

Computational Linguistics

LT3233



Jixing Li

Lecture 4: Context-Free Grammars

Slides adapted from Edward Stabler, Michael Collins

Lecture plan

- Syntax 101
- What are context-free grammars?
- Top-down recognition
- Short break (15 mins)
- Hands-on exercises

N-grams review

Bigram model: Condition on the previous word:

$$\mathbf{P}(w_1, w_2, w_3, \dots, w_n) \approx \prod_{i=1}^n \mathbf{P}(w_i | w_{i-1})$$

$$\mathbf{P}(a \text{ pile of shaving cream}) \approx$$

$$\mathbf{P}(a) \times \mathbf{P}(\text{pile} | a) \times \mathbf{P}(\text{of} | \text{pile}) \times \mathbf{P}(\text{shaving} | \text{of}) \times \mathbf{P}(\text{cream} | \text{shaving})$$

MLE ↓

$$\frac{\text{Count}(a)}{\# \text{ of words}}$$

$$\frac{\text{Count}(a \text{ pile})}{\text{Count}(a)}$$

$$\frac{\text{Count}(\text{pile of})}{\text{Count}(\text{pile})}$$

$$\frac{\text{Count}(\text{of shaving})}{\text{Count}(\text{of})}$$

$$\frac{\text{Count}(\text{shaving cream})}{\text{Count}(\text{shaving})}$$

→ Can extend to trigrams, 4-grams, 5-grams, etc.

Languages are not Markovian

N-gram models are insufficient models of language



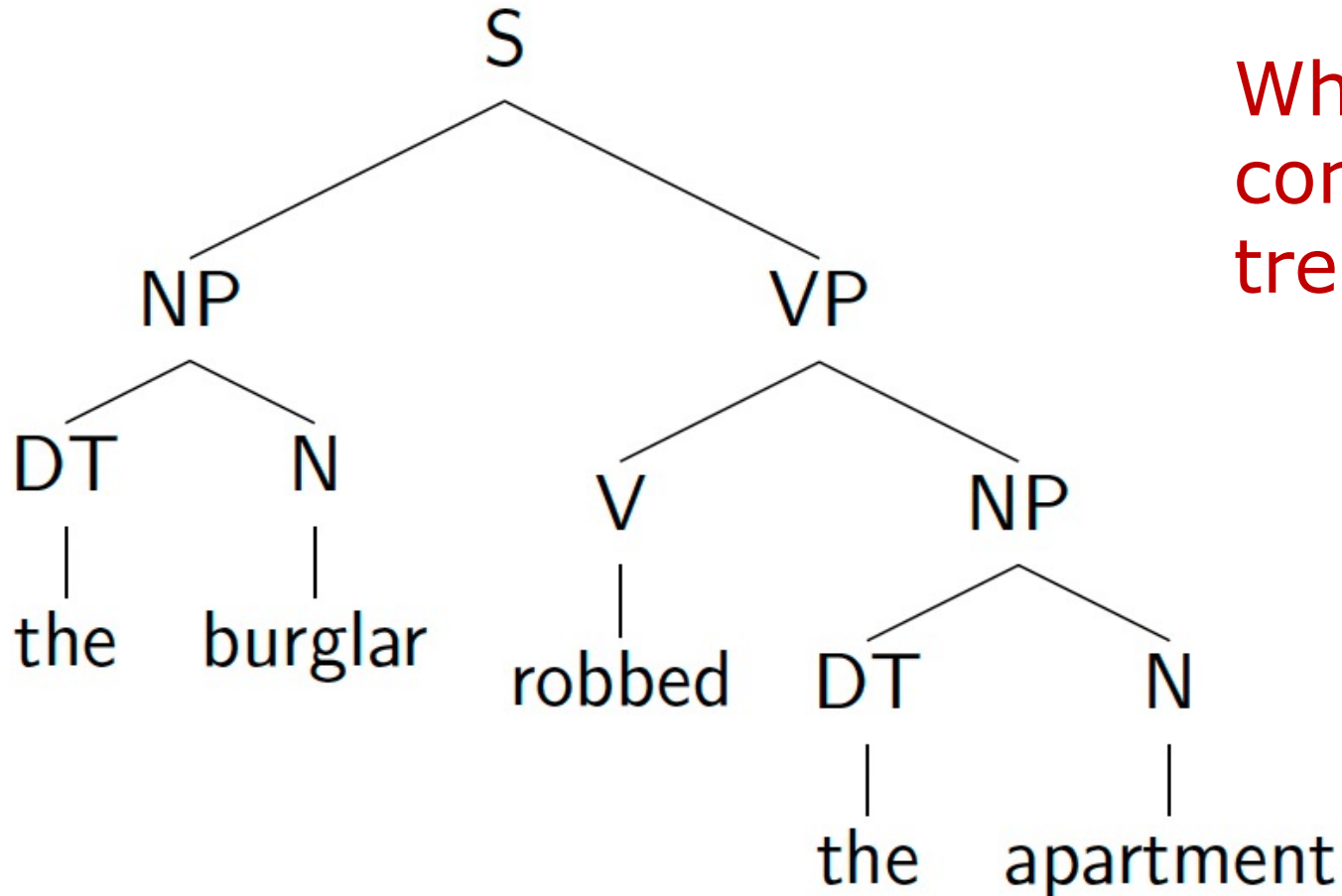
Bulldogs bulldogs bulldogs fight fight fight.

The bulldogs that are fought by the bulldogs that the bulldogs fight, fight.

