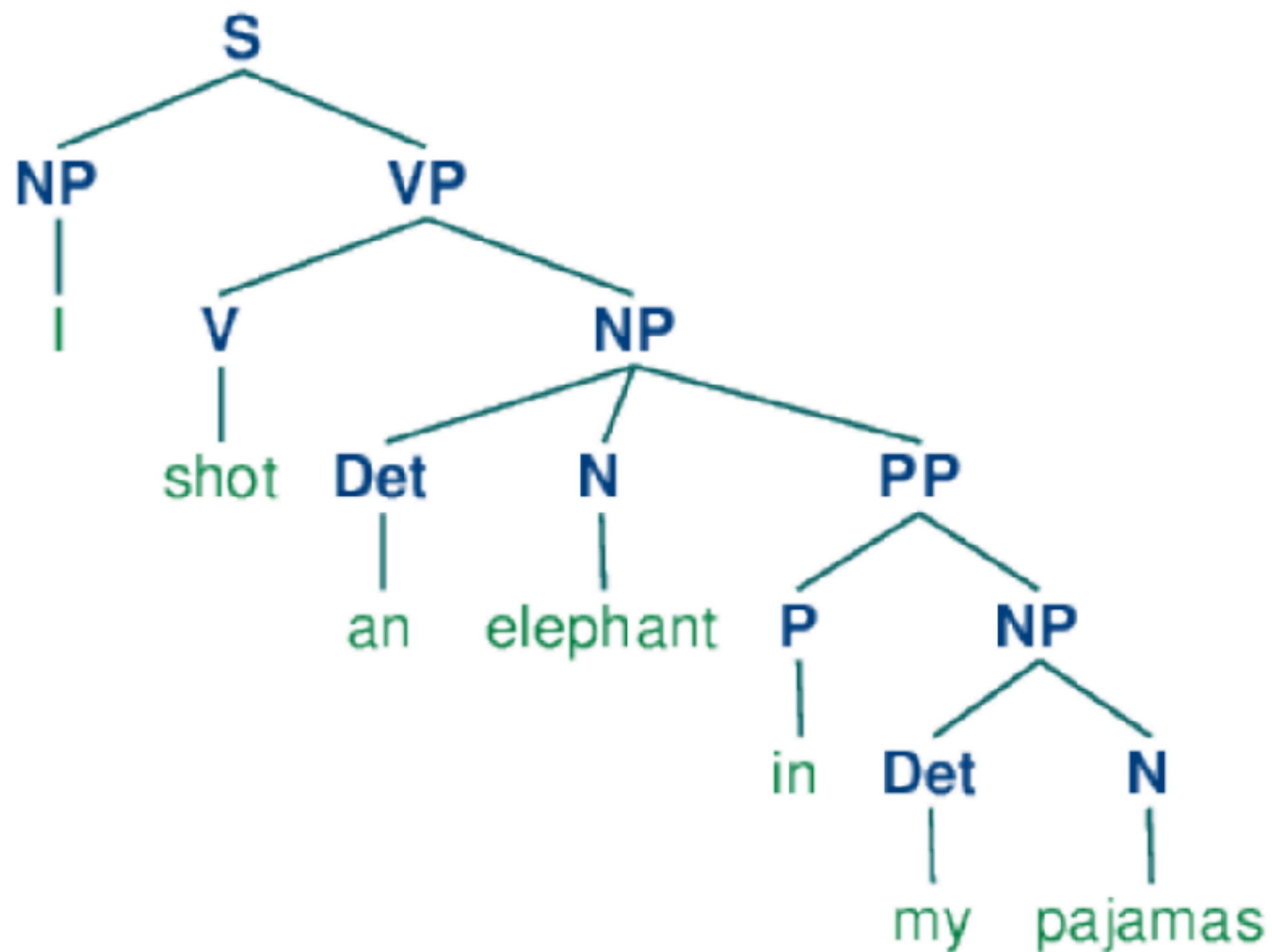
I prefer an afternoon lecture

A (Constituency) Parse Tree



Ambiguity

- ▶ Structural ambiguity occurs when the grammar can assign more than one parse to a sentence



Cocke-Kasami-Younger (CKY) algorithm

- ▶ Bottom-up parsing
 - Start with words
- ▶ Dynamic programming
 - save the results in a table/chart
 - re-use these results in finding larger constituents
- ▶ Presumes a CFG in Chomsky Normal Form

