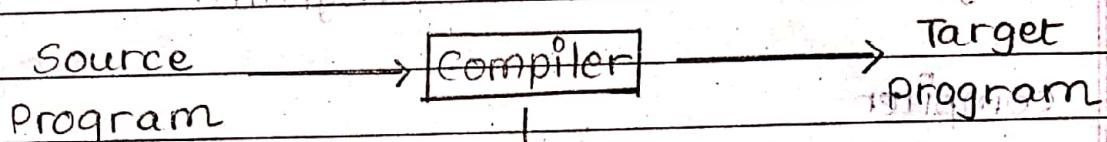


## What is a compiler? Design of compiler

It is a program that reads a program written in one language [source language] and translates it into an equivalent target program.



### → Classification of Compiler

- 1. Single Pass
  - 2. Multi Pass
  - 3. Load and go
  - 4. Debugging
  - 5. Optimized
- Debugging is like interpreters.

### → Analysis/Synthesis model of communication

Analysis



Synthesis

### → Analysis Tools.

- 1. Structure Editors
  - 2. Pretty Printers
  - 3. Static checkers
  - 4. Interpreters
- It is actually command based.
- To discuss potential bug.

- Following are additional analysis tools
  - ex. Text formattor silicon compiler
  - query interpreter
  - Output is in circuit form.
  - Output is in circuit diagram.
- Analysis of the source program.

compiler

Classification of compiler

Analysis and synthesis model of compiler

Analysis tools of compiler

## 1. Linear Analysis / Lexical Analysis / Scanning

The stream of character making up the source program is to read from left to right and grouped into tokens there are sequences of character having a collective meaning.

- Token.

Operators

Keywords

Datatype

Signs

## 2. Hierarchical Analysis / Syntax Analysis / Parsing.

All inner nodes are always operators and leaves are number or identifier.

In this phase of analysis characters or tokens are grouped hierarchically into nested collections with collective meaning.

### 3. Semantic Analysis (Type - checking)

In this analysis certain checks are performed to ensure that the components of a program fit together meaningfully.

→ compilers. [Language Processing System]

Skeletal Source Program

With an addition of comments to input.

PreProcessor

Machine

Source Program

Compiler

Target Assembly-Program

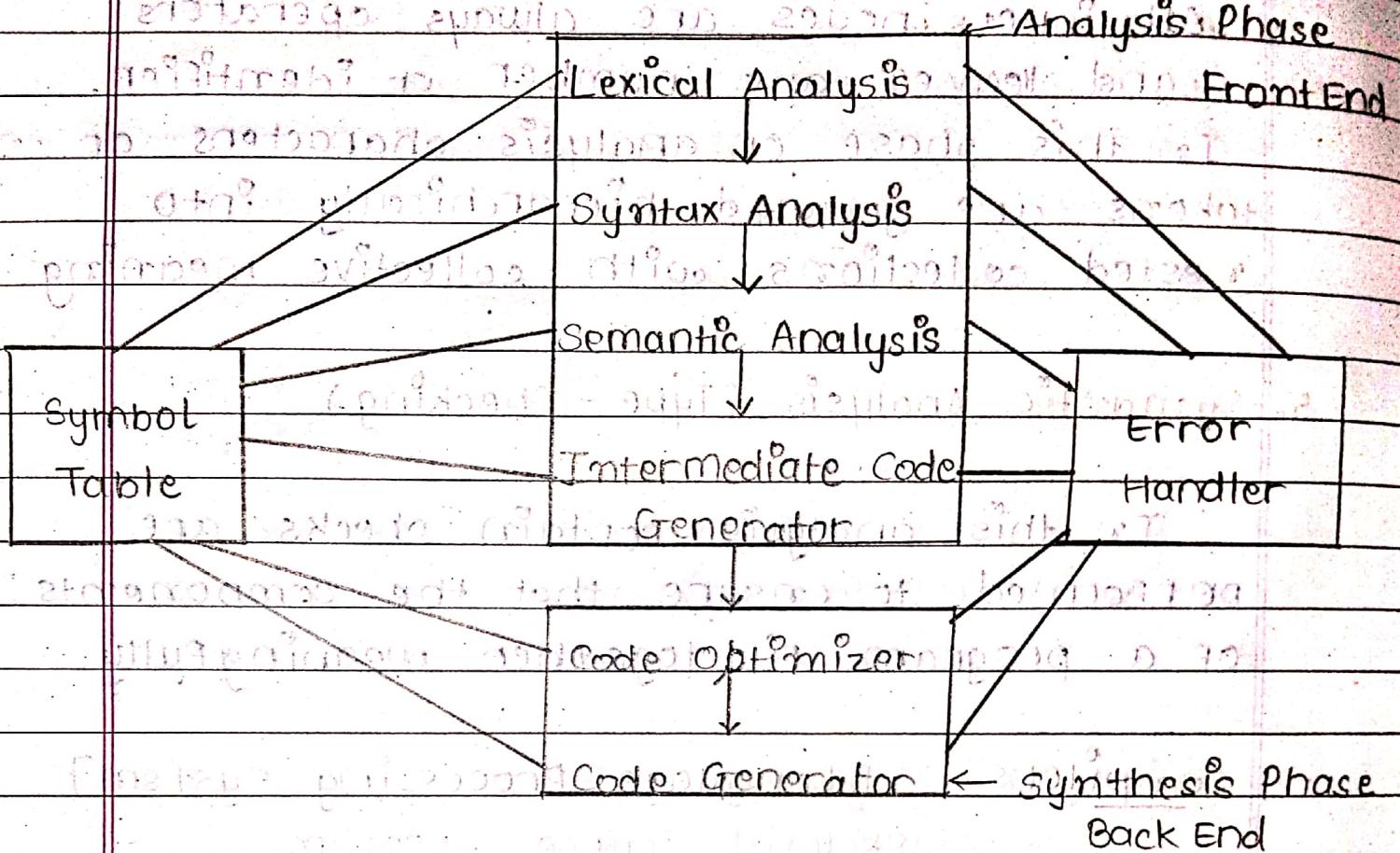
Assembler

Relocatable m/c code

Linker / Loader

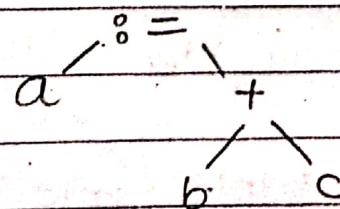
Target. m/c code.

→ Phases of compiler



- Input of one phase is output of other phase.
- Output of Lexical Analysis is sequence of Tokens.
- smallest Individuals of program are called Tokens.
- operators are always first

Ex :-  $a = b + c$



sum :-  $a + b + c$

$t1 \% = b * c$

$t2 \% = a + t1$

sum \% = t2.

## → Types of Error.

1. Lexical Error % - For unknown symbol

2. Syntax Error % - Checking of syntax

3. Semantic Errors % Focus is on type checking,

4. Logical Errors % - bad functionality

## → Cousins of Compiler.

1. Preprocessors

→ Macro processing

→ File inclusion

→ Rational Preprocessor

→ Language Extension

2. Assemblers

3. Two-Pass Assembly

4. Loaders and Linkers (Link-Editors).

## → The grouping of phase.

1. Front-End and Back-End

2. Passes

3. Reducing the number of passes.

→ Compiler construction Tools.

Following are systems for compiler writing process:

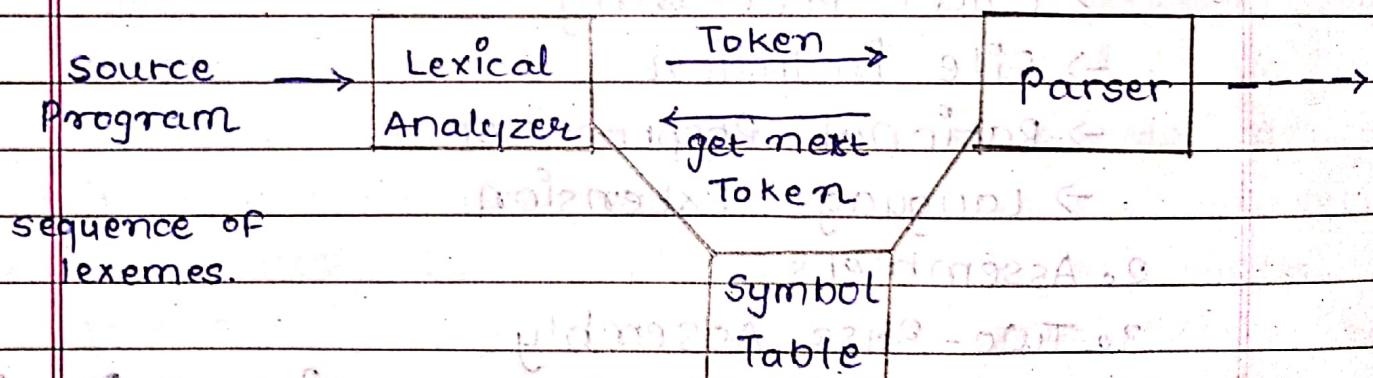
1. compiler - compilers
2. compiler generators
3. Translator writing systems

Following are compiler construction tools:

1. Parser Generators
2. Scanner Generators
3. Syntax directed translation engines
4. Automatic code Generators
5. Data Flow Engines.

Unit 2

→ The Role of the Lexical Analyzer.



- In compiler white space includes tab, newline character, blank.
- Comparison of lexemes → Lexical Analyzer

Lexical Analyzer divided into two phases:

1. Scanning - Simple task
2. Lexical Analysis - Complex Operation

It is the first phase of compiler.

Its main task is to read the input character and produce as output a sequence of tokens that the parser uses for syntax analysis.

Depending upon receiving a get next token command from the parser, the lexical analyzer reads input character until it can identify the next token.

Additional task is stripping out from the source program, comments and white spaces like blank, newline and tab characters.

