

# Context Free Languages

Last Updated: March 20th, 2025

## 1 Introduction

Context free languages are foundational for defining several types of computing languages that occur in practice; including programming languages, markup languages, and languages for communication protocols. CFL's were first studied in relation to natural-language processing in the 1950's.

Their importance stems from the following.

1. CFL's are significantly more expressive than regular languages in that they are capable of defining recursive languages that may have unlimited recursive depth.
2. a CFL can be recognized by a pushdown automaton (PDA). Unlike DFA's a PDA has unlimited memory, albeit in the form of a stack whose access is limited to i) reading the top of the stack, ii) pushing on to the stack, and iii) popping the top of the stack. PDA's are also interesting because their nondeterministic counterparts (NPDA's) are more powerful than PDA's, and the set of languages accepted by an NPDA is equal to the set of CFL languages.
3. Letting **CFL** denote the class of all context-free languages, it can be shown that  $\text{CFL} \subseteq \text{P}$ . This is due to the fact that there exists a dynamic-programming algorithm that can decide any **CFL** language in  $O(n^3)$  steps, where  $n$  is the length of the input word. See Sipser, page 290.

Although every regular language is also a CFL (see the exercises), the converse is not true. For example, it can be proved that the language

$$L = \{a^n b^n | n \geq 0\}$$

is a CFL but is *not* regular.

**Proposition 1.** The language

$$L = \{a^n b^n | n \geq 0\}$$

is not regular.

**Proof.**

1. Suppose that  $L$  is regular and let  $M$  be a DFA that accepts  $L$  and has initial state  $q_0$ .
2. Since  $M$  has a finite number of states and there are infinitely many  $n$  values, there must be at least one state  $q$  for which there is an  $n_1 > 0$  and  $n_2 > 0$  for which
  - (a)  $n_1 \neq n_2$
  - (b) When either  $a^{n_1}$  or  $a^{n_2}$  is read starting in state  $q_0$ , the computation ends at  $q$ .
  - (c) When reading  $b^{n_1}$  starting in state  $q$ , the computation ends in an accepting state  $q_a$ .
3. But then the facts stated in 2) imply that, when starting in state  $q_0$  and reading  $a^{n_2} b^{n_1}$ , the computation moves to state  $q$  after reading  $a^{n_2}$ , followed by moving to state  $q_a$  after reading  $b^{n_1}$ .
4. Then by 3), it follows that  $M$  accepts  $a^{n_2} b^{n_1}$  which is a contradiction, since  $n_1 \neq n_2$ . □

## 2 Context-Free Grammars

A **Context-Free Grammar (CFG)** is a 4-tuple  $(V, \Sigma, R, S)$ , where

1.  $V$  is a finite set of variables
2.  $\Sigma$  is a finite set that is disjoint from  $V$ , called the **terminal set**
3.  $R$  is a finite set of rules where each **rule** has the form

$$A \rightarrow s,$$

where  $A \in V$  and  $s \in (V \cup \Sigma)^*$ . Variable  $A$  is referred to as the **head** of the rule, while  $s$  is referred to its **body**.

4.  $S \in V$  is the **start variable**

**Example 1.** Consider the set of rules

$$R = \{S \rightarrow SS, S \rightarrow aSb, S \rightarrow \varepsilon\}.$$

Then we may use this set of rules to define a CFG  $G = (V, \Sigma, R, S)$ , where

$$V = \{S\},$$

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

and variable  $S$  is the start variable.

For brevity we may list together rules having the same head as follows.

$$S \rightarrow SS \mid aSb \mid \varepsilon.$$

Here, each of the rule bodies is separated by a vertical bar. □

**Example 2.** One common use of CFG's is to provide grammatical formalism for natural languages. For example, consider the set of rules  $R$ :

