

## UNIT-I INTRODUCTION TO COMPILER AND AUTOMATA.

1. Compilers - Analysis of the source program.
2. Phases of Compiler - cousins of the compiler
3. Grouping of phases - Compiler construction tools.
4. Lexical analysis - Role of Lexical Analyzer
5. Input buffering - Specification of tokens - Design of lexical analysis(LEX)
6. Finite automaton (Deterministic and non-deterministic)  
Conversion of regular expression to NDFA - Thompson's.
7. Conversion of NDFA to DFA - Minimization of NDFA.
8. Derivation - parse tree - ambiguity.

### Compilers:-

⇒ Compiler is a program which takes one language (source program) as input and translates it into an equivalent another language (Target program)



- ⇒ Compiler converts high level language to machine level language.
- ⇒ During this process of translation if some errors are encountered then compiler displays them as error messages.
- ⇒ The compiler takes a source program as high level language such as C, PASCAL, FORTRAN and converts it into low level language or a machine level language such as assembly language.

## Analysis of the source program :-

⇒ The compilation can be done in two parts.

1. Analysis part - the source program is read and broken down into constituent pieces.

- The syntax and the meaning of the source string is determined and then an intermediate code is created from the input source program.

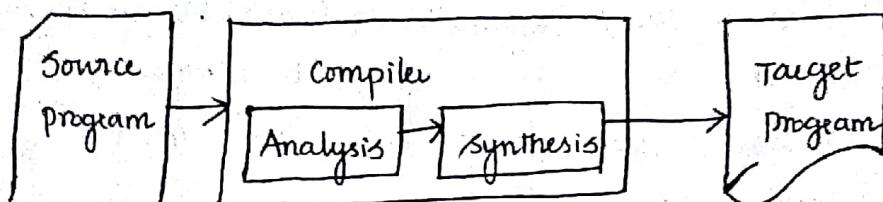


Fig: Analysis and synthesis model.

2. Synthesis part - Intermediate form of the source language is taken and converted into an equivalent target program.

## Execution of program

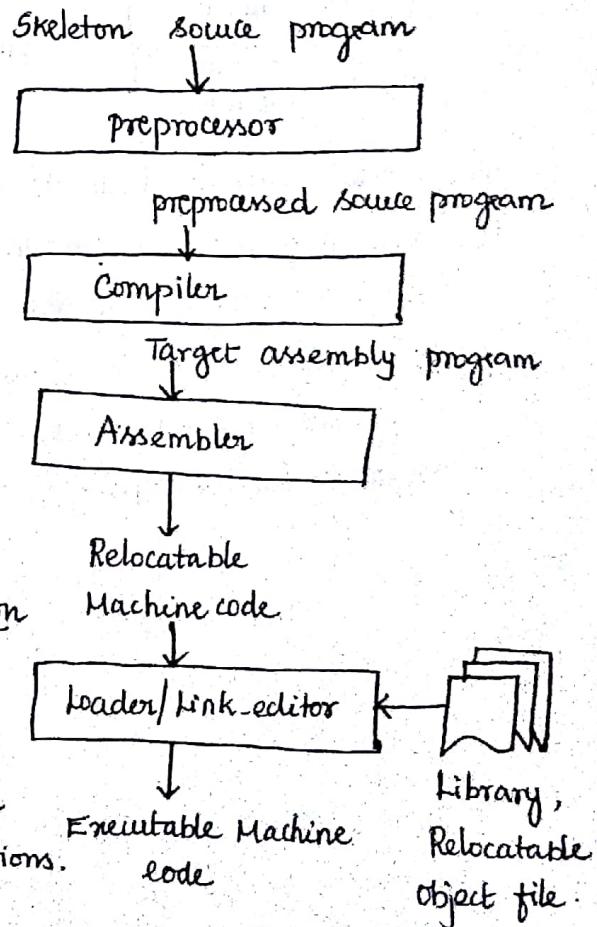
⇒ compiler takes a source program written in high level language as an input and converts it into a target assembly language.

⇒ The assembler then takes this target assembly code as input and produces a relocatable machine code as an input.

⇒ Then a program loader is called for performing the task of loading and link editing.

⇒ Task of loader is to perform the relocation of an object code.

Allocation and placement of load time address and data in memory at proper locations.



- (2)
- ⇒ The link editor links the object modules and prepares a single module from several files of relocatable object modules to resolve the mutual references.
  - ⇒ These files may be library files and these library files may be referred by any program.

### Properties of compiler

1. The compiler itself must be bug-free
2. It must generate correct machine code.
3. The generated machine code must run fast
4. The compiler must be portable
5. It must give good diagnostics and error messages.
6. It must have consistent optimization.

### Phases of compilers :-

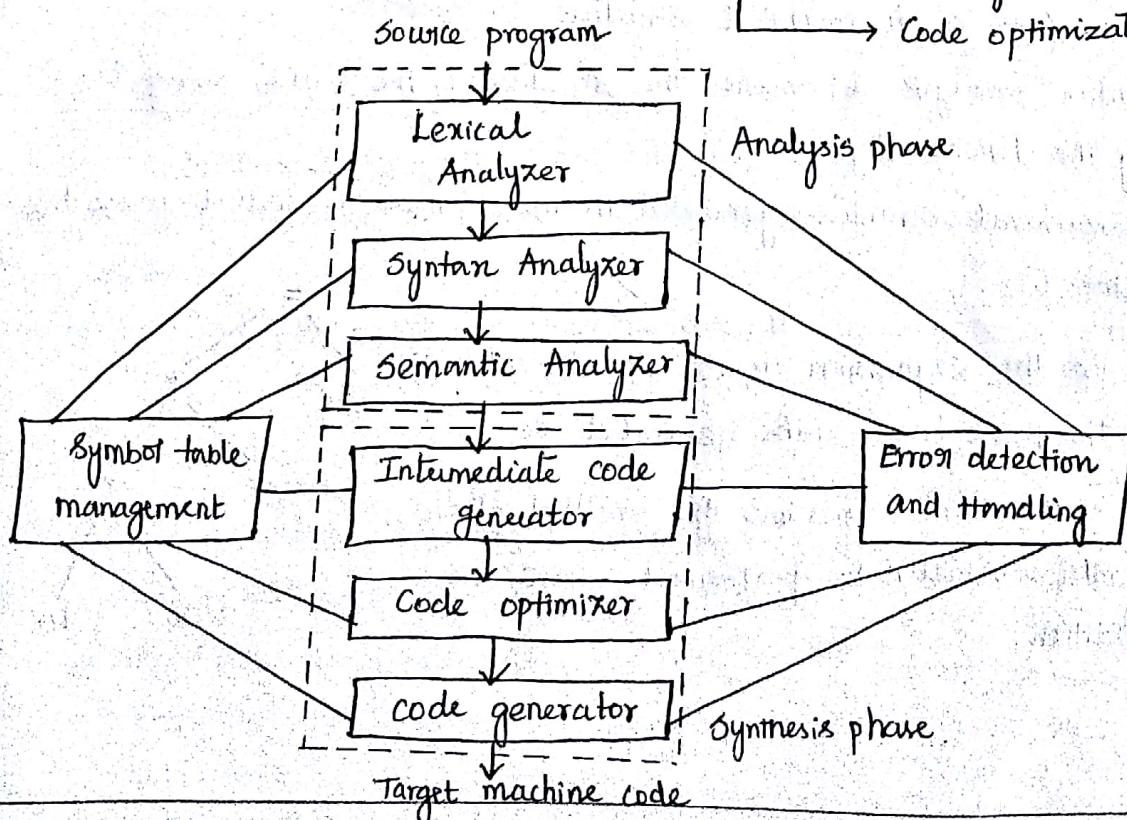
⇒ The process of compilation is carried out in two parts analysis and synthesis

⇒ Analysis is carried out in three phases

→ Lexical analysis  
→ Syntax analysis  
→ Semantic analysis.

⇒ Synthesis is carried out in three phases

→ Intermediate code generation  
→ Code generation  
→ Code optimization.



## Lexical Analysis:

- ⇒ The lexical analysis is also called scanning.
- ⇒ complete source code is scanned and broken up into group of strings called tokens.
- ⇒ A token is a sequence of characters having a collective meaning.

Example: consider the assignment statement

total = count + rate \* 10

Then in lexical analysis phase this statement is broken up into series of tokens as follows,

1. The identifier total
2. The assignment symbol
3. The identifier count
4. The plus sign.
5. The identifier rate
6. The multiplication sign
7. The constant number 10.

⇒ The blank characters which are used in the programming statement are eliminated during lexical analysis phase.

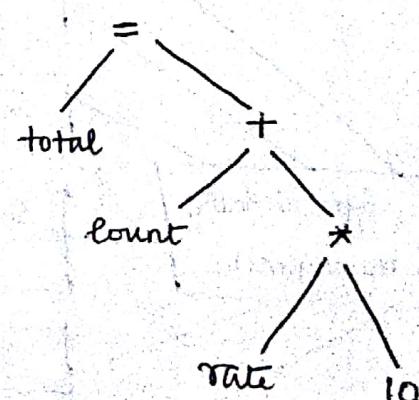
## Syntactic Analysis:

- ⇒ The syntactic analysis is also called parsing.
- ⇒ In this phase the tokens generated by the lexical analyzer are grouped together to form a hierarchical structure.
- ⇒ The syntactic analysis determines the structure of the source string by grouping the tokens together.
- ⇒ The hierarchical structure generated in this phase is called parse tree or syntax tree.

