

# Automata Theory :: Context-Free Languages

Jörg Endrullis

Vrije Universiteit Amsterdam

# Context-Free Languages

A grammar  $G = (V, T, S, P)$  is called **context-free** if all production rules are of the form

$$A \rightarrow u$$

where  $A \in V$  and  $u \in (V \cup T)^*$ .

That is, the left-hand side of all rules is a single variable.

A language  $L$  is **context-free** if there is a context-free grammar with  $L = L(G)$ .

Context-free grammars are used to

- describe the **syntax of programming languages**, and
- form the basis of **parsing algorithms**.

## Exercises (1)

**Show that the following language is context-free:**

$$\{a^n b^m c^{2n} \mid n \geq 0, m \geq 0\}$$

We give a context-free grammar for the language:

$$S \rightarrow aScc \mid T$$

$$T \rightarrow bT \mid \lambda$$

The start variable is  $S$ .

## Exercises (2)

Show that the following language is context-free:

$$L = \Sigma^* \setminus \{ww \mid w \in \Sigma^*\}$$

where  $\Sigma = \{a, b\}$ . What shape do words  $w \in L$  have?

- $w$  has **odd length**;
- $w$  has **even length** such that  $w = uv$  with  $|u| = |v|$  and

$$u(m) \neq v(m)$$

for some  $0 \leq m < |u|$ . Let  $n = |u| - m - 1$ . Then  $w$  is in

$$\begin{aligned} & \Sigma^m \{a\} \Sigma^n \Sigma^m \{b\} \Sigma^n \\ \cup & \underbrace{\Sigma^m \{b\} \Sigma^n}_u \underbrace{\Sigma^m \{a\} \Sigma^n}_v \end{aligned} = \begin{aligned} & \Sigma^m \{a\} \Sigma^m \Sigma^n \{b\} \Sigma^n \\ \cup & \underbrace{\Sigma^m \{b\} \Sigma^m}_u \underbrace{\Sigma^n \{a\} \Sigma^n}_v \end{aligned}$$

We give a context-free grammar for the language:

$$\begin{aligned} S &\rightarrow AB \mid BA \mid A \mid B & A &\rightarrow XAX \mid a & X &\rightarrow a \mid b \\ B &\rightarrow XBX \mid b \end{aligned}$$

## Exercises (3)

**Is the following language context-free?**

$$L = \{ ww \mid w \in \Sigma^* \}$$

where  $\Sigma = \{a, b\}$ .

On the previous slide, we have shown that  $\bar{L}$  is context-free.

This language is not context-free.

We will prove this using the context-free pumping lemma.

The class of context-free languages is not **closed under complement**.

## Derivation Trees and Ambiguity

# Rightmost and Leftmost

Let  $G$  be a grammar, and consider a derivation  $S \Rightarrow^* w$ .

- If  $G$  is **(right) linear**,  $w$  contains **at most one variable**.
- If  $G$  is **context-free**,  $w$  can contain **multiple variables**.

Two strategies to choose which variable to expand:

- **leftmost**: always the leftmost variable
- **rightmost**: always the rightmost variable

$$S \rightarrow SaS \mid b$$

Two derivations of  $bab$ :

**leftmost:**  $S \Rightarrow SaS \Rightarrow baS \Rightarrow bab$

**rightmost:**  $S \Rightarrow SaS \Rightarrow Sab \Rightarrow bab$

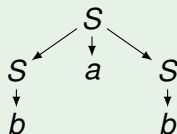
Result depends not on the strategy, but the choice of the rules.

# Derivation Trees

A **derivation tree** for a context-free grammar  $G = (V, T, S, P)$ :

- Nodes have labels from  $V \cup T \cup \{\lambda\}$ . The **root** has label  $S$ .
- A node with label  $A \in V$  can have children labelled
  - $x_1, \dots, x_n$  if there is a rule  $A \rightarrow x_1 \cdots x_n$  with  $n \geq 1$ , or
  - $\lambda$  (single child) if there is a rule  $A \rightarrow \lambda$ .
- Every node with a label from  $T \cup \{\lambda\}$  is a **leaf**.

In the grammar  $S \rightarrow SaS \mid b$ , a **derivation tree** for  $bab$  is:



The labels of the leaves of a derivation tree, read from left to right (skipping  $\lambda$ ), form a word in  $L(G)$ .



Ambiguity

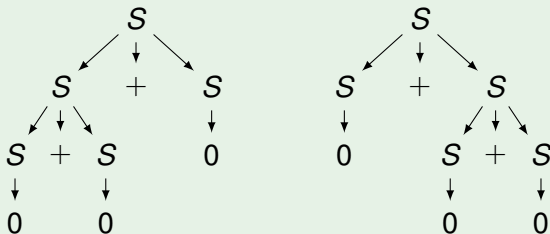
# Ambiguous Grammars

A context-free grammar  $G$  is **ambiguous** if there exists a word  $w \in L(G)$  for which there are multiple derivation trees.

Is the following grammar ambiguous?

$$S \rightarrow S + S \mid 0$$

Yes, there are two derivations trees for  $0 + 0 + 0$ :



# Ambiguous Grammars

Is the following grammar ambiguous?

$$S \rightarrow S + 0 \mid 0$$

No, every word has precisely one derivation tree.

Note that both grammars

$$S \rightarrow S + S \mid 0$$

$$S \rightarrow S + 0 \mid 0$$

generate the same language:

$$\{0(+0)^n \mid n \geq 0\}$$

# Ambiguity in Programming Languages

The following production rules (from ALGOL 60) became known as **dangling else problem**:

*Statement*  $\rightarrow$  if *Condition* then *Statement* |  
if *Condition* then *Statement* else *Statement* |  
...  
  
*Condition*  $\rightarrow$  ...

There are two derivation trees for

if ... then if ... then ... else ...

The production rules are **ambiguous**.

# Ambiguity Typically Unwanted

Ambiguity is typically unwanted:

- derivation trees often used to assign meaning to words,
- multiple derivation tree may result in double meaning.

In practice, ambiguity is often resolved outside of the grammar. For example, by a precedence on the rules:

- For example,  $0 + 2 * 1$  is parsed as  $0 + (2 * 1)$ .

**Ambiguity is undecidable.** That is, there exists no algorithm that decides whether a context-free grammar is ambiguous.

# Inherently Ambiguous Languages

A language  $L$  is **inherently ambiguous** if

- $L$  is context-free, and
- **every** grammar  $G$  with  $L(G) = L$  is ambiguous.

The following context-free language is inherently ambiguous

$$\{ a^n b^m c^\ell \mid n = m \vee m = \ell \}$$

It is generated by the following (ambiguous) grammar

$$\begin{array}{lll} S \rightarrow A \mid B & A \rightarrow Ac \mid C & B \rightarrow aB \mid D \\ & C \rightarrow aCb \mid \lambda & D \rightarrow bDc \mid \lambda \end{array}$$

## Exercise

Find two derivation trees for the word  $abc$ .

Parsing

# Parsing

**Parsing** is the search for a derivation tree for a given word.

## Theorem

For context-free grammars, parsing is possible in  $O(|w|^3)$  time. (Here  $|w|$  is the length of the input word.)

For every context-free grammar  $G$ :

- there exists an algorithm and  $C \geq 0$  such that
- for every word  $w$ , the algorithm determines in at most

$C \cdot |w|^3$  steps

whether  $w \in L(G)$  and a derivation tree if  $w \in L(G)$ .