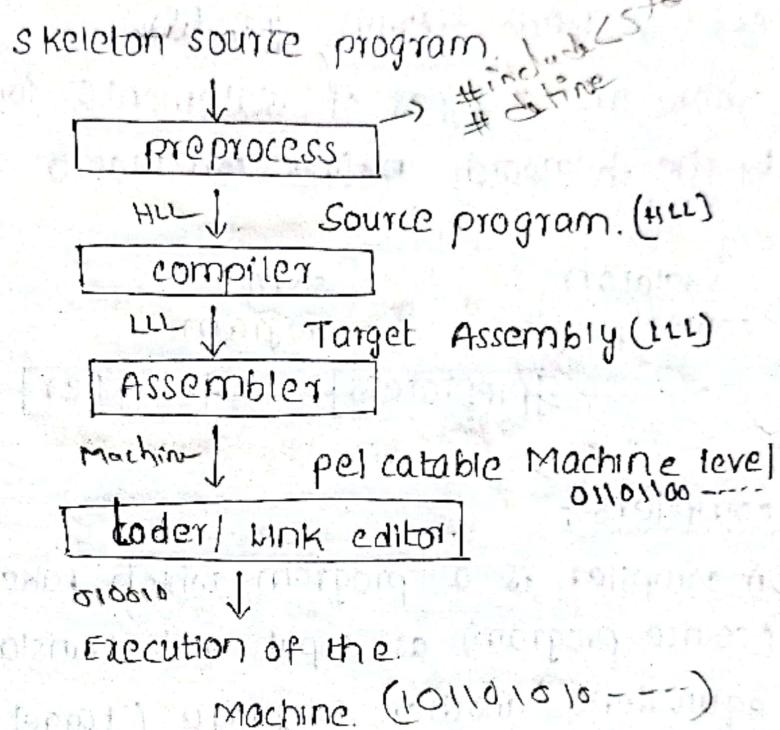


# unit-1

## Introduction to the compilers:

To create a executable form of your source code we may require several programs. These programs work together to produce an executable target program.

The process of execution of program is



The compiler, takes a source program written in a highlevel language as an input and converts it into a target assembler language.

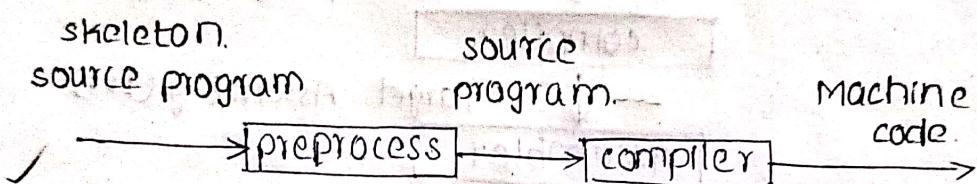
The assembler then takes this target assembler code as input and produces a relocatable machine code as an output then the task of loading and link editing is called to perform the loading of the program into the memory and proper locations and link editor links the object modules and prepares the single module from several modules.

## PREPROCESS:-

- 1) The output of the preprocess may be given as input to compilers.
- 2) preprocess allows users to access Macros in the program. Macro means a set of instructions which can be repeatedly in the program.
- 3) preprocess also allow users to include header file which may required to program.

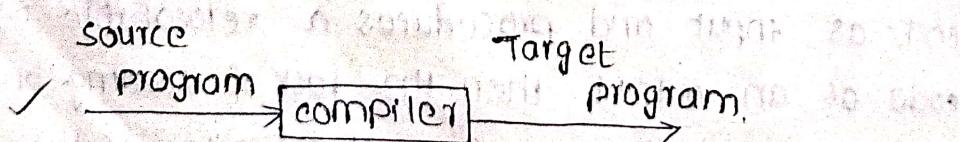
Ex: `#define PI(3.14), #include`

There are 2 types of statements for Macros is given by the keyword define (or) Macro



## COMPILERS:-

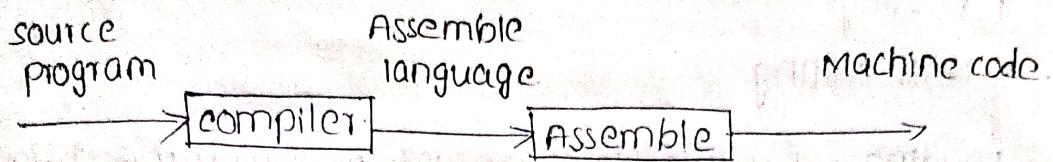
- 1) compiler is a program which takes first language (source program) as input and translates it into an equivalent another language (target program)
- 2) During the process of Translation if some errors are encountered then compiler displays them as error messages.
- 3) The model of Basic compiler represented as



The compiler takes source program as high level language (pascal, c, fortran) and converts into low level (or) Machine level language such as Assembly language. Assembly to Machine level language.

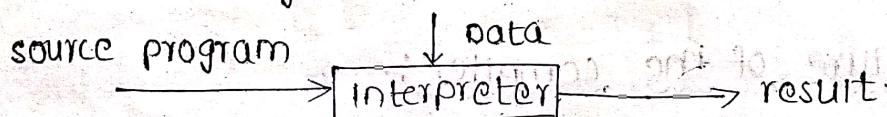
## Assembler :-

Some compilers produce Assembly code as output which is given to Assembler as input. The Assembler is a kind of translator which takes the Assembly program as input and produces:



## Interpreter :-

An interpreter is a kind of translator which produces the result directly when the source program and data is given to it as input. It does not produce objective code but each time program needs execution language needs SNOBOL, LISP, can be translated using interpreter.



The process of interpretation can be carried out in different phases -  
1. LEX Analysis  
2. Syntax Analysis  
3. Semantics Analysis  
4. Direct execution.

## Advantages :-

- 1) Modification of user program can be easily made and implemented as the execution process.
- 2) The type of object that denotes variable may change dynamically.
- 3) Debugging a program and finding an error is a simplified task.
- 4) The Interpreter for the language makes a Machine independent.

## Disadvantages :

- 1) The execution is slower.
- 2) Memory consumption is more.

## Loaders and Linkers :

Loaders is a program which performs two functions

- 1) Loading
- 2) Link editing

Loading: Loading is a process in which the relocatable machine code is read and the relocatable addresses are altered. Then the code with altered instructions and the data is placed in the memory at proper locations.

Link editing: The job of link editor is to make a single program from several files. If code in one file refers the location in another file then such a reference is called External reference.



## Structure of the compiler:

The compilation can be done in 2 parts.

1) Analysis

2) synthesis part

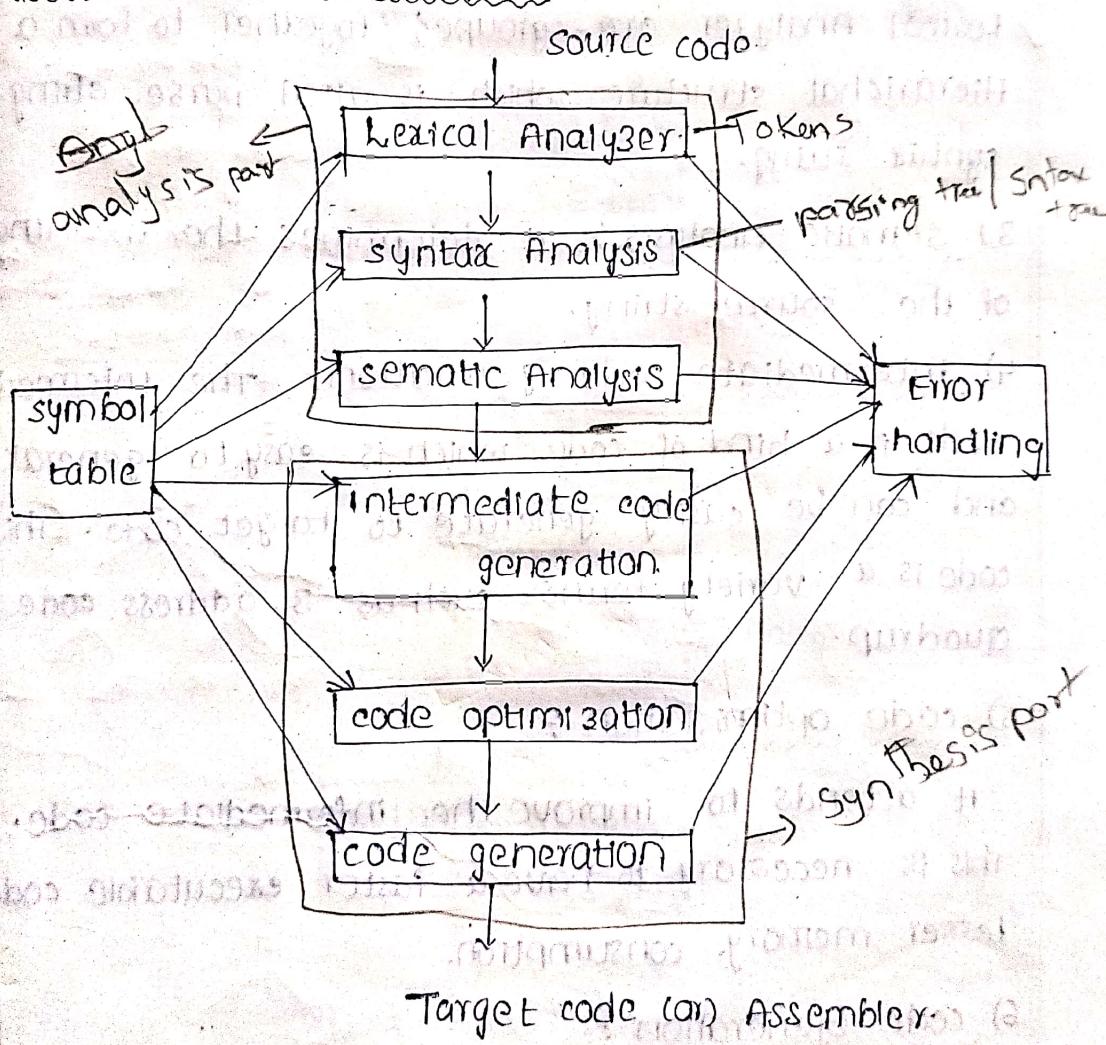
\* Analysis: In Analysis part the source program is read and broken down into pieces. The syntax of the meaning of source string is determined and then an intermediate code is generated from the source code.