Example: For the expression  $\text{total} = \text{count} + \text{rate} * 10$

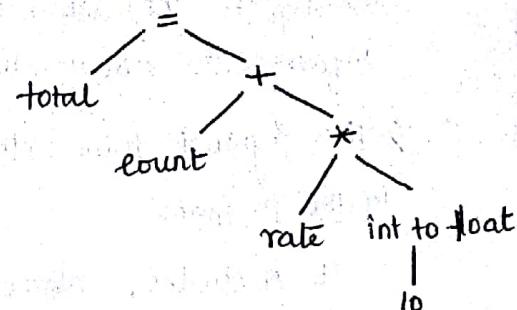
The parse tree can be generated as,

Here, in arithmetic expression the multiplication operation should be performed before the addition.



## Semantic Analysis:-

- ⇒ Semantic analysis determines the meaning of the source string.
- ⇒ For example meaning of source string means matching of parenthesis, in the expression, or matching of if....else statements or performing arithmetic operations of the expressions that are type compatible or checking of scope of operation.
- ⇒ Example: Type compatibility checking of the expression  $\text{total} = \text{count} + \text{rate} * 10$



## Intermediate Code Generation :-

- ⇒ The intermediate code is a kind of code which is easy to generate & it is in variety of forms such as three address code, quadruple, triple, Posix.
- ⇒ The three address code consists of instructions each of which has at most three operands.
- ⇒ Example:
 

```

t1 := int to float(10)
t2 := rate * t1
t3 := count + t2
total := t3
      
```

## Code Optimization :-

- ⇒ The code optimization phase attempts to improve the intermediate code.
- ⇒ This is necessary to have a faster executing code or less consumption of memory.
- ⇒ Thus by optimizing the code the overall running time of the target program can be improved.

## Code Generation :-

- ⇒ Code generation phase the target code gets generated.
- ⇒ The intermediate code instructions are translated into sequence of machine instructions.

### Example:

```

MOV rate, R1
MUL #10.0, R1
MOV Count, R2
ADD R2, R1
MOV R1, total
  
```

## Symbol table Management :-

- ⇒ To support these phases of compiler a symbol table is maintained.
- The task of symbol table is to store identifiers used in the program.
- ⇒ The symbol table also stores information about attributes of each identifier. The attributes of identifiers are usually its type, its scope, information about the storage allocated for it.
- ⇒ The symbol table also stores information about the subroutines used in the program.
  - It includes, Name of the subroutine
  - Number of arguments passed to it
  - Type of these arguments
  - Method of passing these arguments and return type.
- ⇒ Basically symbol table is a data structure used to store the information about identifiers.
- ⇒ The symbol table allows us to find the record for each identifier quickly and to store or retrieve data from that record efficiently.

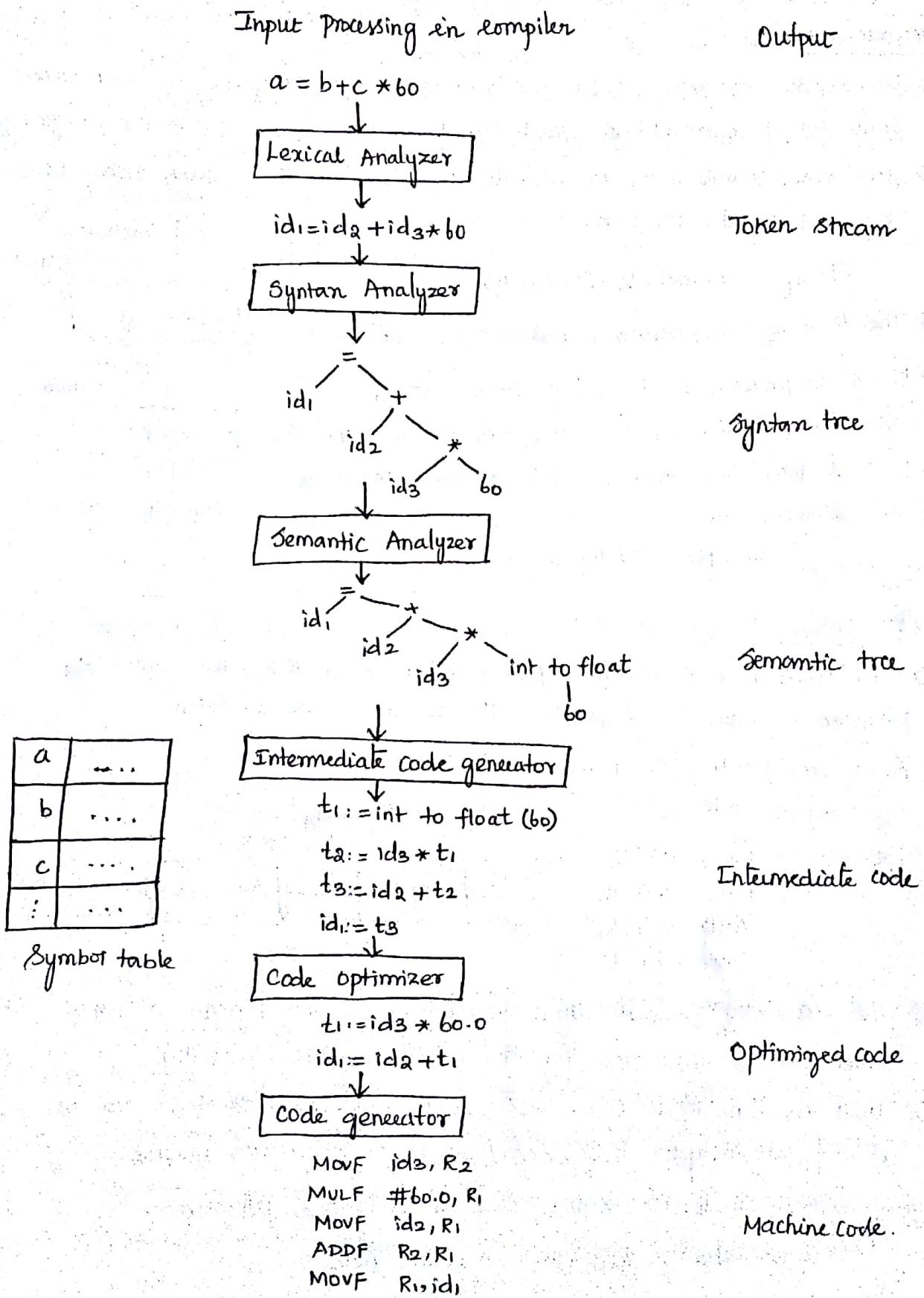
## Error detection and handling:-

- ⇒ In compilation, each phase detects errors. These errors must be reported to error handler whose task is to handle the errors.
- ⇒ Normally, the errors are reported in the form of message.
- ⇒ When the input characters from the input do not form the token, the lexical analyzer detects it as error. (Lexical errors)
- ⇒ Large number of errors can be detected in syntax analysis phase.
  - ✓ Errors in structure
  - ✓ Missing operators
  - ✓ Unbalanced parenthesis

} ⇒ Syntax errors
- ⇒ During semantic analysis; type mismatch kind of error is usually detected.

(1)

Example: Show how an input  $a = b + c * 60$  get processed in compiler.  
 Show that the output at each stage of compiler. Also show the contents of symbol table.



## Cousins of the compiler

Cousins of compiler means the entities in which the compiler typically operates. They are preprocessor, assemblers, loaders and link editors.

### 1. Preprocessor:-

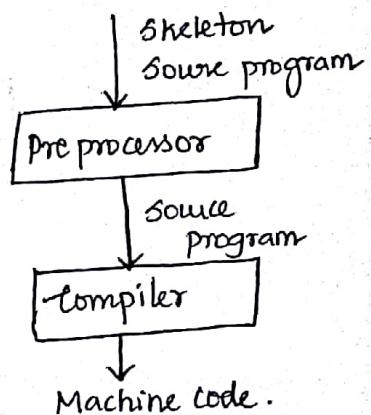
- ⇒ preprocessors allows user to use macros in the program. Macro means some set of instructions which can be used repeatedly in the program.
- ⇒ preprocessor allows users to include the header files which may be required by the program.

Example: `#include <stdio.h>`

⇒ The task of preprocessor is called file inclusion.

⇒ Macro preprocessor: Macro is a small set of instructions. Whenever in a program macro name is identified then that name is replaced by macro definition.

