

Syntax

- Introduction
- Formal Grammars
- Grammars for NLP

Introduction to Syntax₁

- Syntax describes regularity and productivity of a language making explicit the structure of sentences
- Goal of syntactic analysis (parsing):
 - Detect if a sentence is correct
 - Provide a syntactic structure of a sentence.

Introduction to Syntax₂

- It refers to the way words are arranged together
- Basic ideas related to syntax

- **Contiguency**

Groups of words may behave as a single unit or phrase: a constituent. Example: ***Noun phrase***

- **Grammatical relations**

Formalization of ideas from traditional grammar. Example: ***subject*** and ***object***

- **Subcategorization and dependency relations**

Relations between words and phrases

Ex: Verb ***want*** followed by an ***infinitive verb***

Introduction to Syntax₃

- Regular languages and part of speech refers to the way words are arranged together but cannot support easily: **Contiguency, Grammatical relations and Subcategorization and dependency relations**
- They can be modelled by grammars based on **context-free grammars**
- Context-free grammars is a formalism powerful enough to represent complex relations and can be efficiently implemented.
- Context-free grammars are integrated in many language applications.

Introduction to Syntax₄

- A **Context-free grammar** consists of a **set of rules** or productions, each expressing the ways the symbols of the language can be grouped together, and a **lexicon** of words
- An example of a set of rules expressing
 - **NP** (noun phrase) **can be either a ProperNoun or a determiner (Det) folowed by a Nominal**
 - **A nominal can be one or more Nouns**
 - NP → Det Nominal**
 - NP → ProperNoun**
 - Nominal → Noun | Nominal Noun**
 - Other rules can be added:

NLP syntax_ **Det → a** **Det → the** **Noun → flight**

Formal Grammars₁

- Text (string)
 - Free monoid over a vocabulary with the concatenation operation (\cdot)
- Vocabulary (V), set of words (w)
 $w \in V$
- Language Models (LM)
 - Probability distribution over the texts
- Language (L), set of sentences (s)
 $s \in L$
 $L \subset V^*$ usually infinite
- $s = w_1, \dots, w_N$
 $P(s)$ Probability of s

Formal Grammars₂

- Naive Implementation of a LM
 - Enumerate $s \in L$
 - Compute $p(s)$, e.g. counting occurrences
 - Parameters of the model $|L|$
- But ...
 - L is usually not enumerable
 - How to estimate the parameters?

- Simplifications

- History

- $h_i = \{ w_i, \dots, w_{i-1} \}$
- Usually up to 5-grams
- Google's n-grams

- Markov Models

$$p(s) = p(w_1 \dots w_N) = P(w_1^N) = \prod_{i=1}^N p(w_i | h_i)$$

Formal Grammars₃

Language (L) over vocabulary V $L \subset V^*$

How to define L, i.e. how to decide if $s \in L$?

- If we use a LM: **If $p(s) > 0$ then $s \in L$**
- The usual way: Generative Approach
 - A sentence is correct if it is grammatical (according to a grammar)
 - Getting a grammar **G** such that $s \in L_G$ being L_G the language generated (or recognized) by the grammar **G**
 - There are many types of grammars, the most used are the Phrase Structure Grammars (PSG) or Context Free Grammar (CFG).
 - **$G = \langle N, \Sigma, P, S \rangle$**
 - Set of N non terminal symbols
 - Set of Σ terminal symbols
 - Set of P productions or rules
 - $S \in N$, axiom