Synthesis part: In synthesis part the intermediate form of the source code is taken and converted into an equivalent target program during this process if certain code has to be optimized for efficient execution.

The Analysis part is carried out into 3 parts.

- 1) Lexical Analysis: The source program is read and broken into stream of tokens. strings which are called tokens.
- 2) Syntax Analysis: The tokens are arranged in a hierarchical structure that helps for finding the syntax of the source string.
- 3) Semantic Analysis: The meaning of the source string is determined.

## phases of the compiler:



1) Lexical Analysis :- It is also called scanning. The complete source code is scanned and it broken up into group of strings called Tokens. The token is a sequence of characters having a collective meaning.

Ex:-  $\text{total} = \text{count} + \text{rate} * 60$

\* The blank characters which are used in the programming statements are eliminated during lexical analysis.

2) Syntax Analysis :- The syntax analysis is also called parsing. In these phase the tokens are generated by Lexical Analyzer are grouped together to form a Hierarchical structure which is called parse string (or) syntax string.

3) Semantic Analysis :- It determines the meaning of the source string.

4) Intermediate code generation :- The intermediate code is a kind of code which is easy to generate and can be easily generate to target code. (This code is a variety forms such as 3 address code, Quadrup.)

5) Code optimization :-

It attends to improve the intermediate code. This is necessary to have a faster executable code (or) lesser memory consumption.

6) Code Generation :-

In this phase, the target code gets generated the intermediate code is translated to the machine instructions.

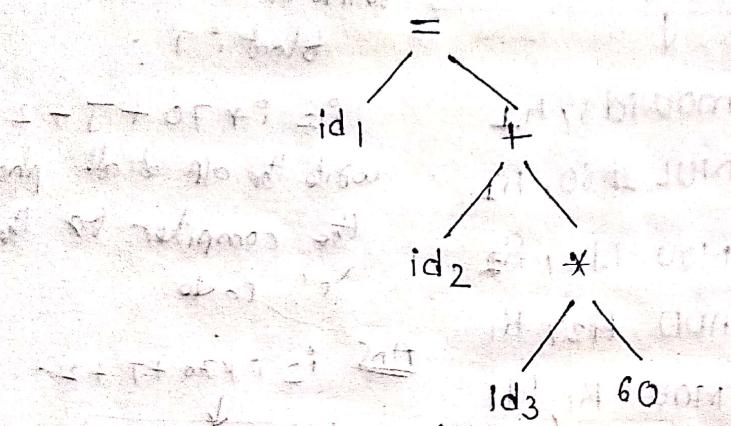
Ex:  $\text{total} = \text{id}_1 + \text{id}_2 * 60$

Identifier  $\rightarrow$  Assignment operator  $\rightarrow$  Int/Float

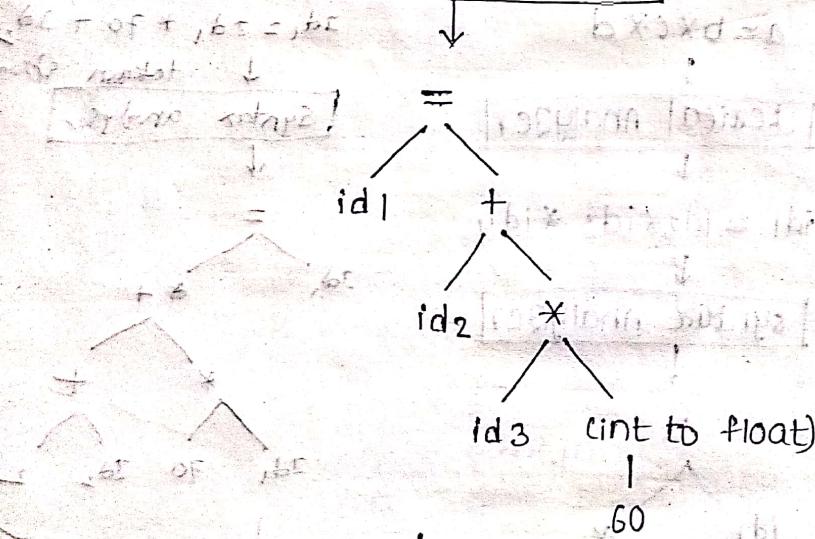
### Lexical Analyzer

$\text{id}_1 = \text{id}_2 + \text{id}_3 * 60$

### Syntax Analyzer



### Semantic Analyzer



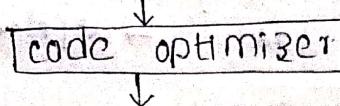
### Intermediate code generation

$t_1 = \text{int to float}(60)$

$t_2 = \text{id}_3 * t_1$

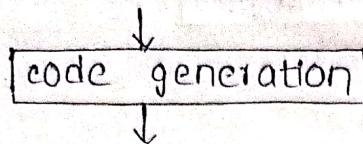
$t_3 = \text{id}_2 + t_2$

$\text{id}_1 = t_3$



$$t_1 = id_3 * 60$$

$$id_1 = id_2 + t_1$$



MOV id<sub>3</sub>, R<sub>1</sub>

MUL #60, R<sub>1</sub>

MOV id<sub>2</sub>, R<sub>2</sub>

ADD R<sub>2</sub>, R<sub>1</sub>

MOV R<sub>1</sub>, id<sub>1</sub>

Ex Describe the O.P. to the various phases of compiler with respect to following statements:

position = initial + rate \* 60

Ex consider the following format of 'C' code:

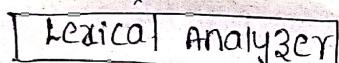
float i, T;

i = i + 70 + T + 2;  
write the O.P. at all phases of the compiler for the above 'C' code

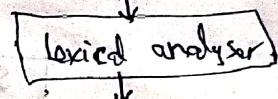
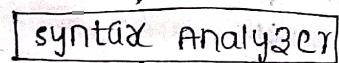
Ans i = i + 70 + T + 2.

Ex-2: a = b \* c \* d.

$$a = b * c * d$$

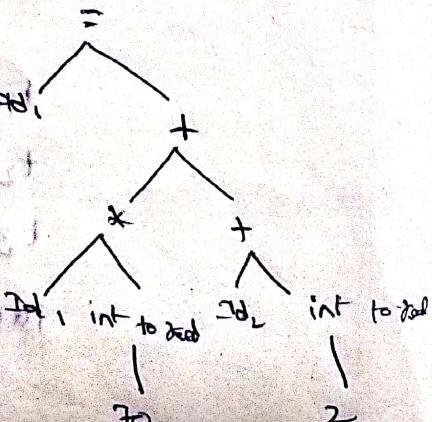
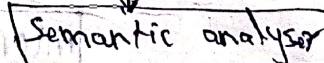
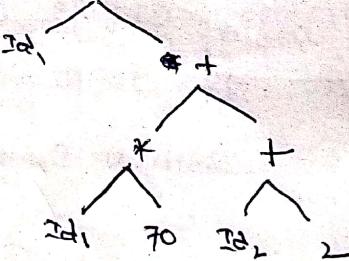
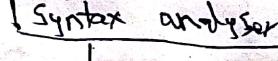


$$id_1 = id_2 * id_3 * id_4$$

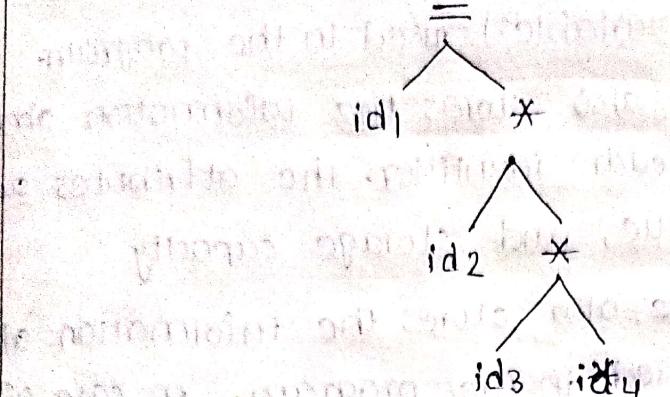


id<sub>1</sub> = id<sub>1</sub> \* id<sub>2</sub> + id<sub>3</sub> + id<sub>4</sub>

↓ token stream



## semantic analysis



## Intermediate code generation

$$E_1 = \text{id}_3 * \text{id}_4$$

$$t_1 = \text{id}_2 * t_1$$

$$\text{id}_1 = t_2$$

## code optimization

$$t_1 = \text{id}_4 * \text{id}_3$$

$$t_2 = \text{id}_2 * t_1$$

## code generation

## Intermediate code generation

$$t_1 = \text{int to float} (70)$$

$$t_2 = \text{id}_1 * t_1$$

$$t_3 = \text{int to float} (2)$$

$$t_4 = \text{id}_2 + t_3$$

$$t_5 = t_2 + t_4$$

$$\text{id}_1 = t_5$$

## code optimizer

$$t_{12} \text{id}_1 * 70.0$$

$$t_{12} = \text{id}_2 + 2.0$$

$$t_3 = t_1 + t_2$$

$$\text{id}_1 = t_3$$

## code generator

mov  $\text{id}_1, R_1$

mul 70.0, R\_1

mov  $\text{id}_2, R_2$

ADD 2.0, R\_2

ADD R\_2, R\_1

mov R\_1,  $\text{id}_1$

Mov  $\text{id}_3, R_2$

MUL  $\text{id}_4, R_2$

Mov  $\text{id}_2, R_1$

MUL  $R_2, R_1$

Mov  $R_1, \text{id}_1$

MOV R\_1,  $\text{id}_1$



## Symbol table Management

To support the phases of the compiler a symbol table is maintained. The task of the symbol table is to store identifiers (variables) used in the program.