Example: `#define PI 3.14`



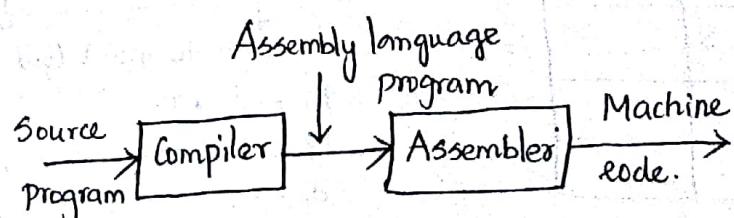
### 2. Assemblers:-

⇒ The assembler is a kind of translator which takes the assembly program as input and produces the machine code as output.

⇒ An assembly code is a mnemonic version of machine code.

Example:

```
MOV a, R1  
MUL #5, R1  
ADD #7, R1  
MOV R1, b
```



⇒ The assembler converts these instructions in the binary language which can be understood by the machine. (Machine code)

⇒ This machine code is a relocatable machine code that can be passed directly to the loader/linker for execution purpose.

⇒ Pass 1 assembler - complete scan of the input program.

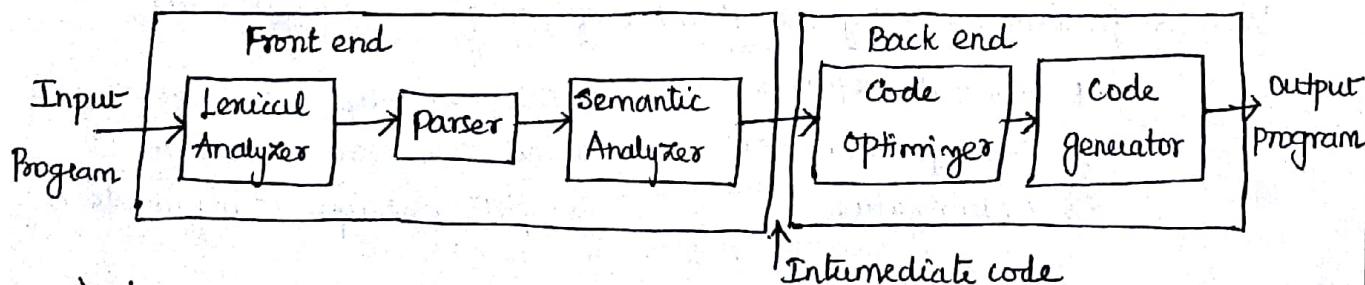
Pass 2 assembler - Relocatable machine code.

### 3. Loaders and Link Editors

- ⇒ Loader is a program which performs two functions loading and link editing.
- ⇒ Loading is the process in which the relocatable machine code is read and the relocatable addresses are altered. Then that code with altered instructions and data is placed in the memory at proper location.
- ⇒ The job of link editor is to make a single program from several files of relocatable machine codes.

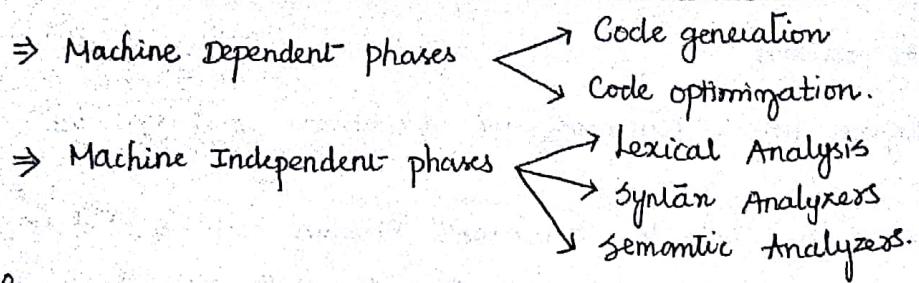
### Grouping of phases :-

- ⇒ Different phases of compiler can be grouped together to form a front end and back end.
- ⇒ The front end consists of those phases that primarily dependent on the source language and independent on the target language.
- ⇒ The front end consists of analysis part. Typically it includes,
  - Lexical Analysis, syntax analysis and semantic analysis.
- ⇒ The back end consists of those phases that are typically dependent upon the target language and independent on the source language.
- ⇒ It includes code generation and code optimization.



⇒ The front end and back end of the compiler is very much advantageous because of following reasons -

1. By keeping the same front end and attaching different back ends one can produce a compiler for some source language on different machines.
2. By keeping different front ends and same back end one can compile several different languages on the same machine.



### Concept of pass:-

- ⇒ One complete scan of the source language is called pass. It includes reading and input file and writing to an output file.
- ⇒ Many phases can be grouped one pass. It is difficult to compile the source program into a single pass, because the program may have some forward reference.
- ⇒ On the other hand, if we group several phases into one pass, we may be forced to keep the entire program in the memory. Therefore memory requirement may be large.
- ⇒ There may be a requirement of more than one pass in compilation process.
  - ⇒ A typical arrangement for optimizing compiler is one pass for scanning and parsing.
  - ⇒ One pass for semantic analysis and third pass for code generation and target code optimization.

### Factors affecting the number of passes:-

Various factors affecting the number of passes in compilers are,

- |                        |                                  |
|------------------------|----------------------------------|
| 1. Forward reference   | * The compiler can require       |
| 2. Storage limitations | more than one passes to          |
| 3. Optimization.       | complete subsequent information. |

### Difference between phase and pass

#### phase

⇒ The process of compilation is carried out in various steps. These steps are referred as phases.

⇒ Phases of compilation are lexical analysis, syntax analysis, semantic analysis, intermediate code generation, code generation, code optimization.

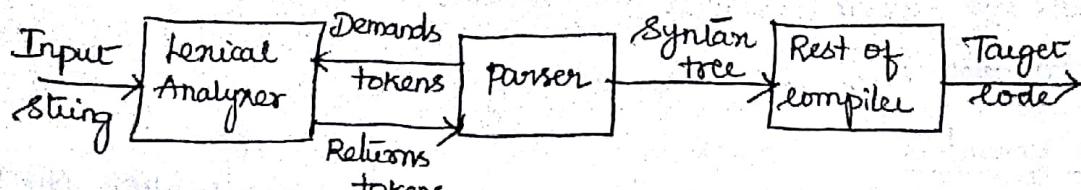
#### pass

⇒ Various phases are logically grouped together to form a pass.

⇒ Process of compilation can be carried out in single pass or in multiple passes.

## Lexical Analysis:-

⇒ Lexical Analyzer is the first phase of compiler. The lexical analyzer reads the input source program from left to right one character at a time and generate the sequence of tokens.



## Role of lexical analyzer:-

1. It produces stream of tokens.
2. It eliminates blank and comments.
3. It generates symbol table which stores the information about identifiers, constants encountered in the input.
4. It keeps track of line numbers.
5. It reports the error encountered while generating the tokens.

## Lexical errors:-

⇒ These types of errors can be detected during lexical analysis phase.

Typical lexical phase errors are

- Spelling errors. Hence get incorrect tokens.
- Exceeding length of identifier or numeric constants.
- Appearance of illegal characters.

## Error recovery action

1. If there are more than one pattern matching, then the lexical analyzer has to choose the longest lexeme matched.
2. Another rule is that the lexical rule written in LEX specification file must be ordered.

## Compiler construction Tools

- ⇒ Writing a compiler is tedious and time consuming task. There are some specialized tools for helping in implementation of various phases of compilers.
- ⇒ These tools are called compiler construction tools. Various compiler construction tools are,

### 1. Scanner generators:

- ⇒ These generators generate lexical analyzers. The specification given to these generators are in the form of regular expressions.
- ⇒ The UNIX has utility for a scanner generator called LEX. The specification given to the LEX consists of regular expressions for representing various tokens.

### 2. PARSER generators:

- ⇒ These produce the syntactic analyzer. The specification given to these generators is given in the form of context free grammar.
- ⇒ Typically UNIX has a tool called YACC which is a parser generator.

### 3. Syntactic-directed translation engines:-

- ⇒ In this tool the parse tree is scanned completely to generate an intermediate code.
- ⇒ The translation is done for each node of the tree.

### 4. Automatic code generators:-

- ⇒ These generators take an intermediate code as input and converts each rule of intermediate language into equivalent machine language.
- ⇒ The template matching technique is used. The intermediate code statements are replaced by templates that represent the corresponding sequence of machine instructions.

### 5. Data flow Engines:-

- ⇒ The data flow analysis is required to perform good code optimization. The data flow engines are basically useful in code optimization.

## Tokens, Patterns, Lenemes.

Tokens  $\Rightarrow$  It describes the class or category of input string.

Example: Identifiers, keywords, constants.

Patterns  $\Rightarrow$  Set of rules that describes the token.

Lenemes  $\Rightarrow$  sequence of characters in the source program that are matched with the pattern of token.

Example: int, i, num, vars etc..

## Issues in Lexical analyzer:-

1. The lexical analysis and syntax analysis are separated out for simplifying the task of one or other phases. This separation reduces the burden on parsing phase.
2. Compiler efficiency gets increased if the lexical analyzer is separated. This is because a large amount of time is spent on reading the source file.
3. Input alphabet peculiarities and other device specific anomalies can be restricted to the lexical analyzer. The representation of non-standard symbol can be isolated in lexical analyzer.

