

Syntax Analysis (parsing)

Syntax analysis or parsing is the second phase of a compiler. In this phase, we shall learn the basic concepts used in the construction of a parser.

We have seen that a lexical analyzer can identify tokens with the help of regular expressions and pattern rules. But a lexical analyzer cannot check the syntax of a given sentence due to the limitations of the regular expressions. Regular expressions cannot check balancing tokens, such as parenthesis. Therefore, this phase uses context-free grammar (CFG), which is recognized by push-down automata.

CFG, on the other hand, is a superset of Regular Grammar, as depicted below:



It implies that every Regular Grammar is also context-free, but there exists some problems, which are beyond the scope of Regular Grammar. CFG is a helpful tool in describing the syntax of programming languages.

Context-Free Grammar

In this section, we will first see the definition of context-free grammar and introduce terminologies used in parsing technology.

A context-free grammar has four components:

- A set of non-terminals (V). Non-terminals are syntactic variables that denote sets of strings. The non-terminals define sets of strings that help define the language generated by the grammar.
- A set of tokens, known as terminal symbols (Σ). Terminals are the basic symbols from which strings are formed.
- A set of productions (P). The productions of a grammar specify the manner in which the terminals and non-terminals can be combined to form strings. Each production consists of a non-terminal called the left side of the

production, an arrow, and a sequence of tokens and/or non-terminals, called the right side of the production.

- One of the non-terminals is designated as the start symbol (S); from where the production begins.

The strings are derived from the start symbol by repeatedly replacing a non-terminal (initially the start symbol) by the right side of a production, for that non-terminal.

Example

We take the problem of palindrome language, which cannot be described by means of Regular Expression. That is, $L = \{ w \mid w = w^R \}$ is not a regular language. But it can be described by means of CFG, as illustrated below:

$G = (V, \Sigma, P, S)$

Where:

$V = \{Q, Z, N\}$

$\Sigma = \{0, 1\}$

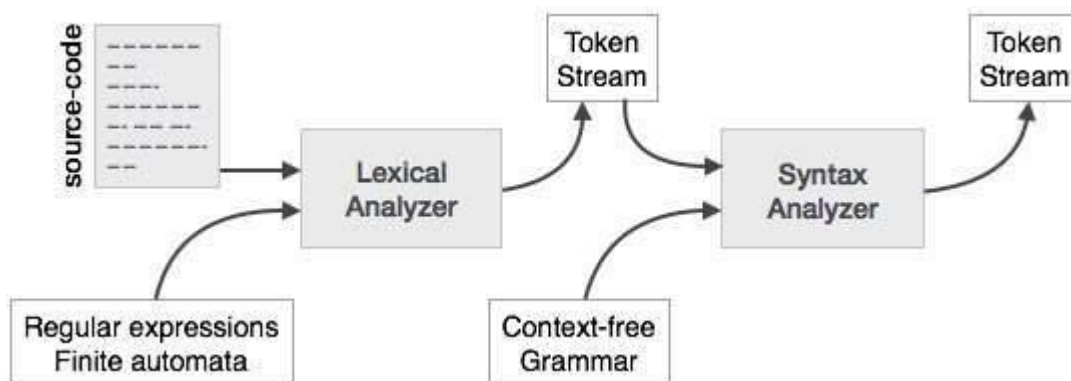
$P = \{ Q \rightarrow Z \mid Q \rightarrow N \mid Q \rightarrow \square \mid Z \rightarrow 0Q0 \mid N \rightarrow 1Q1 \}$

$S = \{ Q \}$

This grammar describes palindrome language, such as: 1001, 11100111, 00100, 1010101, 11111, etc.

Syntax Analyzers

A syntax analyzer or parser takes the input from a lexical analyzer in the form of token streams. The parser analyzes the source code *token stream* against the production rules to detect any errors in the code. The output of this phase is a **parse tree**.