Grammars for NLP₁

Example of CFG

$G_1 = \langle N_1, T_1, P_1, \text{SENTENCE} \rangle$

$N_1 = \{\text{SENTENCE}, \text{NP}, \text{VP}, \text{RNP}, \text{PP}\}$

$T_1 = \{\text{det}, \text{n}, \text{adj}, \text{vi}, \text{vt}, \text{prep}\}$

$P_1 = \{$

1 SENTENCE \rightarrow NP VP.

2 NP \rightarrow det n RNP.

3 NP \rightarrow n RNP.

4 NP \rightarrow np RNP.

5 RNP \rightarrow ϵ .

6 RNP \rightarrow PP RNP.

7 RNP \rightarrow adj RNP.

8 VP \rightarrow vi.

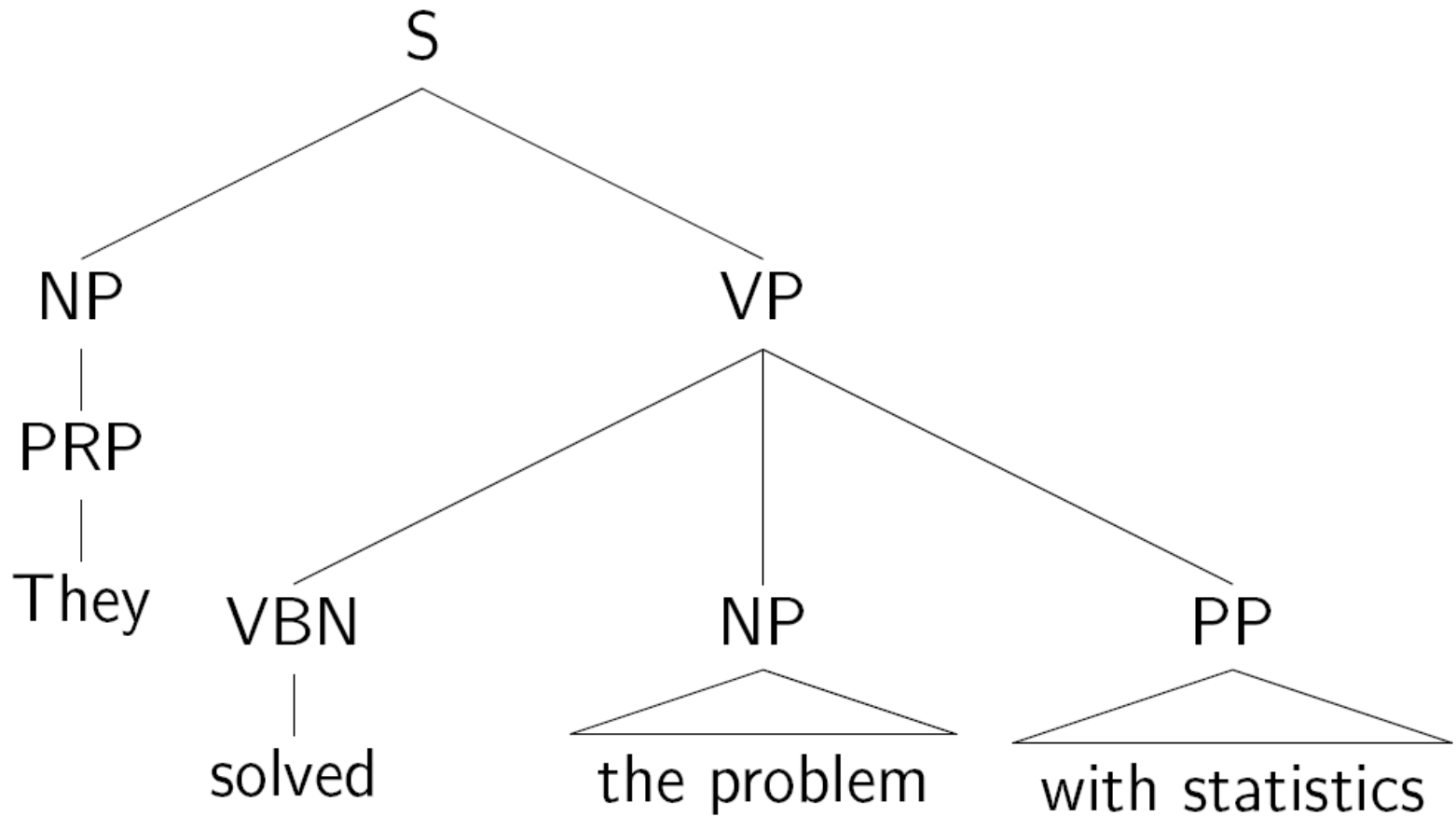
9 VP \rightarrow vt NP.

