

COMPILER DESIGN

MODULE-II

Introduction to Syntax Analysis.

SYLLABUS :

Role of syntax analyzer - Syntax error handling.
Review of Context Free Grammars - Derivations and parse trees, Eliminating Ambiguity. Basic parsing approaches - Eliminating Left Recursion, left factoring. Top-Down parsing - Recursive Descent Parsing, Predictive parsing, LL(1) grammars

Role of Syntax Analyzer :

- Every programming languages has precise rules to describe its syntactic structure.
- The syntax of pgmng lang. constructs can be specified by context-free grammars (CFG) or BNF (Backus-Naur Form) notation.
- The parser obtains a string of tokens from the lexical analyzer and verifies that this string of tokens can be generated by the grammar for the source language..
- The parser also reports any Syntax errors to the user and it also recover from the commonly occurring errors to continue processing the remainder of the program.
- For well-formed programs, the parser constructs a parse tree and passes it to the rest of the compiler for further processing.
- Note : CFG are recognized by push-down automata

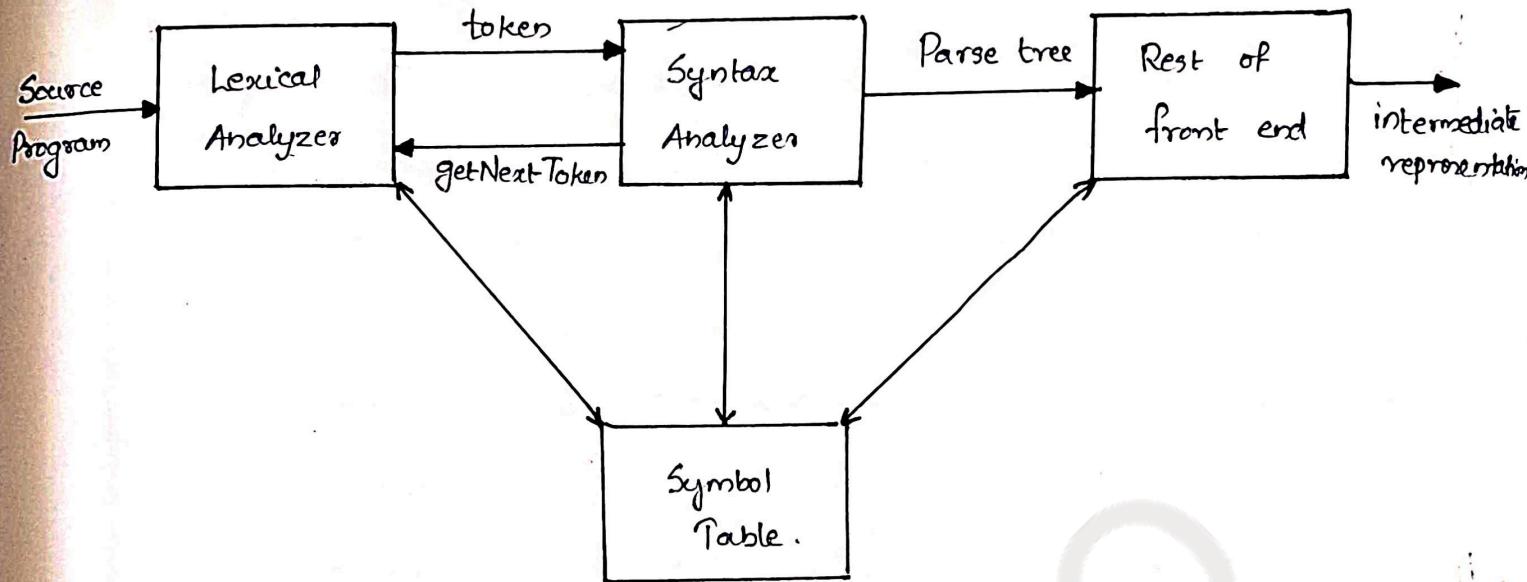


Fig: position of parser in compiler model
 (Syntax analyzer)

There are 3 types of parsers :

- ① Universal parser
- ② Top-down parser
- ③ Bottom-up parser.

- universal parsing methods such as CYK algorithm and Earley's algorithm can parse any grammars. But these methods are inefficient to use.
- Thus the commonly used parsing methods are : top-down and bottom-up approaches.
- Top-down methods build parse trees from the top (ie root) to the bottom (ie leaves), while bottom-up methods start from the leaves and work their way up to the root. In both the cases, the input to the parser is scanned from left-to-right, one symbol at a time.

Syntax Error Handling:

The goals of an error handler in a parser are:

- To report the presence of errors clearly & accurately.
- Recover from each error quickly to detect subsequent errors.
- Add minimal overhead to the processing of correct programs.

- Here, the aim is to detect the syntax errors.

Error-Recovery strategies:-

- ① Panic-mode recovery
- ② Phrase-level recovery
- ③ error-productions
- ④ Global corrections.

1) Panic-mode Recovery:-

In this method, when the parser discovers an error, it discards the input symbols one at a time, until one of a designated set of synchronizing tokens is found.

The synchronizing tokens are usually delimiters, such as semicolon or } , whose role in the source pgm is clear & unambiguous.

Disadvantage: Panic-mode correction often skips a considerable amount of input without checking it for additional errors.

Advantages: - Simplicity.

- It is guaranteed not to go into an infinite loop

2) Phrase-level Recovery :

- On discovering an error, the parser may perform local correction on the remaining input; that is it may replace a prefix of the remaining input by some string that allows the parser to continue.
- Typical local correction is to replace a comma by a semicolon, delete an extraneous semicolon, or insert a missing semicolon.
- The choice of the local correction is left to the compiler designer.
- we must be careful to choose replacements that do not lead to infinite loops.
- This method has been used in several error-repairing compilers, as it can correct any input string.
- Drawback :- Difficulty in coping with situations in which the actual error has occurred before the point of detection.

3) Error-productions :

- We can augment the grammar for the language with productions that generate the common erroneous constructs.
- Such a parser can detect the anticipated errors when an error production is used during parsing.
- The parser can then generate appropriate error diagnostics about the erroneous construct that has been recognized in the input.

4) ~~Global~~ Correction:

- Given an incorrect input string x and grammar G , these algorithms will find a parse tree for a related string y , such that the no: of insertions, deletions and changes of tokens required to transform x into y is as small as possible.
- That means, these algorithms will choose a minimal sequence of changes to obtain a globally least-cost correction.
- Unfortunately, these methods are too costly to implement in terms of time and space, so these techniques are currently of theoretical interest only (not practical).

Context-Free Grammars : (CFG)

A context-free grammar, CFG , consists of terminals (T_0), variables or the non-terminals (V), productions (P) and a start symbol (S). i.e,

$$G = (V, T, P, S).$$

-) Terminals : Basic symbols from which strings are formed.
-) Non-terminals : variables that denote set of strings. They impose a hierarchical structure on the language that is key to syntax analysis and translation.
-) In a grammar, one non-terminal is distinguished as the start symbol, and the set of strings it denotes is the language generated by the grammar. Conventionally, the productions for the start symbol are listed first.
-) The productions of a grammar specifies the manner in which the terminals and non-terminals can be combined to form strings.

Each production consists of:

- A non-terminal called the head or left side of the production.
- The symbol ' \rightarrow '. Sometimes $::=$ has been used in place of the arrow.
- The body or right side consisting of zero or more terminals and non-terminals. The components of the body describe one way in which strings of the non-terminal at the head can be constructed.

Derivations and Parse-trees:

Derivations: Beginning with the start symbol, each rewriting step replaces a non-terminal by the body of one of its productions. This corresponds to the top-down construction of a parse tree.

For example, consider the below grammar:

$$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$$

The string $-(id + id)$ is a sentence of the above grammar because we can derive this string as given below.

$$\begin{aligned} E &\rightarrow -E \\ &\rightarrow -(E) \\ &\rightarrow -(E + E) \\ &\rightarrow -(id + E) \\ &\rightarrow -(id + id) \end{aligned}$$

- This is the left-most derivation.

The same sentence $-(id + id)$ can also be derived using the given below right-most derivation:

$$\begin{aligned} E &\rightarrow -E \\ &\rightarrow -(E) \\ &\rightarrow -(E+E) \\ &\rightarrow -(E+id) \\ &\rightarrow \underline{\underline{-(id+id)}} \end{aligned}$$

Thus, the derivations are of two types:

- ① Left-most derivation (LMD): the left-most non-terminal in each sentential is always chosen. This can be denoted using ' \xrightarrow{lm} '. i.e., $\alpha \xrightarrow{lm} \beta$.
- ② Right-most derivation (RMD): The right-most non-terminal is always chosen. In this case, we write it as, $\alpha \xrightarrow{rm} \beta$.