This way, the parser accomplishes **two tasks**, i.e., **parsing the code, looking for errors and generating a parse tree as the output of the phase.**

Parsers are expected to parse the whole code even if some errors exist in the program. Parsers use error recovering strategies, which we will learn later.

Derivation

A derivation is basically a sequence of production rules, in order to get the input string. During parsing, we take two decisions for some sentential form of input:

- Deciding the non-terminal which is to be replaced.
- Deciding the production rule, by which, the non-terminal will be replaced.

To decide which non-terminal to be replaced with production rule, we can have two options.

Left-most Derivation

If the sentential form of an input is scanned and replaced from **left to right**, it is called **left-most derivation**. The sentential form derived by the left-most derivation is called the **left-sentential form**.

Right-most Derivation

If we scan and replace the input with production rules, from **right to left**, it is known as **right-most derivation**. The sentential form derived from the right-most derivation is called the **right-sentential form**.

Example

Production rules:

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow id$

Input string: $id + id * id$

The left-most derivation is:

$E \rightarrow E * E$

$E \rightarrow E + E * E$

$E \rightarrow id + E * E$

$E \rightarrow id + id * E$

$E \rightarrow id + id * id$

Notice that the left-most side non-terminal is always processed first.

The right-most derivation is:

$E \rightarrow E + E$

$E \rightarrow E + E * E$

$E \rightarrow E + E * id$

$E \rightarrow E + id * id$

$E \rightarrow id + id * id$

Parse Tree

A parse tree is a **graphical depiction of a derivation**. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the parse tree. Let us see this by an example from the last topic.

We take the left-most derivation of $a + b * c$

The left-most derivation is:

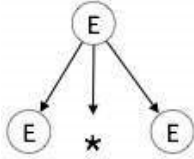
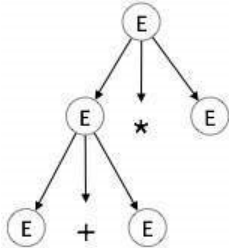
$E \rightarrow E * E$

$E \rightarrow E + E * E$

$E \rightarrow id + E * E$

$E \rightarrow id + id * E$

$E \rightarrow id + id * id$

Step 1:	$E \rightarrow E * E$	 <pre> graph TD E1((E)) --> E2((E)) E1 --> S1[*] E1 --> E3((E)) </pre>
Step 2:	$E \rightarrow E + E * E$	 <pre> graph TD E1((E)) --> E2((E)) E1 --> S1[+] E1 --> E3((E)) E2 --> E4((E)) E2 --> S2[*] E2 --> E5((E)) </pre>

Step 3:	$E \rightarrow id + E * E$	<pre> graph TD E1((E)) --> E2((E)) E1 --> S1[*] E1 --> E3((E)) E2 --> E4((E)) E2 --> S2[+] E2 --> E5((E)) E4 --> id1[id] </pre>
Step 4:	$E \rightarrow id + id * E$	<pre> graph TD E1((E)) --> E2((E)) E1 --> S1[*] E1 --> E3((E)) E2 --> E4((E)) E2 --> S2[+] E2 --> E5((E)) E4 --> id1[id] E5 --> id2[id] </pre>
Step 5:	$E \rightarrow id + id * id$	<pre> graph TD E1((E)) --> E2((E)) E1 --> S1[*] E1 --> E3((E)) E2 --> E4((E)) E2 --> S2[+] E2 --> E5((E)) E4 --> id1[id] E5 --> id2[id] E3 --> id3[id] </pre>

In a parse tree:

- All leaf nodes are terminals.
- All interior nodes are non-terminals.
- In-order traversal gives original input string.

A parse tree depicts associativity and precedence of operators. The deepest sub-tree is traversed first, therefore the operator in that sub-tree gets precedence over the operator which is in the parent nodes.

Ambiguity

A grammar G is said to be ambiguous if it has more than one parse tree (left or right derivation) for at least one string.