## Input Buffering:-

$\Rightarrow$  The lexical analyzer 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 portion of the input scanned.

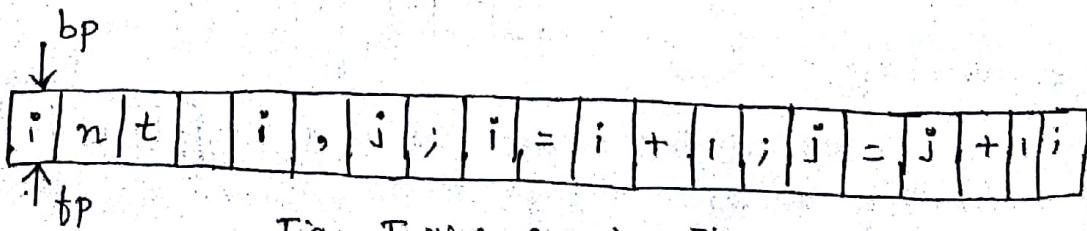
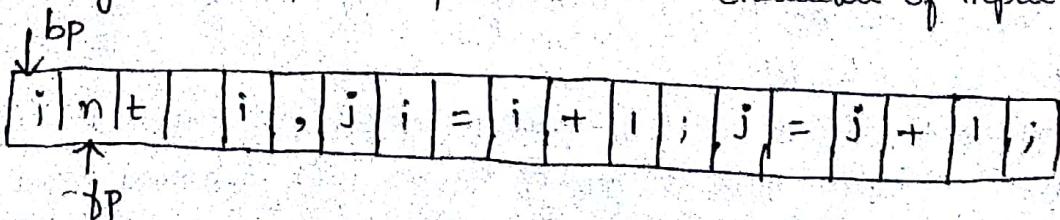
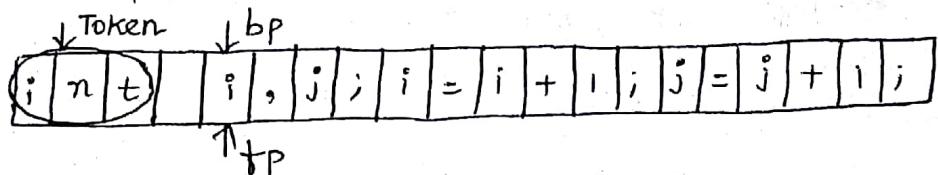


Fig: Initial configuration.

$\Rightarrow$  Initially both the pointers point to the 1st character of input string



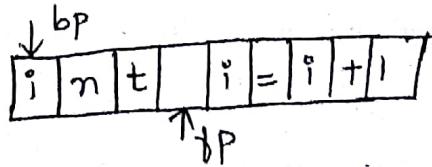
- ⇒ The forward\_ptr moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme.
- ⇒ Here, forward\_ptr(fp) encounters a blank space the lexeme "int" is identified.
- ⇒ The fp will be moved ahead at white space. When fp encounters white space, it ignore and moves ahead. Then both the begin\_ptr(bp) and forward\_ptr(fp) are set at next token i.



- ⇒ The input character is thus read from secondary storage. But reading in this way from secondary storage is costly. Hence buffering technique is used.
- ⇒ A block of data is first read into a buffer and then scanned by lexical analyzer. There are two methods used in this context.
  - one buffer scheme and two buffer scheme.

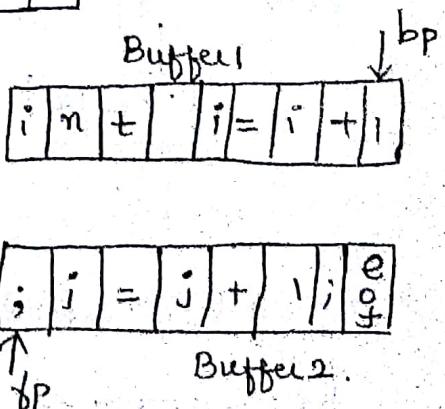
### 1. one buffer scheme:-

- ⇒ In this one buffer scheme, only one buffer is used to store the input string.
- ⇒ But the problem with this scheme is that if lexeme is very long then it crosses the buffer boundary, to scan rest of the lexeme the buffer has to be refilled, that makes overwriting the first part of lexeme.



### 2. Two buffer scheme:-

- ⇒ In this method two buffers are used to store the input string.
- ⇒ The first buffer and second buffer are scanned alternately.
- ⇒ When end of current buffer is reached the other buffer is filled.
- ⇒ The only problem with this method is that if length of the lexeme is longer than length of buffer then scanning input cannot be scanned completely.



## Specification of tokens:

- ⇒ To specify the tokens regular expressions are used. When a pattern is matched by some regular expression then token can be recognized.
- ⇒ Regular expressions are mathematical symbolisms which describe the set of strings of specific language. It provides convenient and useful notation for representing tokens.

## Definition of Regular expression

1.  $\emptyset$  is a regular expression that denotes the set containing empty string.
2. If  $R_1$  and  $R_2$  are RE then  $R = R_1 \cup R_2$  also RE (Union operation)
3. If  $R_1$  and  $R_2$  are RE then  $R = R_1 \cdot R_2$  also RE (Concatenation operation)
4. If  $R_1$  is a RE then  $R = R_1^*$  also regular expression (Kleen closure)

A language denoted by regular expressions is said to be a regular set or a regular language.

## Example:

1. RE for a language containing the strings of length two over  $\Sigma = \{0, 1\}$ .

Solution:  $R.E = (0+1)(0+1)$

2. RE for a language containing strings which end with "abb" over  $\Sigma = \{a, b\}$

Solution:  $(a+b)^* abb$

3. RE for recognizing identifiers

Solution:  $R.E = \text{letter} (\text{letter} + \text{digit})^*$

Where letter =  $(A, B, \dots, Z, a, b, \dots, z)$  and digit =  $(0, 1, 2, \dots, 9)$

4. RE for the language containing all the strings with any number of 'a's and 'b's

Solution:  $R.E = (a+b)^*$

5. RE for the language accepting all the strings which are ending with "00" over the set  $\Sigma = \{0, 1\}$

Solution:  $R.E = (0+1)^* 00$

6. RE to represent the valid date format in mm-dd-yy format

Solution:  $[1-12] - [1-31] - [0-9][0-9][0-9][0-9]$

## Design of lexical analysis (LEX) :

### Introduction - LEX

⇒ LEX is a unix utility which generates lexical analyzer. LEX scans the source program in order to get the stream of tokens and these tokens are related together.

⇒ The LEX specification file can be created using the extension .l. For example, the specification file can be x.l.

⇒ This x.l file is then given to LEX compiler to produce lex.yy.c.

⇒ This lex.yy.c is a C program which is actually a lexical analyzer program.

⇒ The LEX Specification file stores the regular expressions for the tokens and the lex.yy.c file consists of the tabular representation of the transition diagrams constructed for the regular expression.

⇒ Finally the compiler compiles this generated lex.yy.c and produces an object program a.out, when some input stream is given to a.out then sequence of tokens gets generated.

### Structure of LEX

LEX program consists of three parts

#### 1. Declaration section

- ↳ declaration of variable, constants
- ↳ regular definitions can also be written.

% { Declaration section

% }

%/%

#### 2. Rule section

- ↳ consists of regular expressions with associated actions.

Rule section

%/%

#### 3. Procedure section

- ↳ All the required procedures are defined.

Auxiliary procedure section.

Following commands are used to run the lex program x.l

- \$ lex x.l → Generates lex.yy.c
- \$ cc lex.yy.c → compiles lex.yy.c
- \$ ./a.out → Runs an executable file.

Example: Lex program for word counting.

```
% {
    int char_cnt=0, word_cnt=0, line_cnt=0; // Declaration section
}%
word [^ \t\n]+
%%%
{word} {word_cnt++; char_cnt+=yystrlen;} // Rule section
\n {char_cnt++; line_cnt++;}
    . {char_cnt++;}
}%
main()
{
    yylex();
    printf ("\n The character count = %d", char_cnt);
    printf ("\n The word count = %d", word_cnt);
    printf ("\n The line count = %d", line_cnt);
    printf ("\n");
}
int yywrap() // This function is called when end of file is
{
    return 1;
}
```

## Finite Automata

A finite automata is a collection of 5-tuple  $(Q, \Sigma, \delta, q_0, F)$

Where,  $Q$  - finite set of states  $\delta$  - Transition function.

$\Sigma$  - input alphabet  $F$  - set of final states.

$q_0$  - initial state

## Types of Automata

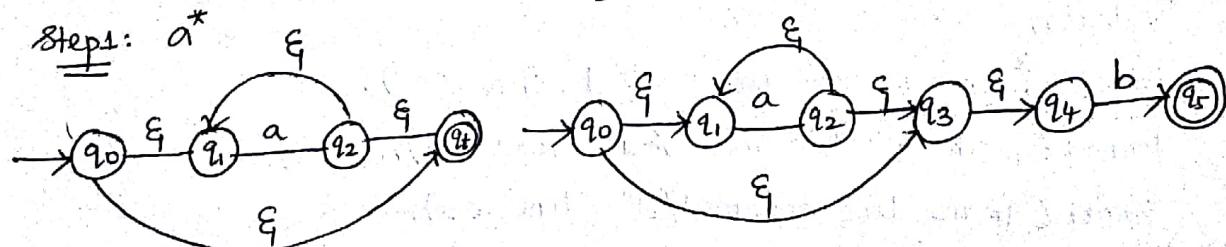
1. Deterministic Finite Automata  $\Rightarrow$  each move is uniquely determined on current input and current state.
2. Non-deterministic Finite Automata  $\Rightarrow$  Two or more path from current state to next state for a input.