Eg: $S \rightarrow AB ; A \rightarrow aaA |\lambda ; B \rightarrow Bb | \lambda$ (' λ ' means ' ϵ ')

Derive the string 'aab' using both LMD & RMD.

Soln:

Left-most derivation :-

$$\begin{aligned} S &\rightarrow AB \\ &\rightarrow aaAB \\ &\rightarrow aaB \\ &\rightarrow aaBb \\ &\rightarrow aab \\ &\underline{\underline{aab}} \end{aligned}$$

This requires 5 steps.

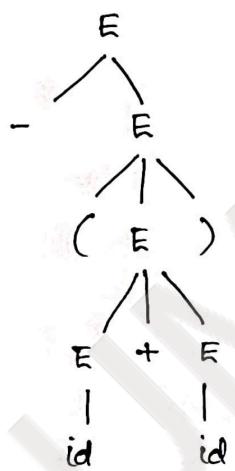
Right-most derivation:

$$\begin{aligned} S &\rightarrow A B \\ &\rightarrow A B b \\ &\rightarrow A b \\ &\rightarrow a a A b \\ &\rightarrow \underline{\underline{aab}}. \end{aligned}$$

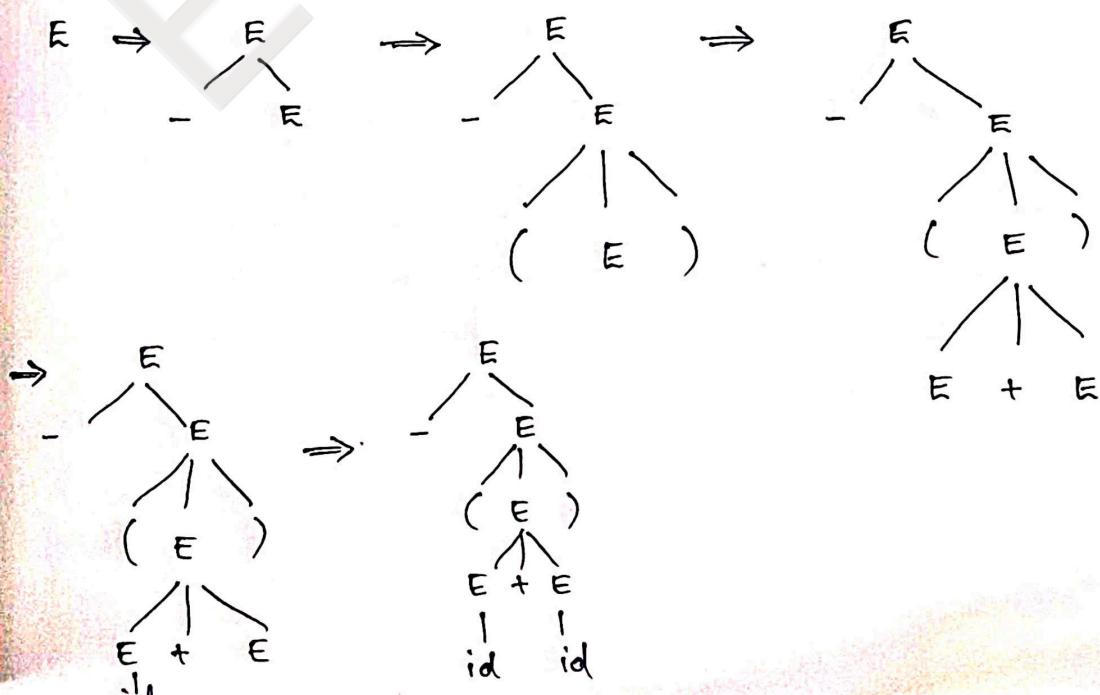
Parse Trees:

Parse tree is a graphical representation of a derivation.

e.g.: For the previous example, -(id+id), the parse tree is,



The sequence for generating this parse tree is, (in case of LRD),



Ambiguity:

A grammar that produces more than one parse tree for some sentence is called ambiguous grammar. That means, an ambiguous grammar is one that produces more than one left-most derivation or more than one right-most derivation for the same sentence.

Eg: The below grammar,

$$E \rightarrow E+E \mid E * E \mid (E) \mid id$$

permits two distinct LMDs for the sentence $id + id * id$.

i,

$$\text{1} \quad E \rightarrow E+E$$

$$\rightarrow id + E$$

$$\rightarrow id + E * E$$

$$\rightarrow id + id * E$$

$$\rightarrow id + id * id$$

$$\text{2} \quad E \rightarrow E * E$$

$$\rightarrow E + E * E$$

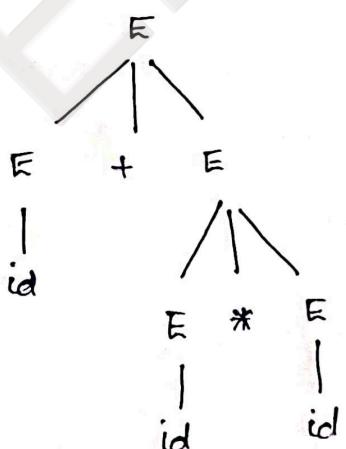
$$\rightarrow id + E * E$$

$$\rightarrow id + id * E$$

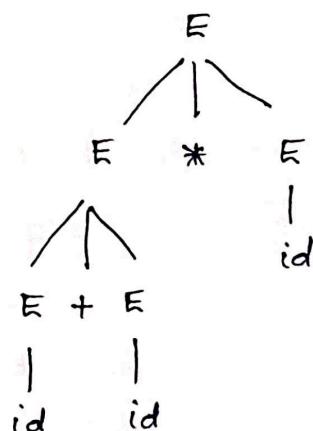
$$\rightarrow id + id * id$$

Thus, the above grammar is considered to be ambiguous.

The corresponding two parse trees are:



and



Elimination of Left-Recursion:

- A grammar is left-recursive if it has a non-terminal 'A' such that there is a derivation $A \xrightarrow{*} A\alpha$ for some string ' α '.
- i.e., A grammar is left-recursive if it is of the form,

$$A \rightarrow A\alpha | \beta$$

- Then, this grammar can be converted into,

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' | \epsilon \end{aligned}$$

[Note: A' can be any non-terminal other than the ones present in the given grammar]

- Now the productions are non-left-recursive productions.
- This also works for any number of A-productions.

i.e., Suppose the grammar is,

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n ,$$

Then we can replace the A production as,

$$\begin{aligned} A &\rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A' \\ A' &\rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \epsilon . \end{aligned}$$

Note:- A grammar of the form $S \rightarrow \alpha S | \beta$ is said to be Right recursive. But right recursion does not create any problem for Top-down parsers. Hence, there is no need of eliminating right recursion.

Examples :-

$$B \rightarrow Bc \mid b$$

Solⁿ:

$$\alpha = c \text{ and } \beta = b.$$

$$\Rightarrow B \rightarrow bB'$$

$$B' \rightarrow cB' \mid \epsilon.$$

=====

Q) Eliminate left-recursion.

1. $A \rightarrow ABd \mid Aa \mid a$

$$B \rightarrow b \mid \epsilon.$$

Solⁿ: $A \rightarrow aA'$

$$A' \rightarrow BdA' \mid aA' \mid \epsilon$$

$$B \rightarrow b \mid \epsilon.$$

2. $E \rightarrow E + E \mid E * E \mid a$

Solⁿ: $E \rightarrow aE'$

$$E' \rightarrow +EE' \mid *EE' \mid \epsilon$$

3. $E \rightarrow E + T \mid T$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow id$$

Solⁿ: $E \rightarrow TE'$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow id$$

4. $A \rightarrow Ac \mid Aad \mid bd \mid cd$

soln: $A \rightarrow bdA' \mid cdA'$

$A' \rightarrow cA' \mid adA' \mid \epsilon$

H:W
5. $S \rightarrow A$

$A \rightarrow Ad \mid Ae \mid aB \mid ac$

$B \rightarrow bBc \mid f$

soln: $S \rightarrow A$

$A \rightarrow aBA' \mid acA'$

$A' \rightarrow dA' \mid eA' \mid \epsilon$

$B \rightarrow bBc \mid f$

Examples of Indirect-left recursion:-

1. $A \rightarrow Ba \mid Aa \mid c$

$B \rightarrow Bb \mid Ab \mid d$

soln: In this example, let's consider the production,

$B \rightarrow Ab$

This seems to be non-left recursive.

