

Syntax Analysis

Syntax analysis or parsing is the second phase of a compiler. In this chapter, 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.

The main goal of syntax analysis is to identify an organization in code. It establishes whether a text adheres to the desired style or not. This phase's primary goal is to determine whether the coder wrote accurate source code or not.

By using tokens to build the parse tree, syntax analysis is based on rules unique to the programming language being used. It also establishes the grammar or syntax of the language and the source tongue's structure.

Here is a summary of the duties carried out during this phase:

- Take the lexical analyzer's symbols.
- determines whether or not the expression is syntactically accurate.
- Send in all grammar mistakes
- Create a parse tree, a type of hierarchical arrangement.

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 **on- 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:

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

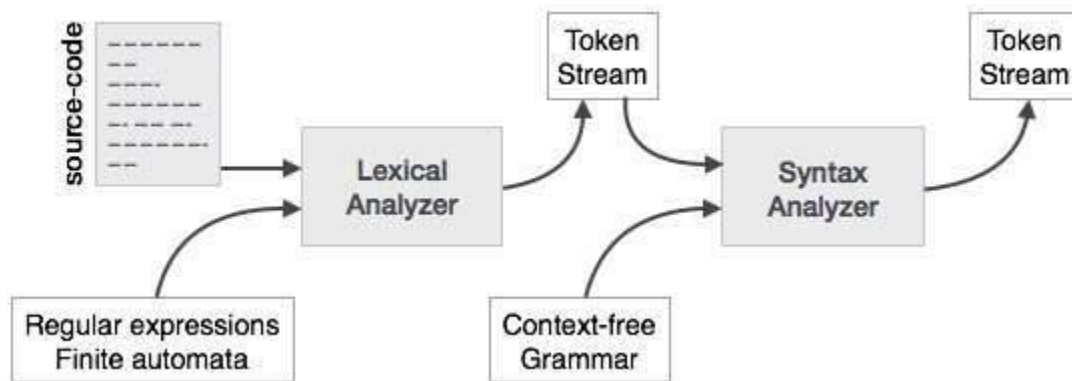
Where:

$$V = \{ Q, Z, N \}$$
$$\Sigma = \{ 0, 1 \}$$
$$P = \{ Q \rightarrow Z \mid Q \rightarrow N \mid Q \rightarrow \varepsilon \mid Z \rightarrow 0Q0 \mid N \rightarrow 1Q1 \}$$
$$S = \{ Q \}$$

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 in this chapter.

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:

$$E \rightarrow E + E$$
$$E \rightarrow E * E$$
$$E \rightarrow \text{id}$$

Input string: $\text{id} + \text{id} * \text{id}$

The left-most derivation is:

$$E \rightarrow E * E$$
$$E \rightarrow E + E * E$$
$$E \rightarrow \text{id} + E * E$$
$$E \rightarrow \text{id} + \text{id} * E$$
$$E \rightarrow \text{id} + \text{id} * \text{id}$$

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

The right-most derivation is:

$$E \rightarrow E + E$$
$$E \rightarrow E + E * E$$
$$E \rightarrow E + E * \text{id}$$
$$E \rightarrow E + \text{id} * \text{id}$$
$$E \rightarrow \text{id} + \text{id} * \text{id}$$

Parse Tree

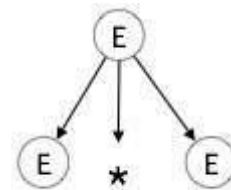
A parse tree is a graphical depiction of a derivation. It is convenient to see how strings are derived from the start symbol. The start symbol of the derivation becomes the root of the parse tree. Let us see this by an example from the last topic.

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

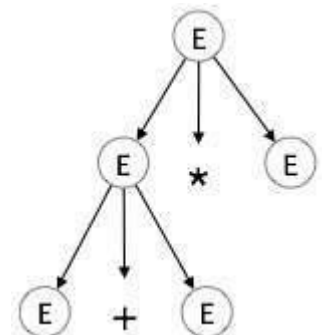
The left-most derivation is:

$$E \rightarrow E * E$$
$$E \rightarrow E + E * E$$
$$E \rightarrow id + E * E$$
$$E \rightarrow id + id * E$$
$$E \rightarrow id + id * id$$

Step 1:

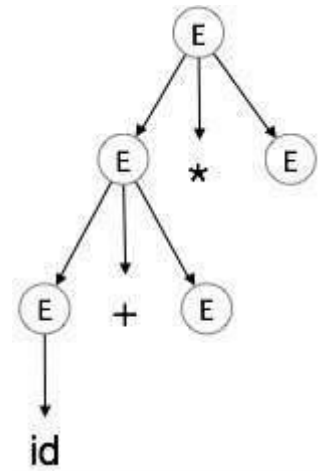
$$E \rightarrow E * E$$


Step 2:

$$E \rightarrow E + E * E$$


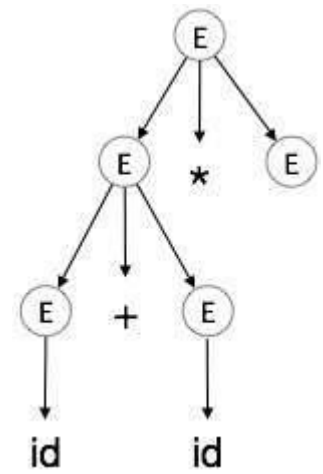
Step 3:

$E \rightarrow id + E * E$



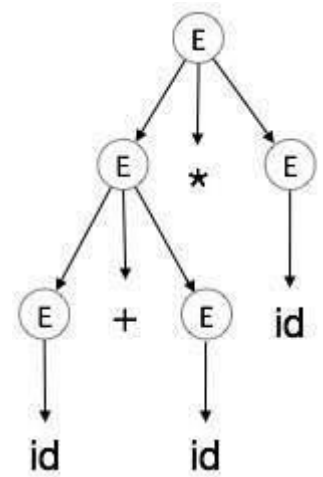
Step 4:

$E \rightarrow id + id * E$



Step 5:

$E \rightarrow id + id * id$



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.

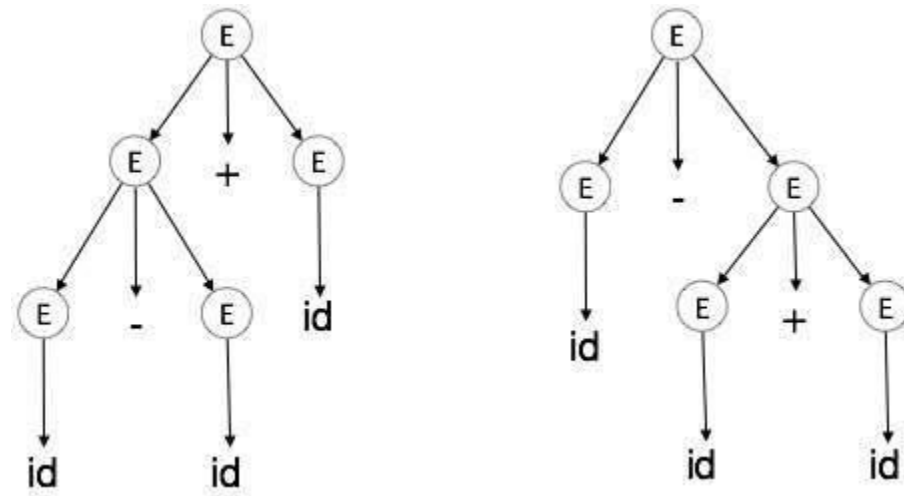
Example

$E \rightarrow E + E$

$E \rightarrow E - E$

$E \rightarrow id$

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, (id + id) + 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, id ^ (id ^ 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.

Left Recursion

A grammar becomes left-recursive if it has any non-terminal 'A' whose derivation contains 'A' itself as the left-most symbol. Left-recursive grammar is considered to be a

problematic situation for top-down parsers. Top-down parsers start parsing from the Start symbol, which in itself is non-terminal. So, when the parser encounters the same non-terminal in its derivation, it becomes hard for it to judge when to stop parsing the left non-terminal and it goes into an infinite loop.

Example:

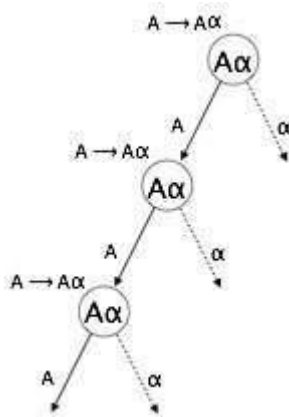
$$(1) A \Rightarrow A\alpha \mid \beta$$

$$(2) S \Rightarrow A\alpha \mid \beta$$

$$A \Rightarrow S\delta$$

(1) is an example of immediate left recursion, where A is any non-terminal symbol and α represents a string of non-terminals.

(2) is an example of indirect-left recursion.



A top-down parser will first parse the A , which in-turn will yield a string consisting of A itself and the parser may go into a loop forever.

Removal of Left Recursion

One way to remove left recursion is to use the following technique:

The production

$$A \Rightarrow A\alpha \mid \beta$$

is converted into following productions

$$A \Rightarrow \beta A'$$

$$A' \Rightarrow \alpha A' \mid \epsilon$$

This does not impact the strings derived from the grammar, but it removes immediate left recursion.

Second method is to use the following algorithm, which should eliminate all direct and indirect left recursions.