Chomsky Normal Form

N	Finite set of non-terminal symbols	NP, VP, S
Σ	Finite alphabet of terminal symbols	the, dog, a
R	Set of production rules, each $A \rightarrow \beta$ $\beta =$ single terminal (from Σ) or two non-terminals (from N)	$S \rightarrow NP VP$ Noun \rightarrow dog
S	Start symbol	

Chomsky Normal Form (CNF)

- ▶ Any CFG can be converted into weakly equivalent CNF
- ▶ In CNF, each non-terminal generates two non-terminals

$$A \rightarrow B C \gamma$$

$$\begin{aligned} A &\rightarrow X1 \gamma \\ X1 &\rightarrow B C \end{aligned}$$

$$S \rightarrow Aux NP VP$$

$$S \rightarrow X1 VP$$

$$X1 \rightarrow Aux NP$$

Chomsky Normal Form (CNF)

▶ Left hand side (LHS) rules

- LHS will have non-terminals

VP → Verb NP PP ❌

VP → Verb NP ✔️

▶ Right hand side (RHS) rules

- Two non-terminals

VP → teaching NP ❌

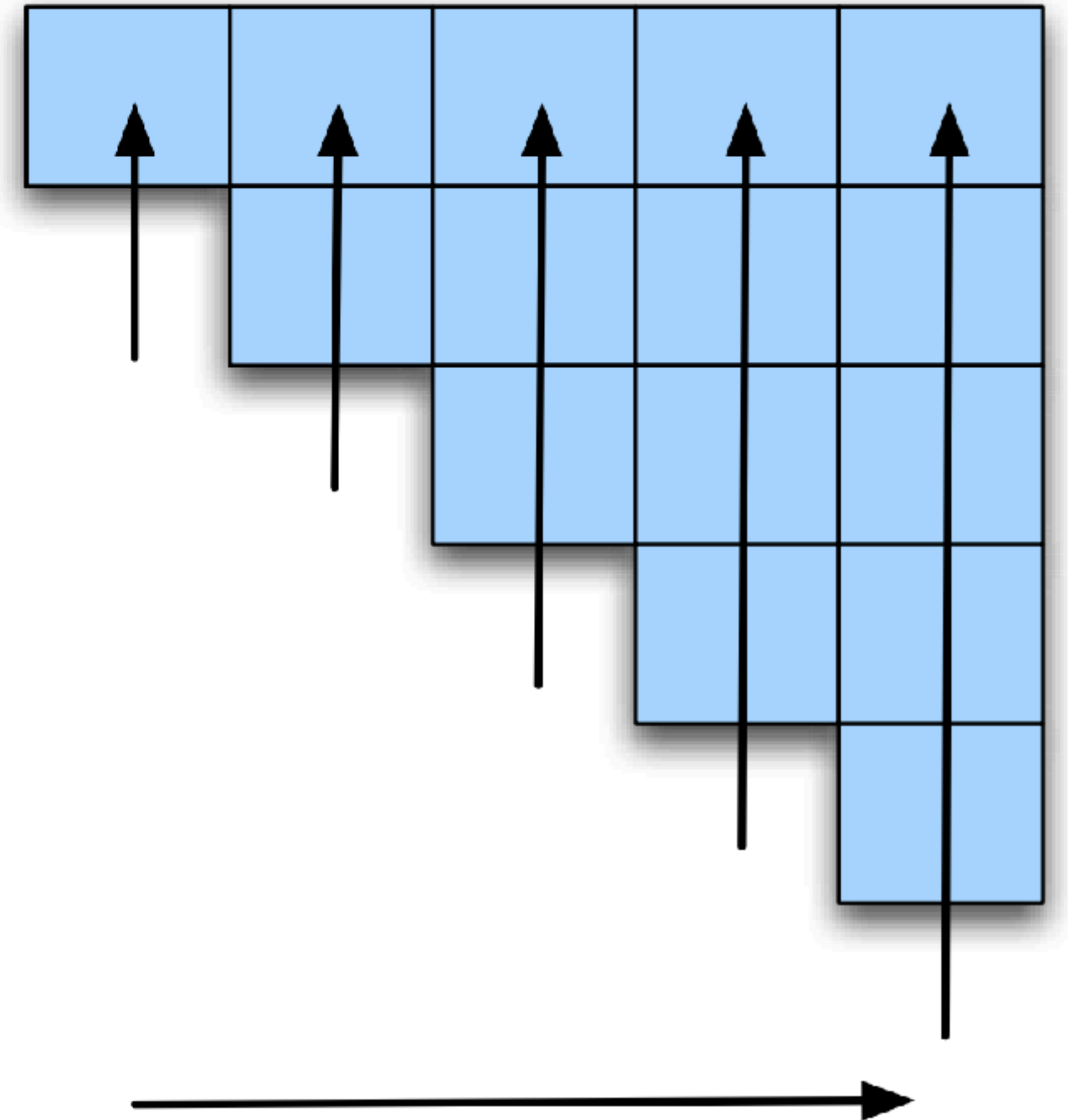
- One terminal

VP → eat ✔️

\mathcal{L}_1 Grammar	\mathcal{L}_1 in CNF
$S \rightarrow NP VP$	$S \rightarrow NP VP$
$S \rightarrow Aux NP VP$	$S \rightarrow X1 VP$
	$X1 \rightarrow Aux NP$
$S \rightarrow VP$	$S \rightarrow book \mid include \mid prefer$
