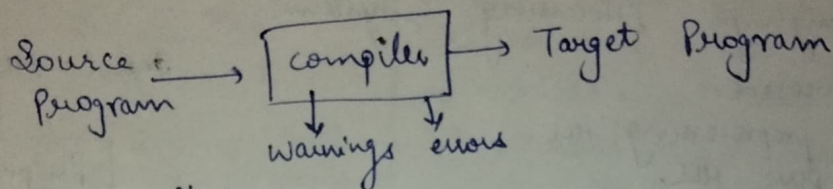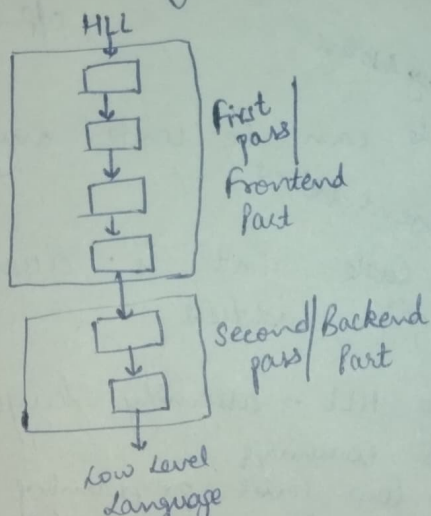# Compiler Design

* Compiler → It is a software/program that converts a program written in HLL ( source language) to a Low level Language ( object / Target language).
  • It also reports errors present in source programs.

Source ——→ | compiler | → Target Program
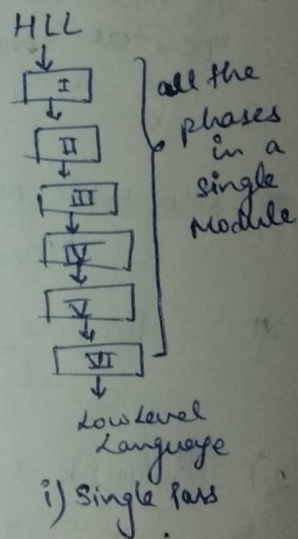Program
              ↓        ↓
          warnings  errors

* Types of compilers
1) Single Pass Compilers → Is a type of compiler that processes the source code only once.

2) Multi Pass Compiler → It is a type of compiler that processes the source code multiple times (to convert HLL → Low Level lang ) i.e to convert source code to target / object code.

HLL

First pass / frontend Part

Second/Backend pass/ Part

Low level Language

ii) Multi Pass (Two Pass)

HLL

I
II
III
IV
V
VI

} all the phases in a single Module

Low Level Language

i) Single Pass

* There are 2 parts of compilation Process/ 2 major phases of Compiler.

Analysis | Synthesis

1) It breaks up the source code / prog. into small parts and creates an Intermediate representation of the source prog.

1) It takes the Intermediate representation of the source prog. as Input and creates the desired target code / program.

**\* Language Processing Systems:**

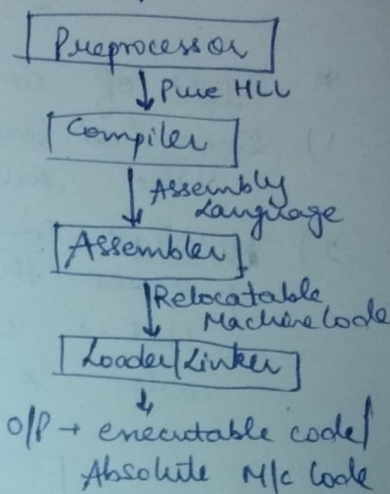We write the programs in a HLL which is easy for us to understand in Compiler design

- These(HLL) programs are fed into a series of tools and OS components to get the desired code that can be used by the Machine. This is called Language Processing System.

1) **Pre - Processor →**
   - In preprocessing, HLL converted to pure HLL.
   - Preprocessor removes the preprocessor directories
     ( # include <stdio.h>
   and will add the respective files
   i.e. <u>file Inclusion</u>
   - It will do <u>macro expansion</u>, operator conversion
     ( eg  a-- ; preprocessing )
              a = a-1

Input = HLL/Source Code
↓
| Preprocessor |
↓ Pure HLL
| Compiler |
↓ Assembly Language
| Assembler |
↓ Relocatable Machine code
| Loader/Linker |
↓
O/P → executable code/ Absolute M/c Code

\# <u>Relocatable M/c code</u> ⎬ <sup>Temporary Address</sup> Code can be loaded anywhere in the Memory.