## Conversion of regular expression of NDFA - Thompson's construction

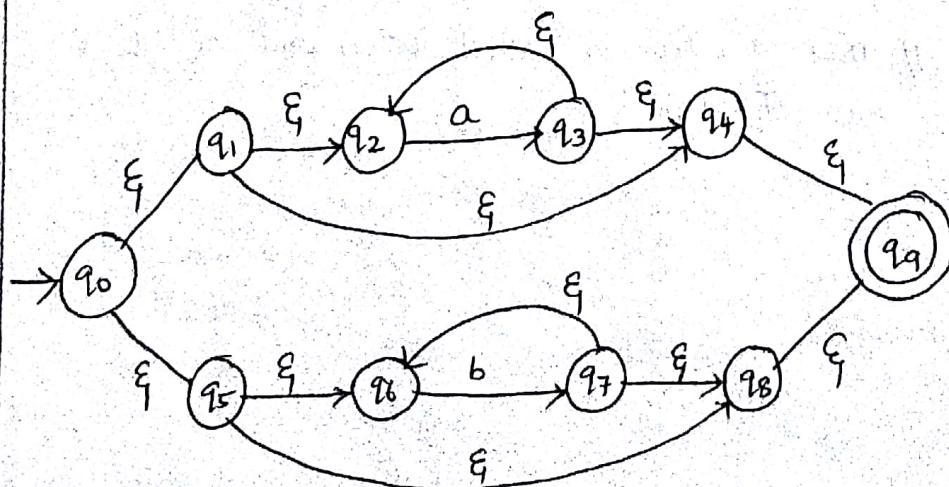
Prob 1: construct NFA equivalent to  $r = a^* b$

Solution:

Step 2:  $r = a^* b$



Prob 2: Draw NFA for  $a^* | b^*$



## Direct method for conversion of RE to FA

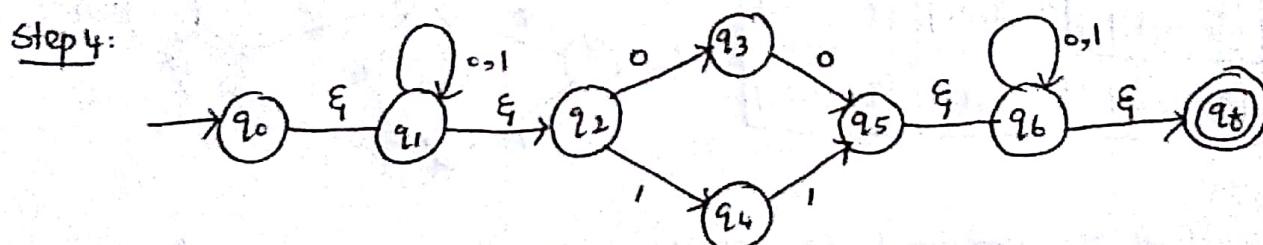
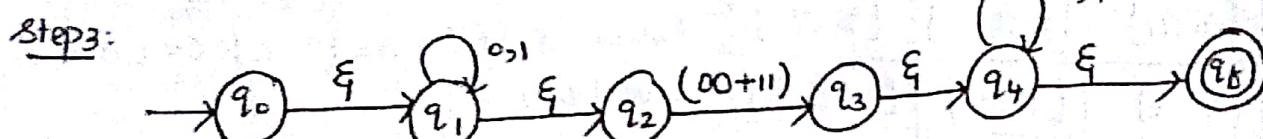
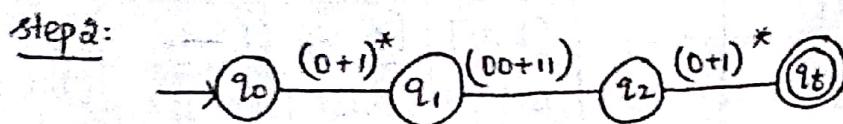
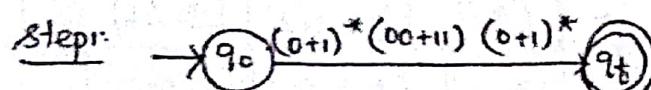
Step 1: Design a transition diagram for given RE using NFA with  $\epsilon$ -moves.

Step 2: convert this NFA with  $\epsilon$  to NFA without  $\epsilon$ .

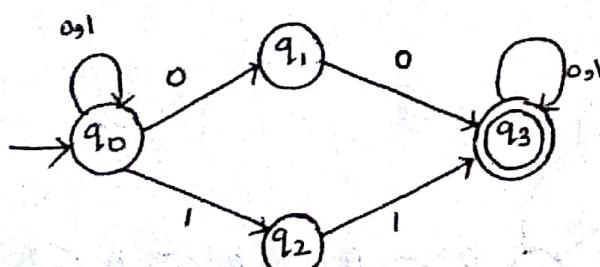
Step 3: convert the obtained NFA to equivalent DFA.

Point: construct DFA to accept the regular expression  $(0+1)^*(00+11)(0+1)^*$

Solution:  $RE = (0+1)^*(00+11)(0+1)^*$



Step 5: Eliminate  $\epsilon$ -transitions,



Step 7: From the transition table

$$\delta([q_0, q_1], 0) = [q_0, q_1, q_3] \rightarrow \text{New State}$$

$$\delta([q_0, q_1], 1) = [q_0, q_2] \rightarrow \text{New State}$$

Now,

$$\delta([q_0, q_1, q_3], 0) = [q_0, q_1, q_3]$$

$$\delta([q_0, q_1, q_3], 1) = [q_0, q_2, q_3] \rightarrow \text{New State}$$

$$\delta([q_0, q_2], 0) = [q_0, q_1]$$

$$\delta([q_0, q_2], 1) = [q_0, q_2, q_3]$$

$\alpha$	$\Sigma$	0	1
$q_0$		$\{q_0, q_1\}$	$\{q_0, q_2\}$
$q_1$		$q_3$	-
$q_2$		-	$q_3$
$q_3$		$q_3$	$q_3$

Step 6: Transition Table

Now,

$$\delta([q_0, q_2, q_3], 0) = [q_0, q_1, q_3]$$

$$\delta([q_0, q_2, q_3], 1) = [q_0, q_1, q_3]$$

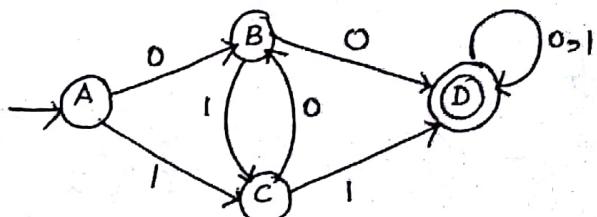
Transition table

$\Sigma$	0	1
$Q$		
[ $q_0$ ]	[ $q_0, q_1$ ]	[ $q_0, q_2$ ]
[ $q_1$ ]	$q_3$	-
[ $q_2$ ]	-	$q_3$
* [ $q_3$ ]	$q_3$	$q_3$
[ $q_0, q_1$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_2$ ]
[ $q_0, q_2$ ]	[ $q_0, q_1$ ]	[ $q_0, q_2, q_3$ ]
* [ $q_0, q_1, q_3$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_2, q_3$ ]
* [ $q_0, q_2, q_3$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_2, q_3$ ]

Minimized DFA

$\Sigma$	0	1
$Q$		
[ $q_0$ ]	[ $q_0, q_1$ ]	[ $q_0, q_2$ ]
[ $q_0, q_1$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_2$ ]
[ $q_0, q_2$ ]	[ $q_0, q_1$ ]	[ $q_0, q_2, q_3$ ]
[ $q_0, q_1, q_3$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_1, q_3$ ]
[ $q_0, q_2, q_3$ ]	[ $q_0, q_1, q_3$ ]	[ $q_0, q_1, q_3$ ]

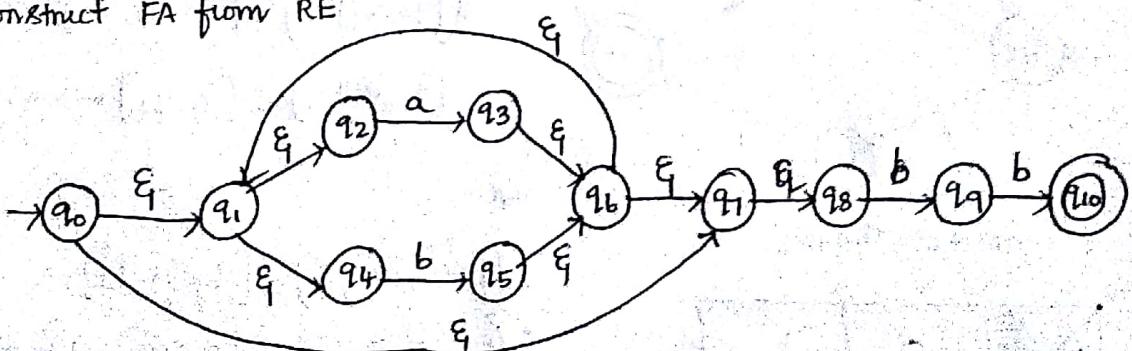
Assume  $A = [q_0]$ ,  $B = [q_0, q_1]$ ,  $C = [q_0, q_2]$ ,  $D = [q_0, q_1, q_3]$



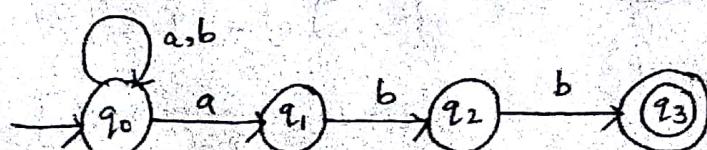
Prob: Give the minimized DFA for the following expression  $(a|b)^*abb$ .