	$S \rightarrow Verb NP$
	$S \rightarrow X2 PP$
	$S \rightarrow Verb PP$
	$S \rightarrow VP PP$
$NP \rightarrow Pronoun$	$NP \rightarrow I \mid she \mid me$
$NP \rightarrow Proper-Noun$	$NP \rightarrow TWA \mid Houston$
$NP \rightarrow Det Nominal$	$NP \rightarrow Det Nominal$
$Nominal \rightarrow Noun$	$Nominal \rightarrow book \mid flight \mid meal \mid money$
$Nominal \rightarrow Nominal Noun$	$Nominal \rightarrow Nominal Noun$
$Nominal \rightarrow Nominal PP$	$Nominal \rightarrow Nominal PP$
$VP \rightarrow Verb$	$VP \rightarrow book \mid include \mid prefer$
$VP \rightarrow Verb NP$	$VP \rightarrow Verb NP$
$VP \rightarrow Verb NP PP$	$VP \rightarrow X2 PP$
	$X2 \rightarrow Verb NP$
$VP \rightarrow Verb PP$	$VP \rightarrow Verb PP$
$VP \rightarrow VP PP$	$VP \rightarrow VP PP$
$PP \rightarrow Preposition NP$	$PP \rightarrow Preposition NP$

CKY algorithm

- ▶ Fills the upper-triangular matrix a column at a time
 - From left to right
 - From bottom to top
- ▶ This scheme guarantees that at each point in time we have all the information we need