\# <u>Absolute M/c Code</u> <sup>Permanent Address</sup> → Code that is assembled to work at one specific address

2) **COMPILER →**
   - convert pure HLL → assembly Language
   - also gives errors $ warnings
   - Assembly Language → low level programing language → not in binary form.

3) **ASSEMBLER →**
   <sub>is a program</sub>
   - for every platform (H/w + OS) we have a assembler.
   - A assembler for one platform will not work for another platform.

   Assembly code --converted--> | executable M/c code |

   H/w → intel
   OS → window

\# <u>Executable M/c code</u> → can be loaded at any time for can be run.

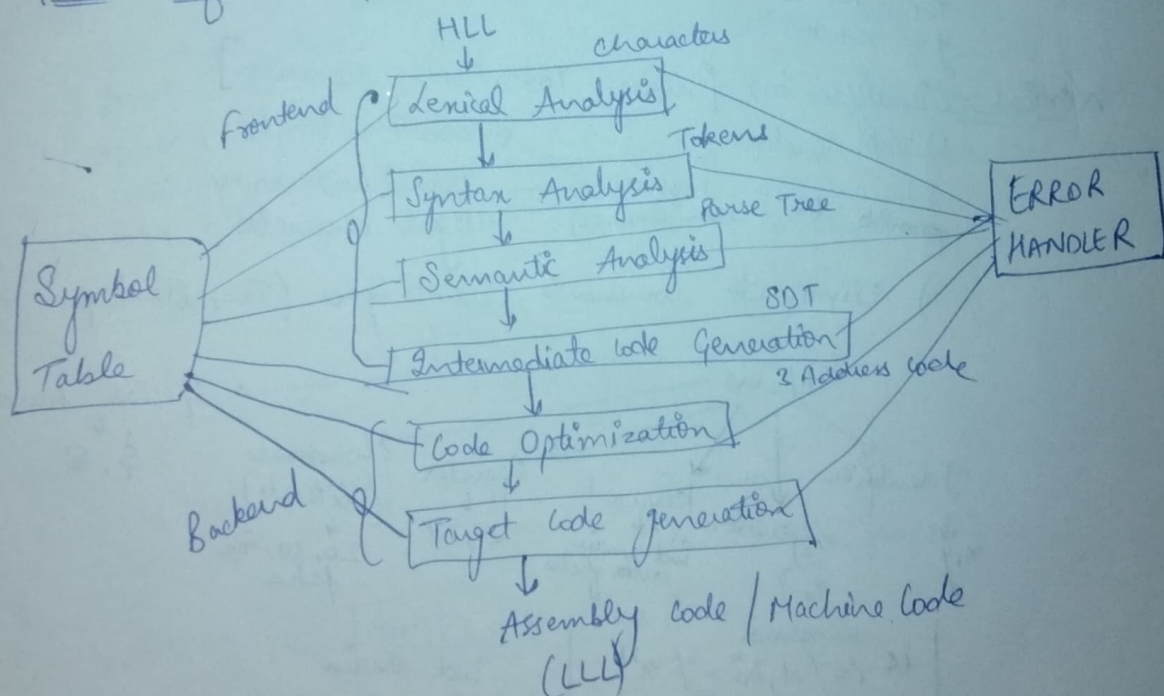4) **LOADER / LINKER →** converts relocatable code → absolute M/c code.
   A | Linker | links diff. obj files into a single file (executable file
   $ then | Loader | — loads that executable file

in the Memory & executes it.

**NOTE** → If we talk of any compiler like Turbo C, gcc etc. all these phases are included to convert Source Code → M/c Code.