But if we replace 'A' by its production $A \rightarrow Ba$, then the above B-production becomes,

$B \rightarrow Bab$

Now this production becomes left-recursive. Hence in such grammars, we have to do one extra step to eliminate left-recursion.

Thus, the solⁿ is:

Step 1: Eliminate left recursions from A-productions.

$$\text{i, } A \rightarrow Ba | Aa | c$$

becomes,

$$A \rightarrow BaA' | cA'$$

$$A' \rightarrow aA' | \epsilon$$

Now, the given grammar becomes,

$$A \rightarrow BaA' | cA'$$

$$A' \rightarrow aA' | \epsilon$$

$$B \rightarrow Bb | Ab | d$$

Step 2:

Substituting the productions of A, in $B \rightarrow Ab$, we get,

$$A \rightarrow BaA' | cA'$$

$$A' \rightarrow aA' | \epsilon$$

$$B \rightarrow Bb | BaA'b | cA'b | d$$

Step 3:

Now eliminate the left recursions from the B-productions.

$$\text{i, } A \rightarrow BaA' | cA'$$

$$A' \rightarrow aA' | \epsilon$$

$$B \rightarrow cA'bB' | dB'$$

$$B' \rightarrow bB' | aA'bB' | \epsilon$$

This is the final grammar after eliminating left recursions.

2) $X \rightarrow XSB | Sa | b$

$S \rightarrow Sb | Xa | a$.

Soln: The above grammar includes indirect left-recursion.

Step 1: Eliminate left-recursion from X .

$$X \rightarrow SaX' | bX'$$

$$X' \rightarrow SbX' | \epsilon$$

$$S \rightarrow Sb | Xa | a$$

Step 2: Substitute X in $S \rightarrow Xa$.

$$X \rightarrow SaX' | bX'$$

$$X' \rightarrow SbX' | \epsilon$$

$$S \rightarrow Sb | SaX'a | bX'a | a$$

Step 3: Eliminate left-recursion from S .

$$X \rightarrow SaX' | bX'$$

$$X' \rightarrow SbX' | \epsilon$$

$$S \rightarrow bX'aS' | aS'$$

$$S' \rightarrow bS' | aX'aS' | \epsilon$$

H:W
3)

$$S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | \epsilon$$

Soln: Step 1: Remove left recursions from A, S .

~~$S \rightarrow Aa | b$~~ . But S is already free from left-recursion.
 ~~$A \rightarrow Ac | Sd | \epsilon$~~

Step 2: Now substitute S in $A \rightarrow Sd$.

$$S \rightarrow Aa | b$$

$$A \rightarrow Ac | Aad | bd | \epsilon$$

Step 3: $S \rightarrow Aa | b$.

$$A \rightarrow bdA' | A'$$

$$A' \rightarrow CA' | adA' | \epsilon$$

Elimination of Left-factoring :

- A grammar contains left-factor if it is of the form:

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | \gamma$$

- To eliminate left-factoring, we have to convert the above grammar into the following form:

$$A \rightarrow \alpha A' | \gamma$$

$$A' \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$$

eg: $S \rightarrow iEts | iEtses | a$
 $E \rightarrow b.$

Here, the S-production contains left-factor. i.e., the 'iEts' (i, common prefix) is common. Hence we convert it to,

$$S \rightarrow iEts s' | a$$

$$s' \rightarrow es | e$$

$$E \rightarrow b.$$

Examples :

Q1) $A \rightarrow ab | ac$

Soln: $A \rightarrow a A'$

$$A' \rightarrow B | c$$

Q2) $S \rightarrow SaB | Sac | ab$

Soln: $S \rightarrow Sas' | ab$

$$s' \rightarrow B | c$$

$$Q3) \quad S \rightarrow aAd | aB$$

$$A \rightarrow a | ab$$

$$B \rightarrow ccd | ddc$$

$$sol^n: \quad S \rightarrow aS'$$

$$S' \rightarrow Ad | B$$

$$A \rightarrow aA'$$

$$A' \rightarrow b | \epsilon$$

$$B \rightarrow ccd | ddc$$

$$Q4) \quad A \rightarrow aAB | aA | a$$

$$B \rightarrow bB | b.$$

$$sol^n: \quad A \rightarrow aA'$$

$A' \rightarrow AB | A | \epsilon \rightarrow$ Here again left factoring, $A' \rightarrow AS | \epsilon$
 $S \rightarrow B | \epsilon.$

$$B \rightarrow bB'$$

$$B' \rightarrow B | \epsilon.$$

$$Q5) \quad A \rightarrow aAB | aBc | aAc$$

$$sol^n: \quad A \rightarrow aA'$$

$$A' \rightarrow AB | Bc | Ac$$

Again, this is a grammar with common prefixes.

$$A \rightarrow aA'$$

$$A' \rightarrow AS | Bc$$

$$S \rightarrow B | c$$

$$Q6) \quad S \rightarrow assbs | asasb | abb | b$$

$$sol^n: \quad S \rightarrow aS' | b.$$

$$S' \rightarrow ssbs | Sasb | bb$$

Again left factoring,

$$S \rightarrow as' | b$$

$$s' \rightarrow SA | bb.$$

$$A \rightarrow Sbs | asb.$$

Q.7)

$$S \rightarrow a | ab | abc | abcd.$$

Soln:

$$S \rightarrow as'$$

$$s' \rightarrow b | bc | bcd | \epsilon.$$

Again taking the common prefix,

$$S \rightarrow as'$$

$$s' \rightarrow bA | \epsilon$$

$$A \rightarrow c | cd | \epsilon$$

Again repeating,

$$S \rightarrow as'$$

$$s' \rightarrow bA | \epsilon$$

$$A \rightarrow cB | \epsilon$$

$$B \rightarrow d | \epsilon.$$

—

H.W.

Q8)

$$A \rightarrow abB | aB | cdg | cdeB | cdfB.$$

Soln:

$$A \rightarrow aA' | cdg | cdeB | cdfB. \quad [\text{Again common prefix : 'cd'}]$$

$$A' \rightarrow bB | B.$$

↓

$$\begin{array}{c} A \rightarrow cA'' | aA' \\ A'' \rightarrow dg | deB | dfB \\ A' \rightarrow bB | B. \end{array}$$

→

$$\begin{array}{c} A \rightarrow cdA'' | aA' \\ A'' \rightarrow g | eB | fB \\ A' \rightarrow bB | B. \end{array}$$

—

TOP-DOWN PARSING:

Top-down parsing can be viewed as the problem of constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in pre-order. Equivalently, top-down parsing can be viewed as finding a left-most derivation for an input string.

Eg: Construct top-down parse tree for the sentence $id + id * id$ for the grammar,

