

# Introduction to Algorithms and Data Structures

## Lecture 22: Context-free languages and grammars

John Longley

School of Informatics  
University of Edinburgh

3 February 2026

# Algorithms and data structures in Language Processing

By 'Language Processing', we mean the processing of both

- ▶ **Artificial** (computer) languages: e.g. Java, Python, HTML, . . .
- ▶ **Natural** (human) languages: e.g. English, Greek, Japanese.

Despite major differences, there's a body of theory common to both kinds: e.g. the formal theory of **generative grammar**.

(Pioneered by Sanskrit scholar **Pāṇini** ~ 500–400 BCE.

Introduced into modern linguistics by **Noam Chomsky** around 1957.)



**Pāṇini**

Photo by Jameela P. – Own work



**Noam Chomsky**

Photo by Andrew Rusk from Toronto, Canada

We'll be looking at this theory (specifically **context-free grammars**) and some of the algorithms/data structures involved.

## Aside: Chomsky in an age of LLMs ...

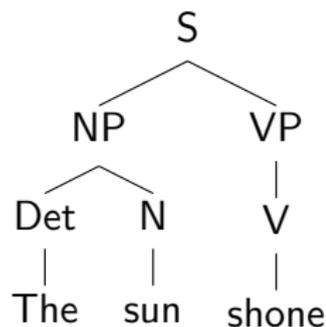
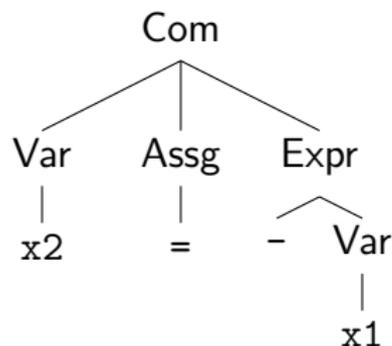
Today's NLP landscape is dominated by **Large Language Models**. These aren't directly based on anything like Chomsky's theories.

However, there are many reasons why Chomsky-style grammars are still important/interesting ...

- ▶ LLMs require huge computing power and massive training. If we're working with a suitably constrained fragment of NL, formal grammars can offer efficient, lightweight algorithms.
- ▶ Unlike LLMs, formal grammars can **guarantee** grammatically correct output. There's a whole area of **text normalization**, where the output of an LLM is run through a grammar checker to enforce certain hard rules.
- ▶ Evidence from probing the attention scores of LLMs suggests that they implicitly 'learn' something like a Chomsky-style grammar. So formal grammars offer a framework for understanding (one aspect of) how LLMs are working.

## Common idea: **syntax trees**

A syntax tree displays the grammatical 'constituent structure' of a language text:

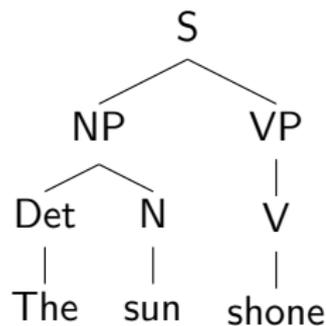
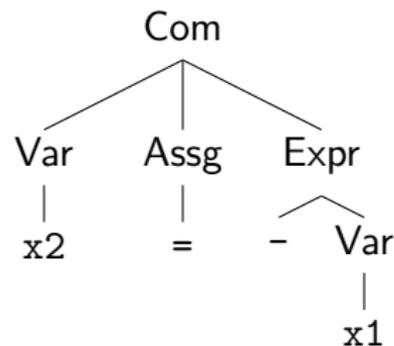


Constructing the tree is an important preliminary to many LP tasks.

A **grammar** for a language is a bunch of rules specifying what syntax trees are possible (and hence what strings are possible).

- ▶ **This lecture:** Defining languages via **context-free grammars**.
- ▶ **Next lecture:** How to find the syntax tree for a given program or sentence (**parsing algorithms**).

## Those trees again



There are two kinds of symbols here:

- ▶ Symbols at the leaves are called **terminals**: these are the basic units from which sentences of the language are built.
- ▶ Symbols at internal nodes are called **non-terminals**: they don't themselves appear in sentences of the language, but name various kinds of 'sub-phrases'.

## Context-free grammars: first example

Let's give a little grammar for **arithmetic expressions**, e.g.

$6 + 7$        $5 * (x + 3)$        $x * ((z * 2) + y)$        $8$        $z$

**Terminals:**  $+, *, (, ), x, y, z, 0, \dots, 9$ .

**Non-terminals:**  $\text{Exp}, \text{Var}, \text{Num}$ .

We designate the non-terminal  $\text{Exp}$  as the **start symbol**.