\* <u>Compiler</u> | <u>Interpreter</u>

| Compiler | Interpreter |
|---|---|
| 1) Takes entire prog at once as Input. | 1) It takes single line of code at a time. |
| 2) Speed – high | 2) Speed – Low |
| 3) generates Intermediate object Code. So Memory requirement is More | 3) Memory requirement is less because no Intermediate code created. |
| 4) Eg → C, C++, Scala uses compiler | 4) Example → perl, python, Ruby, Matlab |
| 5) Errors → all errors are displayed together | 5) Errors → continues translating the prog. until the 1st error is met, in which case it stops. |
| 6) Error detection is difficult. | 6) error detection is easy |
| 7) Compilers are larger in size. | 7) Smaller in size. |
| 8) | 8) |

\* <u>Phases of Compiler</u> →



example → $y = 2 * x$
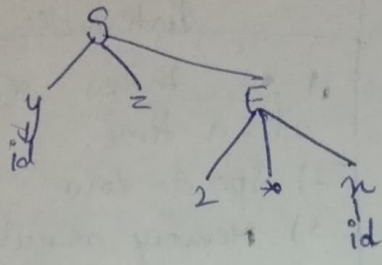
Lexical Analysis to convert into stream of tokens

5 tokens — $y$ = $2$ $*$ $x$ identifier assignment constant operator identifier

$\dfrac{a}{0} \Rightarrow$ Run Time Error

1) **Lexical Analysis** also known as Scanner / Tokenizer

Size, type, scope, reference of variable stored in Symbol Table.
compile Time Error

2) Imp phase → Syntax Analysis act as Parser



3) Semantic Analysis → check logical errors (scoping, properly declared)

SDT → Syntax Directed Translation

4) Imp Phase → Intermediate Code Generation → 3 Address Code → generate Machine Independent Code.

$t_1 = 2 * x$
$y = t_1$

eg $z = a + b * c$
$t_1 = b * c$
$t_2 = a + t_1$
$z = t_2$
} Man. 3 Address } 3 Address Code

5) Code Optimization → Divide Code into Block.
Inside Block → local optimization
Outside Block → Global optimization
$t_1 = 2 * x$
$t_1 = x + x$
$y = x + x$

**Lexical Analysis** → [Lexer, Tokenizer, Scanner]
1) Tokenization
2) Give Error Message — Exceeding length / Unmatched string / Illegal character
provide Lexical error (Lexemes)
3) Eliminate Comments, white Space (Tab, Blank space, Newline)

Tokens ⟶ Special character $\$, @$

Identifier
$x, y$

Separator
$\{, \}$

keyword
int, main, auto, if, else

Operator
$<, >, +$

Constants (literals)
20, 30, True, false

Lexemes
" i = %d, &i = % n "
1
int
main

int main
1    2

Lexeme → Sequence of characters from the Input that match a pattern.

Tokens → Symbolic names for the entities that make up the text of the program. ex: ID, Constant, Keywords, Operators, Punctuations, Literal String.
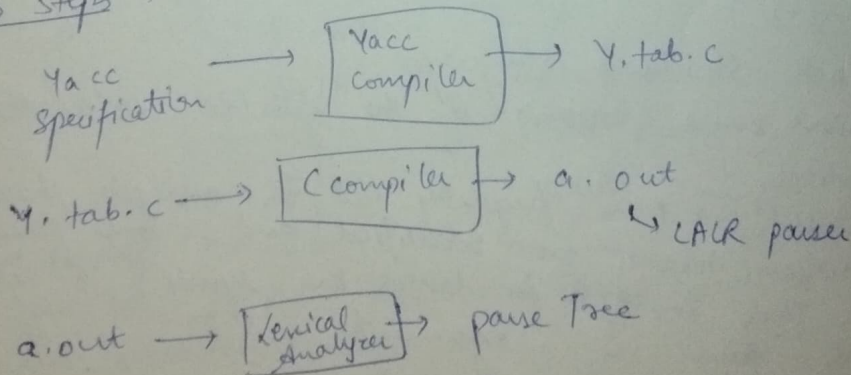
Source Code ⟶ Lexeme ⟶ Token ⟶ Parser

int x = 5;

| Lexeme | Token |
|--------|-------|
| int | Keyword |
| x | id |
| = | Assignment operator |
| 5 | Constant |
| ; | operator |

Keyword id operator constant operator

* **YACC** → It stands for Yet Another Compiler-Compiler (developed by Stephen C Johnson)
- It is a tool for generating Look Ahead Left-to-Right (LALR) parser.
- It takes Input from the Lexical Analyzer and generates parse Tree.
- Syntax Analyzer / Parser is the 2nd phase of the compiler which takes Input as tokens and generates a parse tree.