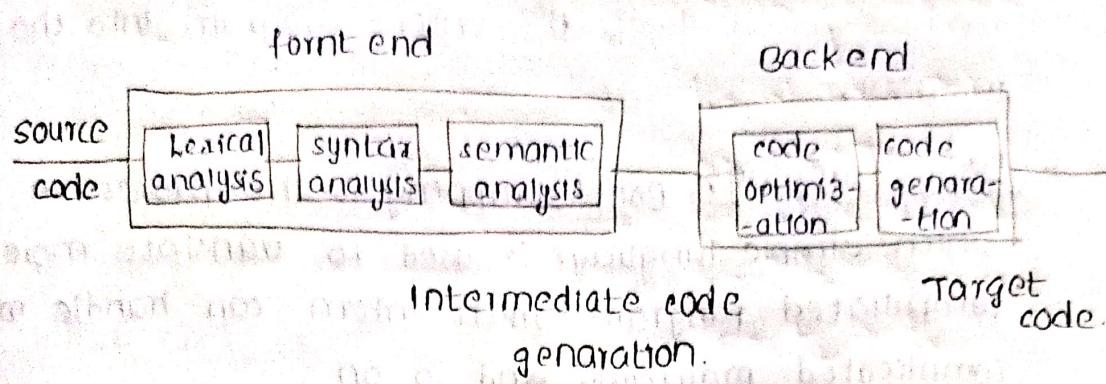
- \* The symbol table also stores the information about the attributes of each identifier. The attributes such as time, scope, value, and storage capacity.
- \* The symbol table also stores the information about the subroutines used in the program. In case of subroutine the symbol table stores the name of the subroutine, number of arguments, type of the arguments, method of passing (call by value, call by reference and return type).
- \* During compilation the lexical analyser detects the identifiers and makes an entry into the symbol table.
- \* Various phases can use the symbol tables in various ways while doing the semantic analysis and intermediate code. We need to know what type of identifiers.

## Error detection and handling :

In compilation each phase detects errors. These errors must be reported to error handler which is used to handle the errors so that the compilation can be produced.

A large number of errors can be detected in syntax analysis such errors are called syntax errors. During semantic analysis type mismatch kind of error is usually detected.

## Grouping of phases :-



Different phases of the compiler can be grouped together to form the front end and backend. The front end consists of phases primarily dependent on the source language and independent on the target language. The front end consists of the analysis part.

The back end consists of phases dependent on the target language and it is independent on the source language.

By keeping same front end and different backends here, produce a compiler for the same source language on different machines.

By keeping different front ends and same Backend -s one can compile different language on the same machines.

### PASS :-

one complete scan of the source language is called pass. It includes reading an input file and writing an output file.

Many phases grouped by using one pass.

It is difficult to compile source code into program.

to separate information into sequential and grouped sets of hardware, often input and output

if we group several phases into one pass we may be forced to keep the entire program into the memory.

\* **Boot strapping:** Boot strapping is the process in which simple language is used to translate more complicated program which intern can handle more complicated programs and so on.

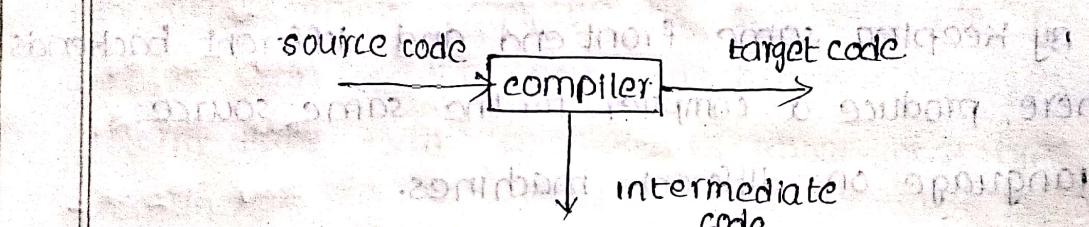
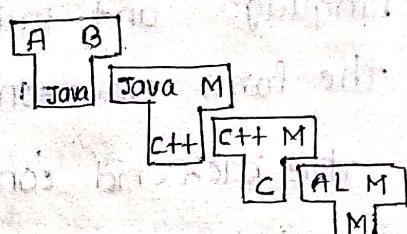
\* writing a highlevel language is a complicated process it takes a lot of time to write a compiler.

\* Boot strapping is represented as b' diagram it consists of '3' languages

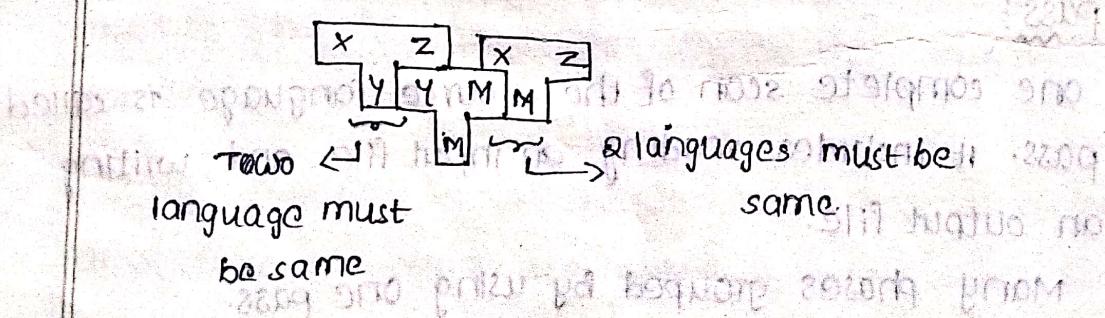
1) source language.

2) Target language.

3) Implementation language.



\* To understand the Boot strapping technique the method is called cross compiler.



suppose we want to write a cross compiler for a new language 'X' the implementation language of this is 'Y' and the target code generated in the language compiler.

is 'z'

Fix. existing compiler 'Y' runs on machine 'M' and generates code for machine 'm'.

If we run  $x, Y, z$  i.e.  $x^z$  using  $Y_M^M$  then we get a compiled  $x_M^z$  that means a compiler a source language 'x' generates the target code in 'z' and which runs on the machine 'M'.

Lexical analysis: It is the first phase of the compiler.

The process of compilation starts first phase called LA

In this phase the input is scanned completely to identify the tokens. So the lexical analysis is also called as scanning. The token structure can be represented with the help of some diagrams. These diagrams are probably known as finite automata.

\* To construct finite automata regular expressions are used. These diagrams translated into programs to identify the tokens

\* The lexical analysis separates the input program into various types of tokens types like, keywords, identifiers and operators and so on.

Ex: int main()

{

int i;

for(i=0; i<10; i++)

{

printf ("Hello world\n");

}

}

Lexemes	Tokens
int	keyword
main	identifier
{	operator
)	closing bracket
{	symbol
int	keyword
i	identifier
for	keyword
printf	identifier
}	symbol
;	symbol

## \* Functions of Lexical analysis

- 1) It is
- 2) It generates symbol table which is a type of data structure to store symbol related information.
- 3) It can count line numbers of the source code.
- 4) It eliminates black spaces, white space and comments.
- 5) It reports the errors encountered while generating the tokens.
- 6) It produces the stream of tokens.

## \* Tokens Lexemes and patterns:

Lexem is a group of characters, token is a class or categorical like keyword, identifiers, operators and some constants.

Pattern is a set of rules to describe the token which is like a A to z, A to Z, 0 to 9 digits.

## \* Lexical errors:

These type of errors can be detected during lexical analysis. Some of the errors are:

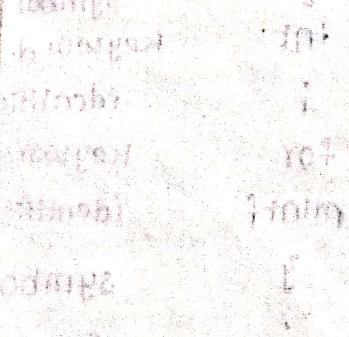
- 1) Exceeding the length of the identifier. For example in the length of the identifier is limited to '10' characters.
- 2) Appearance of illegal characters.