#### → Issues in lexical analysis

1. Simpler Design

2. Compiler Efficiency

3. Compiler Portability.

#### → Define Following Terms :-

1. Lexeme :- It is a sequence of characters in the source program that is matched by the pattern for a token.

2. Pattern :- Set of strings is described by a rule called a pattern associated with the token.

3. Token : smallest individual units of program

(they are terminal symbols produced by lexical analyzer)

Token	Lexeme	Pattern
const	constant	const
function	function	function
if	if	if
relational operator	<, <=, >, >=, =	< or <= or > or >=
num	to	sequence of digit
id.	amount	letter followed by letter(s) or digit(s).

### → Attributes for Token.

When more than one pattern matches a Lexeme, the Lexical analyzer must provide additional information about the particular lexeme that matched to the subsequent phases of the compiler.

For example, the pattern num matches both the strings 0 and 1, but it's essential for the code generator to know what strings was actually matched.

## → Lexical Errors

Few errors not possible to detect during lexical analysis.

For eg. By mistake "if" it is entered "fi" instead "if" keyword then a lexical analyzer cannot feel whether "fi" is a misplacing of the keyword "if".

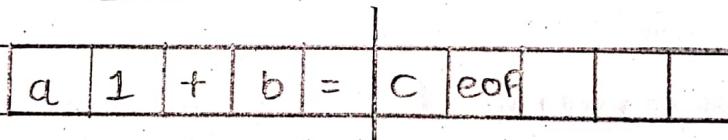
There is simplest recovery strategy called panic mode recovery.

Following are possible error recovery actions.

1. Deleting an extra character.
2. Inserting a missing character.
3. Replacing an incorrect character by a correct character.
4. Transposing two adjacent characters.

## → Input Buffering

### 1. Buffer Pairs



If forward of end of 1st half then begin

    reload second half;

    Forward := Forward + 1

end

else If forward at end of second half then begin

    reload first half;

    move forward to beginning of 1st half.

end

else forward := forward + 1;

- code to advance forward pointer.

This process is for lexical analysis and memory management.

## 2. Sentinels:-

a	=	b	+	eof	c	eof	eof	eof
---	---	---	---	-----	---	-----	-----	-----

forward := forward + 1;

if forward = eof then begin

    if forward at end of first half then begin  
        reload second half;

    forward := forward + 1;

end

else if forward at end of second half then begin

    reload first half;

    move forward to beginning of first half

end

else

    terminate lexical analysis

end

- Lookahead code with sentinels.

## \* Regular Expression

They are important notation for specifying patterns

Each pattern matches a set of strings, so regular expressions will serve as names for sets of strings.

### \* strings and languages

The term alphabet or character class denotes any finite set of symbols.

A string over some alphabet is a finite sequence of symbols, drawn from that alphabet.

In language theory the terms sentence and word are often used as synonymous from the term string.

#### ↳ Length of string :-

e.g.  $s = 'abc'$ ;  
 $s = 'a'$ ;

1st

→ empty string := string with length 0.

→ The term language denotes any set of strings over some fixed alphabet.

→ abstract languages can be presented using  $\{\emptyset\}$ , empty set  $\{\langle \rangle\}$  (no string).

→ exit - if  $x$  and  $y$  are strings then  
concatenation of  $x$  and  $y$  return  
as  $xy$  or  $x.y$ .

$\rightarrow x = 'abc'$   $y = 'xyz'$

$\rightarrow x = 'abc'$   $y = 'xyz'$

$\rightarrow x = 'abc'$   $xy = 'abce'$

$y = e$

$es = se = s$

$\rightarrow$  consider e.g. for concatenation as a product, define  $s, \theta \rightarrow$  length.

$s^0 \rightarrow e$

and for  $i > 0$  define  $s^i$  to be  $s^{i-1} s$ .

$\rightarrow$  Definitions for different parts of string:

1. prefix of  $s$ :

2. suffix of  $s$ :

3. substring of  $s$ :

4. proper prefix, suffix or substring of  $s$

e.g.  $s \neq x$   $x = \text{prefix, suffix, substring}$

5. subsequence of  $s$ :

e.g.  $1 - s = 'abcd' [bd]$

6. operations on language:

3 operation :- (1) union

(2) Concatenation

(3) Closure

union  $\rightarrow$  or, concat  $\rightarrow$  and,

Closure  $\rightarrow$  occurrence of particular symbol.

Consider following eg. Let  $L$  be the set having symbol  $A, B, \dots, z, a, b, \dots, z^2$

$\rightarrow$  New language can be created from  $L$  and  $D$  by applying the operators.

$$1. L \cup D =$$

$$3. L^4 = /LLL \rightarrow (L^+)$$

$$4. L^* = L^0 \rightarrow \epsilon, L^1 \rightarrow L, L^2 \rightarrow LL$$

$*$   $\rightarrow 0$  or more times occurrences.

$$5. D^+ \rightarrow . + \rightarrow 1 \text{ or more times}$$

$$D^1 \quad D^2$$

$$D = DD$$

$$(0 \text{ to } a) \quad (0 \text{ to } 99)$$

$$\rightarrow R.E \rightarrow L(LUD)$$

↳ Definitions of operations on language

Operation

Definition

1.  $L$  union  $M$  ( $L \cup M$ )

$$L \cup M = \{s | s \text{ is in } L \text{ for quantity } (s \in L) \text{ or } s \text{ is in } M\}$$

(Intersection) (Complement)

a. concatenation of languages  $Lm = \{st | s \in L$   
 $t \in m\}$   
 $L$  and  $m$        $L$  and  $t$  is in  $m\}$

b. Kleene closure

of  $L$  written as  $L^*$

$$L^* = \sum_{i=0}^{\infty} L^i$$

$L^*$  denotes 0 or more concatenation of  $L$

c. Positive closure of

$L$  written as  $L^+$

$$L^+ = \sum_{i=1}^{\infty} L^i$$

$\rightarrow L^+$  denotes 1 or more concatenation of  $L$ .

↳ Regular Expressions:

They are kinds of notations which defines different rules.

e.g. letter  $(\text{letter} | \text{digit})^*$

Component 1.

Component 2.

$\rightarrow$  There are rules that define the regular expression over alphabet set associated with each rule is a specification of the language denoted by the regular expression being defined.

(1)  $R \cdot E \rightarrow \epsilon$ ,  $\epsilon$  (input), outcome  $\in E^*$

(Input string)

(Token)

(2) If  $R.E.$  is  $a$ ,

lexeme  $a$

token  $a$

outcome  $\{a\}$

(3) Suppose  $r$  and  $s$ , are regular expression denoting languages, regular sets.

$L(r)$  and  $L(s)$  then

$$3.1 \rightarrow (r) | (s) \rightarrow L(r) \cup L(s)$$

$$3.2 \rightarrow (r).(s) \rightarrow L(r) \cdot L(s)$$

$$3.3 \rightarrow (r)^* \rightarrow (L(r))^*$$

$$3.4 \rightarrow (r) \rightarrow L(r)$$

→ A language denoted by regular expression is said to be a regular set (MCQ)

↳ Priority of operations:-

$$\text{eg. } \Sigma = \{a, b\}$$

R.E.

Sets.

$$1. \underline{a/b} \quad \{a, b\}$$

$$2. \underline{(a/b)} \underline{(a/b)}$$

$$\{\{a, a\}, \{a, b\}, \{b, a\}, \\ \{b, b\}\}$$

$$3. \underline{a^*} \quad \{\epsilon, a, aa, aaa, \dots, \infty\}$$

$$4. \underline{(a/b)^*}$$

$$\{\epsilon, a, b, aa, ab, bb, ba, \dots, \infty\}$$

5.  $a/a^*b$

$\{a, ab, aab, \dots\}$

### ↳ Algebraic properties of regular expressions :-

#### Axiom

#### Description

(Formation logic)

$$(1) rls = slr$$

| is cumulative

$$(2) r(slt) = rs(lt) | is associative$$

$$r(sl)tl = s(ltl) = sltl$$

$$(3) (rs)t = r(st) | concatenation is associative$$

$$(4) r(sl)t = rs|rt | concatenation$$

$$(sl)t = sr|tr | distributes over OR operators.$$

$$(5) \epsilon r = r$$

$\epsilon$  is identity

element for concatenation

$$(6) r^* = (r|\epsilon)^*$$

relation between \* &  $\epsilon$

$$(7) r^{**} = r^*$$

\* is idempotent

→ Regular Definition

For notational convenience we may wish to give names to regular expressions and to define regular expression using some names, if  $\Sigma$  is an alphabet of basic symbols, then a regular definition is a sequence of definitions of other form! [Exn 15-A]

$$d_1 \rightarrow r_1 \quad d \rightarrow \text{definition}$$

$$d_2 \rightarrow r_2 \quad \vdots \quad \epsilon \rightarrow \text{regular expressions}$$

$$d_n \rightarrow r_n$$

$$\text{eg. letter} \rightarrow A|B|C|\dots|z \quad a|b|c|\dots|z$$

$$\text{digit} \rightarrow 0|1|\dots|9$$

$$\text{id} \rightarrow \text{letter} (\text{letter}|\text{digit})^*$$

eg. Regular Definitions for unsigned number

$$\text{digit} \rightarrow 0|1|2|\dots|9$$

$$\text{digits} \rightarrow \text{digit digit}^*$$

$$\text{opt-frac} \rightarrow .\text{digits}|\epsilon$$

$$\text{opt-Expo} \rightarrow C|E (+|-|\cdot|E) \text{ digits}|E$$

$$\Rightarrow \text{num} \rightarrow \text{digits opt-frac opt-Expo}$$

⇒ Notational Shorthands.

1. One or more instance (+)

Regular Expression ( $\sigma$ )

$$\sigma^* = (\sigma^* | \epsilon)$$

$$\sigma^+ = \sigma \sigma^*$$

2. Zero or one instance (?)

(C(C+1-)? digits) ?

3. character classes

Regular Expression [a]

[abc] any character

e.g. 1 (1|d)\*

[A-Z a-z] [A-Z] [a-z] [0-9]\*

→ Language Denoted by Regular Expression  
are called Regular sets.