<u>3 Steps</u> →

Yacc Specification ⟶ [Yacc Compiler] → Y.tab.c

Y.tab.c ⟶ [C compiler] → a.out
                              ↳ LALR parser

a.out ⟶ [Lexical Analyzer] → parse Tree

**LEX Tool** → Lex is a tool / computer program which generates Lexical Analyzer. It is used with YACC parser / generator.
- 1st phase of compiler which converts source code with HLL → stream of Tokens.
- Lex written by Mike Lesk and Eric Schmidt and described in 1975.

Len Source program file.l → | Len compiler | → (lex.yy.c c program) → | Compiler | → a.out (lexical analyser)

I/P Stream → | a.out | → Stream of Tokens

\* **Symbol Table** → Data Structures that are used by compilers to hold Info. about source-program constructs.

• It is used to store Info. about the occurence of various entities such as objects, classes, variable names, functions etc.

• It is used by both analysis phase and synthesis phases. used for following purposes →

1) It is used to store the name of all entities in a structured form at one place.

2) It is used to verify if a variable has been declared

3) It is used to determine the scope of a name.

4) It is used to implement type checking by verifying assignments and expressions in the source code are semantically coueit.

A symbol table can either be linear or hash table

example    int count;
           char x[7]= "NESO ACADEMY";

| Name | Type | Size | Dimension | Line of Declaration | Line of Usage |
|------|------|------|-----------|---------------------|---------------|
| count | int | 2 | 0 | — | 6 |
| x | char | 12 | 1 | — | — |

Operations →

1) Non-Block Structured language.
• Contains single Instance of the variable declaration
• operations:
   a) Insert() → More frequently used in analysis phase (frontend) when tokens are Identified and names are stored in table. The Insert() fn takes the symbol & its value in the form of argument.
           eg    insert (x, int)

b) lookup() → used to search a name & It determines!
   • The existence of symbol in the table
   • Declaration of symbol before It is used
   • check whether the name is used in the scope
   • Initialization of the symbol.

- checking whether the name is declared multiple items.

lookup (symbol)

2) Block Structured languages variable declaration may happen ~~may~~ multiple times.
   a) Insert()
   b) Lookup()
   c) Set() → Specify [The scope of variable]
   d) Reset() → Redefine the scope of variable

Imp → T1: $a?\ (b|c)^* a$

T2: $b?\ (a|c)^* b$

T3: $c?\ (b|a)^* c$  if $c? = 1$ then $(b|a)^* c$ not processed

lexical analyser uses these patterns to recognize 3 ~~typ~~ tokens

T1, T2, T3 over the alphabet a, b & c.

Note that $\boxed{x?}$ specifies 0 or 1. occurence of the symbol $n$

Noted that the analyser outputs the token that matches the longest possible prefix.

If the string "bbaacabc" is processed by the Analyser then which one the following is the sequence of tokens it outputs

a) $T_1 T_2 T_3$    b) $T_1 T_1 T_3$

e) $T_2 T_1 T_3$    d) $T_3 T_3 \frac{bb}{T_2} \frac{aa}{T_1}$

$\frac{bb a}{T_1} \quad \frac{aca}{T_1} \quad \frac{bc}{T3} - (b)$

$\frac{cabc}{T3} - (c)$

→ all options are right (d) is most suitable

d) $\frac{bbaac}{T3} \quad \frac{abc}{T3}$

* **Recursive Grammar:** The grammar G is said to be recursive if in atleast one production there is same variable at both LHS & RHS.

Example:  $S → Sa | b$     Left Recursion
          $S → as | b$      Right Recursion
          $S → a Sb | \epsilon$  Middle Recursion

Types → 1) Left Recursion → The grammar is said to be left recursion if the leftmost variable of RHS is same as the variable at LHS

eg → 1) $S → Sa | b$
     2) $S → AB$
        $A → Aa | b$
        $B → ab | b$

2) Right Recursion → The grammar is said to be right recursion if the rightmost variable of RHS is same as the variable of LHS.

eg → $S → a S | b$

Grammar which is both left & right recursion is ambiguous.

$$eg^n \quad S \to SS | AB$$
$$A \to Aa | a$$
$$B \to Bb | b$$

NOTE → If the grammar is Left Recursive then the parser goto infinite loop so to avoid looping we need to convert the left recursive grammar into right recursive grammar.