10 PP \rightarrow prep NP.

}

Grammars for NLP₂

Parse tree

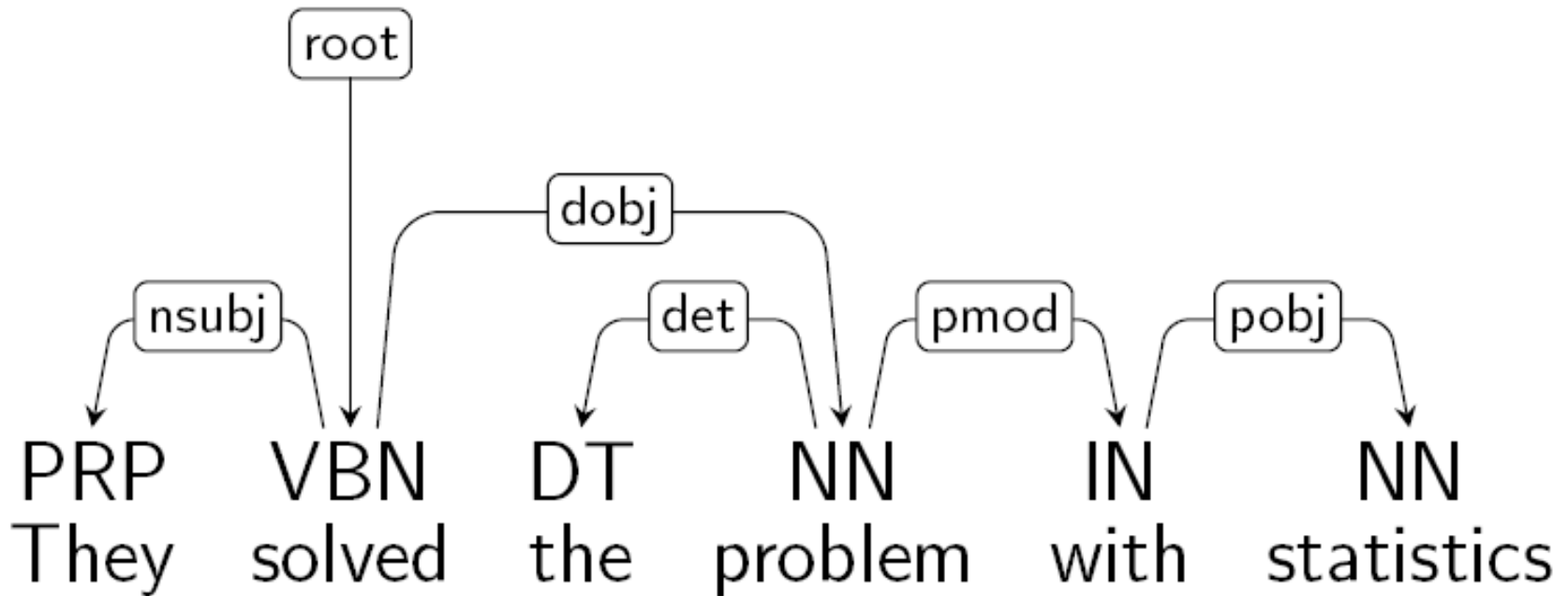


Grammars for NLP₃

- Forms of expressing syntactic structure:
 - Constituent or Phrase Structure (derivation tree = parse tree)
 - Dependency Structure
 - Tree adjoining grammars (TAG)
 - Transformational grammars
 - Systemic grammars
 - Logical Form

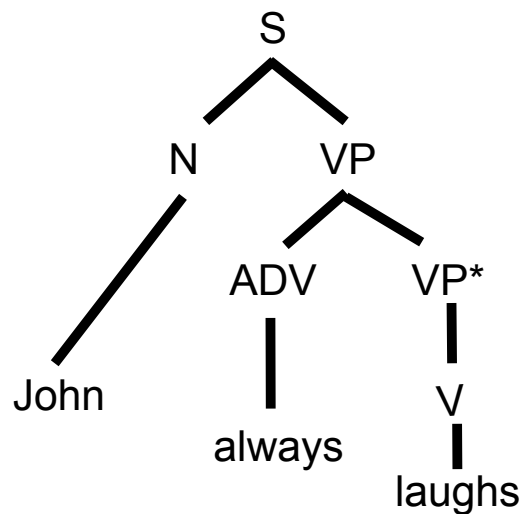
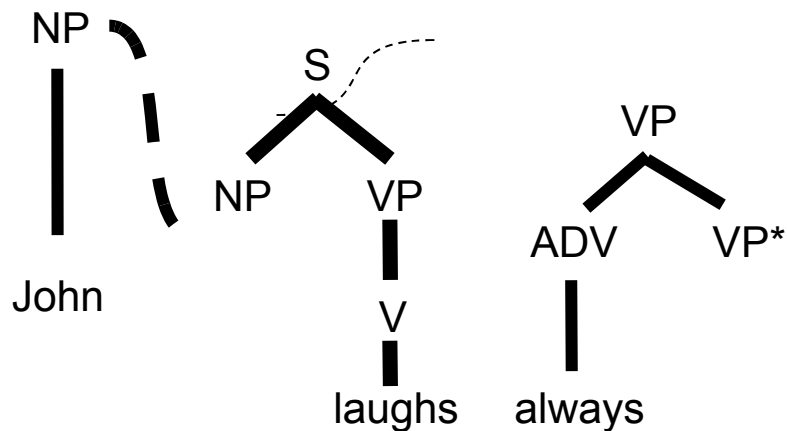
Grammars for NLP₄