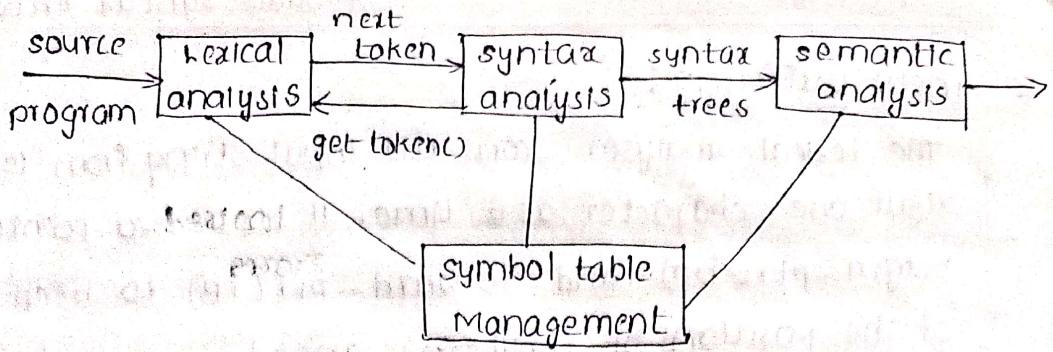
e.g. printf("Hello world"); \$

- 3) unmatched string

e.g. This is a comment \*/

## \* Role of lexical analysis:





The process of compilation uses lexical analysis has the first phase of the compiler. The Lexical analysis brought the request from the syntax analysis as get token(), then syntax analysis gives response to the lexical analysis as next token().

Lexical analysis produce the tokens and syntax analysis consume the tokens so it is called as producer-consumer pair. These two phases also related to the symbol table. in which the maintains about the identifiers to store and recreates.

some times the errors may be generated in the lexical analysis which are called as lexical errors and one handled by the error handlers.

\* Differences between Lexical analysis and syntax analysis

Lexical analysis	syntax analysis
<ol style="list-style-type: none"> <li>1) It is the first phase of the compiler.</li> <li>2) Lexical analysis takes the input as source program and produces output as tokens.</li> <li>3) The errors produced as lexical errors.</li> <li>4) It uses finite automata and regular expression to</li> </ol>	<ol style="list-style-type: none"> <li>1) It is the second phase of the compiler.</li> <li>2) Syntax analysis takes the input as tokens and produce output as syntax trees.</li> <li>3) The errors are produced as syntactical errors.</li> <li>4) It uses the grammars and hierarchical structures to</li> </ol>

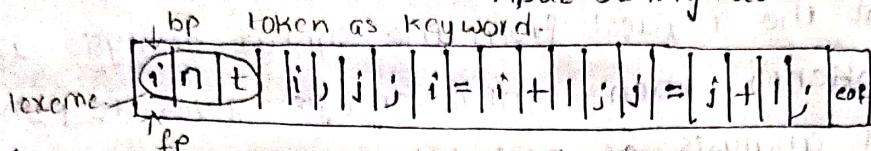
produce tokens

produce syntax trees.

### \* Input Buffering :-

The lexical analyser scans the input string from left to right one character at a time. It uses two pointers begin - ptr (bp) and forward - ptr (fp) to keep track of the position of input i.e. scanned.

Initially both the pointers are pointing to the one(s) first character of the input string as



The forward pointer moves a head to search for the end of the lexeme.

\* As soon as the black space is encountered it indicates the end of the lexeme and lexeme int is identified.

\* The FP will be moved a head at white space when FP encounters a white space. It ignores and move a head.

\* Then both beginning and forward pointer are set at next token i.e. tokens are ready to read.

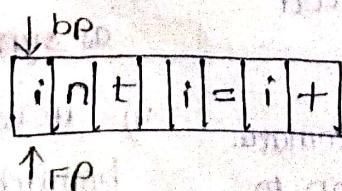
\* The input character is read from the memory is very costly. Hence we use two buffering techniques

1) one buffer scheme.

2) Two buffer scheme.

1) one buffer scheme :-

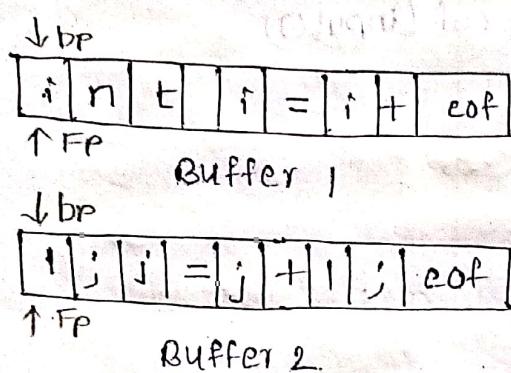
In these one buffer scheme only one buffer is used to store the input string. If lexeme is very long, it crosses the buffer boundary to scan the rest of the lexeme. The buffer has to be refilled it that makes the overwriting the first part of the lexeme.



## 2) Two Buffer scheme:-

To overcome the problem of the one Buffer scheme, two Buffers are used to store the input string. the first buffer and second buffer scan alternatively. when end of the current buffer is reached the other is buffer is filled.

- \* Initially both bp and fp are pointing to the first character of the first buffer.



- \* Then fp moves towards right in search of end of the lexeme, As soon as blank character is recognized. the string between. beginning pointer and fp is recognized.
- \* To identify the boundary of the first buffer end of the buffer (eof) should be placed at the end of the first buffer. Similarly end of the second buffer is also recognized with (eof) at the end of the second buffer.
- \* when FP encounters first eof it recognized end of the first buffer and filling of the second buffer is started.
- \* Alternatively both the buffers can be filled up until end of the input program and stream of tokens is identified.
- \* This eof introduced at the end is called "sentinel"

Code for input Buffering:



```

if FP == eof(buffer)
{
    FP++;
}
else if FP == eof(buffer2)
{
    FP++;
}
else if FP == eof(input())
{
    return;
}
else
{
    FP++;
}

```

\* Recognition of tokens:

For a programming language there are various types of tokens such as identifiers, keywords, constants.

\* The token is usually represented by a pair, token type and token value.

→ Token representation.

Token Type	Token value
------------	-------------

\* The token type tells us the category of token, and the token value gives us the information about the token and it is also called as token attribute.

\* During lexical analysis process the symbol table is maintained. The token value can be a pointer to the symbol table. Because of :-

- In case of identifiers and constants.

\* The lexical analyser reads the input program and generates the symbol table for tokens.

we will consider the encoding of the program.

Token	code	value
if	1	-
else	2	-
while	3	-
for	4	-
identifier	5	ptr to symbol table
constant	6	ptr to symbol table
<	7	1
<=	7	2
>	7	3
>=	7	4
!=	8	5
(	8	6
)	8	7
+	9	1
-	9	2
=	10	1

The corresponding symbol table for identifiers

location	Type	value.
counter		
100	identifier	a
:		
105	constant	10
:		
107	identifier	b
:		
110	constant	2

## program.

\* In the above example scanner scans the input string and recognise the 'if' as keyword and returns token type 'P'. since 'i' indicates keyword 'if'. and Hence 'i' is at the beginning of the token string

\* Next pair (8,1) where '8' indicates parenthesis and '1' indicates opening parenthesis then we can scan the input 'a'. it recognises it has identifier and searches the symbol table to check whether the same entries is present

\* if not it inserts the information about this identifier in symbol table and returns 100.

\* if the same identifier or the variable is already present in the symbol table then lexical analyzer does not insert it into the table instead it returns the location where it is present.

#### \* Specification of tokens :-

For specifying the tokens we use the following :-

1) Alphabets :- a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z

2) strings

3) languages

4) operations the languages

5) Regular expressions

6) Regular grammar

Description :-

1) Alphabet is a. symbol or a set of symbols and it is represented as  $\Sigma$  i.e  $\Sigma = \{0, 1, \dots\}$  (or)  $\Sigma = \{a, b, c, \dots\}$

2) string is defined as a set of alphabets

$$\Sigma = \{a, aa, ab, baa\}$$

i) the types of the strings are empty string, prefix string, string length  $- |w|$ , suffix string, substring, string concatenation.

$$(aa)(a)(aa)(aa)$$

$$|S-i|=i$$

- ii) length of the string means to count number of characters  $w = 000$   
 $|w| = 3$
- iii) Empty string means string does not have any characters  
 $s = \emptyset$
- iv) prefix string means starting characters of the string.
- v) suffix string means last characters of the string.
- vi) substring means removes the both suffix and prefix of the string. which is given.
- vii) string concatenation means concataining two strings

- 3) Language is a finite number of strings.
- Operations on the languages [union, concatenation] i.e u, -
- K en closure (\*) which represents 0 or more
- positive closure (+) which represents 1 or more.
- 5) Regular expressions are used to defined as tokens and also defines the grammar i.e regular which is represented as  $G_1 = (V, T, P, S)$

V - variables

T - terminals

P - productions

S - starting symbol

### \* Finite automata :-

- \* From the regular expressions it is used to (or) possible to build a finite automata.
- \* the finite automata is usually a non-deterministic finite automata and is represented by the 'F' on M. By using the transition diagram.
- The represents as follows ;-

$$\text{FOR } M = \{ Q, \Sigma, \delta, q_0, F \}$$

Q - finite number of states  $\{ q_0, q_1, q_2, \dots \}$

$\Sigma$  - Input alphabet  $\{ 0, 1, \dots \}$

$\delta$  - Transition function. -  $\{\delta = (\theta, \epsilon) \rightarrow Q\}$

$q_0$  - initial state.

$F$  - final state.

