## CONTEXT FREE GRAMMAR:

Context free grammar Is a specification for the syntax of a programming language. Informally a Context Free Grammar (CFG) is defined as a finite set of rules used for deriving or generating strings or sentences in a language, which consists of:

- A set of replacement rules, each having a Left-Hand Side (LHS) and a Right-Hand Side (RHS).

- Two types of symbols called as variables and terminals.

- LHS of each rule is a single variable (non-terminal).

- RHS of each rule is a string of zero or more variables and/or terminals.

CFG is also defined as a set of variables (non-terminals) each of which represents a language. The languages represented by the variables are described recursively in terms of each other. The primitive symbols are called as terminals and the rules relating the variables are called as productions.

Context Free Grammar is formally defined as a 4 tuple (V, T, P, S) where:

- V and T are finite sets of variables and terminals respectively.
- P is the finite set of productions. Each production is of the form A→α where A is a variable and α is a string of symbols from (V U T)*
- S is a special variable called as start symbol.

***Sentential form:*** A sentence that contains variables & terminals is called as a sentential form.

***For example***, If (α, β) ϵ P, then we write the production as α →β where β is the sentential form.

Context-free grammars are useful for describing arithmetic expressions and block structure. Example for context free grammar is given below:

Given a set of productions:

&lt;exp&gt; → &lt;exp&gt; + &lt;exp&gt;

&lt;exp&gt; → &lt;exp&gt; * &lt;exp&gt;

&lt;exp&gt; → (&lt;exp&gt;)

&lt;exp&gt; → id

Example E → E + E, E → E * E, E → (E), E → x, E → y

**LANGUAGE OF CONTEXT FREE GRAMMAR:**

A sequence of tokens is syntactically legal if it can be derived by applying the productions of the CFG. A context-free grammar defines a language which is a set of strings (sequences) of tokens (terminals), each string of tokens is derivable from the production rules of the CFG.

Consider the following simplified grammar for expressions

*expr → expr op expr | ( expr ) | id | num*

*op → + | − | * | /*

- ✓ The string:*(id+num)*id* is syntactically legal and part of the language

- ✓ Similarly: *id*(num-id)* is also part of the language

- ✓ However: *(id+num* is NOT part of the language

- ✓ Similarly: *id*-num+* is NOT part of the language

**DERIVATIONS:**

Derivation is an ordered tree which is defined as sequence of replacements of a substring in a sentential form. To check whether a sequence of tokens is legal or not:

- ✓ We start with a nonterminal called the start symbol

- ✓ We apply productions, rewriting nonterminals, until only terminals remain

- ✓ derivation replaces a nonterminal on LHS of a production with RHS

- ✓ The → symbol denotes a derivation step

*For example*, a derivation for *(id + num)*id* is given below:

```
expr    →  expr op expr
        →  ( expr ) op expr
        →  ( expr op expr ) op expr
        →  ( expr + expr ) op expr
        →  ( expr + expr ) * expr
        →  ( id + expr ) * expr
        →  ( id + num ) * expr
        →  ( id + num ) * id
```

When deriving a sequence of tokens more than one nonterminal may be present and can be expanded. There are two types of derivations 1) left most derivation and 2) right most derivation.