→ **long-distance dependencies**

Syntactic structure

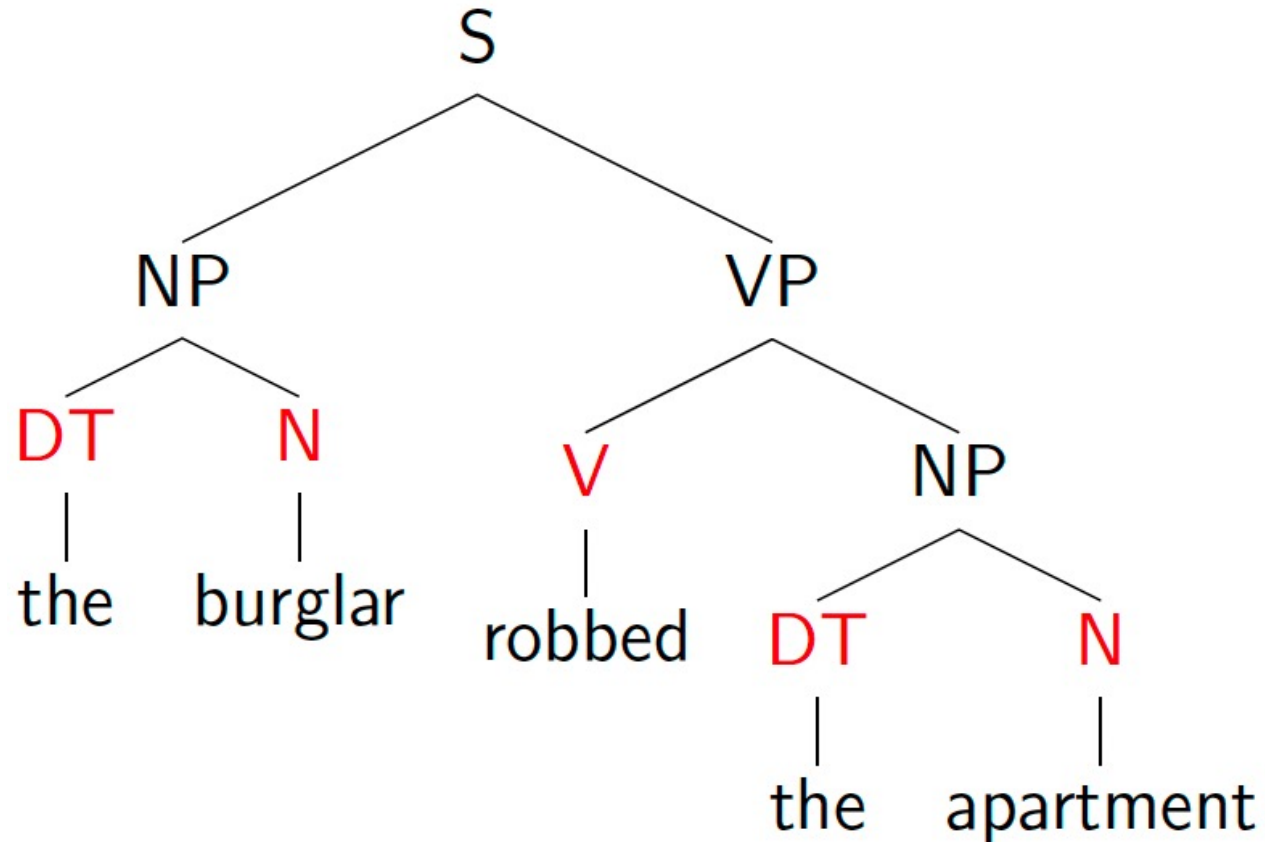
“the burglar robbed the apartment”



What information is conveyed by a syntactic tree?