\* construct an NFA with  $\epsilon$ -moves for the regular expression  $b + ba^*$

Given, regular expression (R.E) =  $b + ba^*$

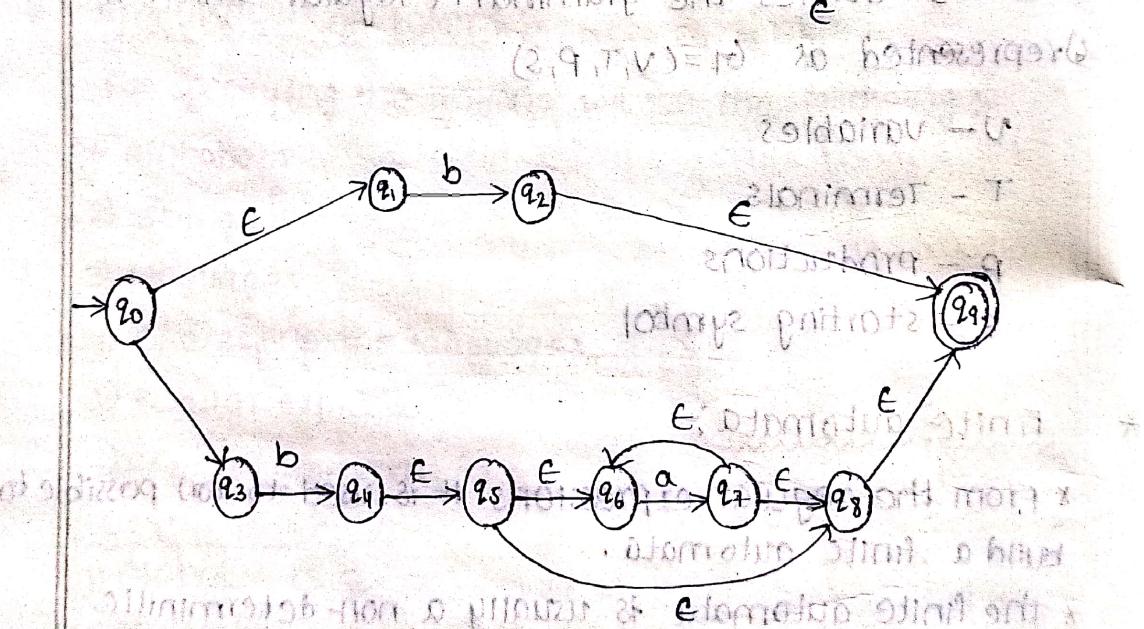
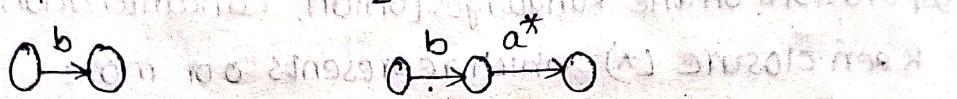
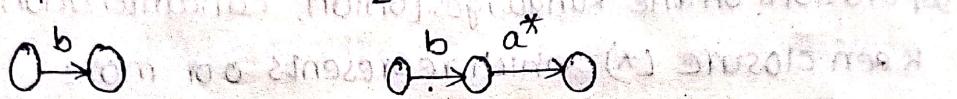
Let,  $r_1 = b$

$r_2 = ba^*$

$r_3 = r_1 + r_2 \equiv b + ba^*$

$r_1 = b$

$r_2 = ba^*$



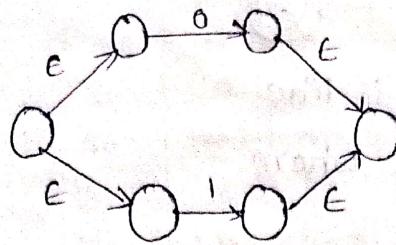
\* construct an NFA with  $\epsilon$ -moves for the R.E is  $(0+1)^*$

Given, R.E =  $(0+1)^*$

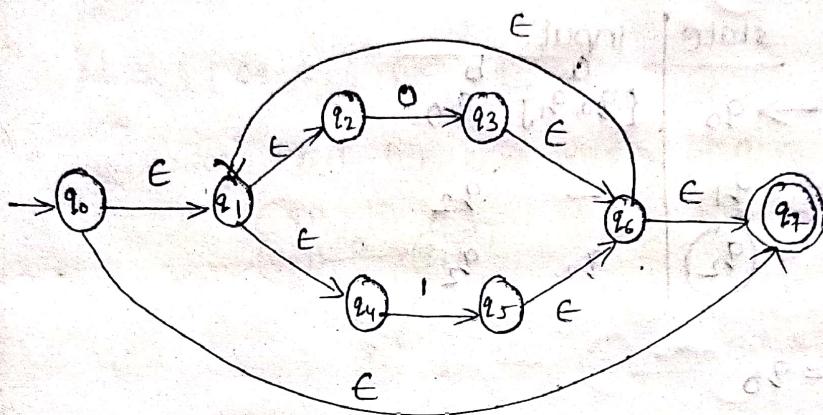
Let,  $r_1 = 0+1$

$r_2 = (0+1)^*$

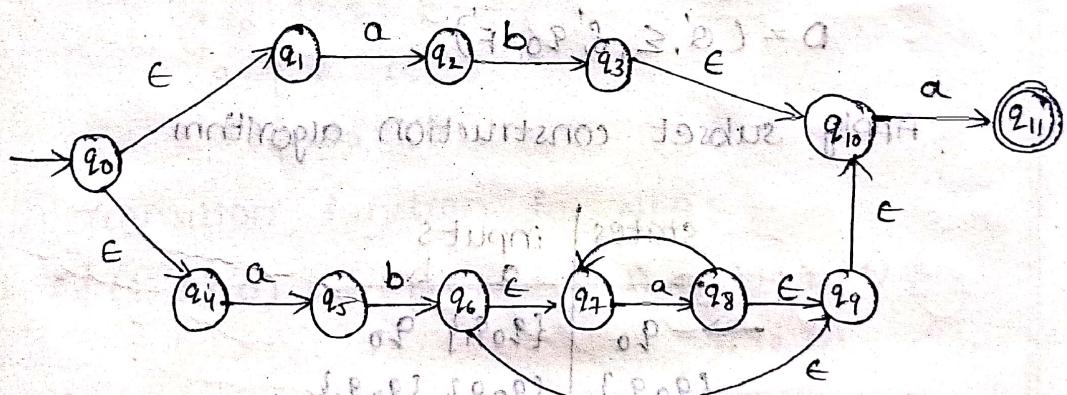
$$r_1 = 0+1$$



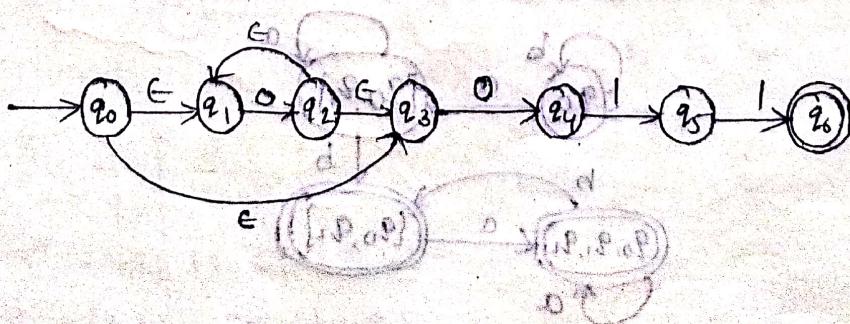
$$r_2 = (0+1)^*$$



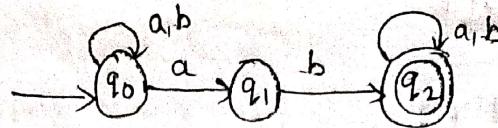
\* construct an NFA with  $\epsilon$  for the R-E is  $(abUaaba^*)^*$



\* construct  $((01)^*01)^*$  for these regular expression draw the NFA with  $\epsilon$ .



\* Construct DFA from the given NFA



Sol: The given NFA  $N'$  is like

$$N = \{Q, \Sigma, \delta, q_0, F\} \text{ where}$$

$$Q = \{q_0, q_1, q_2\}$$

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

$$\delta =$$

state	input	
	a	b
$\rightarrow q_0$	$\{q_0, q_1\}$	$q_0$
$q_1$	-	$q_2$
$q_2$	$q_2$	$q_2$

$$q_0 = q_0$$

$$F = \{q_2\}$$

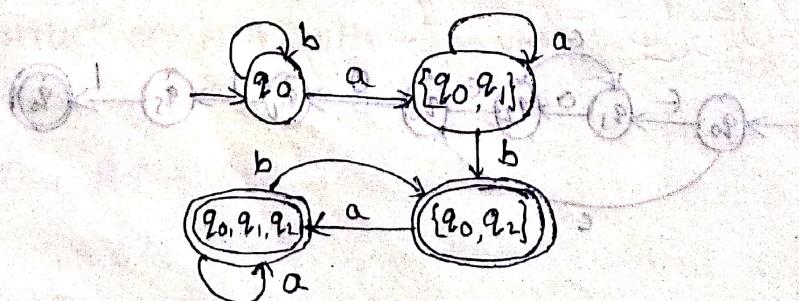
let us consider the DFA 'D' is like

$$D = (Q', \Sigma, \delta', q_0, F')$$

Apply subset construction algorithm

states	inputs
	a      b
$\rightarrow q_0$	$\{q_0\}$ $q_0$
$\{q_0, q_1\}$	$\{q_0, q_1\}$ $\{q_0, q_2\}$
$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$ $\{q_0, q_2\}$
$\{q_0, q_1, q_2\}$	$\{q_0, q_1, q_2\}$ $\{q_0, q_2\}$

$\therefore$  The DFA is



$$Q' = \{q_0, \{q_0, q_1\}, \{q_0, q_2\}, \{q_0, q_1, q_2\}\}$$

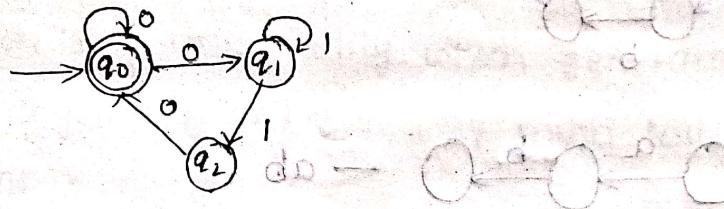
$$\Sigma = \{a, b\}$$

$\delta'$ = states	input	
	a	b
$q_0$	$\{q_0, q_1\}$	$q_0$
$\{q_0, q_1\}$	$\{q_0, q_1\}$	$\{q_0, q_2\}$
$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$
$\{q_0, q_1, q_2\}$	$\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$

$$q_0 = q_0$$

$$f' = \{\{q_0, q_2\}, \{q_0, q_1, q_2\}\}$$

\* construct DFA for the given NFA.



