

Introduction to Parsing

Parsing

ISCL-BA-06

Çağrı Çöltekin
ccolt@infsa.uni-tuebingen.de

University of Tübingen
Seminar für Sprachwissenschaft

Winter Semester 2020/21

version: 2020-09-02 11:30

What is parsing?

- Parsing is the task of assigning a structure to a given sentence
- It is related to recognition: typically we follow the steps taken during derivation to obtain the structure
- From a different perspective, parsing is the inverse of the generation task
- Note: we focus on context-free parsing – the structures we build/recover are trees

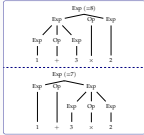
Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 1 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Why do we need parsing?

- The formal approach to languages as sets emphasizes recognition
 - a string is whether in the language or not
- Parsing is in general a step for semantics
 - we cannot assign semantics without structure



Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 2 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Overview

- Representation context-free analyses and parse trees
- Ambiguity
- Top-down parsing
- Bottom-up parsing
- General overview of the parsing methods
- Representing parsing methods: parse forests
- Parsing and semantics

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 3 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Different ways to represent a context-free parse



Semantical 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: $\left[\left[\left[\text{NP } I \right] \right]_{\text{VP}} \left[\text{V saw} \right] \left[\left[\text{NP } \left[\text{Pm}_1 \text{ her} \right] \left[\text{N duck} \right] \right] \right] \right]$

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 4 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

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

Ç. Çöltekin, INF | University of Tübingen

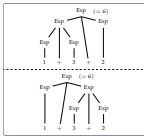
Winter Semester 2020/21 5 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

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*



Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 6 / 20

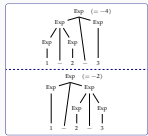
Introduction Ambiguity Top-down parsing Bottom-up parsing

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*



Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 7 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Languages and ambiguity

- A language is ambiguous if there is no unambiguous grammar that can produce it
- For example, the language $a^n b^n c^m \mid a^n b^m c^n$ is ambiguous
 - The strings of the form $a^k b^k c^k$ could be generated by either part of the language definition
- Note: do not confuse ambiguity with different derivations leading to same analysis
 - Ambiguity results in different structures
 - Multiple derivations with the same structure is related to the mechanism used for obtaining the derivations

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 8 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Ambiguity can be removed from a grammar

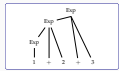
if the language is not ambiguous

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

- This one does not have the ambiguity of

Exp \Rightarrow n
Exp \Rightarrow Exp + Exp

- Both grammars define the same language

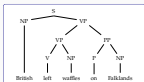
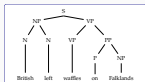


Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 9 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Natural languages are ambiguous



- The grammars we define have to distinguish between two different structures

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 10 / 20

Introduction Ambiguity Top-down parsing Bottom-up parsing

Top-down parsing

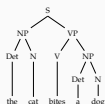
general idea

- 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 the parse grammar, but allow productions that only leads to the input string

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2020/21 11 / 20

Top-down: demonstration



S → NP VP
 NP → Det N
 VP → V NP
 VP → V
 Det → a
 Det → the
 N → cat
 N → dog
 V → bites

From demonstration to parsing

- There may be multiple production applicable
- We need an automatic mechanism to select the correct productions
- We have two actions:
 - predict generate a hypothesis based on the grammar
 - match when a terminal is produced, check if it matches with the terminal in the expected position
 - if matched, continue
 - otherwise, backtrack
- if we eliminate all non terminals, and the complete input string is matched, then parsing successful

Top-down parsing: another demonstration

Introduction Analytic Top-down parsing Bottom-up parsing

the grammar
S → NP VP
NP → Det N
VP → V NP
VP → V
Det → a
Det → the
N → cat
N → dog
V → bites

parse: the cat bites a dog

matched	goal	production
	S	S → NP VP
	NP VP	NP → Det VP
	Det N VP	Det → a ✗
	Det N VP	Det → the ✓
the	N VP	N → cat ✗
the cat	N VP	N → dog ✗
the cat	VP	VP → V
the cat bites	V	V → bites ✓
the cat bites		(rest at the end) ✗
the cat	V NP	VP → V NP
the cat bites	NP	V → bites ✓
the cat bites	Det N	NP → Det N
the cat bites a	N	Det → a ✓
the cat bites a dog	Det	Det → dog ✓

Note that the valid productions yield the parse tree.