Trees visualize parts of speech

DT: determiner; N: noun; V: verb



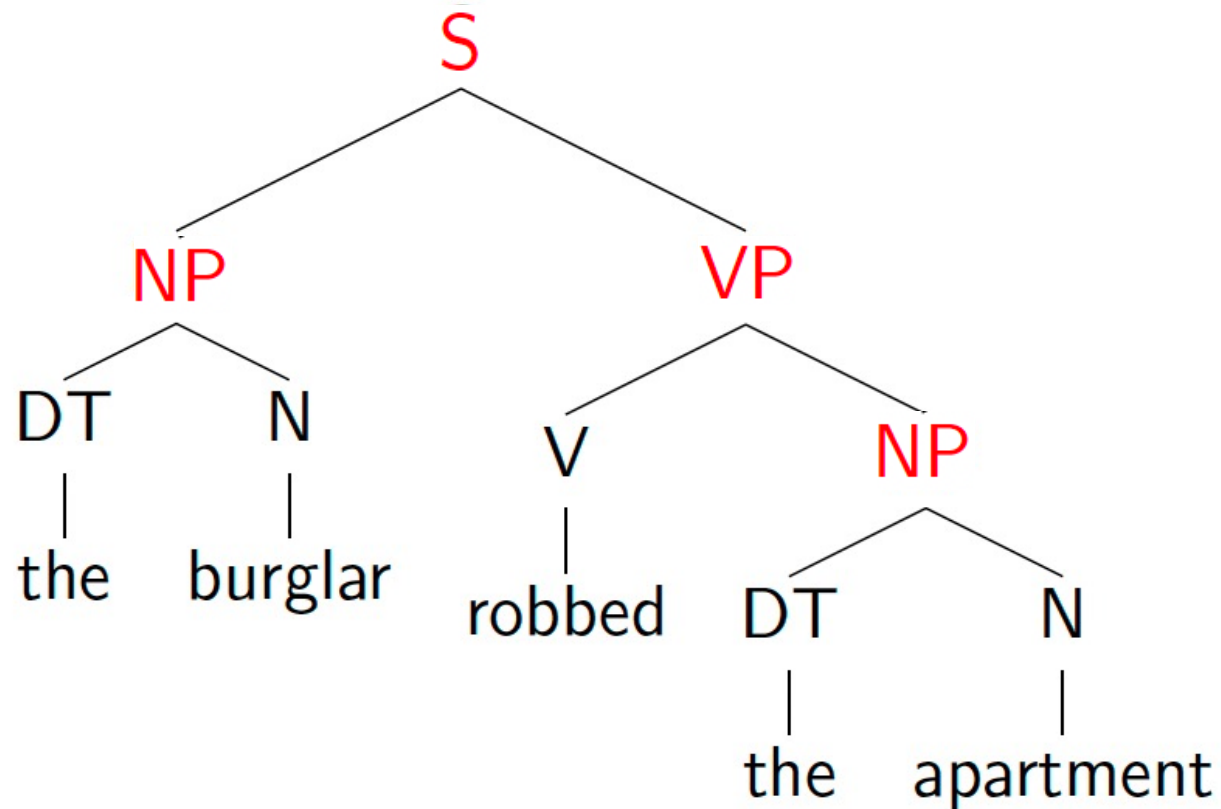
DT → the

N → burglar | apartment

V → robbed

Trees visualize constituency

S: Sentence; NP: noun phrase; VP: verb phrase;



$S \rightarrow NP VP$
 $NP \rightarrow DT N$
 $VP \rightarrow V NP$

Constituency

Example noun phrases:

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

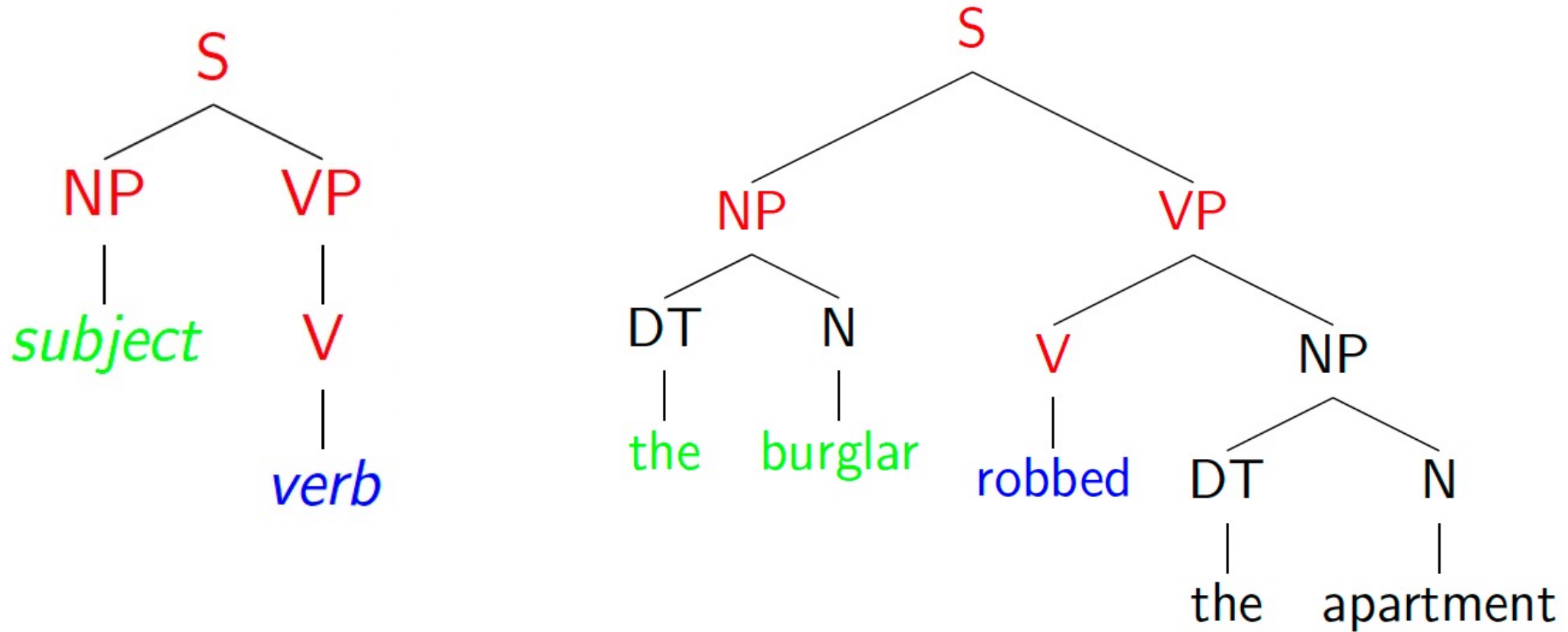
NPs appear in similar syntactic environments, i.e., before a verb:

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

Individual words in an NP may not occur before a verb:

*from *arrive...* *as *attracts...*
*the *is...* *spot *sat...*

Trees visualize grammatical relationships



'the burglar' is the subject of 'robbed'

Context-free grammars

Context-free grammars (CFG) are also called **Phrase-Structure Grammars**. The idea of basing a grammar on constituent structures is formalized by Chomsky (1956)

A CFG consists of **a set of rules** and **a lexicon** of words and symbols

start symbol

S → NP VP

NP → DT N

VP → V NP

DT → the

V → robbed

N → burglar | apartment

non-terminal symbols

terminal symbols

alternate possible expansions

Context-free grammars

The **formal language** defined by a CFG is the set of strings that are derivable from the start symbol.

LO:

S → NP VP

NP → DT N

VP → V NP

DT → the

V → robbed

N → burglar | apartment

Grammatical:

the burglar robbed the apartment

[S [NP [DT the] [NP burglar]] [VP [V robbed] [NP [DT the] [N apartment]]]]

the apartment robbed the burglar

[S [NP [DT the] [NP apartment]] [VP [V robbed] [NP [DT the] [N burglar]]]]

Ungrammatical:

* *robbed the apartment the burglar*

Formal definition of 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

Capital letters like A , B , and S

Non-terminals

S

The start symbol

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

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

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

Strings of terminals

Formal language

Hopcroft and Ullman (1979):

if $A \rightarrow \beta$ is a production of R and α and γ are any strings in the set $(\Sigma \cup N)^*$, then we say that $\alpha A \gamma$ **directly derives** $\alpha \beta \gamma$, or $\alpha A \gamma \Rightarrow \alpha \beta \gamma$.