1st - {palindrome} 2nd - {odd length string}

3rd - {string starting with 010}

4th - {string ending with 010}

5th - {string containing 010}

6th - {string containing 0101}

7th - {string containing 01010}

8th - {string containing 010101}

9th - {string containing 0101010}

10th - {string containing 01010101}

## → Recognition of Tokens :-

Consider the following grammar format :-

stmt → if expr then stmt

  | → if expr then stmt else stmt

expr → term [relOp] term

  | → term [at] ↑ relational operator.

term → id | num |

  | → num

[CFG] → context free grammar

⇒ stmt, expr → Non Terminal

⇒ if, then → Terminals / Direct

⇒ stmt → if expr then stmt ⇒ Production

Label in L.H.S. called start symbol.

Only Non-Terminals can act as label or start symbol.

→ Name given to regular expression called regular definition.

→ \* → closure function

Indicates 0 or more occurrences.

? → 1 or more occurrences.

→ We cannot have Regular Definition for Non-Terminals.

→ Regular Expression Patterns for tokens.

Regular Expression	Related Tokens	Attribute Values
whitespace	-	
if	if	match max 79
then	then	part 100% if
else	else	part - else
id	id	ptr to table entry
num	num	ptr to table entry
<	relap	LT
<=	relap	LE
=	relap	EQ
>	relap	NE
>=	relap	GT
<=	relap	GE

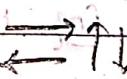
### \* Transition Diagram

- It can be defined as stylish flowchart

Signs :-

(1)  → it is to represent state.

(2)  → Final state (Accepting state)

(3)  → Directed Edges containing some label about workflow.

(4) \* → represent initial state.

(5) other → unknown signs.

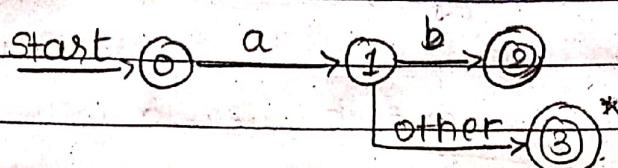
eg.

Regular expressions are converted into  
diagrams.

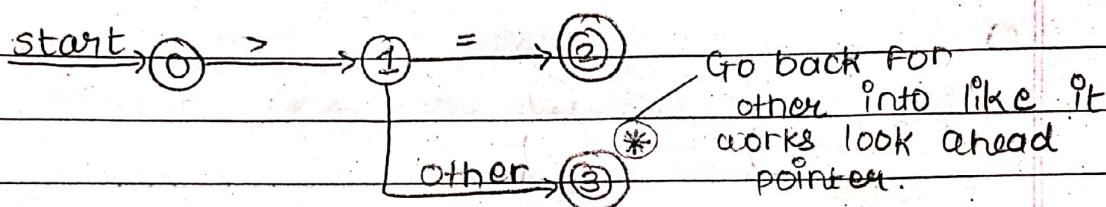
Regular Expressions a label.



(2) R.E.



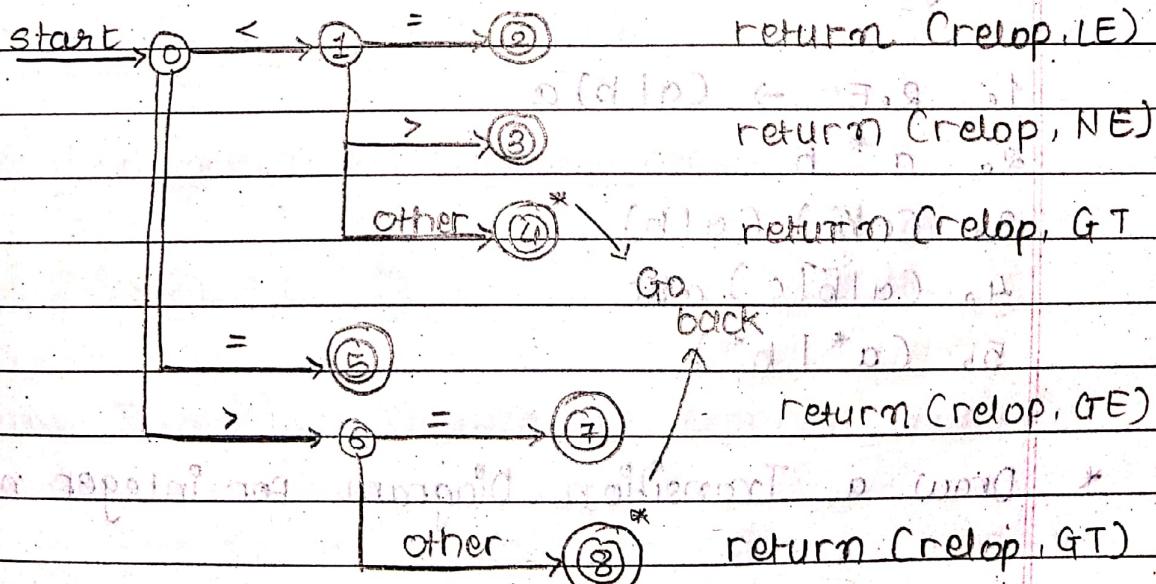
Transition Diagram for relop ( $>=$ ).



→ Transition Diagram for relational operators.

$<$ ,  $<=$ ,  $>$ ,  $>=$ ,  $=$ ,  $<>$

$<$ ,  $<=$ ,  $>$ ,  $>=$ ,  $=$ ,  $<>$

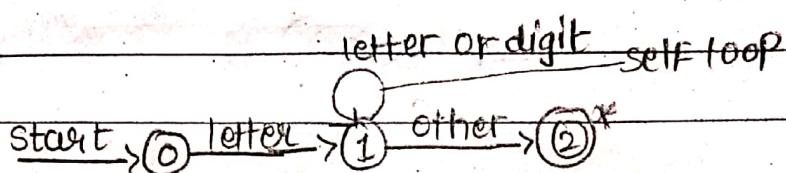


- \* Draw a Transition Diagram for valid identifier.

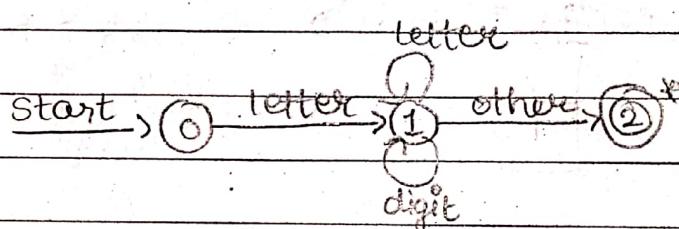
Regular Expression is given as below:-

Letter (Letter | Digit)\*

(1)



(2)



- \* Draw a Transition Diagrams for following regular expressions:-

1. R.E  $\rightarrow (a|b)a$

2.  $a^* b$

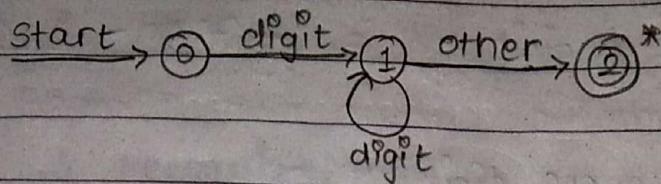
3.  $(a|b)(a|b)$

4.  $(a|b|c)abc$

5.  $(a^* | b^*)$

- \* Draw a Transition Diagram for integer number.

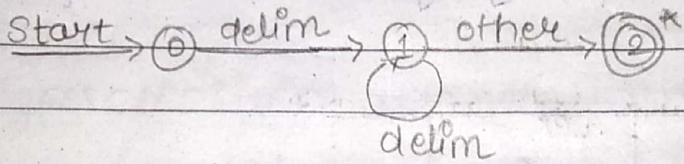
digit\*      or      digit.digit\*



\* Draw a Transition Diagram to recognize [white space. {Blank, Tab, newline}]

$\text{ws}^+$  or  $\text{delim}^+$

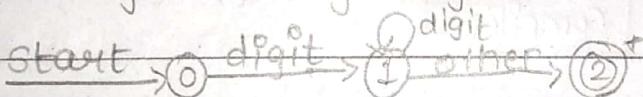
$\text{delim}^+$  or  $\text{delim} \cdot \text{delim}^*$



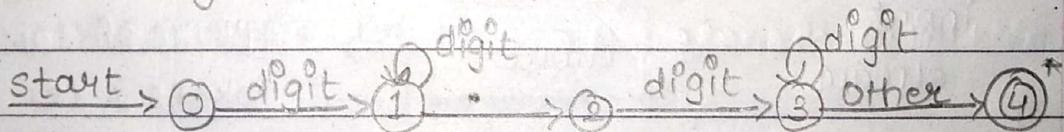
\* Draw Transition Diagram for Real Numbers

e.g. 10.25, 27.1, 123.789

$\text{digit}^+ = \text{digit} \text{digit}^* \leftarrow \text{closure}$



Final Diagram for Real Numbers.

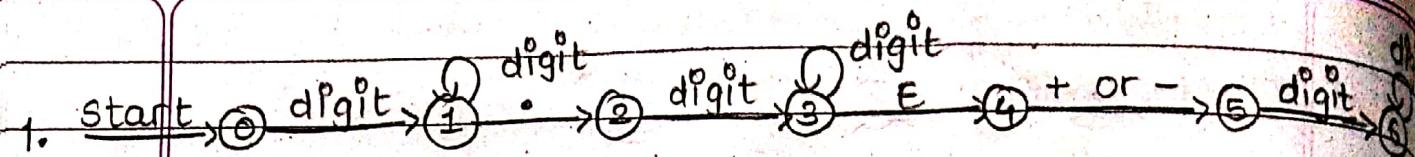


→ Draw Transition Diagram for exponent numeric value.