**Rules:**

$\text{Exp} \rightarrow \text{Exp} + \text{Exp}$

$\text{Exp} \rightarrow \text{Exp} * \text{Exp}$

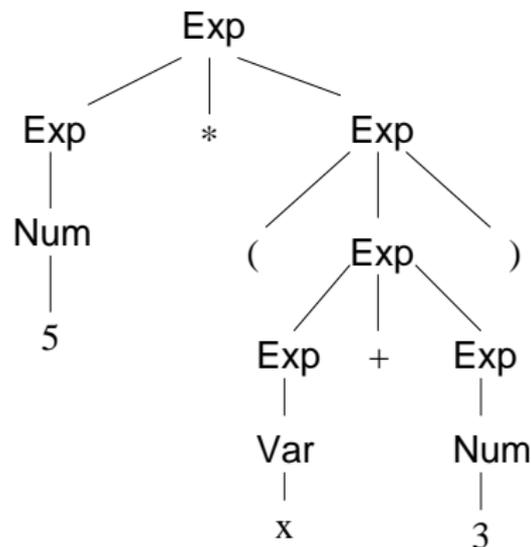
$\text{Exp} \rightarrow \text{Var} \mid \text{Num} \mid (\text{Exp})$

$\text{Var} \rightarrow x \mid y \mid z$

$\text{Num} \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid$   
 $5 \mid 6 \mid 7 \mid 8 \mid 9$

## Generating a syntax tree from a grammar

Beginning with the start symbol, we can grow syntax trees by repeatedly expanding non-terminal symbols using these rules. E.g.:



Exp  $\rightarrow$  Exp + Exp  
Exp  $\rightarrow$  Exp \* Exp  
Exp  $\rightarrow$  Var | Num | ( Exp )  
Var  $\rightarrow$  x | y | z  
Num  $\rightarrow$  0 | ... | 9

This generates  $5 * (x + 3)$  as a legal expression of the language.

# The language defined by a grammar

We can generate **infinitely many** strings from this (finite) grammar!

The **language** defined by the grammar is (by definition) the set of all **strings of terminals** that can be obtained via such a tree.

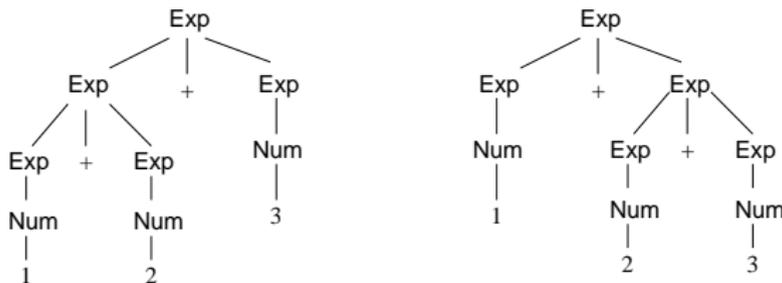
Or as a more ‘machine-oriented’ alternative, we can consider **derivations** involving **sentential forms** (i.e. strings of terminals and non-terminals reachable from the start symbol):

Exp  $\Rightarrow$  Exp \* Exp  
 $\Rightarrow$  Num \* Exp  
 $\Rightarrow$  Num \* (Exp)  
 $\Rightarrow$  Num \* (Exp + Exp)  
 $\Rightarrow$  5 \* (Exp + Exp)  
 $\Rightarrow$  5 \* (Exp + Num)  
 $\Rightarrow$  5 \* (Var + Num)  
 $\Rightarrow$  5 \* (x + Num)  
 $\Rightarrow$  5 \* (x + 3)

Exp  $\rightarrow$  Exp + Exp  
Exp  $\rightarrow$  Exp \* Exp  
Exp  $\rightarrow$  Var | Num | ( Exp )  
Var  $\rightarrow$  x | y | z  
Num  $\rightarrow$  0 |  $\dots$  | 9

## Structural ambiguity

Note that strings such as  $1+2+3$  may be generated by more than one tree (**structural ambiguity**):



This might seem 'harmless' ... but what about  $1+2*3$  ?