Let $\alpha_1, \alpha_2, \dots, \alpha_m$ be strings in $(\Sigma \cup N)^*$, $m \geq 1$, such that

$$\alpha_1 \Rightarrow \alpha_2, \alpha_2 \Rightarrow \alpha_3, \dots, \alpha_{m-1} \Rightarrow \alpha_m$$

We say that α_1 **derives** α_m , or $\alpha_1 \xRightarrow{*} \alpha_m$.

$$\mathcal{L}_G = \{w \mid w \text{ is in } \Sigma^* \text{ and } S \xRightarrow{*} w\}$$

Example

'the dog laughs'

S → NP VP

NP → DT N

DT → the

N → dog

VP → VB

VB → laughs

Derivation

S

NP VP

DT N VP

the N VP

the dog VP

the dog VB

the dog laughs

Rules used

S → NP VP

NP → DT N

DT → the

N → dog

VP → VB

VB → laughs

Problems of CFG

- Cannot express many syntactic phenomena, e.g., subject-verb **agreement**, **wh-movement**, etc.
- Lack **headedness**: The head is the word in the phrase that is grammatically the most important. Heads are passed up the parse tree.
 - N is the head of an NP
 - V is the head of a VP
 - P is the head of PP

Nevertheless, CFG is a good starting point.

Recognizer

The problem of recognition: Given a grammar and a string of words, return *True* if the string has a derivation and *False* otherwise.

'the dog laughs'