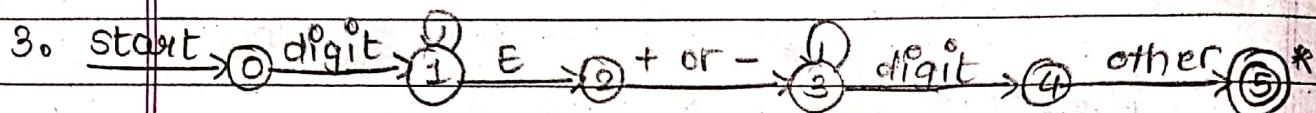
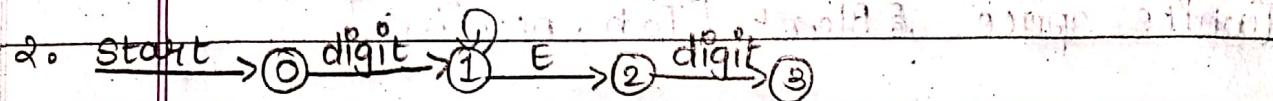
Eg. 23.25 E + 5, 125.1 E - 20, 78E30, 71E + 23

93 E - 235

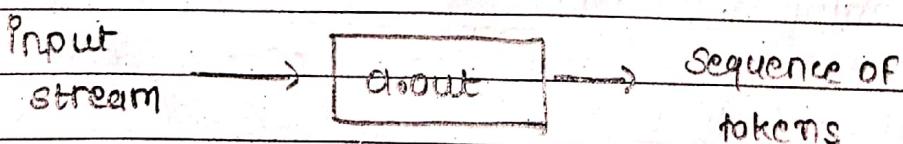
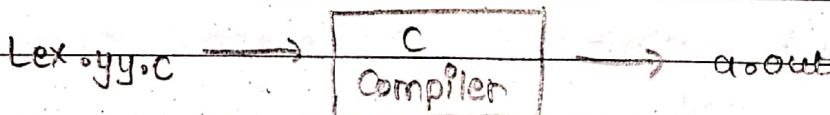
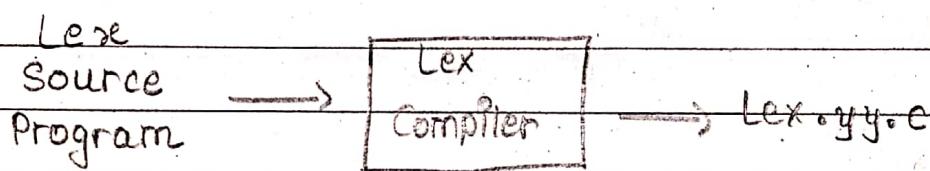
C. ⇒ period in real no.)



[Self loop for more digits].



\* A language for specifying lexical analyzer.



A lex program consists following 3 parts:

1. Declaration

2. Translation Rules

3. Procedures

## → Finite Automata.

A recognizer for a language is a program that takes as input string  $x$  and answers yes if  $x$  is a sentence of the language and no otherwise.

Compile a Regular Expression into a recognizer by constructing a generalized transition diagram called Finite automaton.

A finite automaton can be deterministic or non-deterministic.

### \* Types of Finite Automata

(1) Non-Deterministic Finite Automata (NFA).

(2) Deterministic Finite Automata (DFA).

NFA :- Non-Deterministic Finite Automata.

Non-Deterministic Finite Automata is a mathematical model consisting of

- (1) A set of states  $S$ .
- (2) A set of input symbols  $\Sigma$  (sigma)
- (3) A transition function  $\delta$  that maps state symbol pairs to sets of states
- (4) A state  $s_0$ , that is distinguished as the start or initial state.
- (5) A set of states  $F$  distinguished as accepting or final states.

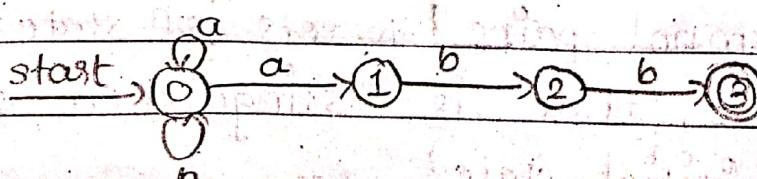
## NFA

## DFA

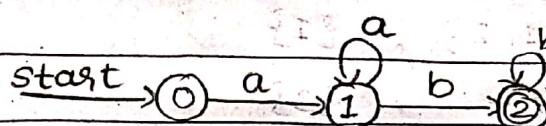
- |   |  |
|---|--|
| 1. Input label $\epsilon$ (sign) is allowed.                              | It is not allowed at all.  |
| 2. NFA is a kind of mathematical model.                                   | It is a special kind of NFA/extension of NFA.                          |
| 3. Edges living from a particular state may have identical input symbols. | edges living from particular state should have distinct input symbols. |
| 4. NFA diagrams are smaller than DFA diagrams.                            | than NFA diagrams.   |
| 5. Backtracking is possible.  | There is no chance of backtracking.                                    |
| 6. NFA works less efficiently than DFA.                                   | NFA.   |

→ Draw NFA diagram for following regular expression and also prepare Transition Table for same.

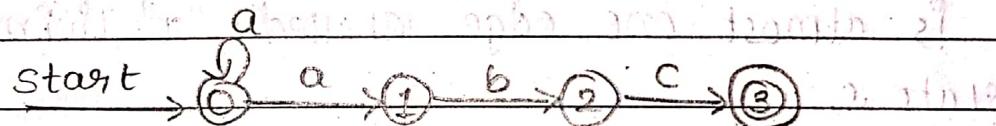
1.  $(a|b)^*abb$



State	Input Symbols	
	a	b
0	{0, 1}	{0}
1	-	{2}
2	-	{3}
3	-	-

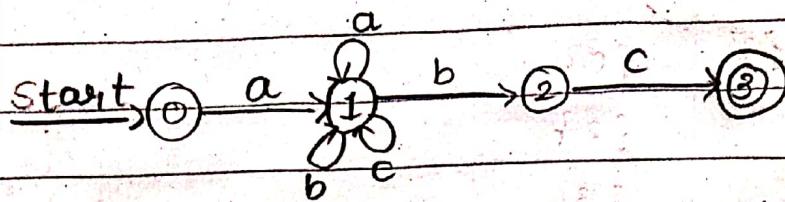
Q.  $aa^*lbb^*$ 

State	Input Symbols	
	a	b
0	{1}	{0} -
1	{1}	-
2	-	{2}

Q.  $a^*abc$ 

State	Input Symbols		
	a	b	c
0	{0, 1}	-	-
1	-	{2}	-
2	-	-	{3}
3	-	-	-

4.  $a(a|b|c)^*bc$



State	Input symbols		
	a	b	c
0	{1}	-	-
1	{1}	{1, 2}	{1}
2	-	{2}	{3}
3	-	-	-

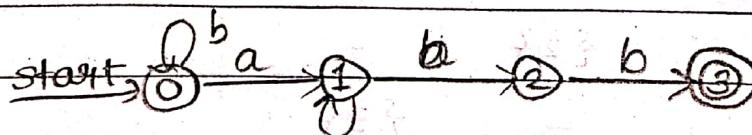
DFA :- Deterministic Finite Automata.

Deterministic Finite Automata, they are special case of NFA in which

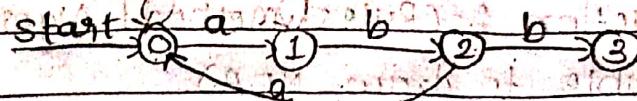
- (1) No state has an Epsilon ( $\epsilon$ ) transition
- (2) For each state  $s$  and input symbol  $A$  there is atmost one edge labelled "a" leaving state  $s$ .

→ Draw DFA diagram for following Regular Expression.

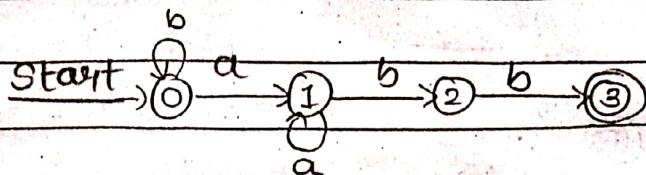
$(a|b)^*abb$



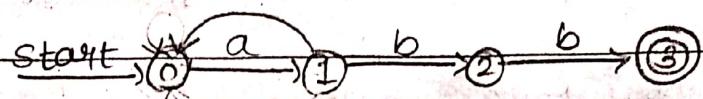
$\rightarrow ababb$



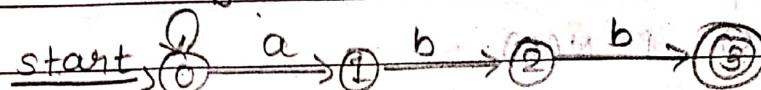
$\rightarrow baabb$



$\rightarrow aaabb$



$\rightarrow bbabb$



### Algorithm Simulating DFA

$s_0 = s_0$ ; initial state

$c_0 = \text{nextchar}$

while  $c \neq \text{eof}$  do

$s_0 = \text{move}(s_0, c_0)$

$c_0 = \text{nextchar}$

end;

if  $s_0$  is in  $E$  then

  return "yes";

else

  return "no";

→ Thompson's Construction :- (Rules / Method / Algorithm)  
[Only Applicable to draw NFA].

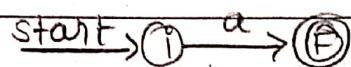
1. Regular Expression =  $\epsilon$

NFA



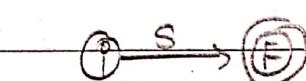
2. Regular Expression =  $a$

NFA



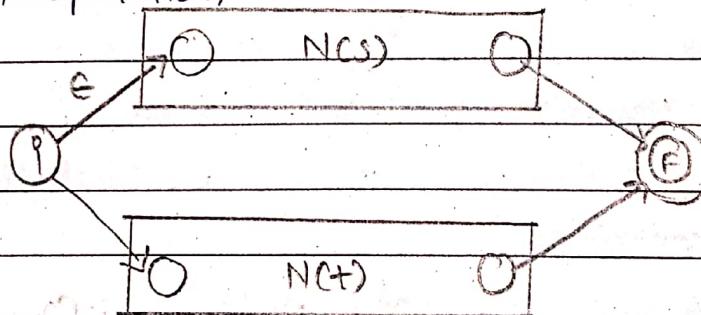
3. Regular Expression  $s \& t$

NFA :  $N(s) \& N(t)$

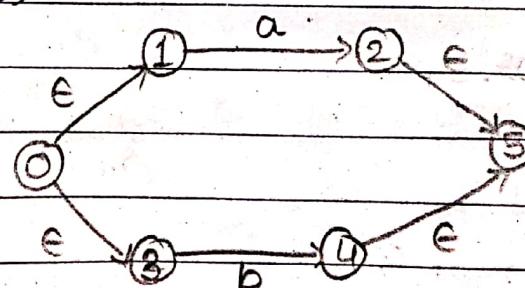


Union.