$$E \rightarrow TE'$$

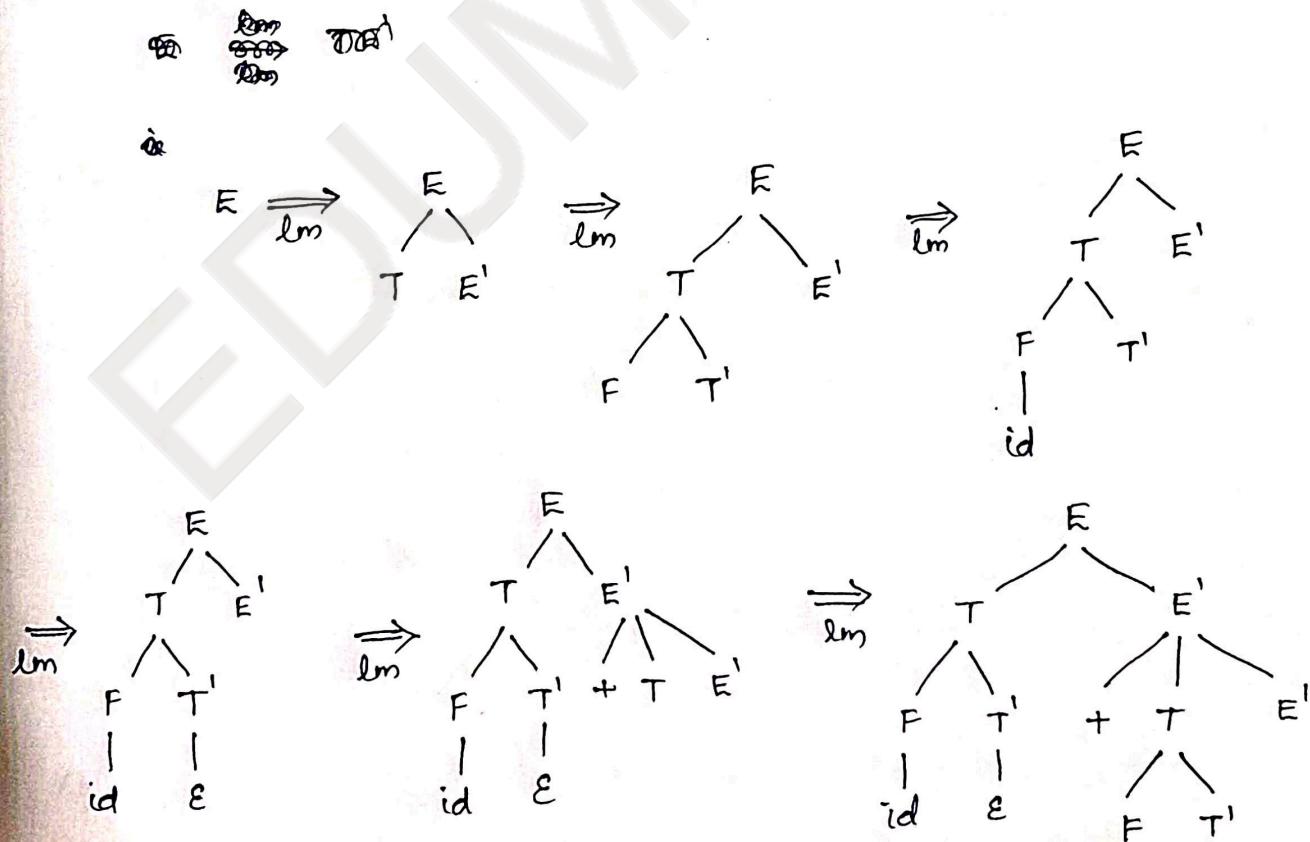
$$E' \rightarrow +TE' \mid \epsilon$$

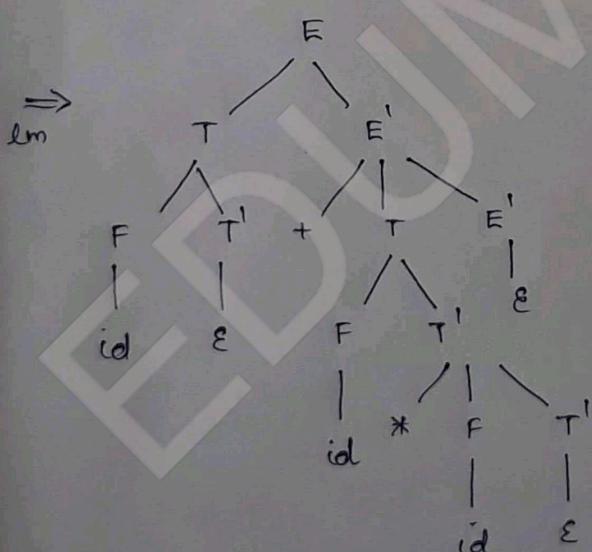
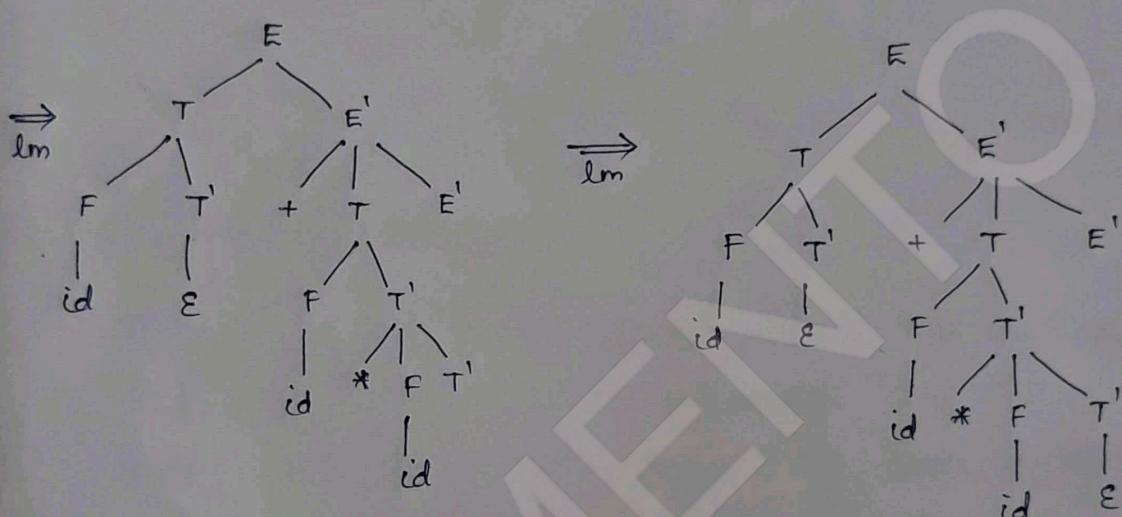
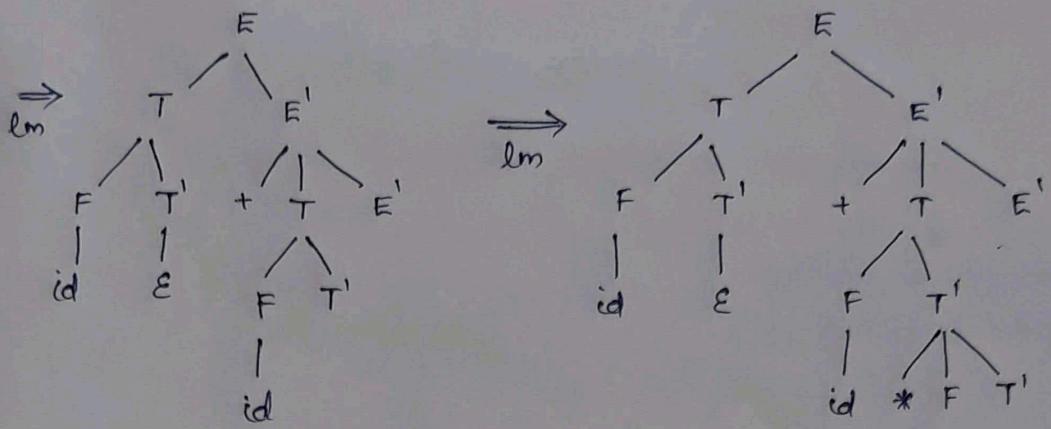
$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