Solution: RE =  $(a|b)^*abb$

Step 1: Construct FA from RE



Step 2: Eliminate  $\epsilon$ -transition.



Transition table:

$\Sigma$	a	b
Q		
$q_0$	$q_0, q_1$	$q_0$
$q_1$	$\phi$	$q_2$
$q_2$	$\phi$	$q_3$
$q_3$	$\phi$	$\phi$

From the table,

$$\delta([q_0, q_1], a) = [q_0, q_1]$$

$\delta([q_0, q_1], b) = [q_0, q_2] \rightarrow \text{new state}$

$$\delta([q_0, q_2], a) = [q_0, q_1]$$

$\delta([q_0, q_2], b) = [q_0, q_3] \rightarrow \text{new state}$

$$\delta([q_0, q_3], a) = [q_0, q_1]$$

$$\delta([q_0, q_3], b) = [q_0]$$

Equivalent DFA table

$\Sigma$	a	b
Q		
$q_0$	$[q_0, q_1]$	$q_0$
$q_1$	$\phi$	$q_2$
$q_2$	$\phi$	$q_3$
$q_3$	$\phi$	$\phi$
$[q_0, q_1]$	$[q_0, q_1]$	$[q_0, q_2]$
$[q_0, q_2]$	$[q_0, q_1]$	$[q_0, q_3]$
$[q_0, q_3]$	$[q_0, q_1]$	$q_0$

Mimimized DFA,

$\Sigma$	a	b
Q		
$q_0$	$[q_0, q_1]$	$q_0$
$[q_0, q_1]$	$[q_0, q_1]$	$[q_0, q_2]$
$[q_0, q_2]$	$[q_0, q_1]$	$[q_0, q_3]$
$[q_0, q_3]$	$[q_0, q_1]$	$q_0$

Assume,

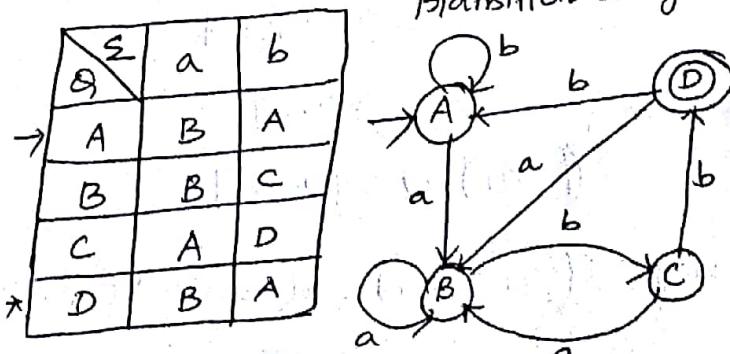
$$q_0 = A$$

$$[q_0, q_1] = B$$

$$[q_0, q_2] = C$$

$$[q_0, q_3] = D$$

Transition diagram,



Prob: Prove that the following two RE are equivalent by showing that the minimum state DFA's are same

$$(i) (a/b)^* \quad (ii) (a^* / b^*)^*$$

Prob: Illustrate Thompson's construction by drawing a NFA for the RE  $(a/b)^*a$ . convert this NFA to DFA using subset construction.

Ans:

Prob: construct the minimized DFA for RE  $(0+1)^* (0+1) 10$

### Conversion of N DFA to DFA.

Prob: Let  $M = (\{q_0, q_1\}, \{0, 1\}, \delta, q_0, \{q_1\})$  be NFA where  $\delta(q_0, 0) = \{q_0, q_1\}$ ,  $\delta(q_0, 1) = \{q_1\}$ ,  $\delta(q_1, 0) = \emptyset$ ,  $\delta(q_1, 1) = \{q_0, q_1\}$  construct its equivalent DFA.

Solution: Given,  $M = (\{q_0, q_1\}, \{0, 1\}, \delta, q_0, \{q_1\})$  be NFA

$$\therefore \text{DFA } M' = (Q', \Sigma, \delta', q_0', F')$$

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

Transition table for NFA

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

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

$$\delta(q_1, 0) = \emptyset$$

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

compute transition from the state  $[q_0, q_1]$

$$\begin{aligned}\delta([q_0, q_1], 0) &= \delta(q_0, 0) \cup \delta(q_1, 0) \\ &= \{q_0, q_1\} \cup \emptyset \\ &= \{q_0, q_1\}\end{aligned}$$

$$\therefore \delta([q_0, q_1], 0) = [q_0, q_1]$$

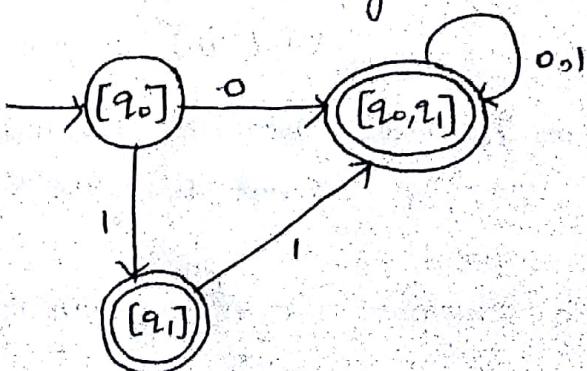
$$\begin{aligned}\text{By } \delta([q_0, q_1], 1) &= \delta(q_0, 1) \cup \delta(q_1, 1) \\ &= \{q_1\} \cup \{q_0, q_1\} \\ &= \{q_0, q_1\}\end{aligned}$$

$$\therefore \delta([q_0, q_1], 1) = [q_0, q_1]$$

$\therefore$  Transition table of DFA

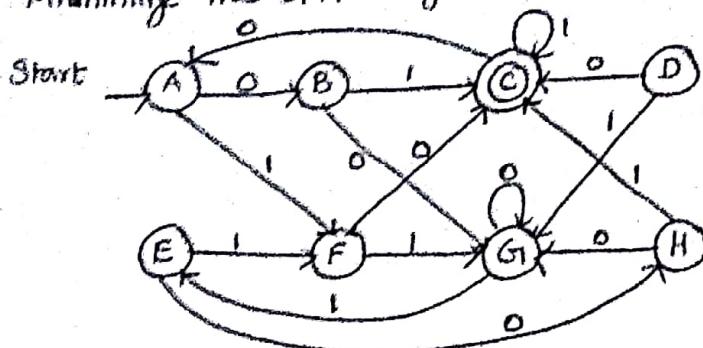
$\Sigma$	0	1
$q_0$	$[q_0, q_1]$	$[q_1]$
$q_1$	$\emptyset$	$[q_0, q_1]$
$[q_0, q_1]$	$[q_0, q_1]$	$[q_0, q_1]$

Transition diagram.



## Minimization of NDFA

Prob 1: Minimize the DFA as given below,



Solution:

Transition table

	0	1
→ A	B	F
B	G <sub>1</sub>	C
(C)	A	C
D	C	G <sub>1</sub>
E	H	F
F	C	G <sub>1</sub>
G <sub>1</sub>	G <sub>1</sub>	E
H	G <sub>1</sub>	C

find the equivalent state

consider pair (B, H)

$$\begin{array}{l|l} \delta(B, 0) = G_1 & \delta(B, 1) = C \\ \delta(H, 0) = G_1 & \delta(H, 1) = C \end{array} \therefore (B, H) \text{ are equivalent}$$

Since (B, H) are equivalent, consider the pair (A, E)

$$\begin{array}{l|l} \delta(A, 0) = B \text{ or } H & \delta(A, 1) = F \\ \delta(E, 0) = H \text{ or } B & \delta(E, 1) = F \end{array} \therefore (A, E) \text{ are equivalent}$$

consider the pair (D, F)

$$\begin{array}{l|l} \delta(D, 0) = C & \delta(D, 1) = G_1 \\ \delta(F, 0) = C & \delta(F, 1) = G_1 \end{array} \therefore (D, F) \text{ are equivalent}$$

∴ the equivalent pairs are (A, E), (B, H) and (D, F)

The minimized DFA will be,

Table filling,

α \ Σ	0	1
→ A	B	D
B	G <sub>1</sub>	C
(C)	A	C
D	C	G <sub>1</sub>
G <sub>1</sub>	G <sub>1</sub>	A

B	x					
C	x	x				
D	x	x	x			
E	(A,E)	x	x	x		
F	x	x	x	(D,F)	x	
G <sub>1</sub>	x	x	x	x	x	x
H	x	(B,H)	x	x	x	x
A		B	C	D	E	F
	<td><td><td><td><td>G<sub>1</sub></td></td></td></td></td>	<td><td><td><td>G<sub>1</sub></td></td></td></td>	<td><td><td>G<sub>1</sub></td></td></td>	<td><td>G<sub>1</sub></td></td>	<td>G<sub>1</sub></td>	G <sub>1</sub>

Pbm 8: Convert the given NFA into its equivalent DFA.