* **Conversion of LR into RR or Removal of LRY**

① $A \to A\alpha | \beta$     LR

    ↓↓ removed

$$A \to \beta A'$$
$$A' \to \alpha A' | \varepsilon$$

② $A \to A\alpha_1 | A\alpha_2 | -- | A\alpha_n$

$A \to \beta_1 | \beta_2 | -- | \beta_n$

    ↓↓

$$A \to \beta_1 A' | \beta_2 A' | \beta_3 A' | -- | \beta_n A'$$
$$A' \to \alpha_1 A' | \alpha_2 A' | -- | \alpha_n A' | \varepsilon$$

**Example**

$$E \to E+T | T \qquad LR$$
$$T \to T*f | f \qquad LR$$
$$F \to id | (E) \qquad \text{No Recursion}$$

**Solⁿ**

$$E \to E+T | T$$

$A \to A\alpha | \beta$

    ↓↓

$A \to \beta A'$        $E \to TE'$

$A' \to \alpha A' | \varepsilon$     $E' \to +TE' | \varepsilon$

$$T \to T*f | f$$
$$T \to FT'$$
$$T' \to *FT' | \varepsilon$$
$$f \to id | (E)$$

**Eg 2 →** $S \to \underset{\alpha_1}{\underline{SaSb}} | \underset{\alpha_2}{\underline{SbSa}} | \underset{\beta}{\varepsilon}$

**Eg 3 →** $S \to \dfrac{SSS}{\alpha} | \dfrac{0}{\beta}$

$$S \to \varepsilon S'$$
$$S' \to aSbS' | bSaS' | \varepsilon$$

$$S \to 0S'$$
$$S' \to SSS'$$

**Eg 4 →** $S \to \underset{\alpha}{\underline{SOS}} | \underset{\beta}{\underline{\varepsilon}}$

$$S \to \varepsilon S'$$
$$S' \to OSIS' | \varepsilon$$

\* **Grammars :-** It is basically set of rules that defines the valid structure of a particular language.

$$G\ (V, T, P, S)\ \text{— Start symbol}$$

- Variables
- Terminals
- Productions

i) Set of terminals (i.e. that terminate, they are not replaced by any other thing further)

ii) Set of non terminals (values/variables that are replaced by terminals). L.H.S

iii) Set of productions :-
on LHS → Non Terminal followed by arrow & on RHS we can have T or V or combo of both

iv) Start Symbol → One of the non-terminal is designated as the start symbol from where the production begins.

## Classification of grammar

1) **Deterministic & Non-Deterministic →**

Non-deterministic → for a given grammar if we have many options on single symbol / variable than that grammar is called Non deterministic.
In NDG we have problem of back tracking.
Example A → $\alpha\beta_1 / \alpha\beta_2 / \alpha\beta_3 / \alpha\beta_4$
we need string $\alpha\beta_4$



Deterministic → Backtracking problem will remove

A → $\alpha A^1$
$A^1$ → $\beta_1 / \beta_2 / \beta_3 / \beta_4$

This process is known as left factoring or eliminating non-deterministic.

2) **Recursive and Non-Recursive →**

Left Recursive → if the leftmost symbol in RHS is equal to symbol of LHS.
$$A → A\alpha / \beta$$

Right Recursive → if the rightmost symbol in RHS is equal to symbol of LHS.
$$A → aA / \beta$$

3) **Ambiguous and Non-Ambiguous →**
A G is said to be ambiguous if there exists more than one derivation tree for the given i/p string. or
More than one leftmost derivation.
or " " " rightmost "

Example:- $E \rightarrow E+E / E*E / id$

we have to derive $id+id*id$

Sol^n  LMD  $E \rightarrow E+E$
$\phantom{E \rightarrow} id+E$
$\phantom{E \rightarrow} id+E*E$
$\phantom{E \rightarrow} id+id*E$
$\phantom{E \rightarrow} id+id*id$

RMD  $E \rightarrow E+E$
$\phantom{E \rightarrow} E+E*E$
$\phantom{E \rightarrow} E+E*id$
$\phantom{E \rightarrow} E+id*id$
$\phantom{E \rightarrow} id+id*id$

RMD  $E \rightarrow E*E$
$\phantom{E \rightarrow} E*id$
$\phantom{E \rightarrow} E+E*id$
$\phantom{E \rightarrow} E+id*id$
$\phantom{E \rightarrow} id+id*id$

Parse Tree-1 (Is Valid)



Parse Tree-2



So we have more than 1 parse tree possible ∴ Grammar is ambiguous.