A left most derivation chooses the leftmost nonterminal to expand and a right most derivation chooses the rightmost nonterminal to expand.
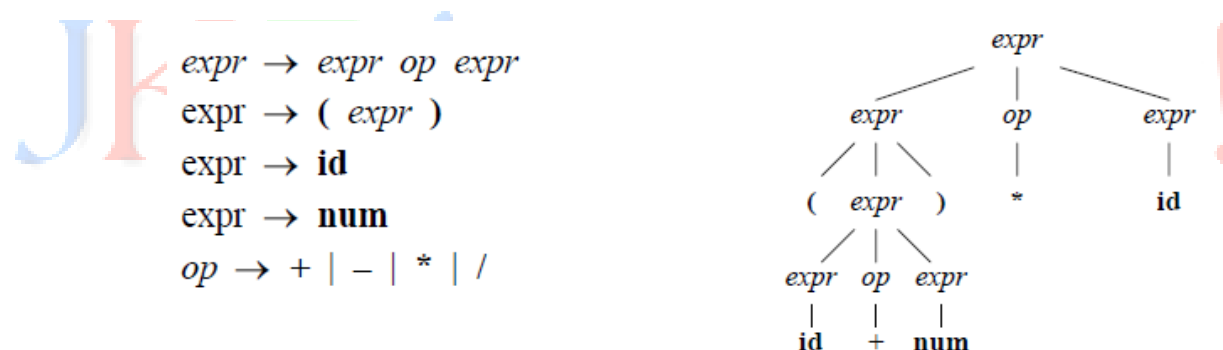
A leftmost derivation for    `(id + num)*id`

$$expr \Rightarrow_{lm} expr\ op\ expr \Rightarrow_{lm} (\ expr\ )\ op\ expr \Rightarrow_{lm} (\ expr\ op\ expr\ )\ op\ expr$$
$$\Rightarrow_{lm} (\ id\ op\ expr\ )\ op\ expr \Rightarrow_{lm} (\ id + expr\ )\ op\ expr$$
$$\Rightarrow_{lm} (\ id + num\ )\ op\ expr \Rightarrow_{lm} (\ id + num\ )\ *\ expr \Rightarrow_{lm} (\ id + num\ )\ *\ id$$

A rightmost derivation for `(id + num)*id`

$$expr \Rightarrow_{rm} expr\ op\ expr \Rightarrow_{rm} expr\ op\ id \Rightarrow_{rm} expr\ *\ id \Rightarrow_{rm} (\ expr\ )\ *\ id$$
$$\Rightarrow_{rm} (\ expr\ op\ expr\ )\ *\ id \Rightarrow_{rm} (\ expr\ op\ num\ )\ *\ id$$
$$\Rightarrow_{rm} (\ expr + num\ )\ *\ id \Rightarrow_{rm} (\ id + num\ )\ *\ id$$

The graphical representation of a derivation is called parse tree, which filters out choice regarding replacement order and it is rooted by the start symbol S. Interior nodes represent nonterminals, Leaf nodes are terminals. The following is a parse tree for ( id + num ) * id

$$expr \rightarrow expr\ op\ expr$$
$$expr \rightarrow (\ expr\ )$$
$$expr \rightarrow id$$
$$expr \rightarrow num$$
$$op \rightarrow +\ |\ -\ |\ *\ |\ /$$



## AMBIGUITY:

A Context Free Grammar is ambiguous if some string w ∈ L(G) has two or more leftmost or rightmost derivations. Sometimes the language generated by an ambiguous grammar has an equivalent unambiguous grammar. Languages that can only be generated by ambiguous grammars are called inherently ambiguous. It can be solved in two ways 1) solved by precedence 2) solved by associativity

## PARSING:

The term parsing is the process of analyzing a string of symbols, either in natural language or in computer languages, according to the rules of a formal grammar.

Programming language that follows some syntax to develop the code has to check for any syntax errors in the code, this will be done in parsing.

Parsing is the process of the generation of a parse tree (or its equivalents) which is a graphical representation of derivation that filters out the choice regarding replacement order; parse tree is the internal structure of compiler or interpreter; corresponding to a given input string w and grammar G. There are two types of parsing techniques. 1) Top Down Parsing and 2) Bottom Up parsing

## TOP DOWN PARSING:

Top down Parsing begin with start symbol of grammar as root of tree and grow it towards leaves. Construction of parse tree in top down parsing is constructed from the top and from left to right. It is classified in to two types based on two forms (i.e. prediction and backtracking)

1) Recursive descent parsing

2) Predictive parsing without backtracking.