$$\begin{aligned}
\langle \text{SENTENCE} \rangle &\rightarrow \langle \text{NOUN-PHRASE} \rangle \langle \text{VERB-PHRASE} \rangle \\
\langle \text{NOUN-PHRASE} \rangle &\rightarrow \langle \text{COMPLEX-NOUN} \rangle \mid \langle \text{COMPLEX-NOUN} \rangle \langle \text{PREP-PHRASE} \rangle \\
\langle \text{VERB-PHRASE} \rangle &\rightarrow \langle \text{COMPLEX-VERB} \rangle \mid \langle \text{COMPLEX-VERB} \rangle \langle \text{PREP-PHRASE} \rangle \\
\langle \text{PREP-PHRASE} \rangle &\rightarrow \langle \text{PREP} \rangle \langle \text{COMPLEX-NOUN} \rangle \\
\langle \text{COMPLEX-NOUN} \rangle &\rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle \\
\langle \text{COMPLEX-VERB} \rangle &\rightarrow \langle \text{VERB} \rangle \mid \langle \text{VERB} \rangle \langle \text{NOUN-PHRASE} \rangle \\
\langle \text{ARTICLE} \rangle &\rightarrow \text{a} \mid \text{the} \\
\langle \text{NOUN} \rangle &\rightarrow \text{trainer} \mid \text{dog} \mid \text{whistle} \\
\langle \text{VERB} \rangle &\rightarrow \text{calls} \mid \text{pets} \mid \text{sees} \\
\langle \text{PREP} \rangle &\rightarrow \text{with} \mid \text{in}
\end{aligned}$$

Here, the variables are the ten parts of speech delimited by  $\langle \rangle$ ,  $\Sigma$  is the lowercase English alphabet, including the space character, and  $\langle \text{SENTENCE} \rangle$  is the start variable.

**Example 3.** A CFG may also be used to define the syntax of a programming language. One fundamental language component to any programming language is that of an *expression*. The following rules imply a CFG for defining expressions formed by a single terminal **a**, parentheses, and the two arithmetic operations  $+$  and  $\times$ . Here  $E$  stands for **expression**,  $T$  for **term**, and  $F$  for **factor**.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow (E) \mid a$$

We have  $V = \{E, T, F\}$ ,  $\Sigma = \{+, \times, a, (, )\}$ , and  $E$  is the start variable.

## 2.1 Grammar derivations

Let  $G = (V, \Sigma, R, S)$  be a CFG, then the language  $D(G) \in (V \cup \Sigma)^*$  of **derived words** is structurally defined as follows.

**Atom**  $S \in D(G)$ .

**Compound Rule** Suppose  $s \in D(G)$ ,  $s$  is of the form  $uAv$  for some  $u, v \in (V \cup \Sigma)^*$ ,  $A \in V$ , and  $A \rightarrow \gamma$  is a rule of  $G$ , then

$$u\gamma v \in D(G).$$

In this case we write  $s \Rightarrow u\gamma v$ , and say that  $s$  **yields**  $u\gamma v$ . In words, to get a new derived word, take an existing derived word and replace one of its variables  $A$  with the body of a rule whose head is  $A$ .

The subset  $L(G)$  of derived words  $w \in D(G)$  for which  $w \in \Sigma^*$  is called the **context-free language (CFL)** associated with  $G$ . Thus, the words of  $L(G)$  consist only of terminal symbols.

## 2.2 The Derivation relation

Let  $u$  and  $v$  be words in  $(V \cup \Sigma)^*$ . We say that  $u$  **derives**  $v$ , written  $u \Rightarrow^* v$  if and only if either  $u = v$  or there is a sequence of words  $w_1, w_2, \dots, w_n$  such that

$$u = w_1 \Rightarrow w_2 \Rightarrow w_3 \Rightarrow \dots \Rightarrow w_n = v.$$

Such a sequence is called a **derivation sequence** from  $u$  to  $v$ .

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}.$$



**Example 4.** Use the CFG from Example 1 to derives the word **aabbaababb**.

$$S \rightarrow SS \mid aSb \mid \varepsilon.$$

**Solution.**

## 2.3 Derivation parse trees

Determining if an arbitrary word belongs to  $L(G)$  is of fundamental importance. But in addition, it is sometimes important to know the structure of the grammar's derivation of the word. For example, if a CFG generates arithmetic expressions, then knowing the structure of the derivation allows one to readily evaluate the expression (assuming the expression terminals have assigned values and the expression operations are properly defined). A **parse tree** for a word  $w \in L(G)$  is a tree whose structure and node labels reflect the derivation  $w$ , where, from left to right, the leaves of the tree are labeled with the letters of  $w$ . Indeed, consider the derivation sequence

$$S = w_1 \Rightarrow w_2 \Rightarrow \cdots \Rightarrow w_n = w.$$

Then the parse tree for  $w$  can be defined in a step-by-step manner. To begin the parse tree  $T_1$  for  $S = w_1$  consists of a single node labeled with  $S$ .