Solution:

Step 1: Find  $\epsilon$ -closure of each state

$$\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\}$$

$$\epsilon\text{-closure}(q_1) = \{q_1, q_2\}$$

$$\epsilon\text{-closure}(q_2) = \{q_2\}$$

Step 2: Find  $\delta'$  transition

Let  $\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\}$  be state A.

$$\begin{aligned}\delta'(A, 0) &= \epsilon\text{-closure} \{ \delta(q_0, 0), q_1, q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_0, 0) \cup \delta(q_1, 0) \cup \delta(q_2, 0) \} \\ &= \epsilon\text{-closure} \{ q_0 \} \\ &= \{q_0, q_1, q_2\} \Rightarrow \text{state A}\end{aligned}$$

$$\begin{aligned}\delta'(A, 1) &= \epsilon\text{-closure} \{ \delta(q_0, 1), q_1, q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_0, 1) \cup \delta(q_1, 1) \cup \delta(q_2, 1) \} \\ &= \epsilon\text{-closure} \{ q_1 \} \\ &= \{q_1, q_2\} \Rightarrow \text{state B}\end{aligned}$$

$$\begin{aligned}\delta'(A, 2) &= \epsilon\text{-closure} \{ \delta(q_0, 2), q_1, q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_0, 2) \cup \delta(q_1, 2) \cup \delta(q_2, 2) \} \\ &= \epsilon\text{-closure} \{ q_2 \} \\ &= \{q_2\} \Rightarrow \text{state C}\end{aligned}$$

Now,

$$\begin{aligned}\delta'(B, 0) &= \epsilon\text{-closure} \{ \delta(q_1, 0), q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_1, 0) \cup \delta(q_2, 0) \} \\ &= \epsilon\text{-closure} \{ \phi \} \\ &= \phi\end{aligned}$$

$$\begin{aligned}\delta'(B, 1) &= \epsilon\text{-closure} \{ \delta(q_1, 1), q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_1, 1) \cup \delta(q_2, 1) \} \\ &= \epsilon\text{-closure} \{ q_1 \} \\ &= \{q_1, q_2\} \Rightarrow \text{state B}\end{aligned}$$

$$\begin{aligned}\delta'(B, 2) &= \epsilon\text{-closure} \{ \delta(q_1, 2), q_2 \} \\ &= \epsilon\text{-closure} \{ \delta(q_1, 2) \cup \delta(q_2, 2) \} \\ &= \epsilon\text{-closure} \{ q_2 \} \\ &= \{q_2\} \Rightarrow \text{state C}\end{aligned}$$

Now,

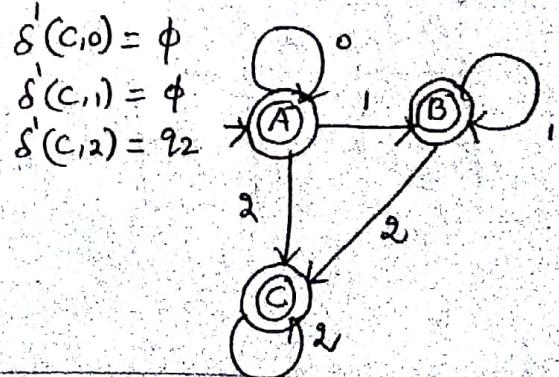
$$\delta'(C, 0) = \epsilon\text{-closure} \{ \delta(q_2, 0) \}$$

$$\begin{aligned}&= \epsilon\text{-closure} \{ \phi \} \\ &= \phi\end{aligned}$$

$$\begin{aligned}\delta'(C, 1) &= \epsilon\text{-closure} \{ \delta(q_2, 1) \} \\ &= \epsilon\text{-closure} (\phi) \\ &= \phi\end{aligned}$$

$$\begin{aligned}\delta'(C, 2) &= \epsilon\text{-closure} \{ \delta(q_2, 2) \} \\ &= q_2\end{aligned}$$

$$\begin{array}{ll} \delta'(A, 0) = A & \delta'(B, 0) = \phi \\ \delta'(A, 1) = B & \delta'(B, 1) = B \\ \delta'(A, 2) = C & \delta'(B, 2) = C \end{array}$$



Ques 2: Using table filling method minimize the following DFA. Draw the transition diagram of resulting DFA.

	O	I
A	B	E
B	C	F
*	D	H
D	E	H
*	F	I
F	G <sub>1</sub>	B
G <sub>1</sub>	H	B
H	I	C
*	A	E

Solution: Here final states are C, F, I

since C, F, I are final states

consider the pair (E, H)

$$\begin{array}{|l|l|} \hline \delta(E, O) = F & \delta(E, I) = I \\ \hline \delta(H, O) = I & \delta(H, I) = C \\ \hline \end{array}$$

∴ (E, H) are equivalent state

consider the pair (B, H)

$$\begin{array}{|l|l|} \hline \delta(B, O) = C & \delta(B, I) = F \\ \hline \delta(H, O) = I & \delta(H, I) = C \\ \hline \end{array}$$

∴ (B, H) are equivalent state.

consider the pair (B, E)

$$\begin{array}{|l|l|} \hline \delta(B, O) = C & \delta(B, I) = F \\ \hline \delta(E, O) = F & \delta(E, I) = E \\ \hline \end{array}$$

(B, E) are equivalent state

since (B, H) and (B, E) are equivalent

consider the pair (A, G<sub>1</sub>)

$$\begin{array}{|l|l|} \hline \delta(A, O) = B \text{ or } H & \delta(A, I) = E \text{ or } B \\ \hline \delta(G_1, O) = H \text{ or } B & \delta(G_1, I) = B \text{ or } E \\ \hline \end{array}$$

∴ (A, G<sub>1</sub>) is equivalent state.

Since (B, E) and (E, H) are equivalent

Since (A, G<sub>1</sub>) and (B, E) are equivalent

consider the pair (A, D)

$$\begin{array}{|l|l|} \hline \delta(A, O) = B \text{ or } E & \delta(A, I) = E \text{ or } H \\ \hline \delta(D, O) = E \text{ or } B & \delta(D, I) = H \text{ or } E \\ \hline \end{array}$$

∴ (A, D) is equivalent state

consider the pair (F, I)

$$\begin{array}{|l|l|} \hline \delta(F, O) = G_1 \text{ or } A & \delta(F, I) = B \text{ or } E \\ \hline \delta(I, O) = A \text{ or } G_1 & \delta(I, I) = E \text{ or } B \\ \hline \end{array}$$

∴ (F, I) is equivalent state

Since (E, H) and (B, H) are equivalent

consider the pair (D, G<sub>1</sub>)

$$\begin{array}{|l|l|} \hline \delta(D, O) = E \text{ or } H & \delta(D, I) = H \text{ or } B \\ \hline \delta(G_1, O) = H \text{ or } E & \delta(G_1, I) = B \text{ or } H \\ \hline \end{array}$$

∴ (D, G<sub>1</sub>) is equivalent state

Since

(A, D) and (E, H) are equivalent

consider the pair (C, I)

$$\begin{array}{|l|l|} \hline \delta(C, O) = D \text{ or } A & \delta(C, I) = H \text{ or } E \\ \hline \delta(I, O) = A \text{ or } B & \delta(I, I) = E \text{ or } H \\ \hline \end{array}$$

∴ (C, I) is equivalent.

Since (D, G<sub>1</sub>) and (B, H) are equivalent

consider the pair (C, F)

$$\begin{array}{|l|l|} \hline \delta(C, O) = D \text{ or } G_1 & \delta(C, I) = H \text{ or } B \\ \hline \delta(F, O) = G_1 \text{ or } D & \delta(F, I) = B \text{ or } H \\ \hline \end{array}$$

∴ (C, F) is equivalent.

Table filling

B	X						
C	X	X					
D	(A, D)	X	X				
E	X	(B, E)	X	X			
F	X	X	(C, F)	X	X		
G <sub>1</sub>	X	X	X	(D, G <sub>1</sub> )	X	X	
H	X	X	X	X	(E, H)	X	X
I	X	X	X	X	(F, I)	X	X

Here

A = G<sub>1</sub> = D

B = H = E

C = I = F

Mimized DFA

Σ	O	I
Q	A	B
A	B	B
B	C	C
C	A	B

## Derivation, parse tree and Ambiguity

Derivations  $\Rightarrow$  Use the production from head to body (i) from start symbol expanding till reaches the given string.