- ▶ In computer languages, structural ambiguity is typically avoided by careful design of the grammar (e.g. enforcing that  $*$  takes precedence over  $+$ ).
- ▶ In natural languages, structural ambiguity is a **fact of life**. E.g. *I saw a man with a telescope.*

## Puzzle (to try later)

Grammar rules again:

$$\text{Exp} \rightarrow \text{Exp} + \text{Exp}$$
$$\text{Exp} \rightarrow \text{Exp} * \text{Exp}$$
$$\text{Exp} \rightarrow \text{Var} \mid \text{Num} \mid ( \text{Exp} )$$
$$\text{Var} \rightarrow x \mid y \mid z$$
$$\text{Num} \rightarrow 0 \mid \dots \mid 9$$

How many possible syntax trees are there for the string below?

$$1 + 2 + 3 + 4$$

## Second example: comma-separated lists

Consider lists of (zero or more) alphabetic characters, separated by commas. E.g.:

$\epsilon$        $a$        $e,d$        $q,w,e,r,t,y$

These can be generated by the following grammar.

**Terminals:**       $a, \dots, z, ,$   
**Non-terminals:** List, Char, Tail  
**Start symbol:** List

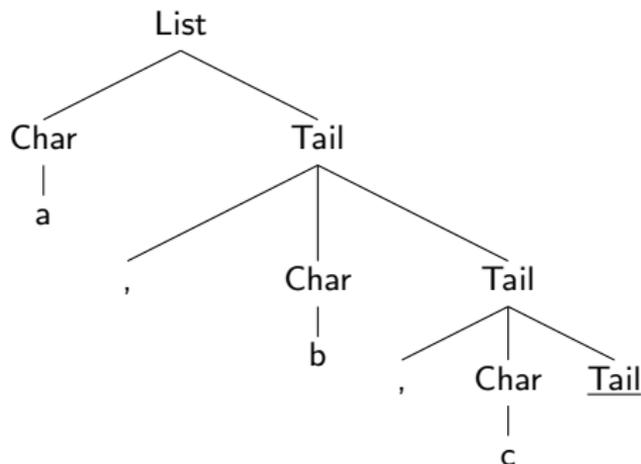
List  $\rightarrow \epsilon \mid \text{Char Tail}$   
Char  $\rightarrow a \mid \dots \mid z$   
Tail  $\rightarrow \epsilon \mid , \text{Char Tail}$

(Note the rules with **empty** right hand side, indicated by  $\epsilon$ .)

## Syntax trees for comma-separated lists

List  $\rightarrow \epsilon \mid \text{Char Tail}$   
Char  $\rightarrow a \mid \dots \mid z$   
Tail  $\rightarrow \epsilon \mid , \text{Char Tail}$

Here's the syntax tree for the list *a, b, c*:



(Note how we've indicated the application of an ' $\epsilon$ -rule'.)

## Context-free grammars: formal definition

A **context-free grammar** (CFG)  $\mathcal{G}$  consists of

- ▶ a finite set  $\Sigma$  of **terminals**,
- ▶ a finite set  $N$  of **non-terminals**, disjoint from  $\Sigma$ ,
- ▶ a choice of **start symbol**  $S \in N$ ,
- ▶ a finite set  $P$  of **productions** of the form  $X \rightarrow \alpha$ , where  $X \in N$ ,  $\alpha \in (\Sigma \cup N)^*$ .

## The language arising from a CFG (optional slide)

A **sentential form** is any sequence of terminals and nonterminals that can appear in a derivation starting from the start symbol.

**Formal definition:** The set of **sentential forms** derivable from  $\mathcal{G}$  is the smallest set  $\mathcal{S}(\mathcal{G}) \subseteq (N \cup \Sigma)^*$  such that

- ▶  $S \in \mathcal{S}(\mathcal{G})$
- ▶ if  $\alpha X \beta \in \mathcal{S}(\mathcal{G})$  and  $X \rightarrow \gamma \in P$ , then  $\alpha \gamma \beta \in \mathcal{S}(\mathcal{G})$ .

The **language** associated with grammar is the set of sentential forms that contain only terminals.

**Formal definition:** The **language** associated with  $\mathcal{G}$  is defined by  $\mathbf{L}(\mathcal{G}) = \mathcal{S}(\mathcal{G}) \cap \Sigma^*$ .

A language  $L \subseteq \Sigma^*$  is defined to be **context-free** if there exists some CFG  $\mathcal{G}$  such that  $L = \mathbf{L}(\mathcal{G})$ .