Now suppose a parse tree  $T_k$  has been associated with  $w_k$ , the  $k$ th word of the derivation. Moreover, assume that, from left to right, the leaves of  $T_k$  are labeled in one-to-one correspondence with the symbols of  $w_k$ . Moreover, assume that  $w_k$  has the form  $w_k = uAv$ , where  $A$  is substituted for a word  $\gamma$ , so that  $w_{k+1} = u\gamma v$ . Then  $T_{k+1}$  is obtained from  $T_k$  by assigning the leaf node labeled with  $A$  a number of children equal to the length of  $\gamma$  and for which, from left to right, the  $i$ th child is labeled with the  $i$ th symbol of  $\gamma$ .

**Example 5.** Use the CFG from Example 3 to derive the expression  $a \times (a + a)$ , and provide the parse tree associated with the derivation.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T \times F \mid F$$

$$F \rightarrow (E) \mid a$$

**Solution.**

## 2.4 Ambiguity

Given a CFG  $G$ , and a word  $w \in L(G)$ , there may be several different derivations of  $w$  from start symbol  $S$ . Many of these derivations however will yield identical parse trees. But in the event that two different derivation sequences of  $w$  from  $S$  yield two different parse trees, then we call  $G$  **ambiguous**. It turns out that an easy way to check for ambiguity is to check that no word  $w$  has more than one *leftmost derivation*.

Given grammar  $G = (V, \Sigma, R, S)$  and word  $w \in L(G)$ , a derivation sequence  $S = w_0 \Rightarrow w_1 \Rightarrow \cdots \Rightarrow w_{n-1} \Rightarrow w_n = w$  is called a **leftmost derivation** of  $w$  provided that, for every  $0 \leq i \leq n-1$ , the yielding of  $w_i$  from  $w_{i-1}$  was obtained by replacing the leftmost variable  $A$  of  $w_{i-1}$  with the body of one of a rule whose head is  $A$ . Therefore, if  $w$  has more than one leftmost derivation, it must be the case that a different sequence of rules were used to derive  $w$ . When this happens we call  $G$  ambiguous, since some words in the grammar have more than one parsing structure.

**Example 6.** Show that the grammar defined by the following rules is ambiguous.

$$\begin{aligned}\langle \text{SENTENCE} \rangle &\rightarrow \langle \text{NOUN-PHRASE} \rangle \langle \text{VERB-PHRASE} \rangle \\ \langle \text{NOUN-PHRASE} \rangle &\rightarrow \langle \text{COMPLEX-NOUN} \rangle \mid \langle \text{COMPLEX-NOUN} \rangle \langle \text{PREP-PHRASE} \rangle \\ \langle \text{VERB-PHRASE} \rangle &\rightarrow \langle \text{COMPLEX-VERB} \rangle \mid \langle \text{COMPLEX-VERB} \rangle \langle \text{PREP-PHRASE} \rangle \\ \langle \text{PREP-PHRASE} \rangle &\rightarrow \langle \text{PREP} \rangle \langle \text{COMPLEX-NOUN} \rangle \\ \langle \text{COMPLEX-NOUN} \rangle &\rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle \\ \langle \text{COMPLEX-VERB} \rangle &\rightarrow \langle \text{VERB} \rangle \mid \langle \text{VERB} \rangle \langle \text{NOUN-PHRASE} \rangle \\ \langle \text{ARTICLE} \rangle &\rightarrow \text{a} \mid \text{the} \\ \langle \text{NOUN} \rangle &\rightarrow \text{trainer} \mid \text{dog} \mid \text{whistle} \\ \langle \text{VERB} \rangle &\rightarrow \text{calls} \mid \text{pets} \mid \text{sees} \\ \langle \text{PREP} \rangle &\rightarrow \text{with} \mid \text{in}\end{aligned}$$

**Solution.**

**Solution Continued.**

**Example 7.** Provide a CFG  $G$  for which

$$L(G) = \{a^n b^n | n \geq 0\}.$$

Provide a derivation of  $a^3 b^3$  and draw its parse tree.