S → NP VP

NP → DT N

DT → the

N → dog

VP → VB

VB → laughs

Derivation

S

NP VP

DT N VP

→ **True**

the N VP

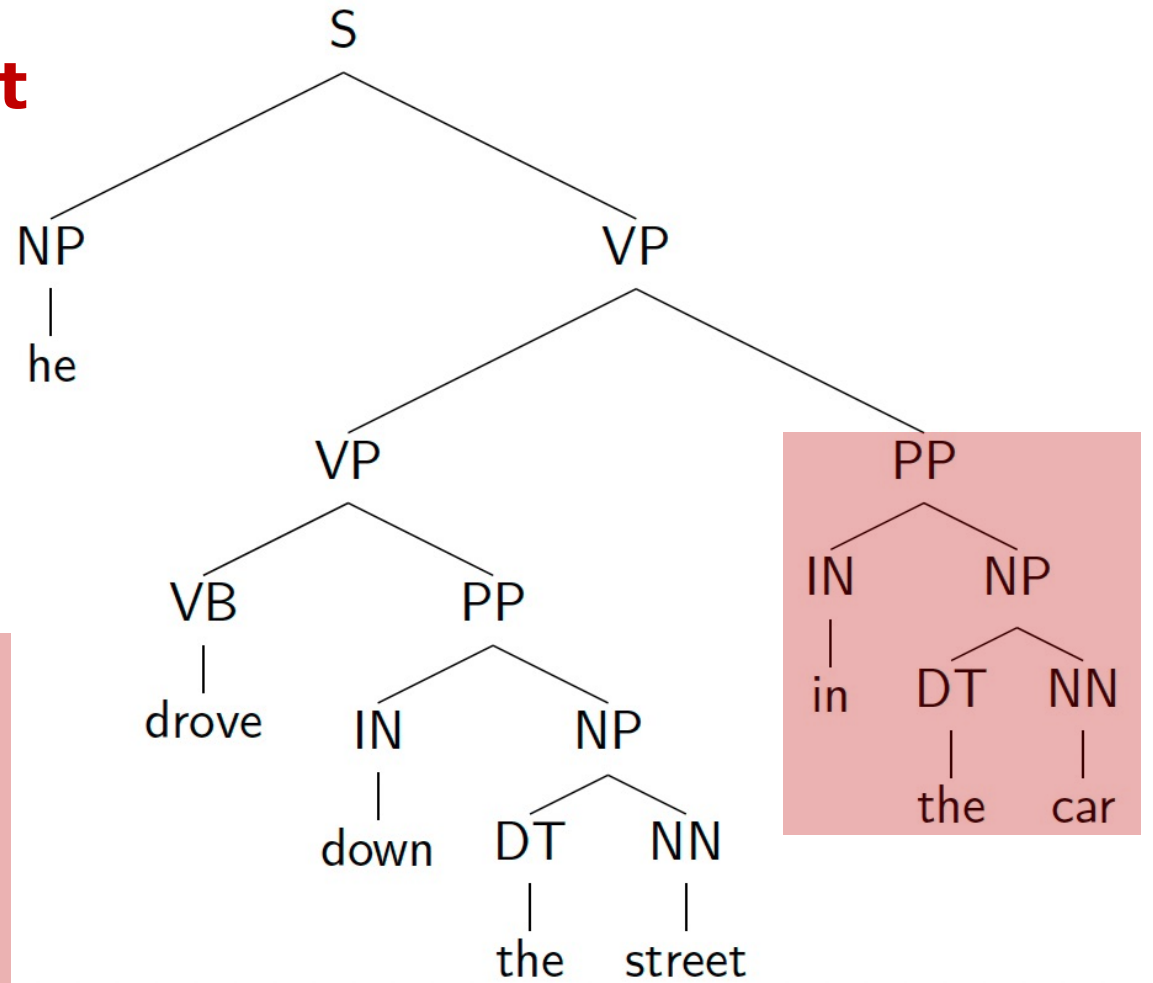
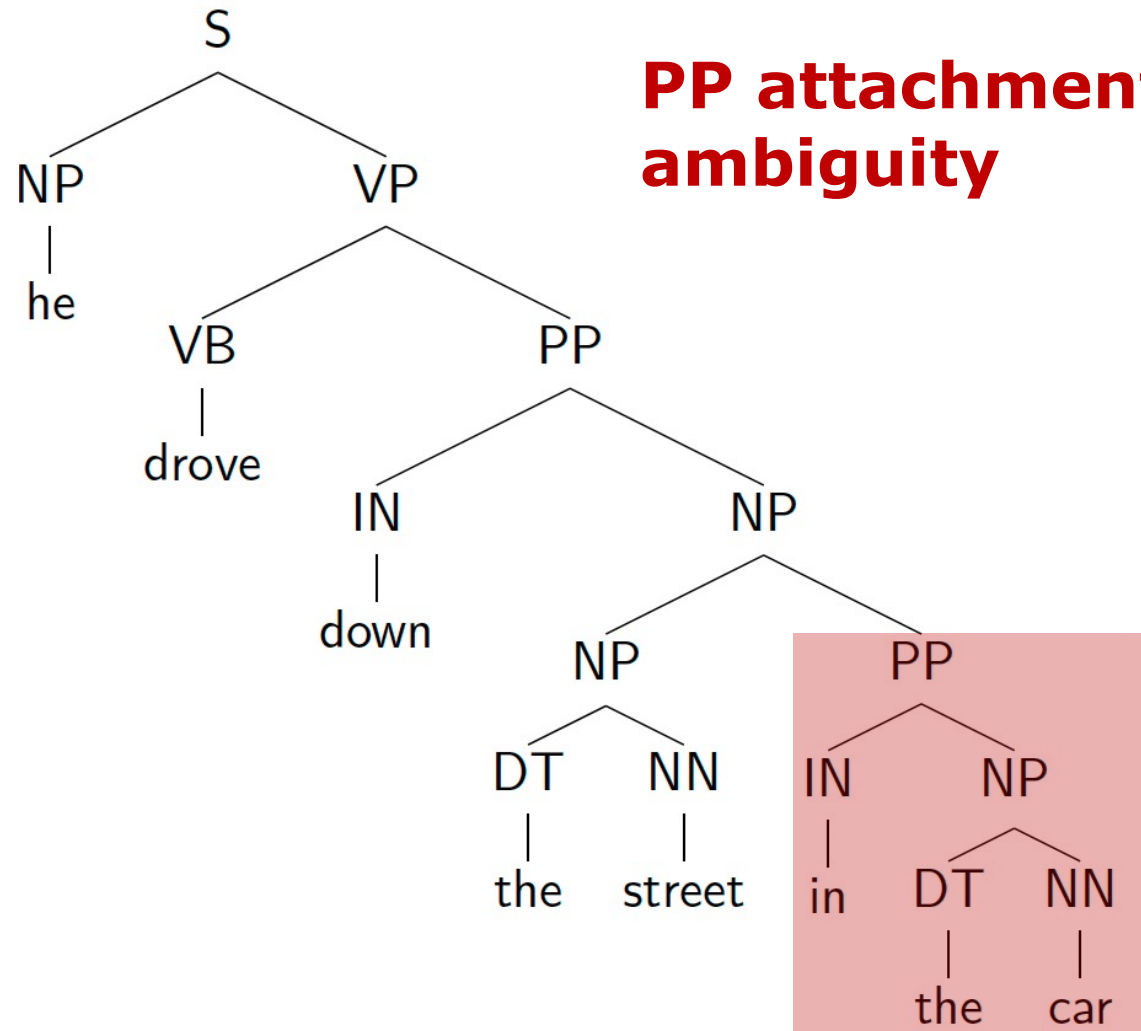
the dog VP