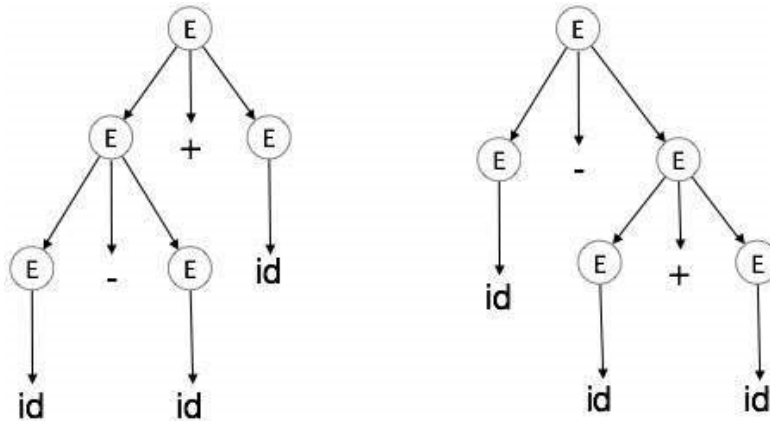
Example

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow \text{id}$

For the string **id + id - id**, the above grammar generates two parse trees:



The language generated by an ambiguous grammar is said to be **inherently ambiguous**. Ambiguity in grammar is not good for a compiler construction. No method can detect and remove ambiguity automatically, but it can be removed by either re-writing the whole grammar without ambiguity, or by setting and following **associativity and precedence** constraints.

Associativity

If an operand has operators on both sides, the side on which the operator takes this operand is decided by the associativity of those operators. If the operation is left-associative, then the operand will be taken by the left operator or if the operation is right-associative, the right operator will take the operand.

Example

Operations such as **Addition, Multiplication, Subtraction, and Division** are **left associative**. If the expression contains:

id op id op id

it will be evaluated as:

(id op id) op id

For example, $(\text{id} + \text{id}) + \text{id}$

Operations like **Exponentiation** are **right associative**, i.e., the order of evaluation in the same expression will be:

id op (id op id)

For example, $\text{id} ^ (\text{id} ^ \text{id})$

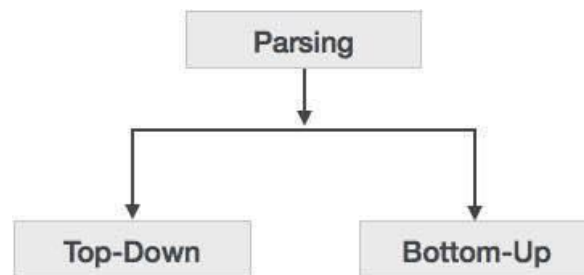
Precedence

If **two different operators share a common operand**, the precedence of operators decides which will take the operand. That is, $2+3*4$ can have **two different parse trees**, one corresponding to $(2+3)*4$ and another corresponding to $2+(3*4)$. By setting precedence among operators, this problem can be easily removed. As in the previous example, mathematically *** (multiplication) has precedence over + (addition)**, so the expression $2+3*4$ will always be interpreted as:

$2 + (3 * 4)$

These methods decrease the chances of ambiguity in a language or its grammar.

Syntax analyzers follow production rules defined by means of context-free grammar. The way the production rules are implemented (derivation) divides parsing into two types: top-down parsing and bottom-up parsing.



Top-down Parsing

When the parser starts constructing the parse tree from the start symbol and then tries to transform the start symbol to the input, it is called top-down parsing.

- **Recursive descent parsing:** It is a common form of top-down parsing. It is called recursive as it uses recursive procedures to process the input. Recursive descent parsing suffers from backtracking.
- **Backtracking:** It means, if one derivation of a production fails, the syntax analyzer restarts the process using different rules of same production. This technique may process the input string more than once to determine the right production.

Bottom-up Parsing

As the name suggests, bottom-up parsing starts with the input symbols and tries to construct the parse tree up to the start symbol.