transition table for NFA:-

	0	1
$q_0$	$\{q_0, q_1\}$	$\emptyset$
$q_1$	$\emptyset$	$\{q_1, q_2\}$
$q_2$	$q_0$	$\emptyset$

transition function for NFA:-

$$\delta(\{q_0, q_1\}, 0) = \delta(q_0, 0) \cup \delta(q_1, 0) = \{q_0, q_1\} \cup \emptyset = \{q_0, q_1\}$$

$$\delta(\{q_0, q_1\}, 1) = \delta(q_0, 1) \cup \delta(q_1, 1) = \emptyset \cup \{q_1, q_2\} = \{q_1, q_2\}$$

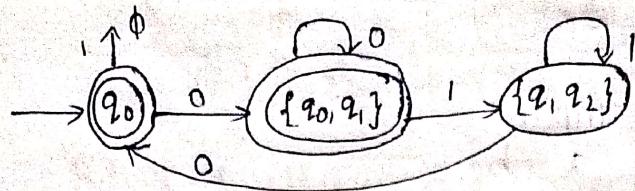
$$\delta(\{q_1, q_2\}, 0) = \delta(q_1, 0) \cup \delta(q_2, 0) = \emptyset \cup q_0 = \{q_0\}$$

$$\delta(\{q_1, q_2\}, 1) = \delta(q_1, 1) \cup \delta(q_2, 1) = \delta\{q_1, q_2\} \cup \emptyset = \{q_1, q_2\}$$

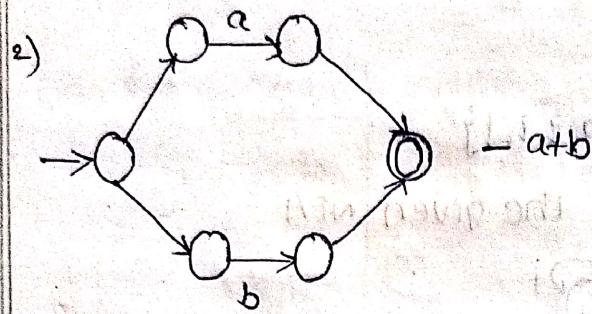
transition table for DFA:-

	0	1
$q_0$	$\{q_0, q_1\}$	$\emptyset$
$\{q_0, q_1\}$	$\{q_0, q_1\}$	$\{q_1, q_2\}$
$\{q_1, q_2\}$	$q_0$	$\{q_1, q_2\}$

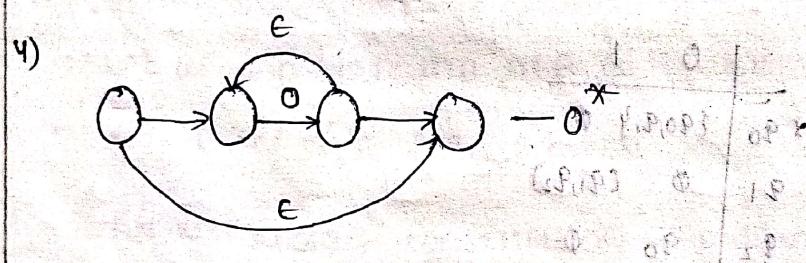
# Transition Diagram for DFA:



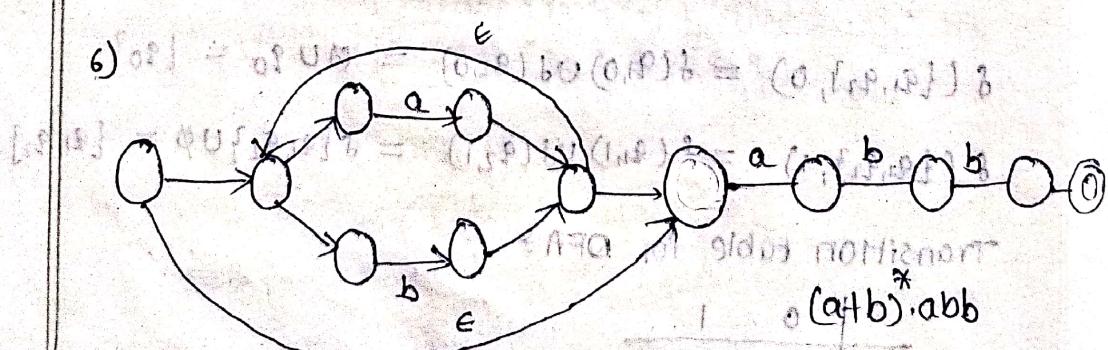
1)  $\epsilon \in Q_0 \cdot Q_0$



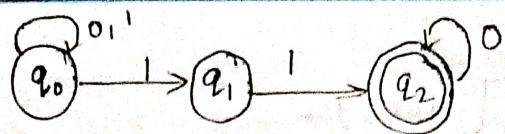
3)  $a \cdot b = ab$



5)  $\epsilon \cdot ab = ab$



$$\begin{array}{c|cc|cc} & 1 & 0 & * & \\ \hline \epsilon & (ab) & ab & ab & \\ (ab) & ab & ab & ab & \\ (ab)^* & ab & ab & ab & \\ (ab)^* \cdot ab & ab & ab & ab & \end{array}$$



so if grammar =  $(V, T, P, S)$

$$q_0 \rightarrow 1q_1 \quad q_1 \rightarrow 1q_2 \quad q_2 \rightarrow 0q_2$$

$$q_0 \rightarrow 0q_0 \quad q_1 \rightarrow 1 \quad q_2 \rightarrow 0$$

$$q_0 \rightarrow 1q_0$$

so there are 'r' grammar.

- \* lexical analyzer generator :- LEX  
various tools have been built for lexical analyzers <sup>constructing</sup> using special purpose notations called regular expressions which are used in recognizing the tokens. A tool called LEX which is a unix utility which generates the lexical analyzer. A LEX is a very much faster in finding the tokens.

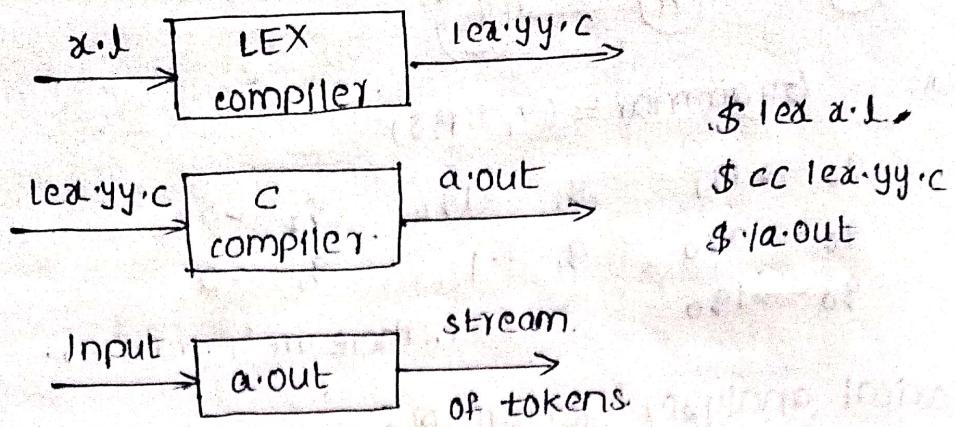
- \* the LEX specification file can be created using the extension (.L).

For example, the LEX specification file can be (x.L) which is given to the lex compiler which gives `lex.c` which is a C program and is also a lexical analyzer program.

- \* the LEX specification file shows the regular expressions for the tokens and the `lex.yyy.c` file consists of tabular representation of the transition diagram constructed for a regular expressions.

- \* In regular expression lexical LEX specification file actions of the pieces of C code which is directly carried over to the `lex.yyy.c`. Finally the compiler compiles the `lex.yyy.c` and produces an object program `l.out`.

- \* when some input is given to the `l.out` then some sequence of tokens generated.



## \* Design of lexical analyzer LEX :-

The LEX specification file consists of three different sections.

- 1) Declaration section.
- 2) Transition rule section.
- 3) Auxiliary procedure section.

### 1) Declaration section:-

In this section declaration of variables constants can be done.

Some regular expression can also be done in the form of regular expression.

• 1. { int a, b; }

float count=0;

### 2) Transition rule section:-

The rule section consists of regular expressions with associated actions.

These translation can be given as

• 1. 1.

Rule1      Action1

Rule2      Action2

Rule3      Action3

• 1. 1.

### 3) Auxiliary procedure section :

The third section is the auxiliary procedure section in which all require procedures are defined.

\* [ . { . . . } ]

. . . . .

"RAM"

"RAV1" | printf ("In noun");

"sing" |

"dance" | printf ("In verb");

. . . . .

main()

{

yylex();

y

int yywrap()

{

return;

y

In the above program, the third section consists of two functions

i) the main function and the yywrap function.

ii) the mainfunction called the 'yy' which is the function

that defined by the yylex i.e LEX.yy.c

\* First we compile the above using LEX compiler and then the LEX compiler will generate a ready 'c' program named LEX.yy.c and yy.c

\* This LEX.yy.c makes use of regular expressions and corresponding actions defining the [lex.i]. To run the LEX program LEX.i i.e LEX.yy.c

[ localhost@cse539] \$ LEX a.i

[localhost@cse539] \$ lex.yy.c



[localhost@CSE539]\$ a.out.