Union Operation structure from thompson's rule.



Ex:  $(a|b)$



only the way for  
one time  
occurrence.

Note :- For union Diamond-Shape is compulsory, and also having six states including initial state and final state.

4. S & T

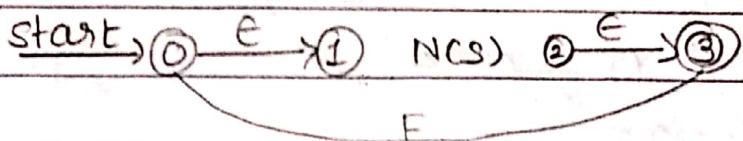
N(S) N(T)

Ex: ab



5. Closure

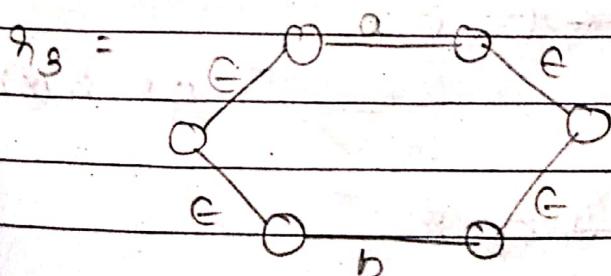
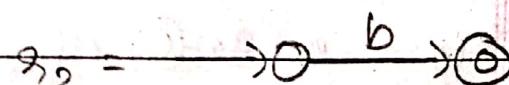
Consider Regular Expression, N(S)



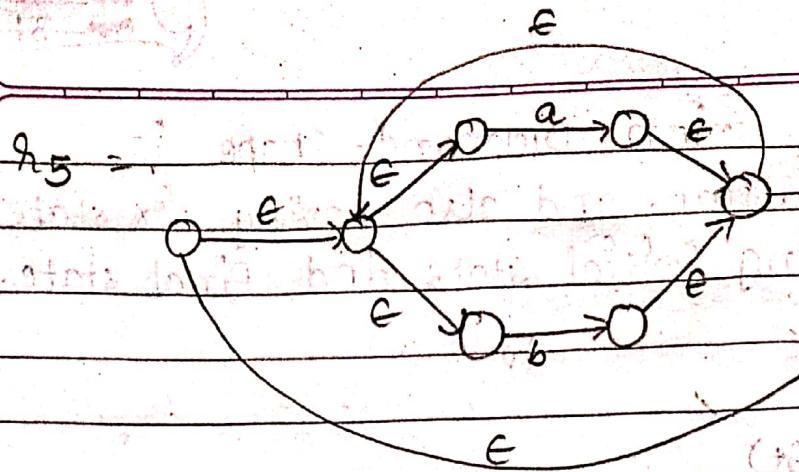
→ As per Thompson construction we cannot have self-loop.

Q. Draw NFA diagram for the following given Regular Expression using Thompson construction Rules. Also provide decomposition structure for the same.

→ (a|b)\*abb.