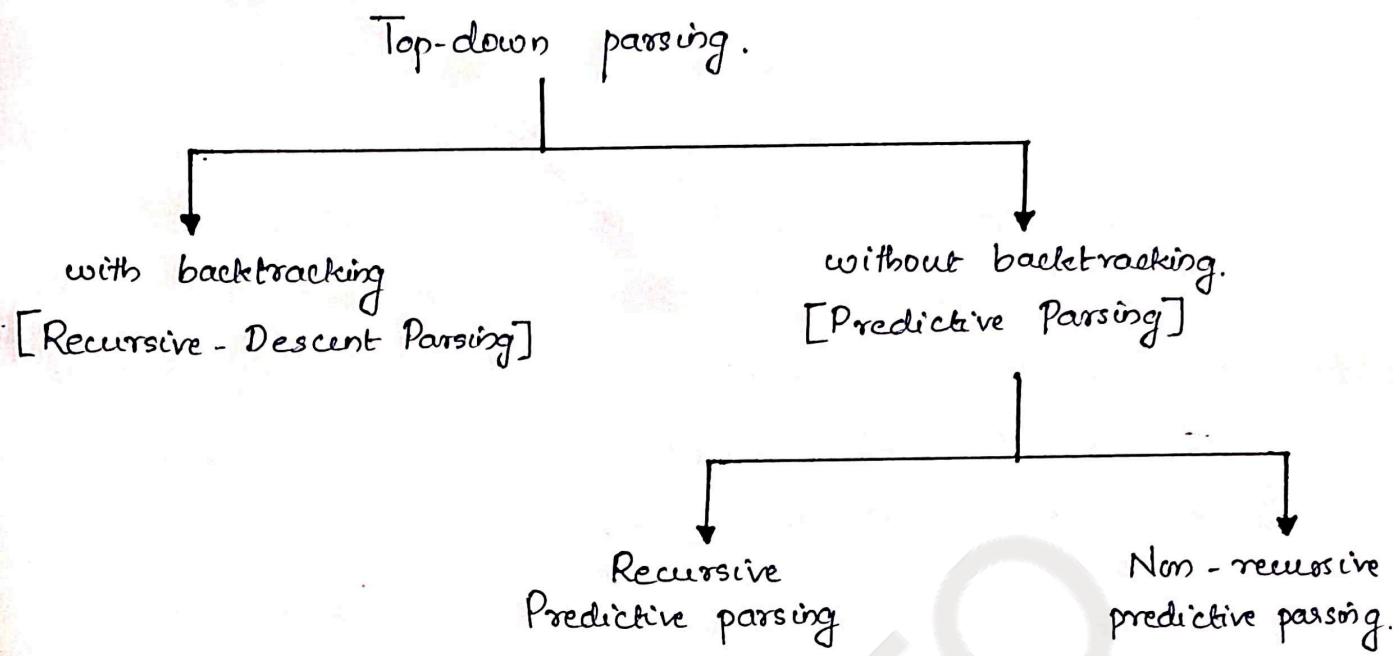
$$F \rightarrow (E) \mid id$$

Soln:





At each step of a top-down parse, the key problem is that of determining the production to be applied for a non-terminal. Either, it sometimes requires backtracking to find the correct ~~or~~ productions, or it needs to look-ahead the input a fixed number of symbols (typically one symbol).



Recursive - Descent Parsing :

- A recursive-descent parsing program consists of a set of procedures, one for each non-terminal.
- Execution begins with the procedure for start symbol, which halts and announces success if its procedure body scans the entire input string.
- Pseudocode for a typical non-terminal is given below.

```

void A()
{
    choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;
    for ( $i=1$  to  $k$ )
    {
        if ( $X_i$  is a non-terminal)
            call procedure  $X_i()$ ;
        else if ( $X_i$  equals the current i/p symbol  $a$ )
            advance the i/p to the next symbol;
        else /* an error has occurred */;
    }
}

```

- General recursive-descent may require backtracking; that is, it may require repeated scans over the input. But, backtracking is not very efficient.
- To allow backtracking, the failure at the final 'else' portion of the above pseudocode is not ultimate failure, but suggests only that we need to return to the first line and try another A-production. Only if there are no more A-productions to try, we declare than an error has been found.

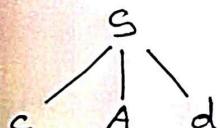
eg: Consider the grammar,

$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a$$

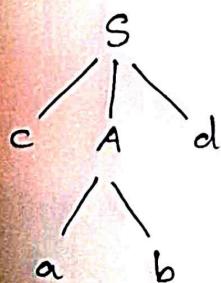
construct a parse-tree for the i/p string, $w = cad$.

Sol:



- i/p pointer will be first pointing to the symbol 'c' of w.

- Now the leftmost leaf, labeled 'c', matches the 1st input symbol of 'w'. So we advance the i/p pointer to 'a'.
- Now consider the next leaf, A. (It's a non-terminal). Hence expand A using the first production $A \rightarrow ab$.

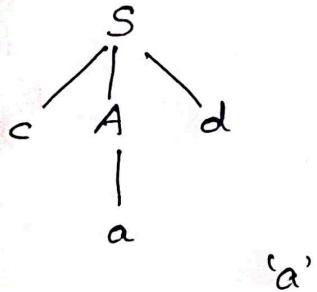


Now the leaf labeled 'a' matches the 2nd i/p symbol of w.

Hence we advance the i/p pointer to the next symbol 'd', and compare it against the next leaf node 'b'. Since,

'b' does not match 'd', we report failure and backtrack to A to see whether there is another alternative for A-production that has not been tried, and produce a match.

- Thus we consider the A-production $A \rightarrow a$.



- Now the leaf 'a' matches the 2nd i/p symbol of w, and the leaf 'd' matches the 3rd i/p symbol. Since we have produced a parse tree for w, we halt and announce successful completion of parsing.

Note :- A left-recursive grammar can cause a recursive-descent parser (even the one with backtracking) to go into an infinite loop.

FIRST & FOLLOW :

- The construction of both top-down and bottom-up parsers is aided by two functions:
 - ① FIRST
 - ② FOLLOW
- During top-down parsing, FIRST & FOLLOW allow us to choose which production to apply, based on the next i/p symbol.
- During panic-mode error recovery, sets of tokens produced by FOLLOW can be used as synchronizing tokens.

FIRST:-

Define ~~FIRST~~ FIRST(α) , where α is any string of grammar symbols, to be the set of terminals that begin strings derived from α .

- If $\alpha \xrightarrow{*} \epsilon$, then ϵ is also in FIRST(α).

- To compute FIRST(x) ,

1) If x is a terminal, then FIRST(x) = { x }

i.e FIRST(terminal) = {that terminal itself}.

2) If $x \rightarrow \epsilon$ is a production, then add ϵ to the FIRST(x).

3) If x is a non-terminal and $x \rightarrow y_1 y_2 \dots y_k$ is a product for some $k \geq 1$,

→ then place ' a ' in FIRST(x) if for some i , ' a ' is in FIRST(y_i), and ϵ is in all of FIRST(y_1), ..., FIRST(y_{i-1}); that is, $y_1 \dots y_{i-1} \xrightarrow{*} \epsilon$.

→ If ϵ is in FIRST(y_j) for all $j = 1, 2, \dots, k$, then add ϵ to FIRST(x).

eg: Consider the grammar,

$$E \rightarrow TE'$$

$$E' \rightarrow +TE'| \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT'| \epsilon$$

$$F \rightarrow (\epsilon)| id.$$

Find the FIRST() of all the non-terminals.

Soln: FIRST(E) = FIRST(T) = FIRST(F) = { (, id) }

$$\text{FIRST}(\epsilon') = \{\epsilon, \epsilon\}$$

$$\text{FIRST}(T) = \text{FIRST}(F) = \{C, \text{id}\}$$

$$\text{FIRST}(T') = \{\ast, \epsilon\}$$

$$\text{FIRST}(F) = \{C, \text{id}\}$$

Examples :-

Q1) $S \rightarrow abc \mid def \mid ghi$. find $\text{FIRST}(S)$.

Solⁿ: $\text{FIRST}(S) = \{a, d, g\}$

Q2) $S \rightarrow ABC \mid ghi \mid jkl$

$$A \rightarrow a \mid b \mid c$$

$$B \rightarrow b$$

$$D \rightarrow d$$

Solⁿ: $\text{FIRST}(D) = d$

$$\text{FIRST}(B) = b.$$

$$\text{FIRST}(A) = \{a, b, c\}$$

$$\begin{aligned}\text{FIRST}(S) &= \{\text{FIRST}(A), g, j\} \\ &= \{a, b, c, g, j\}\end{aligned}$$

Q3) $S \rightarrow ABC$

$$A \rightarrow a \mid b \mid \epsilon$$

$$B \rightarrow c \mid d \mid \epsilon$$

$$C \rightarrow e \mid f \mid \epsilon$$

Solⁿ: $\text{FIRST}(C) = \{e, f, \epsilon\}$

$$\text{FIRST}(B) = \{c, d, \epsilon\}$$

$$\text{FIRST}(A) = \{a, b, \epsilon\}$$

$$\begin{aligned}
 \text{FIRST}(S) &= \text{FIRST}(A) \\
 &\quad \bullet \{a, b, \text{FIRST}(B)\} \\
 &= \{a, b, c, d, \text{FIRST}(C)\} \\
 &= \{a, b, c, d, e, f, \underline{\underline{\varepsilon}}\}
 \end{aligned}$$

Q4) $S \rightarrow aBDh$

$$B \rightarrow cC$$

$$C \rightarrow bc | \varepsilon$$

$$D \rightarrow EF$$

$$E \rightarrow g | \varepsilon$$

$$F \rightarrow f | \varepsilon$$

$$\text{FIRST}(F) = \{f, \varepsilon\}$$

$$\text{FIRST}(E) = \{g, \varepsilon\}$$

$$\text{FIRST}(D) = \text{FIRST}(E)$$

$$= \{g, \text{FIRST}(F)\}$$

$$= \{g, f, \underline{\underline{\varepsilon}}\}$$

$$\text{FIRST}(C) = \{b, \varepsilon\}$$

$$\text{FIRST}(B) = \{c\}$$

$$\text{FIRST}(S) = \{a\}$$

MP
Q5)

$S \rightarrow A$

$$A \rightarrow aB | Ad$$

$$B \rightarrow b$$

$$C \rightarrow g.$$

Soln: The above grammar is left-recursive. Hence first remove left recursion from the grammar and then find the $\text{FIRST}(S)$.