Dependency tree



Grammars for NLP₅

Tree Adjoining Grammars (TAG)



Grammars for NLP₆

Logical Form (close to semantics)

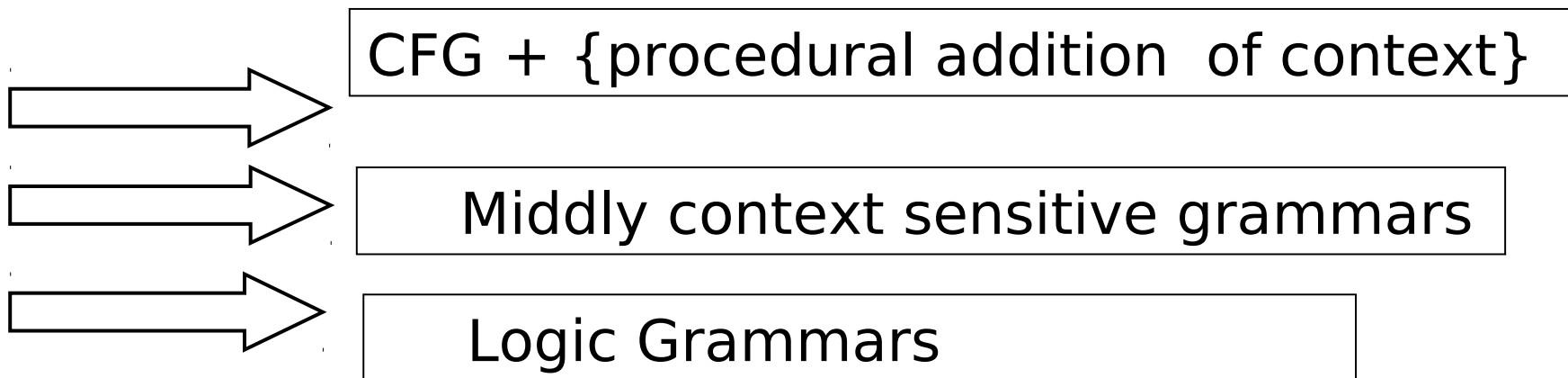
The cat eats fish

$(\exists X \text{ and } (\text{cat } (X),$
 $\quad (\exists Y \text{ and } (\text{fish } (Y),$
 $\quad \quad \text{eat}(X, Y))))))$

Grammars for NLP₇

Expressivity of the grammar

- Minimum: Context free grammar (CFG)
- ¿Is NL context free?
- ¿Sufficient? NO (usually)
- Solution



Grammars for NLP₈

- For practical application in NLP
 - Regular grammars (RG)
 - Equivalent to Finite State Automata and regular expressions
 - Parsed in quadratic time $O(n^2)$
 - Context free grammars (CFG)
 - Equivalent to push-down automata
 - Parsed in cubic time $O(n^3)$
 - Mildly context sensitive grammars
 - Tree Adjoining Grammars (TAG)
 - Linear Context Free Rewriting Systems

Grammars for NLP₉

Obtaining the grammar

- Definition of **T** from the tagset
- Definition of **V**
- Grammar rules **P**
 - Manually
 - Automatically
 - Grammatical inference
 - Semi- automatically

Grammars for NLP₁₀

Improving the grammar

- Transformations of grammars for getting equivalent ones :
 - Removing symbols and productions that cannot be reached
 - Removing unary productions
 - Removing ϵ productions
 - Lexicalization
- Approximation of CFG by RG

Grammars for NLP₁₁

Properties of parsers

- Soundness
 - The output of the parsing is correct according to the grammar
- Termination
 - Any parsing process terminates
- Completeness
 - A parser is complete if given a grammar and a sentence it is sound, produces all the correct parse trees and terminates