

- *Parsing* is the task of analyzing a string of symbols to discover its (inherent) structure
- Typically, the structure (and the valid strings in the language) is defined by a *grammar*
- The output of a parser is a structured representation of the input string, often a *tree*
- *Recognition* is an intimately related task which determines whether a given string is in a language

Ingredients of a parser

(for natural language parsing)

- A formal grammar defining a language of interest
- An algorithm that (efficiently) verifies whether a given string is in the language (*recognizer*) and enumerates the grammar rules used for verification (*parser*)
- A system for ambiguity resolution (not in this course)

Grammars

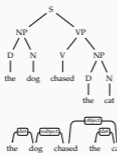
- A grammar is a finite specification of a possibly infinite language
- The most commonly studied type of grammars are *phrase structure grammars*
- Analysis using context-free grammars result in *constituency* or *phrase structure trees*



$S \rightarrow NP VP$ $NP \rightarrow D N$ $VP \rightarrow V NP$
 $V \rightarrow$ chased $D \rightarrow$ the $N \rightarrow$ cat $N \rightarrow$ dog

Why study parsing?

- In general, it is an intermediate step for interpreting sentences
- Applications include:
 - Compiler construction
 - Grammar checking
 - Sentiment analysis
 - Information (e.g., relation) extraction
 - Argument mining
 - ...



Different ways to represent a context-free parse



Sentential form	derivation
S	(start)
NP VP	$S \rightarrow NP VP$
Pm VP	$NP \rightarrow Pm$
I VP	$Pm \rightarrow I$
I V NP	$VP \rightarrow V NP$
I saw NP	$V \rightarrow$ saw
I saw Pm ₁ N	$NP \rightarrow Pm_1 N$
I saw her N	$Pm_1 \rightarrow$ her
I saw her duck	$N \rightarrow$ duck

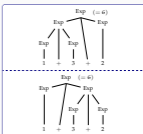
(Labelled) brackets: $\left[\left[\left[\text{I saw} \right] \text{ her duck} \right] \right]$

Relation between different representations

- The parse tree and the bracket representation is equivalent
 - parse trees are easier to read by humans
 - brackets are easier for computers
 - brackets are the typical representation for treebanks
- A parse tree (or bracket representation) can be obtained with a different order of production rules

Grammars and ambiguity

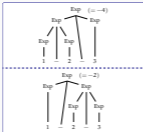
$Exp \rightarrow n$
 $Exp \rightarrow Exp + Exp$
(terminal symbol 'n' stands for any number)



- If a grammar is ambiguous, some sentences produce multiple analyses
- If the resulting analysis lead to the same semantics, the ambiguity is *spurious*

Grammars and ambiguity

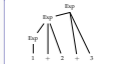
$Exp \rightarrow n$
 $Exp \rightarrow Exp - Exp$



- (terminal symbol 'n' stands for any number)
- Is this ambiguity spurious?
 - If different structures yield different semantics, the ambiguity is *essential*

Ambiguity can be removed from a grammar

$Exp \rightarrow n$
 $Exp \rightarrow Exp + n$
(terminal symbol 'n' stands for any number)

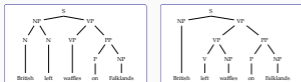


- The grammar above does not have the ambiguity of

$Exp \rightarrow n$
 $Exp \rightarrow Exp + Exp$

- Both grammars define the same language

Natural languages are ambiguous



- The grammars we define have to distinguish between two different structures
- We need methods for ranking analyses

Top-down parsing

general idea

- Start from S, find a sequence of derivations that yield the sentence
- This is simply the same as the generation procedure we discussed earlier
- Attempt to generate all strings from a grammar, but allow only the productions that 'produce' the input string