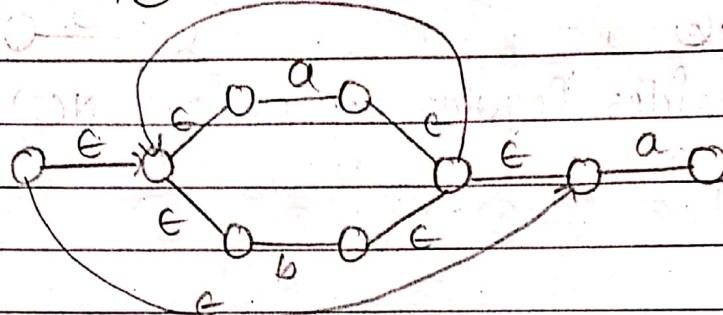
$$g_4 = (g_3)$$



$q_6 =$



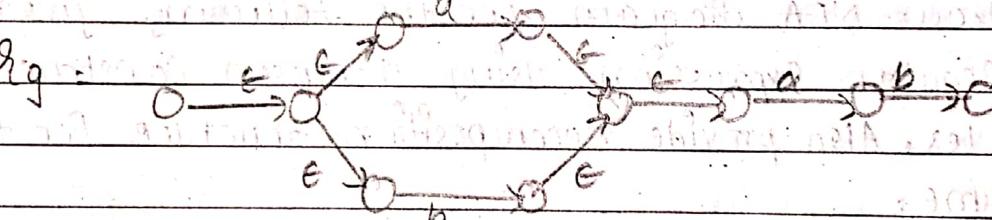
$q_7 =$



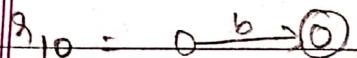
$q_8 =$



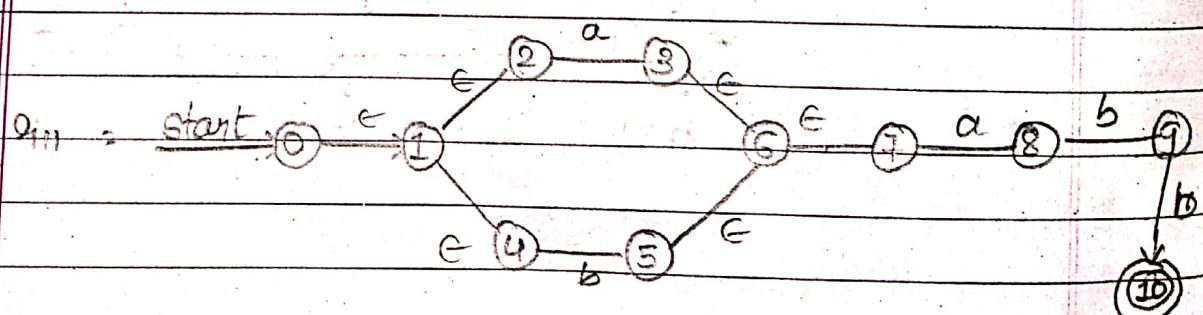
$q_9 =$

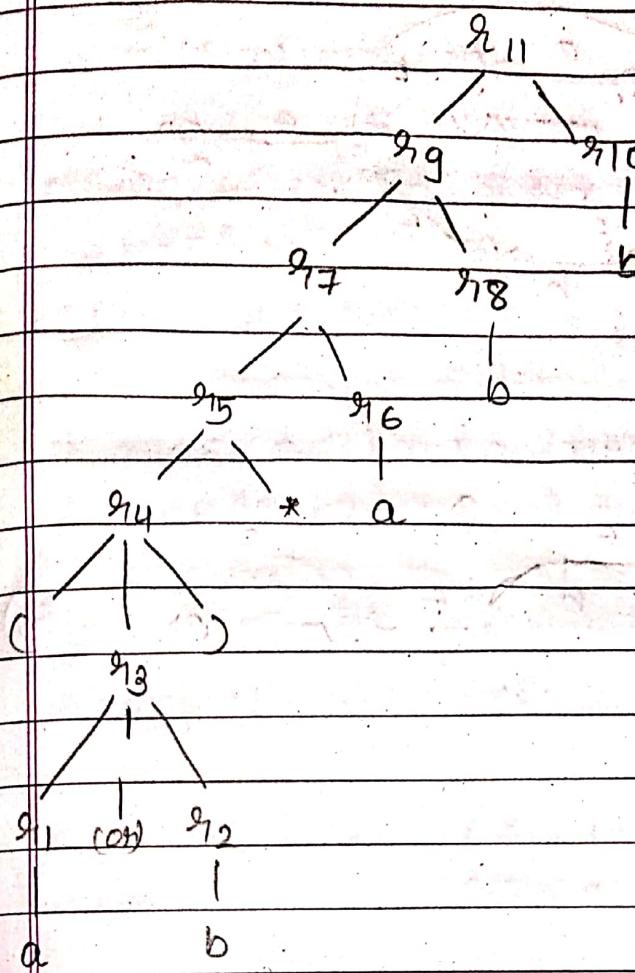


$q_{10} =$



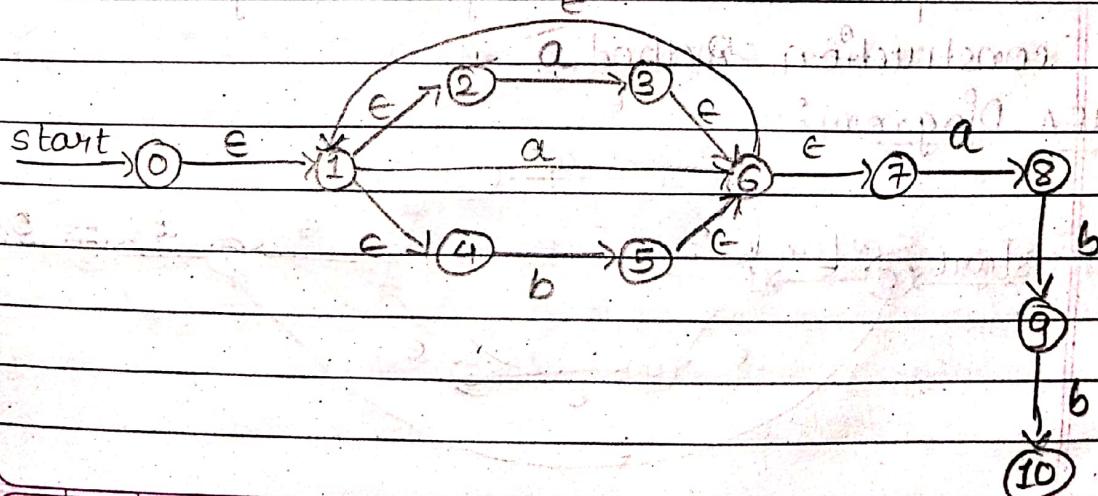
$q_{11} =$



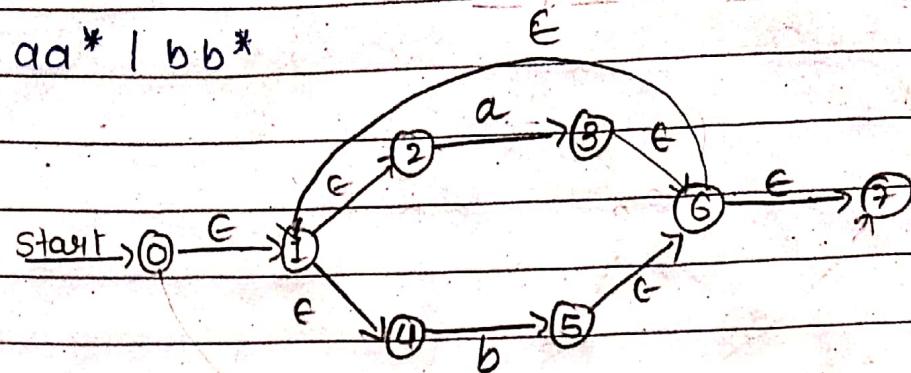
Decomposition Diagram

Q. Draw NFA diagram using Thompson's construction rules for the following Regular Expression.

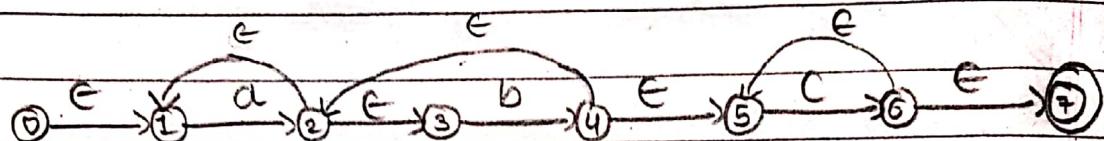
1.  $a(a|b)^*abb$



2.  $aa^* \mid bb^*$



3.  $a^* ab^* bc^* c$



→ Subset Construction.

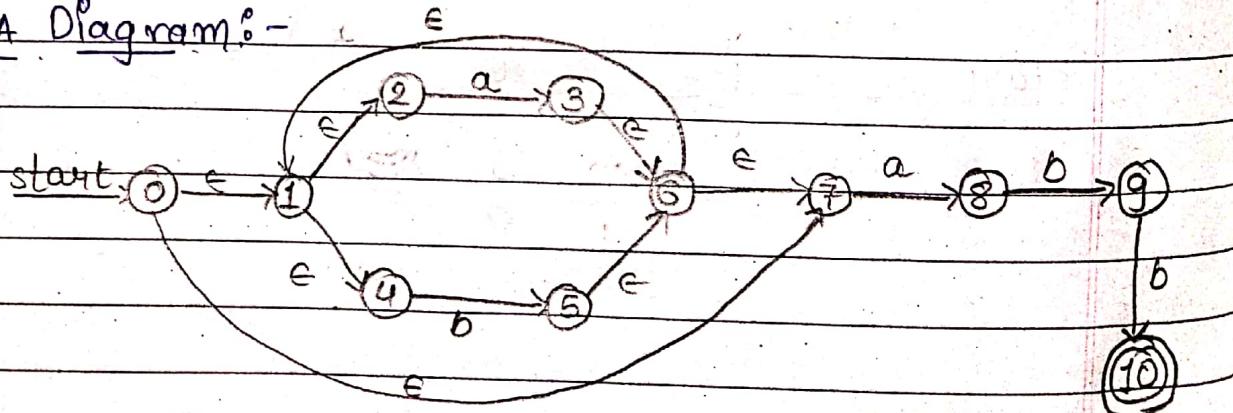
1.  $\epsilon$  - closure (S)

2.  $\epsilon$  - closure (T)

3. more (T, a)

Q. Draw NFA diagram using Thompson's Construction rules for the regular expression. Convert it into its equivalent DFA diagram using subset construction method.

NFA Diagram:-



$\epsilon$ -closure (o)

$$A = \{0, 1, 7, 2, 4\}$$

$$\text{mov}(A, a) = \{3, 8\}$$

$$\epsilon(\text{mov}(A, a)) = \{3, 8, 6, 7, 1, 2, 4\} = B$$

$$A = \{5\}$$

$$\epsilon(\text{mov}(A, b)) = \{5, 6, 7, 1, 2, 4\} = C$$

$$\epsilon(\text{mov}(B, a)) = \{3, 8\}$$

$$\epsilon(\text{mov}(B, b)) = \{5, 9, 6, 7, 1, 2, 4\} = D$$

$$\epsilon(\text{mov}(C, a)) = \{3, 8\} = \{3, 8, 6, 7, 1, 2, 4\} = B$$

$$\epsilon(\text{mov}(C, b)) = \{5\}$$

$$\epsilon(\text{mov}(D, a)) = \{3, 8\} = B$$

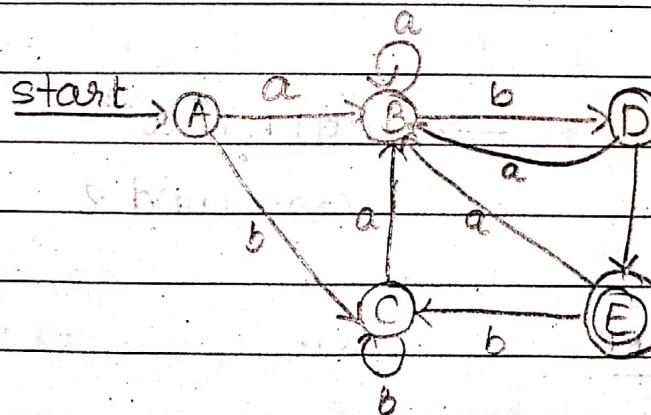
$$\epsilon(\text{mov}(D, b)) = \{10, 5\} = E$$

$$E = \{5, 10, 6, 1, 7, 2, 4\}$$

$$\epsilon(\text{mov}(E, a)) = \{8, 3\} = B$$

$$\epsilon(\text{mov}(E, b)) = \{5\} = C$$

DFA :-



## Algorithm

Initially,  $\epsilon$ -closure( $S_0$ ) is the only state in Dstates and it is unmarked;

while there is an unmarked state  $T$  in Dstates  
do begin

mark  $T$ ;

for each input symbol  $a$  do begin

$U = \epsilon$ -closure(move( $T, a$ ));

if  $U$  is not in Dstates then

add  $U$  as an unmarked state to Dstates

$D_{tran}[T, a] = U$

end

end

## - Regular expression to DFA :-

### Conversion of Regular expression to DFA

Cast without using NFA

→ R.E.  $(a|b)^* a$

$(a|b)^* a$

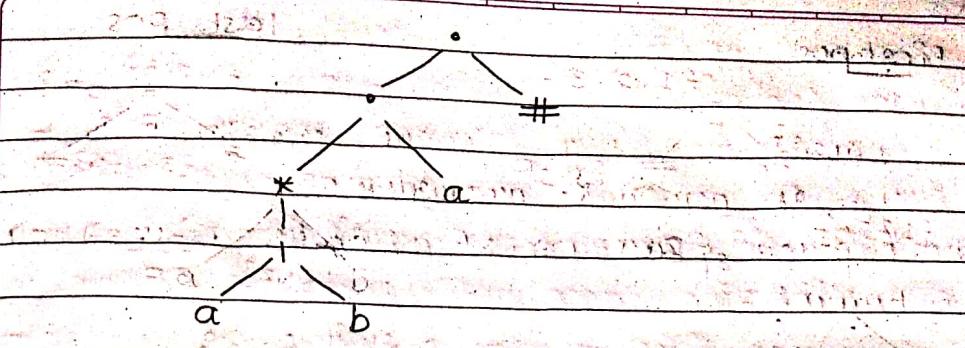
compound 1

$\rightarrow (a|b)^* a^{\#}$

compound 2

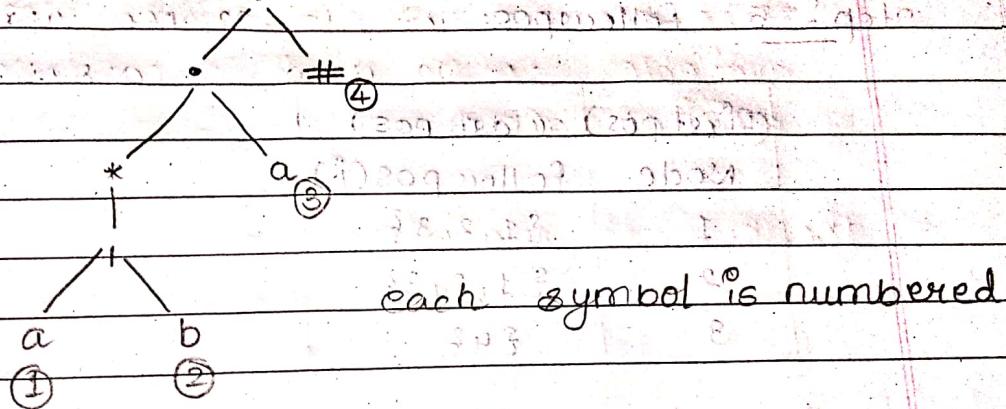
Augmented

## step 2) Syntax Tree

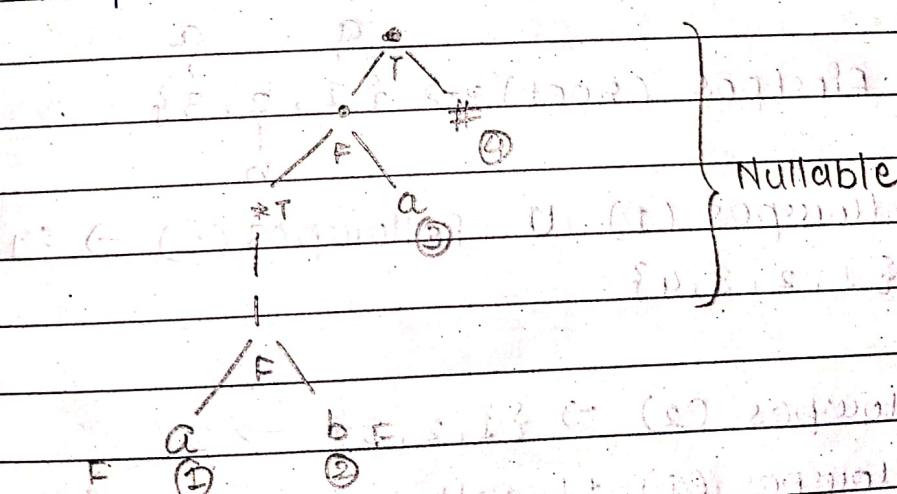


step(3) Numbering of Nodes. (except  $\epsilon$ )