Thus, the grammar changes to,

$$S \rightarrow A$$

$$A \rightarrow aBA'$$

$$A' \rightarrow dA'|\epsilon$$

$$B \rightarrow b$$

$$C \rightarrow g$$

Now find the $\text{FIRST}(.)$.

$$\text{FIRST}(C) = \{g\}$$

$$\text{FIRST}(B) = \{b\}$$

$$\text{FIRST}(A') = \{d, \epsilon\}$$

$$\text{FIRST}(A) = \{a\}$$

$$\text{FIRST}(S) = \text{FIRST}(A) = \{a\}$$

FOLLOW(.) :-

$\text{Follow}(A)$, for a non-terminal A , is the set of terminals that can appear immediately to the right of A in some sentential forms; that is, the set of terminals ' a ' such that there exists a derivation of the form $S \xrightarrow{*} \alpha A \beta$, for some α and β .

Note: If ' A ' is the rightmost symbol in some sentential form, then ' $\$$ ' is in $\text{Follow}(A)$; ' $\$$ ' is a special "endmarker" symbol that is assumed not to be a symbol of any grammar.

To compute $\text{Follow}(A)$, apply the following rule:

- 1) Place $\$$ in ~~follow(A)~~ $\text{Follow}(S)$, where S is the start symbol, and $\$$ is the input right end marker.

- 2) If there is a production $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(B)$ except ' ϵ ' is in $\text{FOLLOW}(B)$.
- 3) If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where $\text{FIRST}(B)$ contains ' ϵ ', then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

eg:- Consider the grammars,

$$E \rightarrow TE'$$

$$E' \rightarrow + TE' | \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' | \epsilon$$

$$F \rightarrow (E) | id$$

find the $\text{FOLLOW}()$ of all non-terminals.

Soln:- $\text{FOLLOW}(E) = \{ \underline{)}, \$ \}$. [since E is the start symbol, we have to add $\$$ to the $\text{FOLLOW}(E)$]

$$\begin{aligned} \text{FOLLOW}(E') &= \{ \text{FOLLOW}(E) \} \\ &= \{ \underline{)}, \$ \} . \end{aligned} \quad // \text{since } E' \text{ is the rightmost symbol.}$$

$$\text{FOLLOW}(T) = \{ \text{FIRST}(E') \}$$

we know that $\text{FIRST}(E) = \{ +, \epsilon \}$. Thus, everything except ' ϵ ' is in $\text{FOLLOW}(T)$; that is, the symbol ' $+$ ' is in $\text{FOLLOW}(T)$.

However, since $\text{FIRST}(E')$ contains ' ϵ ', the productions that contain ' T ' on the right hand side, i.e.,

$$E \rightarrow TE' \text{ and } E' \rightarrow + TE' \text{ becomes,}$$

$$E \rightarrow T \text{ and } E' \rightarrow + T .$$

Now T is the rightmost symbol, and hence,

$\text{FOLLOW}(T)$ includes $\text{FOLLOW}(\epsilon)$ and $\text{FOLLOW}(E')$. (But note that $\text{FOLLOW}(\epsilon) = \text{FOLLOW}(E')$.

Thus,

$$\begin{aligned}\text{FOLLOW}(T) &= \{ +, \text{ FOLLOW}(E) \} \\ &= \{ +,), \$ \}.\end{aligned}$$

Similarly,

$$\text{FOLLOW}(T') = \text{FOLLOW}(T) = \{ +,), \$ \}$$

and,

$$\begin{aligned}
 \text{FOLLOW}(F) &= \{ \text{FIRST}(T') \} \\
 &= \{ *, \text{ FOLLOW}(T) \} \\
 &= \{ *, +,), \$ \}.
 \end{aligned}$$

Imp: Note that follow will never contain ' ϵ '.

(1) Find the FOLLOW(A) and FOLLOW(B).

$$S \rightarrow AaAb \mid BbBa$$

$$A \rightarrow e$$

B → E

5017.

$$\text{FOLLOW}(A) = \{ a, b \}$$

$$\text{follow}(B) = \{ b, a \}$$

$$\text{FOLLOW}(S) = \{ \$ \}$$

୧୨୯

$$S \rightarrow Bb | Cd$$

$$\beta \rightarrow aB|\varepsilon$$

$$c \rightarrow ccl\bar{e}$$

50 | P:

$\text{FOLLOW}(S) = \{\$\}$

$$\text{FOLLOW}(B) = \{ b \}$$

$$\text{FOLLOW}(C) = \{ d \}$$

Q3) $S \rightarrow ACB \mid CbB \mid Ba$

$$A \rightarrow da \mid BC$$

$$B \rightarrow g \mid \epsilon$$

$$C \rightarrow b \mid \epsilon$$

Solⁿ: $\text{FOLLOW}(S) = \{ \$ \}$

$$\text{FOLLOW}(A) = \{ \text{FIRST}(C) \}$$

$$= \{ b, \text{FIRST}(B) \}$$

$$= \{ b, g, \text{FOLLOW}(S) \}$$

$$= \{ b, g, \$ \}$$

$$\text{FOLLOW}(B) = \{ a, \text{FOLLOW}(S), \text{FIRST}(C) \}$$

$$= \{ a, \$, b, \text{FOLLOW}(A) \}$$

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

$$\text{FOLLOW}(C) = \{ \text{FIRST}(B), b, \text{FOLLOW}(A) \}$$

$$= \{ g, \text{FOLLOW}(S), b, h, g, \$ \}$$

$$= \{ g, \$, b, h \}$$

Q4) $S \rightarrow aABb$

$$A \rightarrow c \mid \epsilon$$

$$B \rightarrow d \mid \epsilon$$

Solⁿ: $\text{FOLLOW}(S) = \{ \$ \}$

$$\text{FOLLOW}(A) = \{\text{FIRST}(B)\}$$

$$= \{\underline{d, b}\}$$

$$\text{FOLLOW}(B) = \{\underline{b}\}.$$

Q5) $S \rightarrow ABC$

$$A \rightarrow DEF$$

$$B \rightarrow \epsilon$$

$$C \rightarrow \epsilon$$

$$D \rightarrow \epsilon$$

$$E \rightarrow \epsilon$$

$$F \rightarrow \epsilon$$

soln: $\text{FOLLOW}(S) = \{\$\}$

$$\begin{aligned} \text{FOLLOW}(A) &= \{\text{FIRST}(B)\} \\ &= \{\text{FIRST}(C)\} \quad (\text{since } B \rightarrow \epsilon) \\ &= \{\text{FOLLOW}(S)\} \quad (\text{since } C \rightarrow \epsilon) \\ &= \{\underline{\$}\} \end{aligned}$$

$$\begin{aligned} \text{FOLLOW}(B) &= \{\text{FIRST}(C)\} \\ &= \{\text{FOLLOW}(S)\} \\ &= \{\underline{\$}\} \end{aligned}$$

$$\text{FOLLOW}(C) = \text{FOLLOW}(S) = \{\$\}$$

$$\text{FOLLOW}(D) = \text{FIRST}(E) = \text{FIRST}(F) = \text{FOLLOW}(A) = \{\$\}$$

$$\text{FOLLOW}(E) = \text{FIRST}(F) = \text{FOLLOW}(A) = \{\underline{\$}\}$$

$$\text{FOLLOW}(F) = \text{FOLLOW}(A) = \{\underline{\$}\}$$

Q6)

$$S \rightarrow ABCDE$$

$$A \rightarrow a|\epsilon$$

$$B \rightarrow b|\epsilon$$

$$C \rightarrow c$$

$$D \rightarrow d|\epsilon$$

$$E \rightarrow e|\epsilon$$

solⁿ:

$$\text{FOLLOW}(S) = \{ \$ \}$$

$$\text{FOLLOW}(A) = \{ \text{FIRST}(B) \}$$

$$= \{ b, \text{FIRST}(C) \}$$

$$= \{ \underline{b}, \underline{c} \}$$

$$\text{FOLLOW}(B) = \{ \text{FIRST}(C) \}$$

$$= \{ \underline{\underline{c}} \}$$

$$\text{FOLLOW}(C) = \{ \text{FIRST}(D) \}$$

$$= \{ d, \text{FIRST}(E) \}$$

$$= \{ d, e, \text{FOLLOW}(S) \}$$

$$= \{ \underline{d}, \underline{e}, \$ \}$$

$$\text{FOLLOW}(D) = \{ \text{FIRST}(E) \}$$

$$= \{ e, \text{FOLLOW}(S) \}$$

$$= \{ \underline{e}, \$ \}$$

$$\text{FOLLOW}(E) = \text{FOLLOW}(S)$$

$$= \{ \$ \}$$

LL(1) Grammars :-

Predictive parsers, that is, recursive-descent parsers without backtracking, can be constructed for a class of grammars called LL(1). The first "L" in LL(1) stands for scanning the input from left-to-right, the second "L" for producing a left-most derivation, and the "1" for using one lookahead symbol at each step to make parsing action decisions.