Ambiguous in the sense that " if we have 2 parse trees possible then parser will get confused about which one to generate" or which one is correct.

Ques. check Grammar is ambiguous or not?

i). $S \rightarrow aS / Sa / a$

String — aa

Sol^n



Yes

ii) $E \rightarrow E+E$
$E \rightarrow E-E$
$E \rightarrow id$

$id+id-id$

Sol^n



Yes

iii) $S \to aSb \mid SS$
$S \to \epsilon$      $aabb$
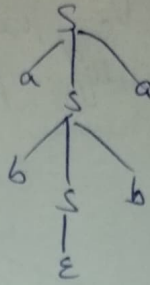


Yes

iv) $S \to aSa \mid bSb \mid \epsilon$      $abba$
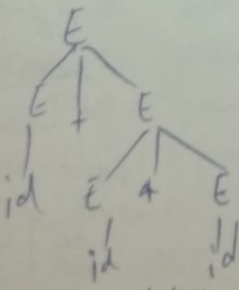


Unambiguous.

v) $S \to SaS \mid b$
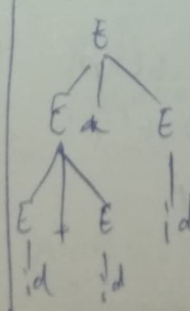


babab         babab

Yes

* Converting Ambiguous → Unambiguous.

$E \to E+E \mid E*E$
$E \to id$

① $id + id + id$   Grammar failed because rules of associativity failed

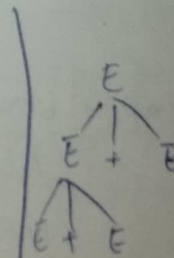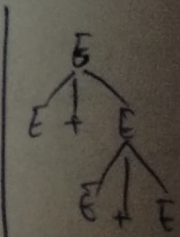② $id + id * id$   Grammar failed ∵ precedence is not taken care of



* evaluated 1st so wrong

+ is evaluated yet ∴ wrong
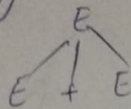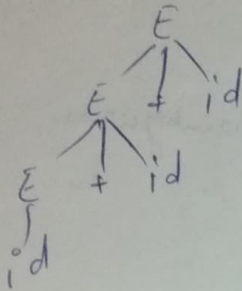
Left Associative | Right Associative

**Getting Both associativities :-** we are defining the GRAMMAR without any order. So we can grow in any direction.

$$E \to E + E$$

To achieve left associativity, we have to grow in left direction only. So now we will restrict the growth of parse tree.

$$E \to E + id \,|\, id$$

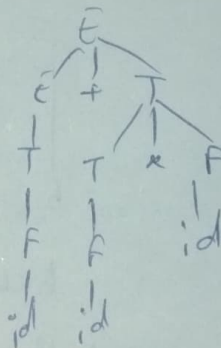**left Recursive**
So, there is no possibility of getting diff. parse Tree.

**for precedance Problems**

we should take care that highest precedance operator should be at the least level.

$+ < *$

$$E \to E + T \,|\, T$$
$$T \to T * F \,|\, F$$
$$F \to id$$

$$id + id * id$$

**Example 2:-** $2 \uparrow 3 \uparrow 2$    $2^{3^2}$    **Right Associative**

$$F \to G \uparrow F \,|\, G$$
$$G \to id$$

**eg 3:-**
$$bExp \to bExp \; OR \; bExp$$
$$\,|\, bExp \; AND \; bExp$$
$$\,|\, NOT \; bExp$$
$$\,|\, True$$
$$\,|\, False$$