Dstates:

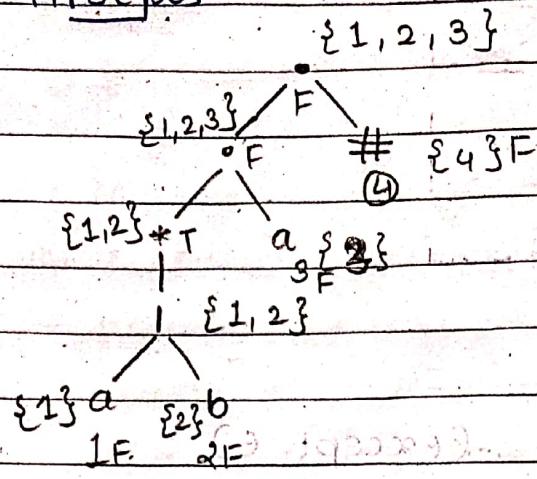


Step(4) Firstpos, \$astpos, nullable

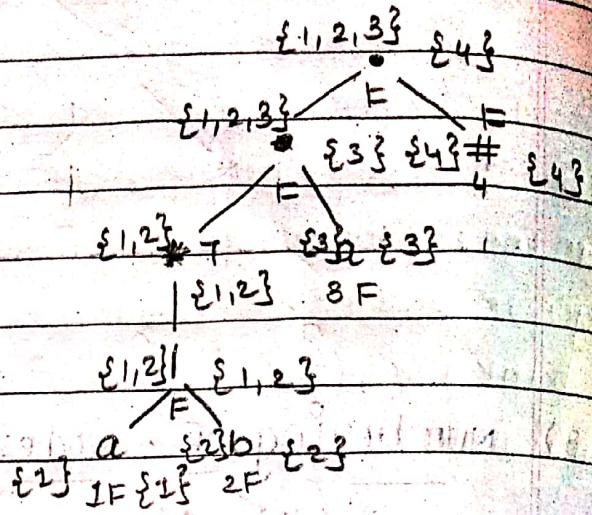


Firstpos, \$astpos

Firstpos



Last pos



Step 5 Followpos

(Firstpos)	(Lastpos)
Node	Followpos( $i^*$ )
1	{1, 2, 3}
2	{1, 2, 3}
3	{4}

Step 6 State Generation

$$S_1 = \text{Firstpos}(\text{root}) \Rightarrow \{ \overset{a}{1}, \overset{a}{2}, \overset{b}{3} \}$$

$$- a : \text{followpos}(1) \cup \text{followpos}(2) \rightarrow \{1, 2, 3\}$$

$$S_1 = \{1, 2, 3, 4\}$$

$$- b : \text{followpos}(2) \rightarrow \{1, 2, 3\} \rightarrow S_1$$

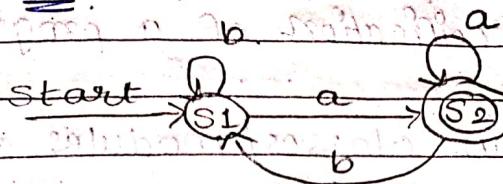
$$\left. \begin{array}{l} - a : \text{followpos}(1) \cup \text{followpos}(3) \rightarrow \{1, 2, 3\} \\ - b : \text{followpos}(2) \rightarrow \{1, 2, 3\} \end{array} \right\} \rightarrow S_1$$

Step 7 T.O.T. (Transition Table)

state	Input symbol	
	a	b
s <sub>1</sub>	s <sub>2</sub>	s <sub>1</sub>
s <sub>2</sub>	s <sub>1</sub>	s <sub>1</sub>

step 8-8

DFA :-



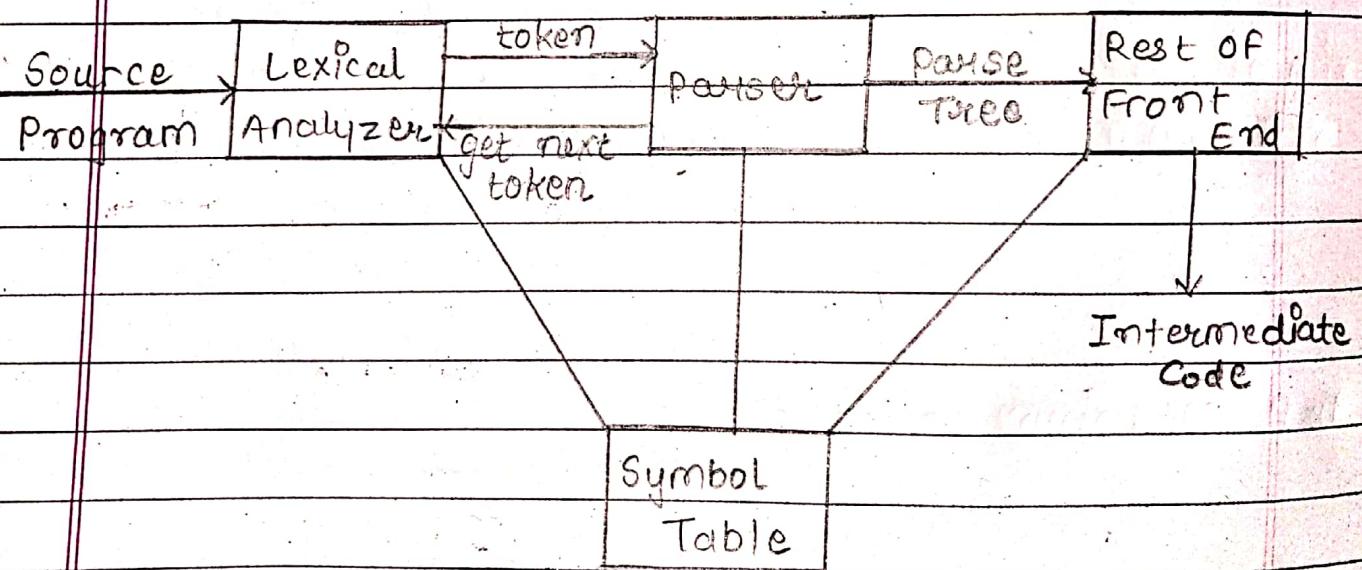
## Syntax Analysis

Every programming language has rule that prescribe the syntactic structure of well formed programs.

Following are some advantages of syntax analysis:

- (1) Grammar gives a precise, easy to understand and syntactic specification of a programming language.
- (2) There are certain classes/modules of grammar available can be used to construct an efficient parser.

=> The Role of the Parser:-



The parser obtains a string of tokens from the lexical analyzer and verifies that the string can be generated by the grammar for the source language.

It is also to get outcome in the form of syntax errors.

### Types of Parser:

1. Universal Parsing Method
2. Top-down Parser
3. Bottom-up Parser

### Types of error

1. Lexical
2. Syntactic
3. Semantic
4. Logical

### - Goals of Error Handler

1. It should report the presence of errors clearly and accurately.
2. It should recover from each error quickly enough to be able to detect subsequent errors.
3. It should not significantly slow down the processing of correct programs.

### - Error Recovery strategies

1. Panic Mode

2. Phrase Level Recovery

3. Error Productions

4. Global Corrections

## - Context-Free Grammar [CFG]

Many programming languages constructs have an inherently recursive structure that can be defined by context free grammar.

e.g. if  $s_1$  and  $s_2$  are statements and  $E$  is expressions

$\Rightarrow$  if  $E$  then  $s_1$  else  $s_2$ .

CFG:  $\text{stmt} \rightarrow \frac{\text{if expr then stmt}}{\overline{T} \quad \overline{NT}} \frac{\text{else stmt}}{\overline{T} \quad \overline{NT}}$

## - Components of CFG

### 1. Terminal

Terminals are the basic symbols from which strings are formed.

The word Token is a synonym for Terminal when we are talking about grammar for programming language.

### 2. Non-Terminal

They are syntactic variables that denote sets of strings.

### 3. Start-Symbol

One Non-Terminal is distinguished as the start symbol and the set of strings it denotes is the language defined by the grammar.

[Always find on left hand side of Regular Expression]

## 4. Production.

The production of grammar specify the manner in which the terminals & non-terminal can be combined to form strings.

Each production consist of a non-terminal followed by an arrow, followed by string of non-terminals and terminals.

Example:-  $S \rightarrow aSb \cup bSb \cup ab$

### Notational Conventions

#### 1. Terminals

- small case letters [a, b, c, ...]
- digits
- operators
- punctuation marks / separators
- bold face strings

#### 2. Non-Terminals

- uppercase letter [A, B, C, ...]
- lowercase italic string
- S

#### 3. Uppercase characters like U, V, W, X, Y, Z [late in series and can act as terminals and non-terminals]

#### 4. Lowercase characters like u, v, w, x, y, z [late in series and can act only as terminals]

5. Greek letters [Arbitrary strings]

$\alpha, \beta, \gamma$  - here, T/NT/TNT production rule

6.  $A \Rightarrow \alpha_1, A \Rightarrow \alpha_2, A \Rightarrow \alpha_3, \dots, A \Rightarrow \alpha_n$

7. whatever symbol available on LHS of arrow sign represents start symbols.