**Example 8.** Provide a CFG  $G$  for which

$$L(G) = \{x^i y^j z^k \mid i, j, k \geq 0 \text{ and } j = 2i \text{ or } k = 2j\}.$$



**Example 9.** Use the CFG from the previous example to provide a derivation of  $x^2y^3z^4$  and draw its parse tree.

## Exercises

1. For the CFG defined in Example 1, provide a derivation for the following words.
  - a. ababab
  - b. aaababbbab
  - c. aababaabbbaabb
2. For the CFG defined in Example 3, provide a derivation and parse tree for the following expressions.
  - a.  $a$
  - b.  $a + a$
  - c.  $a \times (a \times a)$
  - d.  $((a))$

## Exercise Solutions

1. We have the following derivations.

- a. ababab

$$S \Rightarrow SS \Rightarrow aSbS \Rightarrow abS \Rightarrow abSS \Rightarrow abaSbS \Rightarrow ababS \Rightarrow ababaSb \Rightarrow ababab.$$

- b. aaababbbab

$$\begin{aligned} S \Rightarrow SS \Rightarrow aSbS \Rightarrow aaSbbS \Rightarrow aaSSbbS \Rightarrow aaSbSbbS \Rightarrow aaabSbbS \Rightarrow aaabaSbbbS. \\ \Rightarrow aaababbbbS \Rightarrow aaababbbbaSb \Rightarrow aaababbbab. \end{aligned}$$

- c. aababaabbbaabb

$$\begin{aligned} S \Rightarrow SS \Rightarrow aSbS \Rightarrow aSSbS \Rightarrow aaSbSbS \Rightarrow aabSbS \Rightarrow aabSSbS \\ \Rightarrow aabaSbSbS \Rightarrow aababSbS \Rightarrow aababaSbbS \Rightarrow aababaaSbbbS \Rightarrow aababaabbbbS \\ \Rightarrow aababaabbbaSb \Rightarrow aababaabbbbaaSbb \Rightarrow aababaabbbaabb. \end{aligned}$$

2. For the CFG defined in Example 3, provide a derivation and parse tree for the following expressions.

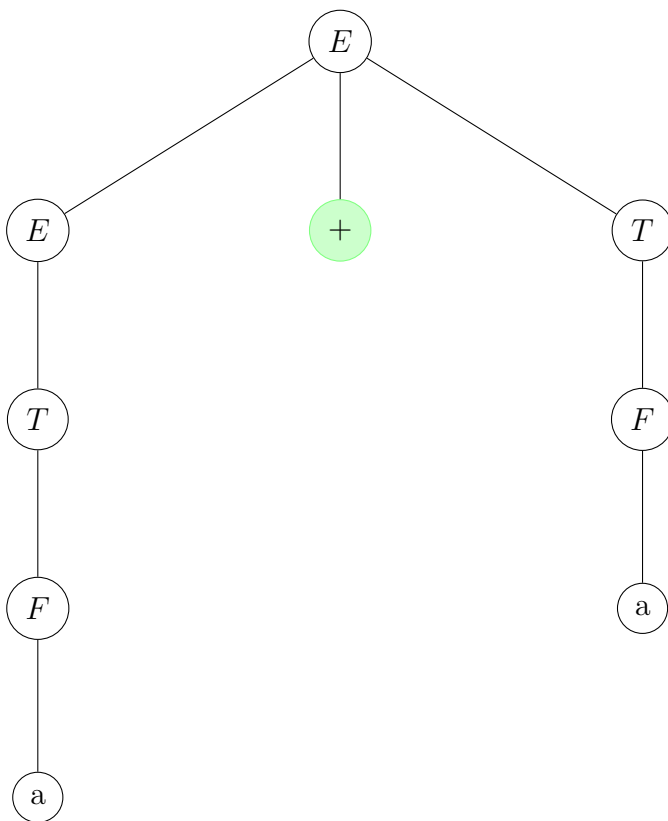
- a.  $a$

$$E \Rightarrow T \Rightarrow F \Rightarrow a.$$



b.  $a + a$

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \Rightarrow a + T \Rightarrow a + F \Rightarrow a + a.$$



c.  $a \times (a \times a)$

d.  $((a))$

$$E \Rightarrow T \Rightarrow F \Rightarrow (E) \Rightarrow (T) \Rightarrow (F) \Rightarrow ((E)) \Rightarrow T \Rightarrow F \Rightarrow ((a)).$$