the dog VB

the dog laughs

Ambiguity

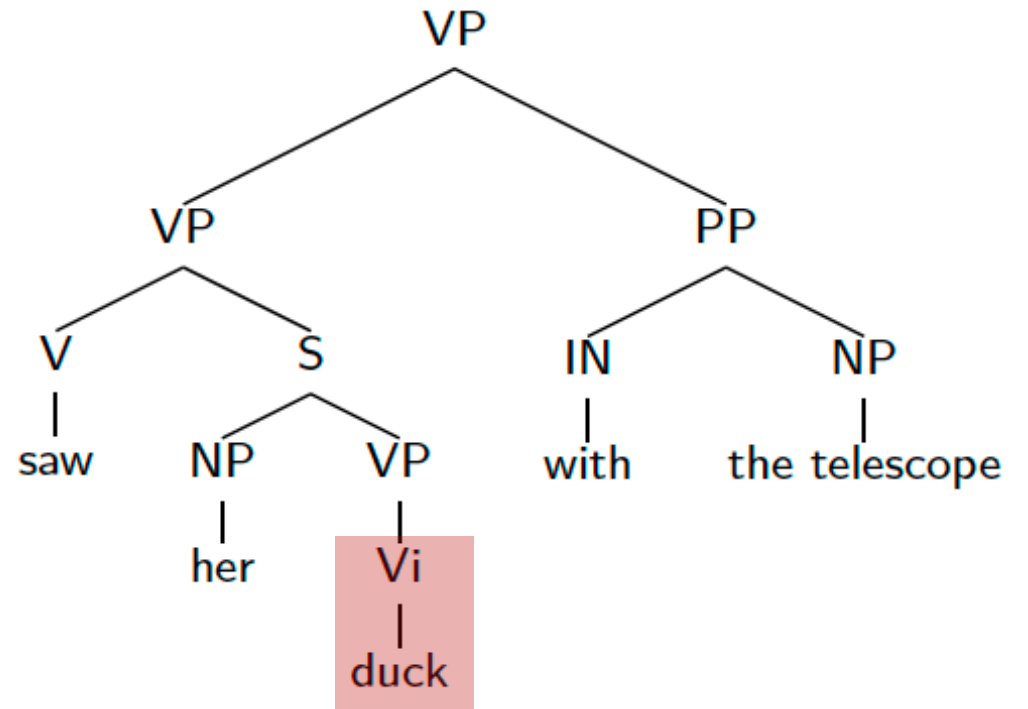
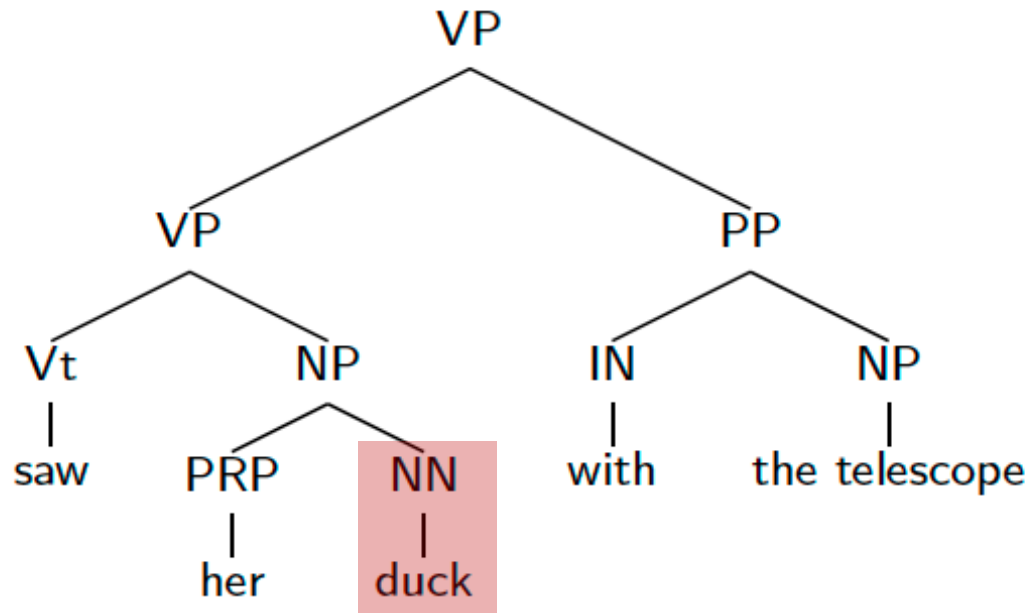
'He drove down the street in the car.'



Source of ambiguity

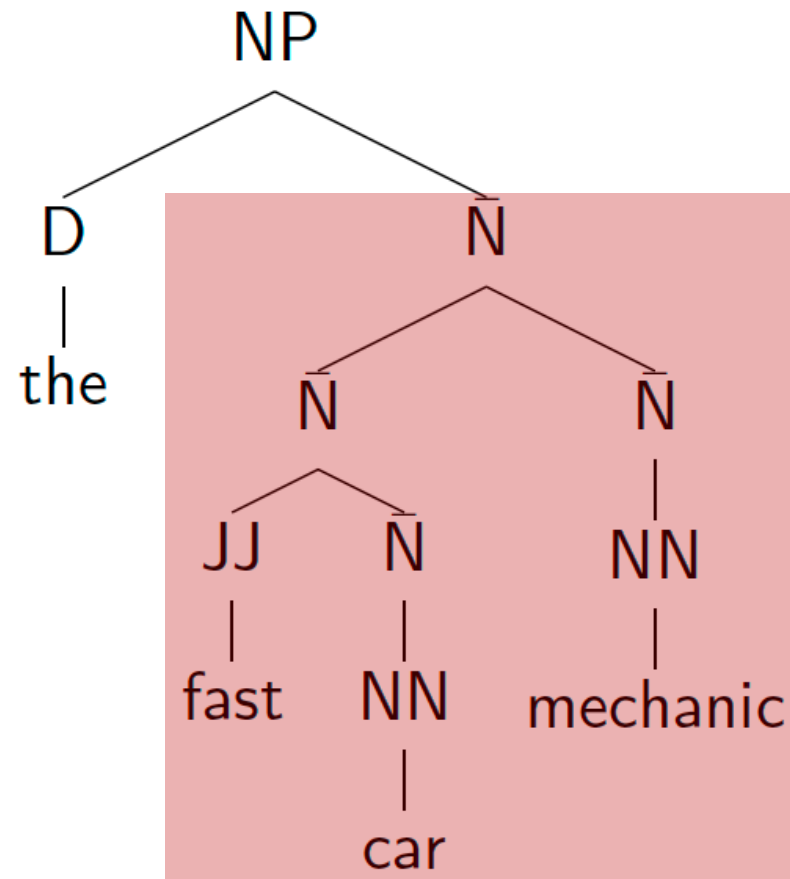
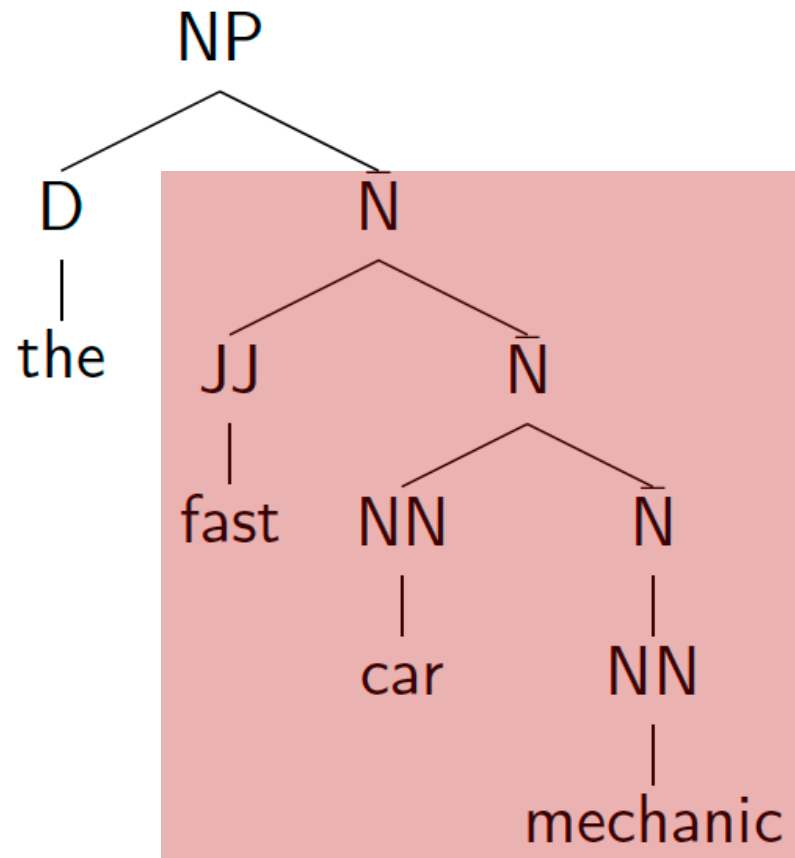
Part-of-speech ambiguity

