

→ Principles of compiler design by Aho & Ullman → Vauosa publication.

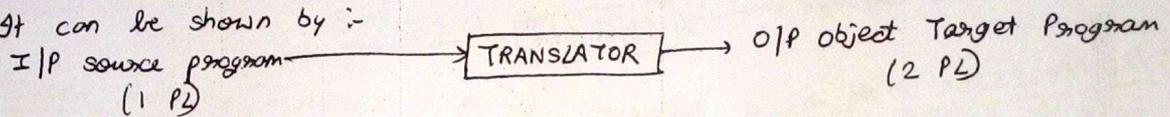
### INTRODUCTION TO COMPILER :-

- TRANSLATOR :- types

- compiler
- interpreter
- assembler
- preprocessor.

### - STRUCTURE OF COMPILER

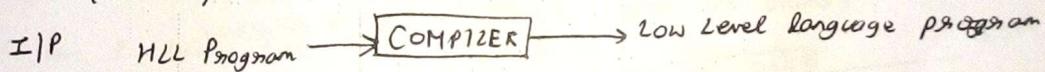
TRANSLATOR :- It converts 1 language to other language. It is a program that takes a program written in 1 programming language as i/p & produces an o/p in a program written in another language. Input program is called SOURCE PROGRAM & output program is called OBJECT or target program. It can be shown by :-



Translator is a generic (general) term.

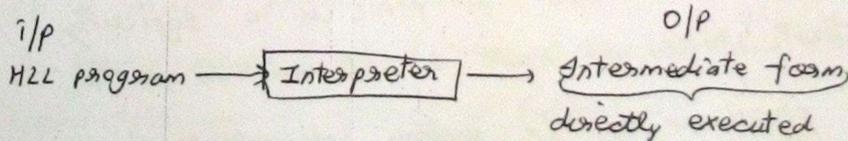
### DIFFERENT TYPES OF TRANSLATOR :-

1) COMPILER :- It is a specific translator which takes a source program written in HLL (C, Fortran) & produces as o/p a program written LLL (assembly or machine) - o/p



Source program is first compiled & converted into object program. This object program is linked with different libraries (LINKER) & then loaded in memory (LOADER) for the execution.

2) INTERPRETER :- It is a translator that transfers a source program written in HLL into an intermediate form (post-fix notation or 3-address code) which can be directly executed by the machines.



## Diff b/w INTERPRETER & COMPILER :-

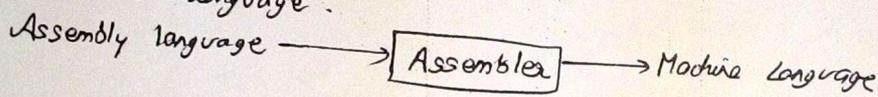
### INTERPRETER

- 1) HLL  $\rightarrow$  Intermediate
- 2) It performs line by line translation of code
- 3) Execution time is slow
- 4) Smaller in size ~~time~~
- 5) Better for implementation of complex programming language constructs

### COMPILER

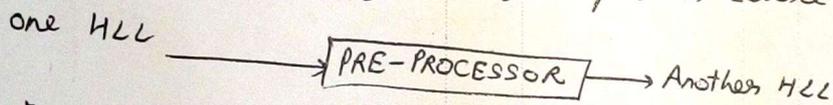
- 1) HLL  $\rightarrow$  LLL
- 2) It takes complete program in 1 go.
- 3) It is comparatively fast, bigger in size.
- 4) Bigger in size
- 5) It is not better than interpreter in case of complex constructs

3) ASSEMBLER :- It is a translator which converts assembly language to machine language.



4) PRE PROCESSOR :- It is a translator that takes a HLL program as input & produces a program in another HLL (changes the representation of code). It performs :-

- $\rightarrow$  include header files with source program.
- $\rightarrow$  copies values of macros in the source program.
- $\rightarrow$  remove the comments
- $\rightarrow$  removes the extra space in source program.