## Assorted remarks

- ▶  $X \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  is simply an **abbreviation** for a bunch of productions  $X \rightarrow \alpha_1, X \rightarrow \alpha_2, \dots, X \rightarrow \alpha_n$ .
- ▶ These grammars are called **context-free** because a rule  $X \rightarrow \alpha$  says that an  $X$  can *always* be expanded to  $\alpha$ , no matter where the  $X$  occurs.  
This contrasts with **context-sensitive** rules, which might allow us to expand  $X$  only in certain contexts, e.g.  $bXc \rightarrow bac$ .
- ▶ Broad intuition: context-free languages allow **nesting of phrase structures to arbitrary depth**. E.g. brackets, begin-end blocks, for-loops, subordinate clauses in English, ...

## A programming language example

Some context-free rules for a little programming language.

```
stmt   → if-stmt | while-stmt | begin-stmt | assg-stmt
if-stmt → if bool-expr then stmt else stmt
while-stmt → while bool-expr do stmt
begin-stmt → begin stmt-list end
stmt-list → stmt | stmt ; stmt-list
assg-stmt → var := arith-expr
bool-expr → arith-expr compare-op arith-expr
compare-op → < | > | <= | >= | == | !=
```

Also need rules for arith-expr and var (in style of earlier grammar).

Grammars like this (often with ::= or : in place of →) are standard in computer language reference manuals. This notation is often called **BNF** (Backus-Naur Form).

## A natural language example

Consider the following lexical classes ('parts of speech') in English:

N	nouns	( <i>alien, cat, dog, house, malt, owl, rat, table</i> )
Name	proper names	( <i>Jack, Susan</i> )
TrV	transitive verbs	( <i>admired, ate, built, chased, killed</i> )
LocV	locative verbs	( <i>is, lives, lay</i> )
Prep	prepositions	( <i>in, on, by, under</i> )
Det	determiners	( <i>the, my, some</i> )

Now consider the following productions (start symbol S):

S	→	NP VP
NP	→	<i>this</i>   Name   Det N   Det N RelCl
RelCl	→	<i>that</i> VP   <i>that</i> NP TrV   NP TrV   NP LocV Prep
VP	→	<i>is</i> NP   TrV NP   LocV Prep NP

## Natural language example in action

Even this modest bunch of rules can generate a rich multitude of English sentences, for example:

- ▶ *this is Jack*
- ▶ *some alien ate my owl*
- ▶ *Susan admired the dog that lay under my table*
- ▶ *this is the dog that chased the cat that killed the rat that ate the malt that lay in the house that Jack built*

## Nesting in natural language

Excerpt from Jane Austen, *Mansfield Park*.

Whatever effect Sir Thomas's little harangue might really produce on Mr. Crawford, it raised some awkward sensations in two of the others, two of his most attentive listeners — Miss Crawford and Fanny. One of whom, having never before understood that Thornton was so soon and so completely to be his home, was pondering with downcast eyes on what it would be *not* to see Edmund every day; and the other, startled from the agreeable fancies she had been previously indulging on the strength of her brother's description, no longer able, in the picture she had been forming of a future Thornton, to shut out the church, sink the clergyman, and see only the respectable, elegant, modernized and occasional residence of a man of independent fortune, was considering Sir Thomas, with decided ill-will, as the destroyer of all this, and suffering the more from . . .

## Reading/browsing

Reading on context-free grammars/languages (sadly not in CLRS):

- ▶ [https://en.wikipedia.org/wiki/Context-free\\_grammar](https://en.wikipedia.org/wiki/Context-free_grammar)  
Aligns quite well with our treatment. Many examples.
- ▶ M. Sipser, *Introduction to the Theory of Computation* (3rd ed), Section 2.1. Online access via UoE library.
- ▶ <https://docs.python.org/3/reference/grammar.html>
- ▶ <https://docs.oracle.com/javase/specs/jls/se7/html/jls-18.html>
- ▶ Treebank Semantics Parsed Corpus (large browsable collection of syntax trees for a variety of English texts):  
<http://www.compling.jp/ajb129/tspc.html>

**Next time:** The **parsing** problem – given a program / expression / sentence, construct its syntax tree. Algorithms for this.

And finally . . .

A musical tribute to Chomsky's groundbreaking book,  
*Syntactic Structures* (1957),  
which introduced the theory of context-free grammars  
and ushered in a new era of linguistics.

'The theory of *Syntactic S.*'

[https://media.ed.ac.uk/playlist/  
dedicated/229853033/1\\_94t0ofs7/1\\_kzedrxoz](https://media.ed.ac.uk/playlist/dedicated/229853033/1_94t0ofs7/1_kzedrxoz)  
(Or follow link from Semester 2 schedule on  
OpenCourse.)