- \* Syntax analysis: The second phase of compiler is syntax analysis. It is also called as parsing (or) parser. It checks the syntax of the language.
- \* The syntax analyzer takes the tokens from the lexical analyzer and groups them in such a way that some programming structure can be recognized.
- \* After grouping the tokens syntax can't be recognized, then syntax error will be generated.

parser:

A parser is a process which takes the input string and produces a parse tree (or) generates the syntax errors.

\* There are two important issues in parsing

1) specification of the syntax:

- It means how to write any programming statement.

\* specifications should be precise and ambiguous should be in detail, should cover all the details of the programming languages, should complete such specification is called free grammar.

2) Representation of the input:

\* After passing, all the subsequent phases of the compiler takes the information from the parse tree being generated.

\* There are two types of algorithms based on the parsing

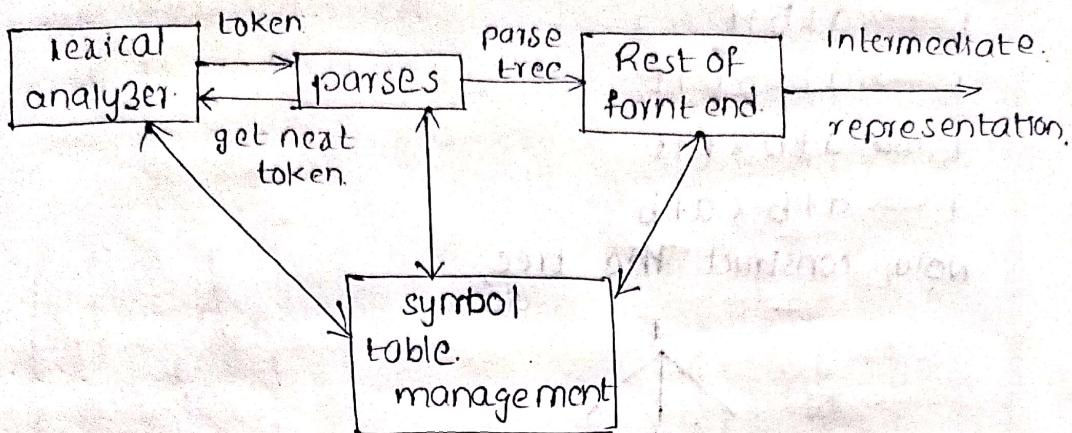
1) Top-down parsing or bottom-up parsing

2) Bottom-up parsing or top-down parsing

3) Role of parser:

In the process of compilation, the parser and lexical analyzer work together, when parser requires string

of tokens it involves lexical analyzer. the lexical analyzer supplies tokens to the syntax analyzer.



A grammar  $G_1 = (V, T, P, S)$  is CFG where the production rules can be formed.  $V \rightarrow (VUT)^*$

Derivational parse trees:-

Derivation of the variables means generating strings. They are two types of derivations.

1) Left most derivation,

2) Right most derivation,

1) Left most derivation:-

Replacing the left most side of the variable in the left hand side of tree is called left most tree.

2) Right most derivation:-

Replacing the right most side of the variable in the right hand side of tree is called right most tree.

\* Examples :- obtain the left most and right most tree.

derivation of a string  $a+b * a+b$

$E \rightarrow E+E | E \times E | E-E | E/a/b$  is the string.

Given string and grammar is

$E \rightarrow E+E | E-E | E \times E | E | a | b$  and,

the string is  $a+b * a+b$ .

so/it

$a+b * a+b$

$E \rightarrow E+E | E-E | E\times E | E/a/b$

$E \rightarrow E+E$

$E \rightarrow E+E+E$

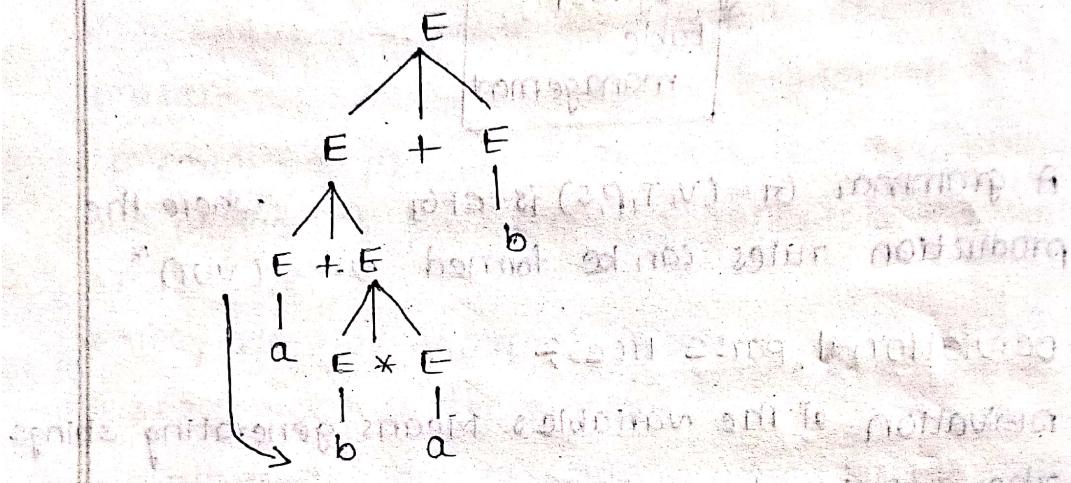
$E \rightarrow a+b+E$

$E \rightarrow a+E\times E+E$

$E \rightarrow a+b\times E+E$

$E \rightarrow a+b\times a+b$

Now, construct the tree,



- 2)  $E \rightarrow E+E | E\times E | E/a/b$  and the string id+id\xid  
both left and right most derivation.

- 3)  $S \rightarrow (L)/a$

$L \rightarrow L, S / S$  obtain left most and right most  
derivation. i) (a, a)

ii) (a, (a, a))

iii) (a, ((a, a), (a, a))) soft parenthesis

- 4) Consider the grammar:

$S \rightarrow aS | aSbS | \epsilon$  obtain the string aaabaab

- 5)  $S \rightarrow aAB | bB A/E$

$A \rightarrow aAb/E$  obtain the string baabb

$B \rightarrow bB | E$

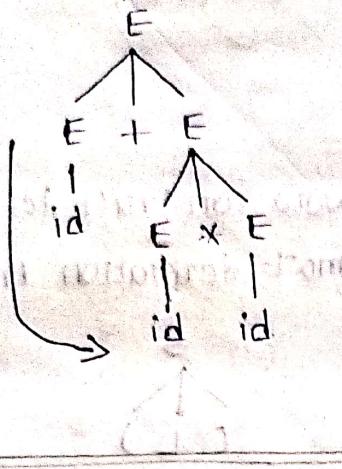
Given,

$$E \rightarrow E+E$$

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

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

Now, construct the left most derivation tree.



$$E \rightarrow E+E$$

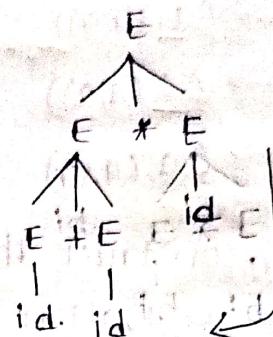
$$E \rightarrow E+E * E$$

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

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

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

Now, construct the right most derivation tree.



3

i) Given,

$$S \rightarrow (L)$$

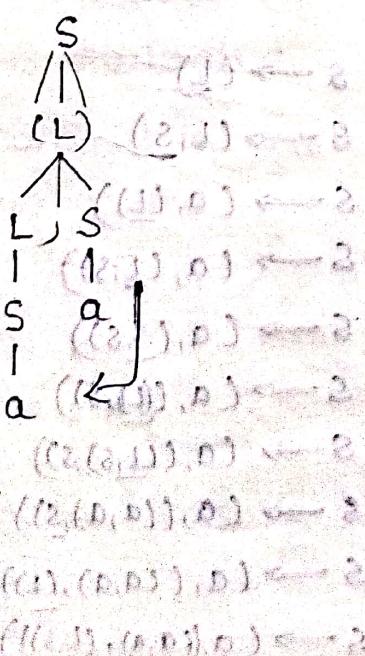
$$S \rightarrow (\underline{L}, S)$$

$$S \rightarrow (\underline{L}, a)$$

$$S \rightarrow (\underline{s}, a)$$

$$S \rightarrow (q, a)$$

Now, construct the right most derivation tree.