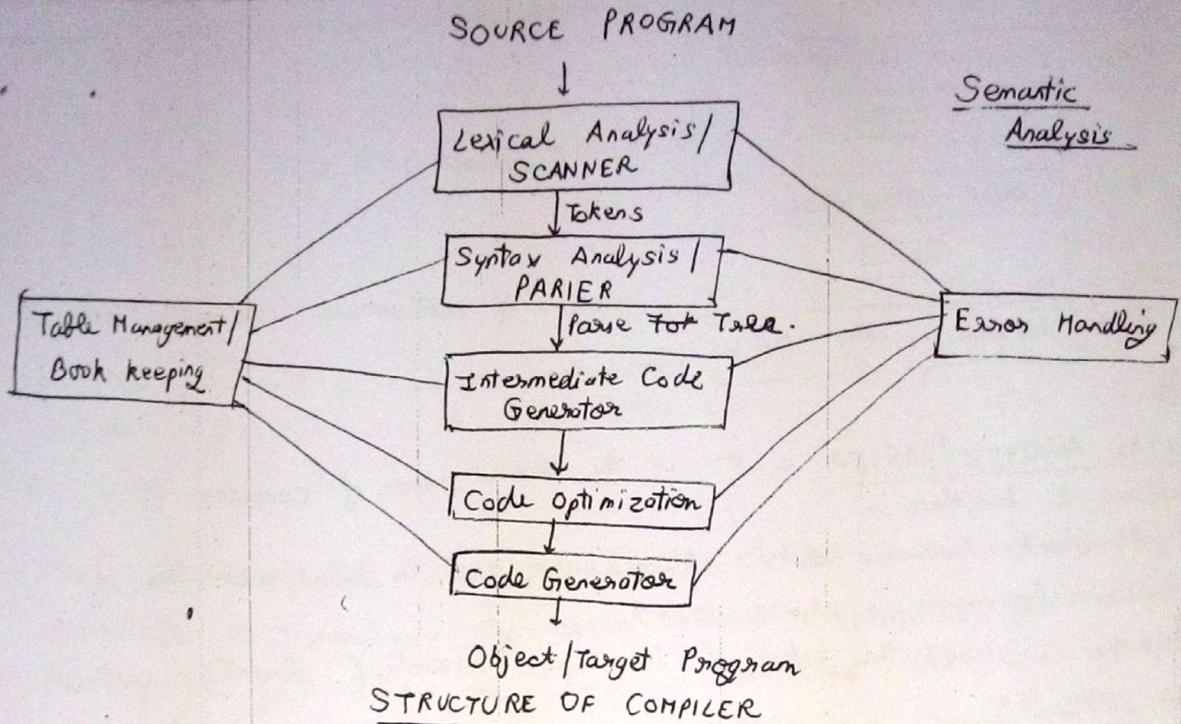


Ex: Include  $\rightarrow$  preprocessor directive

11/2

## STRUCTURE OF COMPILER

The process of compilation is a complex one. It is not possible to complete it in a single step. So it is divided into a series of sub-processes called as PHASES. A phase is a logically cohesive operation that takes as input 1 representation of the source program & produces as o/p another representation of it.



1) LEXICAL ANALYSIS / SCANNER :- It is the first phase of compiler. It reads the source program character by character & groups them into tokens. A token is a sequence of characters that is meaningful & cannot be further broken down. There are generally 5 types of tokens - Keywords, Identifier, Constant, Operator & special symbols.

Ex - There is a statement in FORTRAN language

EQ  
if (5 EQ MAX) GOTO 100

These are 8 tokens here if, (, 5, EQ, MAX, ), GOTO, 100

The output of LA is a stream of tokens that is passed on to the next phase. The first 2 phases - LA & SA are often group together in some pass i.e. LA operates either under the control of parser or as a co-routine with the parser. The parser asks LA for the next token whenever it needs it. LA returns a code for the token it has found.

The LA calls a book keeping routine which installs the token in the symbol table if it is not present there. A symbol table is the data structure used to keep a record of all the tokens. For the same FORTRAN statement, the symbol table can be formed as follows :-

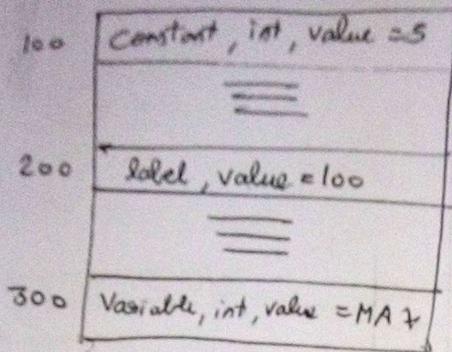
if (5 EQ MAX) GOTO 100

if ( [const, 100] EQ [id, 300] ) GOTO [LABEL, 200]