Example:

Input string : **a + b * c**

Production rules:

$S \rightarrow E$

$E \rightarrow E + T$

$E \rightarrow E * T$

$E \rightarrow T$

$T \rightarrow \text{id}$

Let us start bottom-up parsing

a + b * c

Read the input and check if any production matches with the input:

a + b * c

T + b * c

E + b * c

E + T * c

E * c

E * T

E

S

Example: Consider the grammar:-

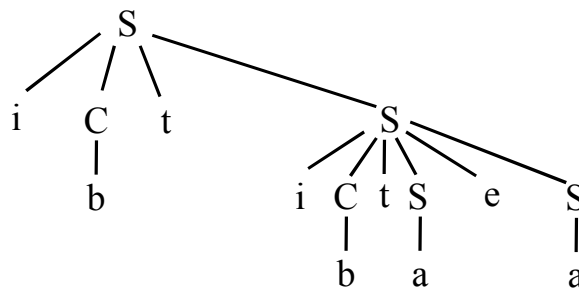
(1) $S \longrightarrow i C t S$

(2) $S \longrightarrow i C t S e S$ ----- (1)

(3) $S \longrightarrow a$

(4) $C \longrightarrow b$

Here i, t and e stand for "if", "then" and "else" respectively and C for "conditional", S for "statement".



Parse Tree T

The left-most derivation corresponding to this parse tree is given by:-

$$\begin{aligned}
S &\xRightarrow{lm} i C t S \\
&\xRightarrow{lm} i b t S \\
&\xRightarrow{lm} i b t i C t S e S \\
&\xRightarrow{lm} i b t i b t S e S \\
&\xRightarrow{lm} i b t i b t a e S \\
&\xRightarrow{lm} i b t i b t a e a
\end{aligned}$$

A right most derivation can be constructed from a parse tree analogously. At each step we replace the right-most non-terminal by the tables of its children. For example, the first two steps of a right most derivation constructed from Figure would be

$$S \Rightarrow i C t S \Rightarrow i C t i C t S e S$$

Examples

1.

$E \rightarrow E + E$

$E \rightarrow E * E$

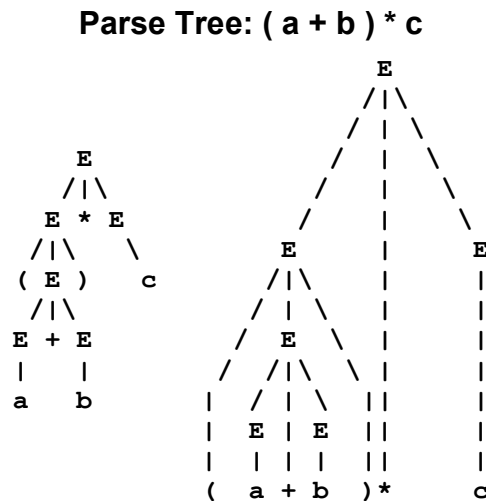
$E \rightarrow (E)$

$E \rightarrow a \mid b \mid c \mid \dots$

The string is $(a + b) * c$

Here is a sequence of replacements that leads to a sequence of terminal symbols.

$E \Rightarrow E * E \Rightarrow (E) * E \Rightarrow (E + E) * E \Rightarrow (a + E) * E \Rightarrow$
 $(a + b) * E \Rightarrow (a + b) * c$



2.

Consider the context-free grammar

$S \rightarrow AB \mid CD$

$A \rightarrow 0A1 \mid 01$

$B \rightarrow 2B \mid 2$

$C \rightarrow 0C \mid 0$

$D \rightarrow 1D2 \mid 12$

and the string **012**. Is the grammar ambiguous or unambiguous? Justify your answer.

$S \Rightarrow AB \Rightarrow 01B \Rightarrow 012$

$S \Rightarrow CD \Rightarrow 0D \Rightarrow 012$

the grammar is ambiguous, two derivation from one side

3.

$E \rightarrow I \mid C \mid E + E \mid E * E$

$I \rightarrow y \mid z$

$C \rightarrow 4$

Where E is start symbol, and the string $y + 4 * z$, is the grammar ambiguous or unambiguous? Justify your answer.