Note:- No left-recursive or ambiguous grammar can be LL(1).

∴ A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha | \beta$ are two distinct productions of G , the following conditions hold :

- 1) For no terminal 'a', do both α and β derive strings beginning with 'a'.
- 2) At most one of α and β can derive the empty string.
- 3) If $\beta \xrightarrow{*} \epsilon$, then α does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$. Likewise, if $\alpha \xrightarrow{*} \epsilon$, then, β does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$.

The first two conditions are equivalent to the statement that $\text{FIRST}(\alpha)$ and $\text{FIRST}(\beta)$ are disjoint sets. The third condition is equivalent to stating that if ϵ is in $\text{FIRST}(\beta)$, then $\text{FIRST}(\alpha)$ and $\text{FOLLOW}(A)$ are disjoint sets, and likewise if ϵ is in $\text{FIRST}(\alpha)$.

The following algorithm collects the informations from FIRST and FOLLOW sets into a predictive parsing table $M[A, a]$, a two-dimensional array, where A is a non-terminal, and ' a ' is a terminal or the symbol '\$', the input endmarker. The algorithm is based on the following idea:

The production $A \rightarrow \alpha$ is chosen if the next input symbol ' a ' is in $\text{FIRST}(\alpha)$. The only complication occurs when $\alpha = \epsilon$ or, more generally, $\alpha \stackrel{*}{\Rightarrow} \epsilon$. In this case, we should again choose $A \rightarrow \alpha$, if the current input symbol is in $\text{FOLLOW}(A)$, or if \$ on the input has been reached and \$ is in $\text{FOLLOW}(A)$.

ALGORITHM :- Construction of a predictive parsing table.

Input : Grammar G

Output : Parsing table M .

Method : For each production $A \rightarrow \alpha$ of the grammar, do the following :

- 1) For each terminal ' a ' in $\text{FIRST}(\alpha)$, add $A \rightarrow \alpha$ to $M[A, a]$.
- 2) If ϵ is in $\text{FIRST}(\alpha)$, then for each terminal ' b ' in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, b]$. If ϵ is in $\text{FIRST}(\alpha)$ and \$ is in $\text{FOLLOW}(A)$, add $A \rightarrow \alpha$ to $M[A, \$]$ as well.

If, after performing the above, there is no production at all in $M[A, a]$, then set $M[A, a]$ to error, which we normally represent by an empty entry in the table.

Eg:- For the below grammar, produce a parsing table.

(Note: For creating the parsing table, you have to find out the $\text{FIRST}()$ and $\text{FOLLOW}()$ of all the non-terminals]

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

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

Sol:- Find out the FIRST() and FOLLOW(). [already done]

$$\text{FIRST}(E) = \{ (, id \} ; \text{ FOLLOW}(E) = \{), \$ \}$$

$$\text{FIRST}(E') = \{ +, \epsilon \} ; \text{ FOLLOW}(E') = \{), \$ \}$$

$$\text{FIRST}(T) = \{ (, id \} ; \text{ FOLLOW}(T) = \{ +,), \$ \}$$

$$\text{FIRST}(T') = \{ *, \epsilon \} ; \text{ FOLLOW}(T') = \{ +,), \$ \}$$

$$\text{FIRST}(F) = \{ (, id \} ; \text{ FOLLOW}(F) = \{ *, +,), \$ \}$$

Now, create the parsing table, [rows \Rightarrow Non-terminals

[rows \Rightarrow non-terminals ; columns \rightarrow all the terminals (except ϵ)
and also add $\$$]]

Non-terminals	Input Symbols (terminals & $\$$)					
	id	+	*	()	$\$$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$				$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$				$F \rightarrow (E)$	

This algorithm can be applied to any grammar G to produce a parsing table M.

•*• IMP :- But, note that for every LL(1) grammar, each parsing-table entry uniquely identifies a production or signals an error. For some grammars, M , (\in the parsing table), may have multiple entries in the same cell. Such grammars are not LL(1). For example, consider the below grammar,

$$S \rightarrow iEtSS' | a$$

$$S' \rightarrow eS | \epsilon$$

$$E \rightarrow b.$$

Construct the parsing table & check if it is LL(1) or not.

Sol:

$$\text{FIRST}(S) = \{ \underline{i}, a \}$$

$$\begin{aligned} \text{FIRST}(S') &= \{ \underline{e}, \underline{\epsilon} \} & \Rightarrow \text{FOLLOW}(S') &= \{ \text{FOLLOW}(S) \} \\ &&&= \{ \$, \text{FIRST}(S') \} \\ &&&= \{ \$, \underline{e} \} \end{aligned}$$

$$\text{FIRST}(E) = \{ \underline{b} \}$$

Parsing table :

	<u>i</u>	<u>t</u>	<u>a</u>	<u>e</u>	<u>b</u>	<u>\$</u>
<u>S</u>	$S \rightarrow iEtSS'$		$S \rightarrow a$			
<u>S'</u>				$S' \rightarrow eS$ $S' \rightarrow \epsilon$		$S' \rightarrow \epsilon$
<u>E</u>					$E \rightarrow b.$	

Here, it is clear that this grammar is not LL(1).

[Because, In the table, $M[S', e]$ contains two productions, $S' \rightarrow eS$ and $S' \rightarrow \epsilon$]

Predictive Parsing :-

The goal of predictive parsing is to construct a top-down parser that never backtracks. To do so, we must transform a grammar in two-ways:

- ① Eliminate left-recursion
- ② Perform left-factoring.

- These rules eliminate the most common causes for backtracking although they do not guarantee to produce LL(1) grammar which is completely back-track free parsing.

Predictive Parsers (i.e. parser w/o backtracking) are of two types :

- ① Recursive Predictive Parser
- ② Non-recursive predictive parser.

Non-Recursive Predictive Parser :-

- It is possible to build a non-recursive predictive parser by maintaining a stack explicitly.
- The key problem during predictive parsing is that of determining the production to be applied for a non-terminal.
- The non-recursive parser looks up the production to be applied for. in a parsing table.

- Requirements :-

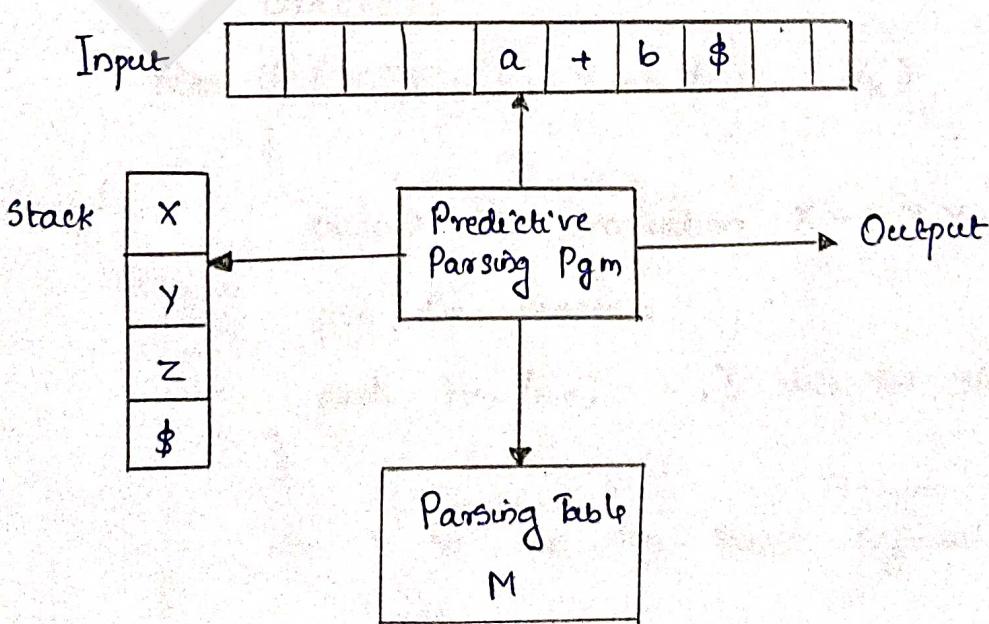
- ① Stack
- ② Parsing Table
- ③ Input Buffer
- ④ Parser.