$$S \rightarrow (L)$$

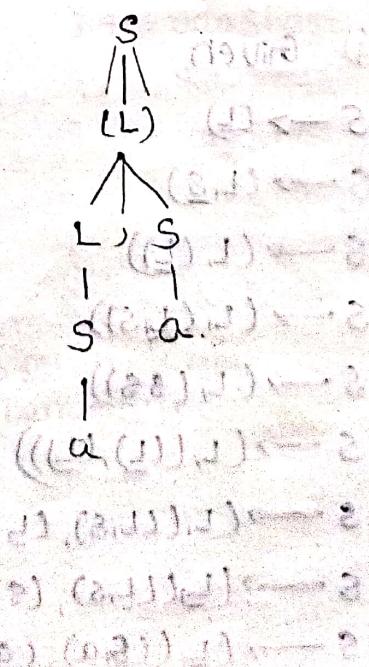
$$S \rightarrow (\underline{L}, S)$$

$$S \rightarrow (\underline{s}, S)$$

$$S \rightarrow (q, \underline{s})$$

$$S \rightarrow (q, a)$$

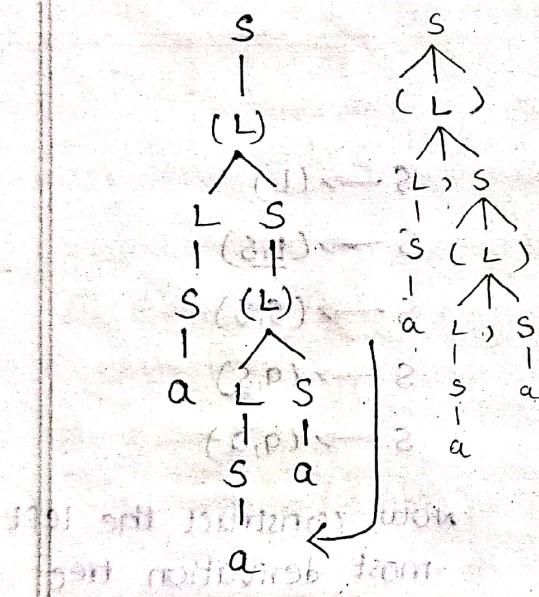
Now, construct the left most derivation tree.



ii) Given,

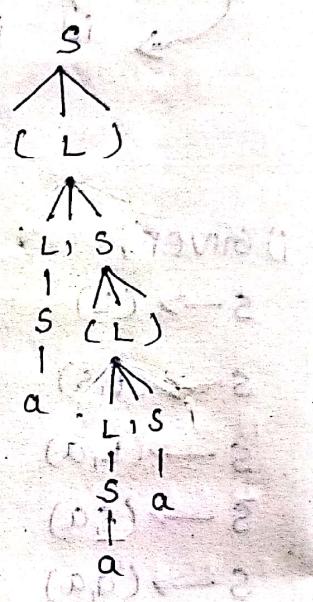
$$\begin{aligned}S &\rightarrow (\underline{L}) \\S &\rightarrow (L, \underline{S}) \\S &\rightarrow (L, (\underline{L})) \\S &\rightarrow (L, (L, \underline{S})) \\S &\rightarrow (L, (L, S)) \\S &\rightarrow (L, (L, a)) \\S &\rightarrow (\underline{L}, (a, a)) \\S &\rightarrow (\underline{S}, (a, a)) \\S &\rightarrow (\underline{a}, (a, a))\end{aligned}$$

Now, construct the right most derivation tree.



$$\begin{aligned}S &\rightarrow (\underline{L}) \\S &\rightarrow (L, \underline{S}) \\S &\rightarrow (\underline{S}, LL) \\S &\rightarrow (a, (\underline{L}, S)) \\S &\rightarrow (a, (\underline{S}, S)) \\S &\rightarrow (a, (a, a))\end{aligned}$$

Now construct left most derivation tree.



iii) Given,

$$\begin{aligned}S &\rightarrow (\underline{L}) \\S &\rightarrow (L, \underline{S}) \\S &\rightarrow (L, (\underline{L})) \\S &\rightarrow (L, (\underline{L}, S)) \\S &\rightarrow (L, (\underline{S}, S)) \\S &\rightarrow (L, ((\underline{L}), (\underline{L}))) \\S &\rightarrow (L, ((\underline{L}, S), (\underline{L}, S))) \\S &\rightarrow (L, ((\underline{L}, S), (S, a))) \\S &\rightarrow (L, ((S, a), (a, a)))\end{aligned}$$

$$\begin{aligned}S &\rightarrow (\underline{L}) \\S &\rightarrow (L, \underline{S}) \\S &\rightarrow (a, (\underline{L})) \\S &\rightarrow (a, (\underline{L}, S)) \\S &\rightarrow (a, (\underline{S}, S)) \\S &\rightarrow (a, ((\underline{L}), S)) \\S &\rightarrow (a, ((\underline{L}, S), S)) \\S &\rightarrow (a, ((a, a), \underline{S})) \\S &\rightarrow (a, ((a, a), (\underline{L}))) \\S &\rightarrow (a, ((a, a), (\underline{L}, S)))\end{aligned}$$

$$S \rightarrow (S, ((a, a), (a, a)))$$

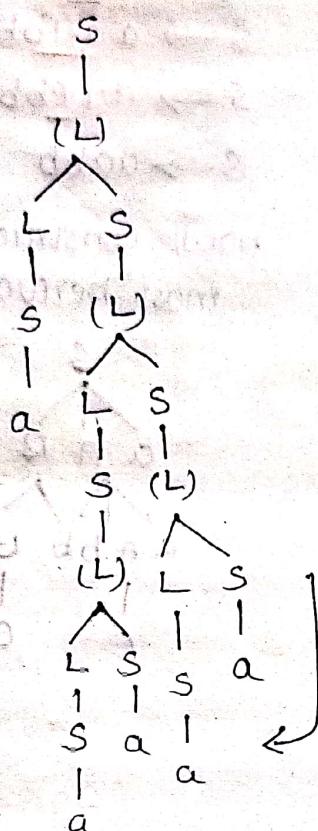
$$S \rightarrow (a, ((a, a), (a, a)))$$

Now, construct the right most derivation tree.

$$S \rightarrow (a, (ca, a), (s, a))$$

$$S \rightarrow (a, ((a, a), (a, a)))$$

left most derivation tree.



5) Given,  $S \rightarrow aAB \mid bBA \mid \epsilon$ ,  $A \rightarrow aAbc$ ,  $b \rightarrow bB/c$ ,  $aabb$

$$S \rightarrow a\cancel{A}B$$

$$S \rightarrow aa\cancel{A}B$$

$$S \rightarrow aa\cancel{c}bB$$

$$S \rightarrow aabb\cancel{B}$$

$$S \rightarrow aabb\cancel{\epsilon}$$

$$S \rightarrow aabb$$

$$S \rightarrow a\cancel{A}B$$

$$S \rightarrow a\cancel{A}b\cancel{B}$$

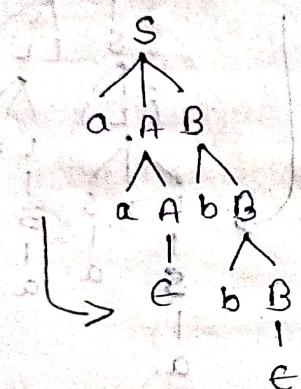
$$S \rightarrow a\cancel{A}b\epsilon$$

$$S \rightarrow aa\cancel{A}bb$$

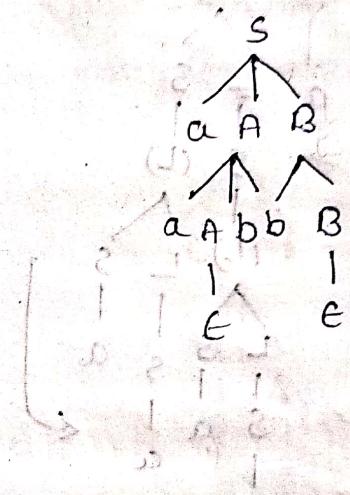
$$S \rightarrow aa\cancel{c}bb$$

$$S \rightarrow aabb$$

NOW, construct left most derivation tree



NOW, construct right most derivation tree



### \* Ambiguous grammar :-

A grammar is said to be ambiguous if it generates more than one parse tree for the language.

Ex:-  $E \rightarrow E+E \mid E * E \mid id$  generate the  $id + id * id$  then find.

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

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

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

$$E \rightarrow E+E$$

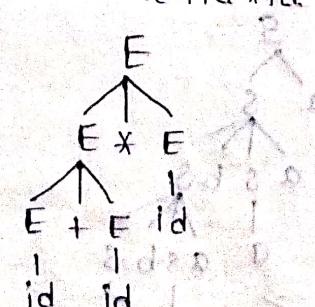
$$E \rightarrow E+E * E$$

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

$$E \rightarrow E+E$$

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

$$\therefore \text{It is ambiguous.}$$



2)  $S \rightarrow asbs$  generate the string abab and prove that ambiguous or not.

$S \rightarrow bsas$

$S \rightarrow \epsilon$

so it  $S \rightarrow \underline{asbs}$

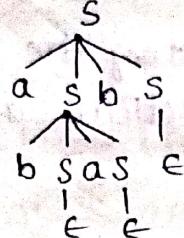
$S \rightarrow \underline{absasbs}$

$S \rightarrow \underline{abeasbs}$

$S \rightarrow ababe$

$S \rightarrow abab$

LMDT



$S \rightarrow \underline{asbs}$

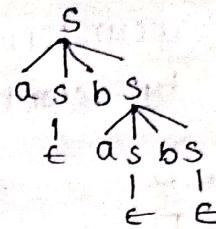
$S \rightarrow \underline{alsbasbs}$

$S \rightarrow \underline{asbaebc}$

$S \rightarrow aebab$

$S \rightarrow abab$

RMDT



3)  $S \rightarrow AB$

$B \rightarrow ab$

$A \rightarrow aa$

$A \rightarrow a$

$B \rightarrow b$

generate the string aab then prove that ambiguous or not.

$S \rightarrow \underline{AB}$

$S \rightarrow \underline{aaB}$

$S \rightarrow aab$

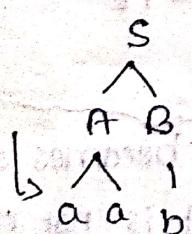
$S \rightarrow \underline{AB}$

$S \rightarrow \underline{Aab}$

$S \rightarrow aab$

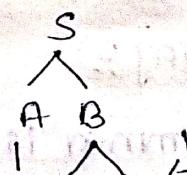
left most derivation

tree



right most derivation

tree



∴ It is a ambiguous.

unit-2