Top-down parsing: problems and possible solutions

Introduction Analytic Top-down parsing Bottom-up parsing

- Trial-and-error procedure leads to exponential time parsing
- But lots of repeated work: dynamic programming may help avoid it
- What happens if we had a rule like NP → NP PP
- some rules may cause infinite loops
- Notice that if we knew which terminals are possible as the initial part of a non-terminal symbol, we can eliminate the unsuccessful matches earlier

Bottom-up parsing
general idea

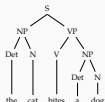
Introduction Analytic Top-down parsing Bottom-up parsing

- Start from the input symbol, and try to *reduce* the input to start symbol
- We need to match parts of the sentential form (starting from the input) to the RHS of the grammar rules
- While top-down process relies on *productions* the bottom-up process relies on *reductions*

production	NP	V	NP	reduction
	Det	N	Det	N
	the	cat	bites	a dog

Bottom-up: demonstration

Introduction Analytic Top-down parsing Bottom-up parsing



S → NP VP
 NP → Det N
 VP → V NP
 VP → V
 Det → a
 Det → the
 N → cat
 N → dog
 V → bites

A (first) introduction to shift-reduce parsing

Introduction Analytic Top-down parsing Bottom-up parsing

- We keep two data structures:
 - a stack for the (partially) reduced sentential form
 - an input queue that contains only terminal symbols

NP V	a dog
------	-------

- We use two operations:

shift shifts a terminal to stack

NP V	a dog	shift	NP V a	dog
------	-------	-------	--------	-----

reduce when top symbols on stack match a RHS, replace them with the LHS of the rule

NP V	a dog	reduce	NP VP	a dog
------	-------	--------	-------	-------

Shift-reduce (bottom-up) parsing a demonstration

Introduction Analytic Top-down parsing Bottom-up parsing

stack	input	rule	stack	input	rule
	the cat bites a dog	shift	NP V	a dog	shift
the	cat bites a dog	Det → the	NP V a	a dog	Det → a
Det	cat bites a dog	shift	NP V Det	dog	shift
Det cat	bites a dog	N → cat	NP V Det dog		N → dog
NP	bites a dog	NP → Det N	NP V Det N		NP → Det N
NP	bites a dog	shift	NP V NP		VP → V NP
NP bites	a dog	V → bites	NP VP		S → NP VP
NP V	a dog	VP → V	S		(done)
NP VP	a dog	S → NP VP			
S	a dog	shift			
S a	dog	shift			
S a dog		Det → A			
S Det dog		N → dog			
S Det N		NP → Det N			
S NP		(stack)			

- All input reduced to S, accept
- Rules form the parse tree

Summary

Introduction Analytic Top-down parsing Bottom-up parsing

- Parsing can be formulated as a top-down or bottom-up search (the search may also be depth-first or breadth first)
- Naive parsing algorithms are inefficient (exponential time complexity)
- There are some directions: dynamic programming, filtering
- Suggested reading for this part: Grune and Jacobs (2007, ch.3)

Next:

- Bottom-up chart parsing: CKY algorithm
- Suggested reading: Grune and Jacobs (2007, section 4.2), Jurafsky and Martin (2009, draft 3rd ed, section 13.2)

Acknowledgments, references, additional reading material

- Please read Grune and Jacobs (2007) chapter 3, a big part part of the lecture follows this chapter

Grune, Dirk and Gerald J. Jacobs (2007). *Parsing Techniques: A Practical Guide*, second. *Monographs in Computer Science*. The first edition is available at http://aladdin.cs.cmu.edu/Books/PTAPG_2nd_Edition/BookIndex.pdf. Springer New York, USA, 9780387086888.

Derivations, David and James H. Martin (2009). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, second. Pearson Education, Inc. ISBN-10: 0-13-035858-1, URL: <http://web.stanford.edu/~jurafsky/slp4/>