START

Arrange non-terminals in some order like $A_1, A_2, A_3, \dots, A_n$

for each i from 1 to n

{

for each j from 1 to $i-1$

{

replace each production of form $A_i \Rightarrow A_j \gamma$

with $A_i \Rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \delta_3 \gamma \mid \dots \mid \gamma$

where $A_j \Rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_n$ are current A_j productions

}

}

eliminate immediate left-recursion

END

Example

The production set

$S \Rightarrow A\alpha \mid \beta$

$A \Rightarrow Sd$

after applying the above algorithm, should become

$S \Rightarrow A\alpha \mid \beta$

$A \Rightarrow A\alpha d \mid \beta d$

and then, remove immediate left recursion using the first technique.

$A \Rightarrow \beta d A'$

$$A' \Rightarrow \alpha d A' \mid \epsilon$$

Now none of the production has either direct or indirect left recursion.

Left Factoring

If more than one grammar production rules has a common prefix string, then the top-down parser cannot make a choice as to which of the production it should take to parse the string in hand.

Example

If a top-down parser encounters a production like

$$A \Rightarrow \alpha \beta \mid \alpha \gamma \mid \dots$$

Then it cannot determine which production to follow to parse the string as both productions are starting from the same terminal (or non-terminal). To remove this confusion, we use a technique called left factoring.

Left factoring transforms the grammar to make it useful for top-down parsers. In this technique, we make one production for each common prefixes and the rest of the derivation is added by new productions.

Example

The above productions can be written as

$$A \Rightarrow \alpha A'$$

$$A' \Rightarrow \beta \mid \gamma \mid \dots$$

Now the parser has only one production per prefix which makes it easier to take decisions.

First and Follow Sets

An important part of parser table construction is to create first and follow sets. These sets can provide the actual position of any terminal in the derivation. This is done to create the parsing table where the decision of replacing $T[A, t] = \alpha$ with some production rule.

First Set

This set is created to know what terminal symbol is derived in the first position by a non-terminal. For example,

$$\alpha \rightarrow t \beta$$

That is α derives t (terminal) in the very first position. So, $t \in \text{FIRST}(\alpha)$.

Algorithm for calculating First set

Look at the definition of $\text{FIRST}(\alpha)$ set:

- if α is a terminal, then $\text{FIRST}(\alpha) = \{ \alpha \}$.
- if α is a non-terminal and $\alpha \rightarrow \epsilon$ is a production, then $\text{FIRST}(\alpha) = \{ \epsilon \}$.
- if α is a non-terminal and $\alpha \rightarrow \gamma_1 \gamma_2 \gamma_3 \dots \gamma_n$ and any $\text{FIRST}(\gamma_i)$ contains t then t is in $\text{FIRST}(\alpha)$.

First set can be seen as:

$$\text{FIRST}(\alpha) = \{ t \mid \alpha \xrightarrow{*} t \beta \} \cup \{ \epsilon \mid \alpha \xrightarrow{*} \epsilon \}$$

Follow Set

Likewise, we calculate what terminal symbol immediately follows a non-terminal α in production rules. We do not consider what the non-terminal can generate but instead, we see what would be the next terminal symbol that follows the productions of a non-terminal.

Algorithm for calculating Follow set:

- if α is a start symbol, then $\text{FOLLOW}(\alpha) = \$$
- if α is a non-terminal and has a production $\alpha \rightarrow AB$, then $\text{FIRST}(B)$ is in $\text{FOLLOW}(A)$ except ϵ .
- if α is a non-terminal and has a production $\alpha \rightarrow AB$, where $B \in \epsilon$, then $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(\alpha)$.

Follow set can be seen as: $\text{FOLLOW}(\alpha) = \{ t \mid S \xrightarrow{*} \alpha t^* \}$

Limitations of Syntax Analyzers

Syntax analyzers receive their inputs, in the form of tokens, from lexical analyzers. Lexical analyzers are responsible for the validity of a token supplied by the syntax analyzer. Syntax analyzers have the following drawbacks -

- it cannot determine if a token is valid,

- it cannot determine if a token is declared before it is being used,
- it cannot determine if a token is initialized before it is being used,
- it cannot determine if an operation performed on a token type is valid or not.

These tasks are accomplished by the semantic analyzer, which we shall study in Semantic Analysis.

Applications of Syntax Analysis

- **Natural language processing:** Syntax analysis is used in natural language processing to analyze and understand the structure of sentences in a language. It helps identify the parts of speech, determine the relationships between the words, and construct a parse tree that represents the hierarchical structure of the sentence.
- **Information extraction:** Syntax analysis can be used to extract structured information from unstructured text, such as identifying names, dates, and locations in a news article or extracting product details from an online shopping website.
- **Machine translation:** Syntax analysis is an important step in the process of machine translation, as it helps to identify the structure and meaning of sentences in the source language and translate them accurately into the target language.
- **Computer science:** In computer science, syntax analysis is an important phase in the process of compiling a program. It checks the source code of a program to ensure that it follows the correct syntax of the programming language in which it is written.
- **Text analytics:** Syntax analysis can be used in text analytics to extract insights and information from large volumes of text data. For example, it can be used to identify common themes or trends in customer reviews or to classify text documents based on their content.