CKY algorithm: a toy example

Rules

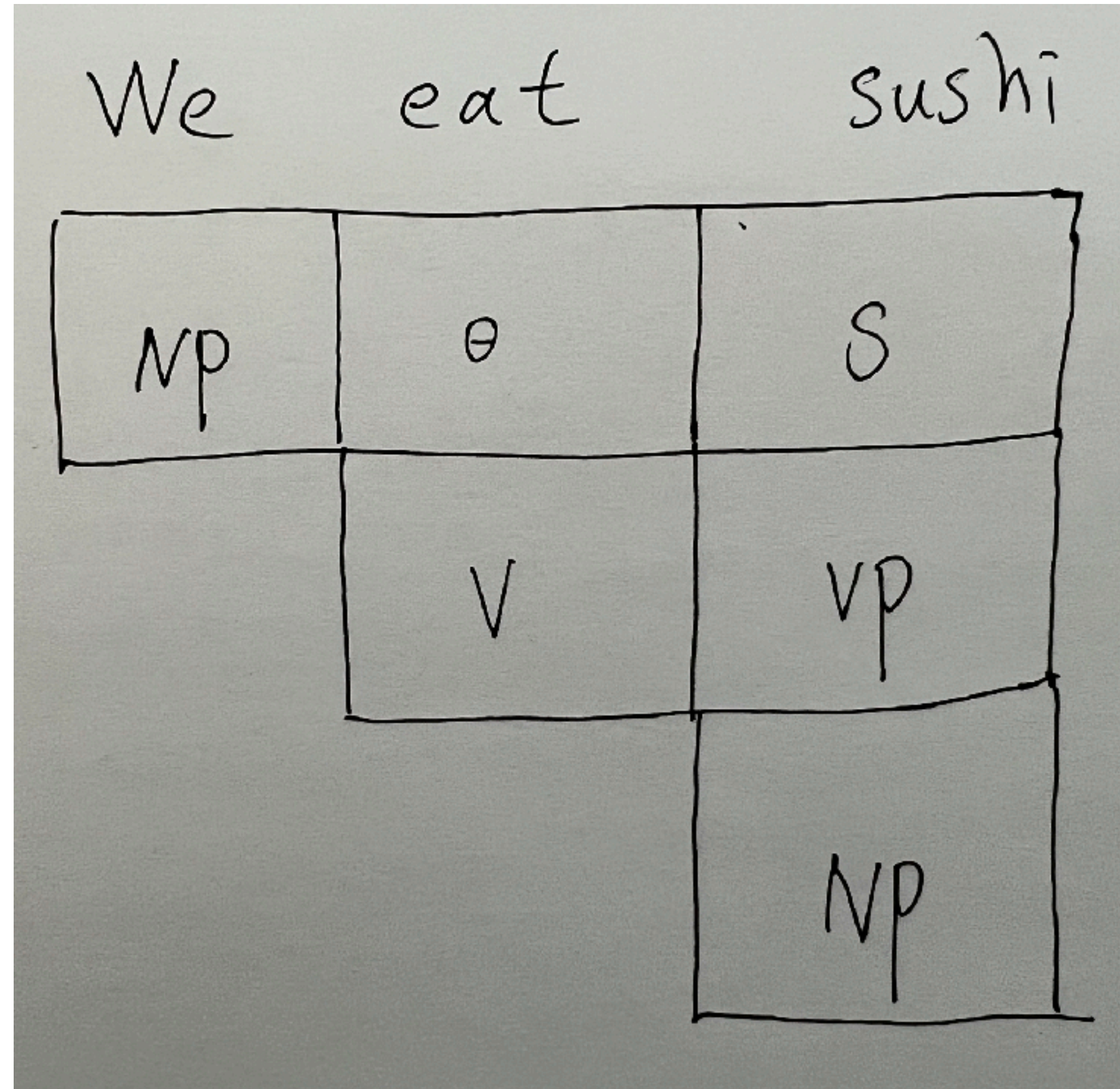
$S \rightarrow NP VP$

$VP \rightarrow V NP$

$V \rightarrow \text{eat}$

$NP \rightarrow \text{we}$

$NP \rightarrow \text{sushi}$



CKY algorithm

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

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

for all $\{A \mid A \rightarrow words[j] \in grammar\}$

$table[j-1, j] \leftarrow table[j-1, j] \cup A$

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

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

for all $\{A \mid A \rightarrow BC \in grammar \text{ and } B \in table[i, k] \text{ and } C \in table[k, j]\}$

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

CKY Example

	<i>Book</i>	<i>the</i>	<i>flight</i>	<i>through</i>	<i>Houston</i>
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]	

\mathcal{L}_1 in CNF

$S \rightarrow NP VP$

$S \rightarrow X1 VP$

$X1 \rightarrow Aux NP$

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

$S \rightarrow Verb NP$

$S \rightarrow X2 PP$

$S \rightarrow Verb PP$

$S \rightarrow VP PP$

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

$NP \rightarrow TWA \mid Houston$

$NP \rightarrow Det Nominal$

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

$Nominal \rightarrow Nominal Noun$

$Nominal \rightarrow Nominal PP$

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

$VP \rightarrow Verb NP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

$VP \rightarrow Verb PP$

$VP \rightarrow VP PP$

$PP \rightarrow Preposition NP$

Exercise

S	→	NP VP
VP	→	VBD NP
VP	→	VP PP
Nominal	→	Nominal PP
Nominal	→	pajamas elephant I
PP	→	IN NP
NP	→	DT NN
NP	→	pajamas elephant I
NP	→	PRP\$ Nominal

VBD	→	shot
DT	→	an my
PRP	→	I
PRP\$	→	my
IN	→	in

I shot an elephant in my pajamas

Exercise

I	shot	an	elephant	in	my	pajamas
---	------	----	----------	----	----	---------

NP, PRP [0,1]						
	VBD [1,2]					
		DT [2,3]				
			NP, NN [3,4]			
				IN [4,5]		
					PRP\$ [5,6]	
						NNS [6,7]

Summary

- ▶ Concept of syntax and constituency
 - Syntax deals with combinations of words
- ▶ Context-free grammar
 - production rules are independent of the context
- ▶ Cocke-Kasami-Younger (CKY) algorithm
 - Bottom-up parsing - start with words
 - Dynamic programming
 - Presumes a CFG in Chomsky Normal Form

Reading

- ▶ Chapter 17: Context-Free Grammars and Constituency Parsing
- ▶ <https://web.stanford.edu/~jurafsky/slp3/17.pdf>