- NN → duck
- Vi → duck



Source of ambiguity

Noun premodifier ambiguity



Top-down backtrack recognition

Rules:

$$\begin{array}{l} S \rightarrow DP \ VP \\ DP \rightarrow \left\{ \begin{array}{l} (D) \ NP \\ \text{Name} \\ \text{Pronoun} \end{array} \right\} \\ NP \rightarrow N \ (PP) \end{array} \quad \begin{array}{l} VP \rightarrow V \ (DP) \ \left(\left\{ \begin{array}{l} PP \\ CP \\ VP \end{array} \right\} \right) \\ PP \rightarrow P \ (DP) \\ AP \rightarrow A \ (PP) \\ CP \rightarrow C \ S \\ AdvP \rightarrow Adv \end{array} \quad \begin{array}{l} NP \rightarrow AP \ NP \\ NP \rightarrow NP \ PP \\ NP \rightarrow NP \ CP \\ VP \rightarrow AdvP \ VP \\ VP \rightarrow VP \ PP \\ AP \rightarrow AdvP \ AP \end{array}$$

$$\alpha \rightarrow \alpha \ \text{Coord} \ \alpha$$

(for $\alpha = D, V, N, A, P, C, Adv, VP, NP, DP, AP, PP, AdvP, S, CP$)

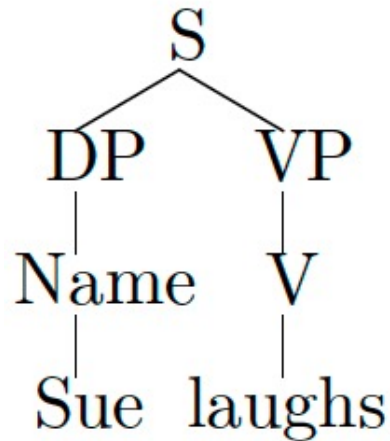
Top-down backtrack recognition

Lexicon:

D →	$\left\{ \begin{array}{c} \text{the} \\ \text{a} \\ \text{some} \\ \text{every} \\ \text{one} \\ \text{two} \end{array} \right\}$	A →	$\left\{ \begin{array}{c} \text{gentle} \\ \text{clear} \\ \text{honest} \\ \text{compassionate} \\ \text{brave} \\ \text{kind} \end{array} \right\}$	N →	$\left\{ \begin{array}{c} \text{student} \\ \text{teacher} \\ \text{city} \\ \text{university} \\ \text{beer} \\ \text{wine} \end{array} \right\}$	V →	$\left\{ \begin{array}{c} \text{laughs} \\ \text{cries} \\ \text{praises} \\ \text{criticizes} \\ \text{says} \\ \text{knows} \end{array} \right\}$	Adv →	$\left\{ \begin{array}{c} \text{happily} \\ \text{sadly} \\ \text{impartially} \\ \text{generously} \end{array} \right\}$
Name →	$\left\{ \begin{array}{c} \text{Bill} \\ \text{Sue} \\ \text{José} \\ \text{Maria} \\ \text{Presidents Day} \\ \text{Tuesday} \end{array} \right\}$	Pronoun →	$\left\{ \begin{array}{c} \text{he} \\ \text{she} \\ \text{it} \\ \text{her} \\ \text{him} \end{array} \right\}$	P →	$\left\{ \begin{array}{c} \text{in} \\ \text{on} \\ \text{with} \\ \text{by} \\ \text{to} \\ \text{from} \end{array} \right\}$	C →	$\left\{ \begin{array}{c} \text{that} \\ \epsilon \\ \text{whether} \end{array} \right\}$	Coord →	$\left\{ \begin{array}{c} \text{and} \\ \text{or} \\ \text{but} \end{array} \right\}$

Top-down backtrack recognition

Derive: '*Sue laughs*'