The parser mimics a left-most derivation. If 'w' is the i/p that has been matched so far, then the stack holds a sequence of grammar symbols α such that,

$$S \xrightarrow[\text{lm}]{*} w\alpha.$$

The table-driven parser has an input buffer, a stack containing a sequence of grammar symbols, a parsing table and an output stream. The i/p buffer contains the string to be parsed, followed by the endmarker \$. We also use \$ symbol to indicate the bottom of the stack, which initially contains the start symbol of the grammar on top of \$.

The parser is controlled by a program that considers X , the symbol on top of the stack, and 'a', the current input symbol. If X is a non-terminal, the parser chooses an X -production by consulting the entry $M[X,a]$ of the parsing table M . Otherwise, it checks for a match between the terminal X and current i/p symbol 'a'.



Algorithm :-

Table-driven predictive parsing.

Input :- A string ' w ' and a parsing table M for grammar G_1 .

Output :- If ' w ' is in $L(G_1)$, a LMD of ' w '; otherwise, an error indication.

Method :- Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol S of G_1 on top of the stack, above $\$$. The given below program uses the predictive parsing table M to produce a predictive parse for the input.

Pgm :-

let ' a ' be the first symbol of ' w ';

let ' X ' be the top stack symbol;

while ($X \neq \$$) // if stack is not empty.

{

if ($X = a$)

pop the stack and let ' a ' be the next symbol of w

else if (X is a terminal)

error;

else if ($M[X, a]$ is an error entry)

error;

else if ($M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$)

{

output the production $X \rightarrow Y_1 Y_2 \dots Y_k$;

pop the stack;

push Y_k, Y_{k-1}, \dots, Y_1 into the stack, with Y_1 on top;

}

let X be the top stack symbol;

}

Example :- Consider the grammar with $w = id + id * id$.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' | \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' | \epsilon$$

$$F \rightarrow (E) | id$$

[Note: we have already created the parsing table M, for this given grammar].

Matched	Stack	Input	Action.
id	E \$	id + id * id \$	
	TE' \$	id + id * id \$	output $E \rightarrow TE'$
	FT'E' \$	id + id * id \$	output $T \rightarrow FT'$
	id T'E' \$	id + id * id \$	output $F \rightarrow id$
	T'E' \$	+ id * id \$	match +
	E' \$	+ id * id \$	output $T' \rightarrow \epsilon$
	+ TE' \$	+ id * id \$	output $E' \rightarrow +TE'$
	TE' \$	id * id \$	match *
	FT'E' \$	id * id \$	output $T \rightarrow FT'$
	id T'E' \$	id * id \$	output $F \rightarrow id$
id + id	T'E' \$	* id \$	match *
	* FT'E' \$	* id \$	output $T' \rightarrow *FT'$
	FT'E' \$	id \$	match *
id + id * id	id T'E' \$	id \$	output $F \rightarrow id$
	T'E' \$	\$	match *
	E' \$	\$	output $T' \rightarrow \epsilon$
	\$	\$	output $E' \rightarrow \epsilon$

Hence the string w is accepted by the parser.

Eliminating Ambiguity :-

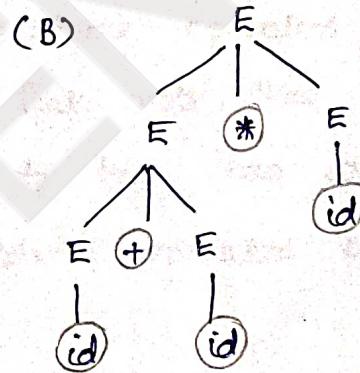
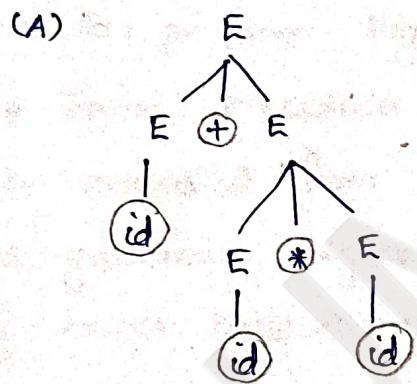
To eliminate the ambiguity from a grammar, you have to re-write the grammar by considering two factors:

- ① Precedence of operators.
- ② Associativity of operators.

eg: Consider the below grammar for the string $id + id * id$.

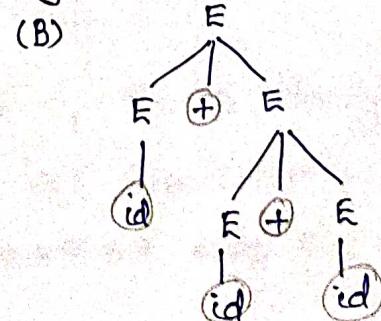
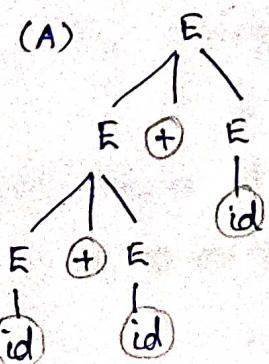
$$E \rightarrow E+E \mid E * E \mid (E) \mid id$$

we can construct two parse trees:



In parse tree (A), $id * id$ is done before '+' operator. But in parse tree (B), '+' is done before *. But we know that * have higher precedence than +. And hence the parse tree (A) is correct.

Similarly, consider the string $id + id + id$.



we know that + is left associative and hence in id + id + id, the first + should be evaluated first followed by the second one. (ie, (id+id)+id). Thus it is clear that the parse tree (A) is correct.

Thus in order to ensure the associativity and precedence of operators, we have to convert the grammars to,

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

$$T \rightarrow T * F \mid F$$

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

i.e., In order to ensure the precedence of operators, we have to perform higher precedence operators first followed by the lower precedence ones. But note that the parse trees are evaluated from bottom to top. Thus while rewriting the grammar, higher precedence operators will be produced by lower productions.

i.e. $E \rightarrow E + E \mid E * E$ is converted like

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

$$T \rightarrow T * F \mid F$$

Now $*$ will be executed before $+$.

Similarly to ensure the associativity of operators, we make the grammar left recursive (for left-associative operators).

i.e. $E \rightarrow E + E$ can be left or right recursive.

By changing it to,

$E \rightarrow E + T \mid T$, this is changed to left recursion only.

and hence it becomes left associative.

a)

Eliminate the ambiguity from the grammar

$$E \rightarrow E * E \mid E - E \mid E \wedge E \mid E / E \mid E + E \mid (E) \mid id$$

The associativity of the operators is as given below. The operators are listed in the decreasing order of precedence.

(i) ()

(ii) / and + are right associative

(iii) ^ is left associative

(iv) * and - are left associative,

Sol:

$$\begin{aligned} E &\rightarrow E * T \mid E - T \mid T && // \text{lowest precedence for } * \text{ & -} \\ T &\rightarrow T \wedge F \mid F \\ F &\rightarrow G \mid F \mid G + F \mid G \\ G &\rightarrow (E) \mid id && // \text{highest precedence for } () \end{aligned}$$

To check the correctness, consider any string.

$$\cancel{(id + id)} * id \Rightarrow (id + id) * id$$

$$E \rightarrow E * T$$

$$\rightarrow T * T$$

$$\rightarrow F * T$$

$$\rightarrow G * T$$

$$\rightarrow (E) * T$$

$$\rightarrow (T) * T$$

$$\rightarrow (F) * T \rightarrow (G + F) * T \rightarrow (id + F) * T$$

$$\rightarrow (id + G) * T \rightarrow (id + id) * T$$

$$\rightarrow (id + id) * F \rightarrow (id + id) * G$$

$$\rightarrow (id + id) * id$$