Two types of derivations are,

i) Left most derivation (LMD)

ii) Right most Derivation (RMD)

$\Rightarrow$  LMD is a derivation in which the leftmost nonterminal is replaced first from the sentential form. (i)  $S \xrightarrow{lm} \alpha$ , then  $\alpha$  is left sentential form.

$\Rightarrow$  RMD is a derivation in which rightmost nonterminal is replaced from the sentential form. (ii)  $S \xrightarrow{rm} \alpha$ , then  $\alpha$  is right sentential form.

## Parse tree (Derivation tree)

$\Rightarrow$  It is a graphical representation for the derivation of the given production rule for a given CFG.

Properties:

$\Rightarrow$  Root node is always a node indicating start symbol.

$\Rightarrow$  Derivation is read from left to right.

$\Rightarrow$  Leaf node is always terminal nodes.

$\Rightarrow$  Interior node are always non-terminal nodes.

Example:

Consider the grammar G has the production.

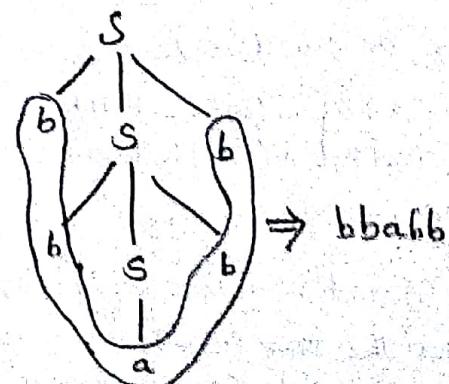
$S \rightarrow bSb \mid a \mid b$  and string "bbabb"

Derivation:

$S \rightarrow bSb$

$\rightarrow bbSbb [S \rightarrow bSb]$

$\rightarrow bbabb [S \rightarrow b]$



Problems: construct the derivation tree for the string "aababbba" using LMD and RMD using  $S \rightarrow aB/bA$ ,  $A \rightarrow a|aS|bAA$ ,  $B \rightarrow b|bS|aBB$ .

Solution:

LMD

$$S \rightarrow aB$$

$$\Rightarrow aaBB \quad [B \rightarrow aBB]$$

$$\Rightarrow aabbBB \quad [B \rightarrow aBB]$$

$$\Rightarrow aaabBB \quad [B \rightarrow b]$$

$$\Rightarrow aaabbB \quad [B \rightarrow b]$$

$$\Rightarrow aaabbaBB \quad [B \rightarrow aBB]$$

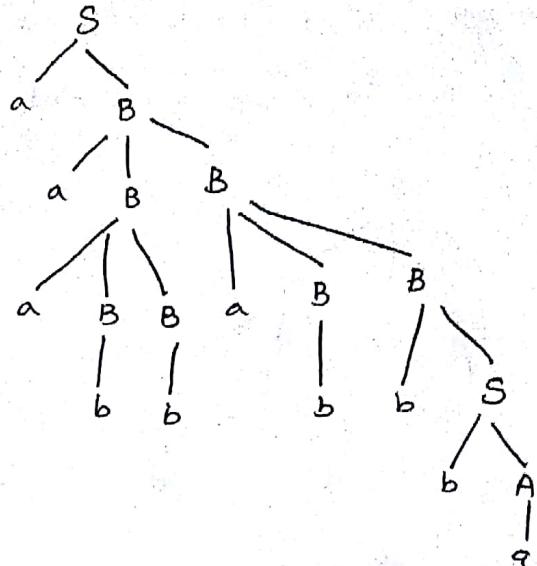
$$\Rightarrow aaabbabB \quad [B \rightarrow b]$$

$$\Rightarrow aaabbabbS \quad [B \rightarrow bS]$$

$$\Rightarrow aaabbabbbA \quad [S \rightarrow bA]$$

$$\Rightarrow aaabbabbba \quad [A \rightarrow a]$$

Parse tree



R.M.D

$$S \rightarrow aB$$

$$\Rightarrow aaBB \quad [B \rightarrow aBB]$$

$$\Rightarrow aabBS \quad [B \rightarrow bS]$$

$$\Rightarrow aaBBba \quad [S \rightarrow bA]$$

$$\Rightarrow aaB bba \quad [A \rightarrow a]$$

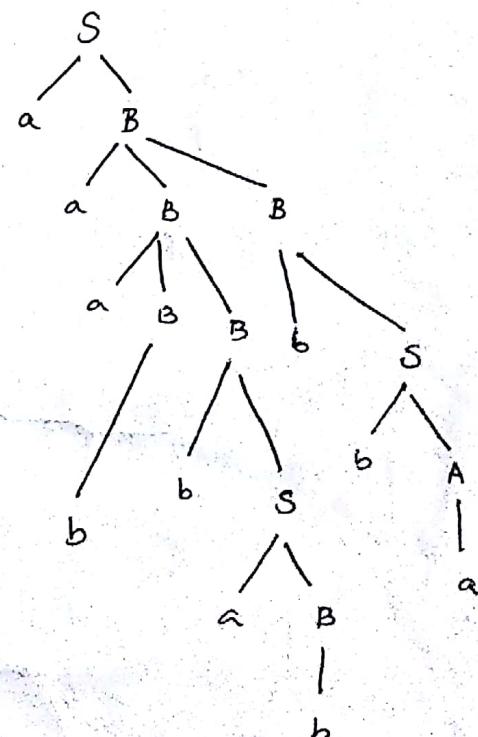
$$\Rightarrow aaaBBbba \quad [B \rightarrow aBB]$$

$$\Rightarrow aaaBbsbba \quad [B \rightarrow bS]$$

$$\Rightarrow aaaBbaBbba \quad [S \rightarrow aB]$$

$$\Rightarrow aaababbba \quad [B \rightarrow b]$$

$$\Rightarrow aaabbabbba \quad [B \rightarrow b]$$

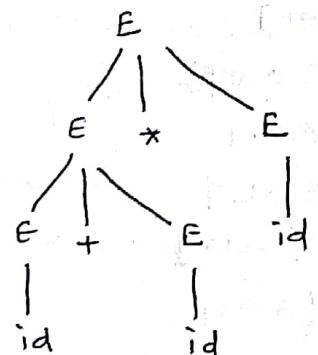
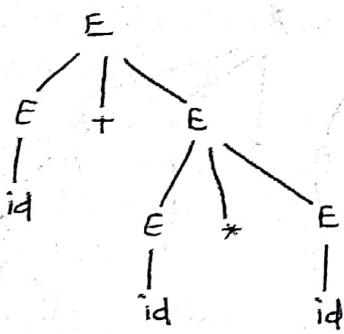


## Ambiguity:

If there exists more than one parse tree for a given grammar, that means there could be more than one leftmost or rightmost derivation. Possible and then that grammar is said to be ambiguous grammar.

Example:  $E \rightarrow E+E \mid E*E \mid id$

Then for string  $id + id * id$



Problem: Show that following grammar is ambiguous.

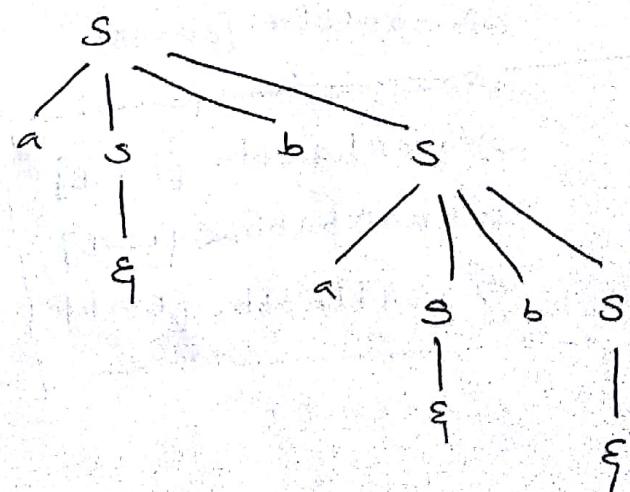
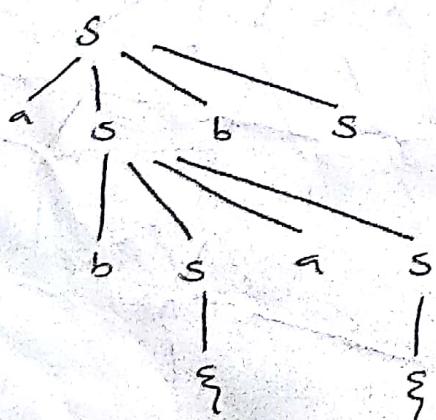
$$S \rightarrow a S b S$$

$$S \rightarrow b S a S$$

$$S \rightarrow \epsilon$$

Solution:

Consider the string "abab"



From the above parse tree, we can conclude  
the given grammar is ambiguous.