<u>step</u>	<u>predicted</u>	<u>input</u>
0.	S	Sue laughs
1.	DP VP	Sue laughs

<u>step</u>	<u>predicted</u>	<u>input</u>
0.	S	Sue laughs
1.	DP VP	Sue laughs
2a.	D NP VP	Sue laughs
2b.	NP VP	Sue laughs
2c.	Name VP	Sue laughs
2d.	Pronoun VP	Sue laughs
2e.	DP Coord DP VP	Sue laughs

expand: DP

Top-down backtrack recognition

<u>step</u>	<u>predicted</u>	<u>input</u>			
0.	S	Sue laughs			
1.	DP VP	Sue laughs			
2a.	D NP VP	Sue laughs			
2b.	NP VP	Sue laughs			
2c.	Name VP	Sue laughs			
2d.	Pronoun VP	Sue laughs			
2e.	DP Coord DP VP	Sue laughs			
3aa.	D Coord D NP VP	Sue laughs			
3ab.	the NP VP	Sue laughs			
3ac.	a NP VP	Sue laughs			
3ad.	some NP VP	Sue laughs			
3ae.	every NP VP	Sue laughs			
3af.	one NP VP	Sue laughs			
3ag.	two NP VP	Sue laughs			
3aa.	D Coord D NP VP		3aa.	D Coord D NP VP	
4aaa.	D Coord D Coord D NP VP		4aaa.	D Coord D Coord D NP VP	Sue laughs
4aab.	the Coord D NP VP		4aab.	the Coord D NP VP	Sue laughs
4aab.	a Coord D NP VP		4aab.	a Coord D NP VP	Sue laughs
4aab.	some Coord D NP VP		4aab.	some Coord D NP VP	Sue laughs
4aab.	every Coord D NP VP		4aab.	every Coord D NP VP	Sue laughs
4aab.	one Coord D NP VP		4aab.	one Coord D NP VP	Sue laughs
4aab.	two Coord D NP VP		4aab.	two Coord D NP VP	Sue laughs

The process will never terminate!

Left recursion: category X can contain another category X as its first element

Top-down backtrack recognition

Remove left recursion:

$$\begin{array}{l} S \rightarrow DP \ VP \\ DP \rightarrow \left\{ \begin{array}{l} (D) \ NP \\ Name \\ Pronoun \end{array} \right\} \\ NP \rightarrow N \ (PP) \end{array} \quad \begin{array}{l} VP \rightarrow V \ (DP) \ \left(\left\{ \begin{array}{l} PP \\ CP \\ VP \end{array} \right\} \right) \\ PP \rightarrow P \ (DP) \\ AP \rightarrow A \ (PP) \\ CP \rightarrow C \ S \\ AdvP \rightarrow Adv \end{array} \quad \begin{array}{l} NP \rightarrow AP \ NP \\ \del{NP \rightarrow NP \ PP} \\ \del{NP \rightarrow NP \ CP} \\ VP \rightarrow AdvP \ VP \\ \del{VP \rightarrow VP \ PP} \\ AP \rightarrow AdvP \ AP \end{array}$$

~~$\alpha \rightarrow \alpha \text{ Coord } \alpha$~~

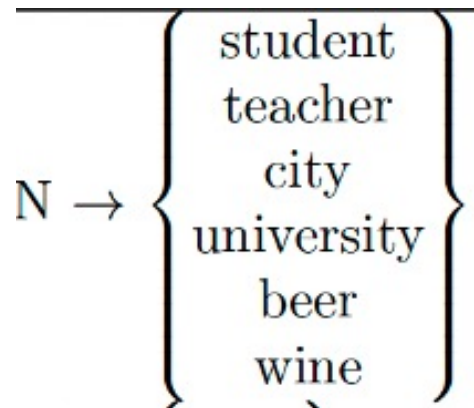
(for $\alpha = D, V, N, A, P, C, Adv, VP, NP, DP, AP, PP, AdvP, S, CP$)

Top-down backtrack recognition

step	predicted	input
0.	S	Sue laughs
1.	DP VP	Sue laughs
2a.	D NP VP	Sue laughs
2b.	NP VP	Sue laughs
2c.	Name VP	Sue laughs
2d.	Pronoun VP	Sue laughs
3aa.	the NP VP	Sue laughs
3ab.	a NP VP	Sue laughs
3ac.	some NP VP	Sue laughs
3ad.	every NP VP	Sue laughs
3ae.	one NP VP	Sue laughs
3af.	two NP VP	Sue laughs

step	predicted	input
0.	S	Sue laughs
1.	DP VP	Sue laughs
2a.	D NP VP	Sue laughs
2b.	NP VP	Sue laughs
2c.	Name VP	Sue laughs
2d.	Pronoun VP	Sue laughs
3ba.	N VP	Sue laughs
3bb.	N PP VP	Sue laughs

step	predicted	input
0.	S	Sue laughs
1.	DP VP	Sue laughs
2a.	D NP VP	Sue laughs
2b.	NP VP	Sue laughs
2c.	Name VP	Sue laughs
2d.	Pronoun VP	Sue laughs
3ca.	Bill VP	Sue laughs
3cb.	Sue VP	Sue laughs
3cc.	José VP	Sue laughs
3cd.	Maria VP	Sue laughs
3ce.	Presidents Day VP	Sue laughs
3cf.	Tuesday VP	Sue laughs



scan: Sue

4cb'	VP	Sue laughs
------	----	------------

Top-down backtrack recognition

4cb'	VP	Sue laughs
4cb'a	V	Sue laughs
4cb'b	V DP	Sue laughs
4cb'c	V PP	Sue laughs
4cb'd	V CP	Sue laughs
4cb'e	V VP	Sue laughs
4cb'f	V DP PP	Sue laughs
4cb'g	V DP CP	Sue laughs
4cb'h	V DP VP	Sue laughs
4cb'aa	laughs	Sue laughs

scan: laughs

4cb' VP Sue laughs

Problems of top-down parsing

- Slow: **expand** requires looping through the whole list of rules
- Non-terminating with left-recursion

To do

- Do HW3
- Optional reading: **NLTK** Ch8:1-3; **SLP** Ch12.1-4;
Stabler Ch1