OR — less
AND
NOT — least level

precedance
NOT > AND > OR
OR ⎱ left Associative
AND ⎰

$$G \to Not \; G \,|\, True \,|\, False$$

**Soln:-**
$$E \to E \; Or \; F \,|\, F$$
$$F \to F \; And \; G \,|\, G$$

Egn:    A → A $ B | B          Tell associativity & precedence

        B → B # C | C

        C → C @ D | D

        D → d

Sol^ns   $ → Left  Associative          Precedence

        # →    "        "    "            @ > # > $

        @ →    "        "    "

Egn:  E → E + F | E * F

        F → id

   Sol^n:    + and *   are   at   same   level

                ∴ same  precedence

                 & both defined as Left Associative


**\* Parsers :→ (Syntax Analyzer)**

                        Parsers

        Top Down Parser ←        → Bottom up Parsers → or Shift Reduce Parsers

TDP with full backtracking        TDP without backtracking          operator Precedence parser

        ↓                                                              ↓
      Brute Force Method                                          only used for operator Grammar

                Recursive Descent Parser        Non-Recursive Descent LL(1) Parser

                                                                                        LR Parsers

                                        LR(0)    SLR(1)    LALR(1)    CLR(1)


        LR parsers can scan string from L→R
        use reverse of RMD

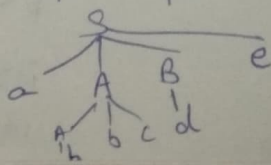\* Bottom up parser :→ generates the parse tree for the given input string with the help of grammar productions by compressing the non-terminals. i.e, it starts from the terminals and ends on the start symbol

        It uses reverse of Right most derivation.

Top Down :→ generates parse tree for the given I/P string with the help of grammar productions by expanding the non-terminals i.e. it starts from the start symbol & ends on the terminals
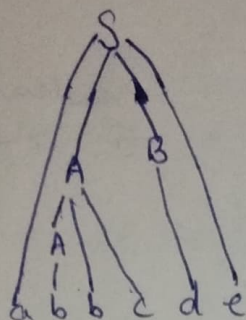
        It uses LMD.

Egn:  S → aAbe      A → Abc | b     B → d

                                                abbcde

Top Down → what is the next production we should use
Bottom Up → when to reduce the given terminal



* **Operator Precedence Parser:-** Ambiguous grammar are not allowed in any parser except operator precedence parser.

Eg:- $E \rightarrow E+E \mid E*E \mid id$

Soln:- operator relation Table / Operator precedence relation table

|     | id  | +   | *   | $   |
| --- | --- | --- | --- | --- |
| id  | —   | >   | >   | >   |
| +   | <   | >   | <   | >   |
| *   | <   | <   | >   | >   |
| $   | <   | <   | <   | —   |

$id > +$

- Two id's will never compared b/c they will never come side by side.
- Identifier will be given highest precedence compared to any other operator.
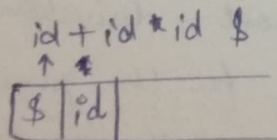- $ has least precedence compared to any other operator
- + is left associative
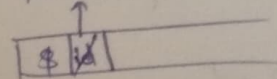$$\frac{1+3+3}{3+3}$$
- * precedence > + precedence

Input → id+id*id

id+id*id $
↑ *

| $ | id |   |
stack

id+id*id $
↑

| $ | id |
id+id*id$
↑
| $ | + |

id+id*id$
↑
| $ | + | id |
id+id*id$

↑ look ahead :- Pointer

$ < id     Push id

id > +     Pop    E
id – E                id+id * id
$ < +     Push

id > +     Push
id > *     Pop     id – E    E    E
                          id+id*id

Top of Stack < look ahead ∴ Push

| $ | id | + | id | * |

id + id * id $
        ↑

id+id*id$

| $ | * | id |

| $ | + | * | id |

| $ | + | * |

*<id Push

id>$ Pop
id – E
*>$ Pop

E       E        E
|        |        |
id + id * id

    E
   /|\
  E  E  E
  |  |  |
  id + id * id

| $ | * |

    E
   /\
  E  E
  |  |
  E  E  E
  |  |  |
  id + id * id

+>$ Pop

**Disadvantage** → No. of entries i.e. if we have 4 operators
we will have 16 entries
for 10 operators we will have 100 entries
for N operator → $O(n^2)$ Size of table will be very big

So to dec. the size of the table, we use operator
for table.