## RECURSIVE DESCENT PARSING:

A recursive descent parser is a kind of top-down parser built from a set of mutually recursive procedures (or a non-recursive equivalent) where each such procedure usually implements one of the production rules of the grammar. Thus the structure of the resulting program closely mirrors that of the grammar it recognizes.
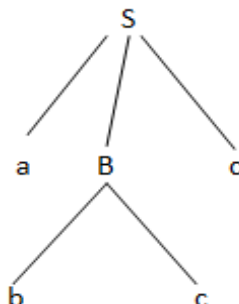
Recursive descent with backtracking is a technique that determines which production to use by trying each production in turn. Recursive descent with backtracking is not limited to LL(k) grammars, but is not guaranteed to terminate unless the grammar is LL(k). Even when they terminate, parsers that use recursive descent with backtracking may require exponential time. It consists of set of procedures one for each non terminal, which uses backtracking.
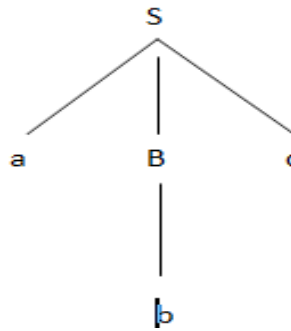
Now let us consider a grammar:

$S \rightarrow aBc$

$B \rightarrow bc/b$

Let the input string be abc then constructing a parse tree

Here it fails so it backtracks



It makes repeated scans of the input. A left recursive grammar can cause a recursive decent parser, even one with backtracking, to go into an infinity loop. A grammar is said to be left recursive if it has a non terminal A such that there is a derivation A→Aα for some string α.

Top down parsing methods cannot handle left recursive grammars, so a transformation that eliminates left recursion is needed. To get non-left recursive productions we rewrite the productions **A → Aα | β** as

A→βA'

A'→αA'|ε

i.e. if A→Aα$_1$|Aα$_2$|....|Aα$_n$ | β$_1$|β$_2$....|β$_m$ then it can be written as follows:

A→β$_1$A' | β$_2$A'....β$_m$A'

A'→α$_1$A'| α$_2$A' | α$_3$A' . . . . | α$_n$A' | ε

Let us consider an example:

E→E+T/T

T→T*F/F

F→(E)/id

The grammar can be written as

E→E+T
E→T
T→T*F
T→F
F→(E)
F→id

---

Now consider E→E+T | T this is in the form of A→Aα, so we can rewrite the grammar as

E→TE'

E'→+TE'/E

Similarly this can be    T→FT'

T'→*FT'/E

F→(E)/id

## PREDICTIVE PARSING:

A predictive parser is a recursive descent parser that does not require backtracking. Predictive parsing is possible only for the class of LL(k) grammars, which are the context-free grammars for which there exists some positive integer k that allows a recursive descent parser to decide which production to use by examining only the next k tokens of input.

The LL(k) grammars therefore exclude all ambiguous grammars, as well as all grammars that contain left recursion. Any context-free grammar can be transformed into an equivalent grammar that has no left recursion, but removal of left recursion does not always yield an LL(k) grammar. A predictive parser runs in linear time.

Predictive parsers can be depicted using transition diagrams for each non-terminal symbol where the edges between the initial and the final states are labeled by the symbols (terminals and non-terminals) of the right side of the production rule.

Left Factoring:

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing. Now we have few steps to follow:

1) Collect all the productions with the same left side and begin with same symbols on right hand side.
     Example:       E→T+E/T

                    T→int/int*T/(E)

                    T→int

                    T→int*T

2) Combine the common strings into a single production and then append a new terminal symbol to the en**d** of this new production.
     Example:       T→int y

3) Create new production using this new non-terminal for each of the suffixes to the common production.

METHODS IN LL:

- FIRST($\alpha$)
- FOLLOW(A)

FIRST ($\alpha$):

It is defined as the set of terminals that begin strings derived from $\alpha$ where $\alpha$ is any grammar symbol.