### \* Derivations [providing alternative information]

Here, production is treated as rewriting rule in which the non-terminal on the left is replaced by the string on the right side of the production.

e.g.  $E \Rightarrow -E$

$E$  "derives"  $-E$

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(\text{id})$

→ In abstract string  $\alpha A \beta \Rightarrow \alpha \gamma \beta$

if  $A \Rightarrow \gamma$  is a production

if  $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \alpha_3 \dots \Rightarrow \alpha_n$

we say  $\alpha_1 \Rightarrow \alpha_n$

1.  $\Rightarrow$  derives in one step

2.  $\xrightarrow{*}$  derives in 0 or more steps

3.  $\xrightarrow{+}$  derives in 1 or more steps

### Purpose of derivations

1.  $\alpha \xrightarrow{*} \alpha$  for any string  $\alpha$ .

2. if  $\alpha \xrightarrow{*} \beta$  and  $\beta \Rightarrow \gamma$  then  $\alpha \xrightarrow{*} \gamma$

→ Define sentence.

$S \Rightarrow w$  must be terminals.

→ Define Sentential Form

$S \Rightarrow x$  can be T / NT / TNT Both

A sentence is a sentential form with no non-terminals.

Derivation

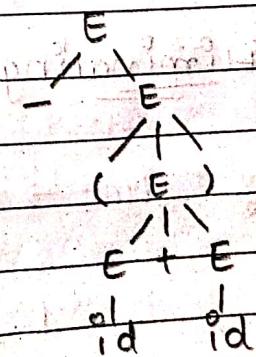
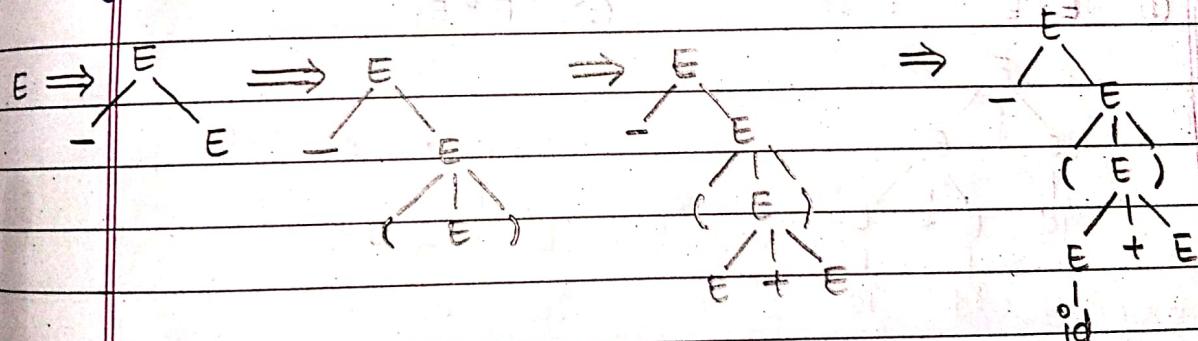
e.g.  $E \rightarrow E + E$  Input :- - (id + id)

$E \rightarrow E * E$

$E \rightarrow (E) \quad E \xrightarrow{lm} - E \xrightarrow{lm} - (E) \xrightarrow{lm} - (E + E)$

$E \rightarrow id \quad - \xrightarrow{lm} - (id + id) \xleftarrow{lm} - (id + E)$

Syntax Tree



Input :-  $id + id * id$

Above input may have double or two possibilities

$$E \Rightarrow E + E$$

$$E \Rightarrow E * E$$

$$\Rightarrow id + E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

$$\Rightarrow id + id * id$$

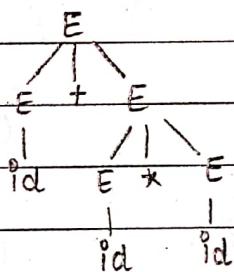
Ambiguity :- A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

OR

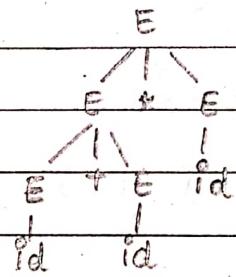
An ambiguous grammar is one that produces more than one left most or more than one right most derivation for the same sentence.

Ex :-

①  $E + E$



②  $E * E$



Ambiguity is not efficient in compilation process.

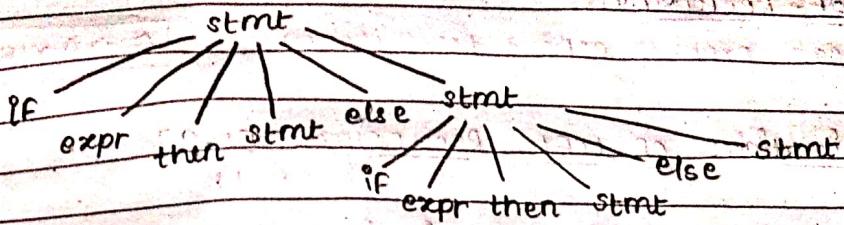
### \* Eliminating Ambiguity

e.g.  $stmt \rightarrow if\ expr\ then\ stmt$

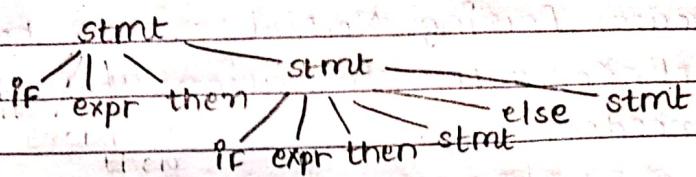
|  $if\ expr\ then\ stmt\ else\ stmt$

| other

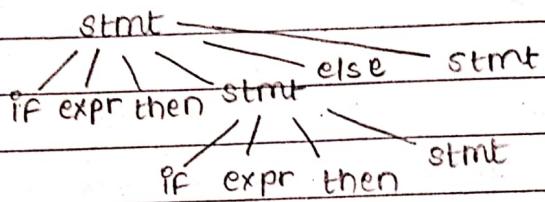
(i)



(ii)



(iii)



→ stmt → matched-stmt  
 | unmatched-stmt

matched-stmt → if expr then matched-stmt else  
 matched-stmt

| other

unmatched-stmt → if expr then stmt  
 | if expr then matched-stmt  
 | else unmatched-stmt

Left Recursion.eg.  $\text{expr} \rightarrow \text{expr} + \text{term}$ 

$$A \rightarrow A\alpha | B$$

\* Elimination of Left Recursion [L.R.]

A Grammar is L.R. if it has a non-terminal  $A$  such that there is a derivation  $A \xrightarrow{*} A\alpha$  for some string  $\alpha$ .

Top-Down Parsing Method cannot handle L.R Grammars. So, a transformation that eliminates L.R. is needed.

eg.  $E \rightarrow E + T | T$        $A \rightarrow A\alpha | B$   
 $T \rightarrow T * F | F$        $A \rightarrow A\alpha | B$   
 $E \rightarrow (E) | id$

(i)  $E \rightarrow TE'$

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

(ii)  $T \rightarrow FT'$

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

(iii)  $F \rightarrow (E)$

$$F \rightarrow id$$

$$\Rightarrow A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | B_1 | B_2 | \dots | B_n$$

$$A \rightarrow B_1 A' | B_2 A' | \dots | B_n A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | E$$

### Algo. Eliminating Left Recursion

1. Arrange the non-terminals in some order  
 $A_1, A_2, \dots, A_n$
2. For  $i := 1$  to  $n$  do begin
  - For  $j := 1$  to  $i-1$  do begin
 

replace each production of the form  $A_i \rightarrow A_j Y$  by the productions

$$A_i \rightarrow S_1 Y | S_2 Y | \dots | S_k Y$$

where  $A_i \rightarrow S_1 | S_2 | \dots | S_k$  are all the current  $A_j$ -productions; end

eliminate the immediate left recursion among the  $A_i$ -productions

end

→ eliminate L.R. 0 in the following given grammar.

$$S \rightarrow Aa/b \quad -①$$

$$A \rightarrow Aa/Sd/e \quad -②$$

**[1]**  $S \rightarrow Aa/b$  using eq. -②.

$$S \rightarrow Sd/a/b \quad -③$$

**[2]**  $A \rightarrow Aa/Sd/e$  using eq. -①

$$A \rightarrow Aa/Aa/bd/e \quad -④$$

$$S \rightarrow Sdab - ③$$

$$A \rightarrow Ac | Aalbd | G - ④$$

$$A \rightarrow A\alpha | \beta - \star \Rightarrow A \Rightarrow BA^! \\ \Rightarrow A^! \Rightarrow \alpha A^! e$$

→ Left Factoring:

It is a grammar transformation that is used for producing a grammar suitable for predictive parsing.

e.g.  $\text{stmt} \rightarrow \text{if expr then stmt else stmt}$   
 $\quad \quad \quad | \text{if expr then stmt}$

Algorithm Left Factoring.

I/P : G (Grammar)

For each N.T.  $A$  find the longest prefix  $\alpha$  common to two or more of its alternative.

If  $\alpha \neq \epsilon$ , that is, there is a non-trivial common prefix replace all the capital  $A$  products

$A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \alpha \beta_n | \gamma$ , where  $\gamma$  represents all alternatives that do not begin with  $\alpha$  by

$$A \rightarrow \alpha A' | \gamma - ①$$

$$A' \rightarrow \beta_1 | \beta_2 | \dots | \beta_n - ②$$

Here,  $A'$  is a new N.T. Repeatedly apply this transformation until no two alternatives for a N.T. have a common prefix.

$$\rightarrow S \rightarrow iEtS | iEtSeS | a - ①$$
$$E \rightarrow b - ②$$

For  $A \rightarrow \alpha A'$

$$A \rightarrow iEtSA' | a$$
$$A' \rightarrow e | es$$
$$E \rightarrow b$$
$$A \rightarrow \alpha B_1 | \beta B_2$$
$$\Downarrow$$
$$A \rightarrow \alpha A' | \gamma$$
$$A' \rightarrow \beta_1 | \beta_2$$

e

ctions