FOLLOW (A):

It is defined as the set of terminals that can appear immediately after A in some sentential form. We have some rules for both FIRST ($\alpha$) and FOLLOW (A). They are:

**FIRST ($\alpha$):**

- If $\alpha$ is a terminal then FIRST ($\alpha$) is $\{\alpha\}$.

- If $\alpha \rightarrow \varepsilon$ is a production then add $\varepsilon$ to FIRST ($\alpha$).

- If $\alpha$ is a non-terminal and $\alpha \rightarrow Y_1, Y_2, Y_3, \dots Y_k$ is a production then place a in First($\alpha$) if for some I, a is in FIRST($Y_i$) and $\varepsilon$ is in all of FIRST($Y_1$)....FIRST($Y_{i-1}$) and if $\varepsilon$ is in FIRST($Y_j$) for all j=1,2,....k then add $\varepsilon$ to FIRST($\alpha$).

**FOLLOW (A):**

- Place \$ in FOLLOW (A) where A is a start symbol and \$ is input right end marker.

- If there is a production A$\rightarrow\alpha$B$\beta$ then everything in FIRST($\beta$) is in FOLLOW($\beta$) except '$\varepsilon$', irrespective of whether $\beta$ may or may not contains $\varepsilon$.

- If there is a production A$\rightarrow\alpha$B$\beta$ (or) A$\rightarrow\alpha$B where FIRST($\beta$) contains $\varepsilon$ then everything in FOLLOW (A) is in FOLLOW (B).

Now let us consider an example grammar.

E$\rightarrow$E+T/T

T$\rightarrow$T*F/F

F$\rightarrow$Id/(E)

---

The above example has lest recursion so we have determinate it.

First consider,

   E➔E+T

   E➔T

  Here α=T,     β1=+E        β=E

So this can be written as

   E➔TE'

   E'➔+E/E

Now consider

  T➔T*F    (or)   T➔F*T

  T➔F              T➔F

Here α=F, β1= *T, β2=E

   T➔FT'

   T'➔*T/E

So now after removing left recursion the productions are:

         E➔TE'

         E'➔+E/E

         T➔FT'

         T'➔*T/E

         F➔id/(E)

  Now let us calculate FIRST (α).

FIRST (E) = FIRST (T) = FIRST (F)

Here FIRST (F) = {id,c}  so from rule (3)

     FIRST (T) = {id,c}

     FIRST (E') = {E, +} from rules (2&3)

FIRST (T') = {E, *}

FIRST (+) = {+}

FIRST (*) = {*}

FIRST (id) = {id)

FIRST (( ) = {(}

FIRST ( )) = {)}

Now let us calculate the FOLLOW (A) for the given grammar:

FOLLOW (E) = {$,)}

FOLLOW (E') = FOLLOW (E)

FOLLOW (T) = {+, $,)}

FOLLOW (T') = FOLLOW (T)

FOLLOW (F) = FOLLOW {*, $, +,)}

Now let us construct a predictive parse table. It has some rules:

- For each production A→α of the grammar do the step 2 & 3

- For each terminal 'a' in FIRST (α) add A→α to M[A, a]

- If E is in FIRST (α) add A→ to M [A, b] for each terminal b in FOLLOW (A). If E is in FIRST (α) and $ is in FOLLOW (A) then add A→ to M [A, $].

-  Make each undefined entry as an error, in the table M.

|     | +      | *     | id     | (      | )      | $      |
|-----|--------|-------|--------|--------|--------|--------|
| E   |        |       | E→TE'  | E→TE'  |        |        |
| E'  | E'→+E  |       |        |        | E'→E   | E'→E   |
| T   |        |       | T→FT'  | T→FT'  |        |        |
| T'  | T'→E   | T'*T  |        |        | T'→E   | T'→E   |
| F   |        |       | F→id   | F→(E)  |        |        |