5
MAX
100

TYPE
Memory

address

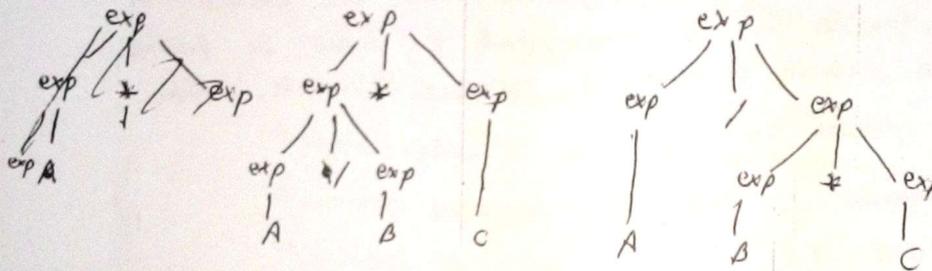


12/1

2) SYNTAX ANALYSIS / PARSER - It is the second phase of compiler. It performs 2 functions :-

- 1) It checks that the tokens obtained are acc. to the specification of source language (grammar) or not.
- 2) If yes it groups the tokens & forms a hierarchical structure called as parse tree.

Ex:- We have an exp.  $A/B * C$ . there can be  $3^2$  possible parse trees for this expressions :-



The correct parse tree depends on syntactic rules defined in grammar.

→ Parse tree of if statement & while loop in BOOK.

3) INTERMEDIATE CODE GENERATION :- This phase uses the parse tree & creates a stream of simple estimation in the form of 3 address code -  
Ex:- For the same exp.  $A/B * C$ , its 3-address code can be written as follows :-

$T_1, T_2$  are temporary variables  $T_1 = A/B, T_2 = T_1 * C$ .

AMP

4) CODE OPTIMIZATION :- This is optional phase designed to improve the (optional)

intermediate code so that the final object program takes less space & less time. there are 2 types of optimization :-

- a) Local optimization
- b) Loop Optimization.

In this method, local transformations are applied on the code. Ex:-

Elimination of common sub-expressions. Suppose :-

$A = B + C + D$   
 $E = B + C + F$   
 We can write  $T_1 = B + C$  ~~repeatedly~~  $\rightarrow A = T_1 + D$   
 $E = T_1 + F$

Another example elimination of jumps over jumps. This is a statement.

if  $A > B$  goto  $L_2$   $\rightarrow$  these instructions can be written as :-  
     goto  $L_3$                       if  $A < B$  goto  $L_3$   
 $L_2$  \_\_\_\_\_                       $L_2$

$\rightarrow$  This is optimization done in the loop. In this case, we move a computation that produces the same result everytime to a point in the program just before the loop is enhanced. Then, this computation is done only once. This type of computation is called LOOP INVARIANT.

5) CODE GENERATION :-

This is the final phase of compiler. It produces the object code by deciding on the memory locations of data, selecting the registers etc.

Ex:  $A = B + C$

Assembly :-   load B  
                   add C  
                   store A

6) TABLE MANAGEMENT / BOOK KEEPING :- This part of compiler keeps a track of the names or tokens used in the program & record essential info about them (datatype etc). There is a data structure used here called SYMBOL TABLE. Most of the info is collected in the early phases i.e. LA & SA.

7) ERROR HANDLING :- This part is invoked whenever there is a flow or error in source program. It warns the programmer by using a diagnostic message. Every phase can detect an error.

## COMPILER WRITING TOOLS

There are the tools that help to generate compilers automatically.

Ex:- Compiler - compilers

Compiler generators

Translator writing system.

These are 2 more tools - LEX & YACC (yet another compiler compiler) -  
LEX is used to generate LA & YACC is used to generate SA.

The input given to these tools is - a description of lexical & syntactic structure of source language, a description of output & a description of target system.

17/1

PASS

A pass is a combination of one or more phases into a single module.

Ex:- The first 2 phases LA & SA work together in the same pass.

A pass reads either the source program or op of previous pass makes the transformation as specified by its phases & writes the op into an intermediate file which is read by the next pass.

The no. of passes depends upon the structure of source language & the environment in which the compiler has to operate.

MULTI-PASS COMPILER :- It uses less space bcz space occupied by 1 pass can be reused by the next pass. It is slower in execution bcz each pass reads & writes into intermediate files.

## BOOTSTRAPPING

It is a process of obtaining a compiler for a complete language by using a compiler for a subset of that language. A compiler is characterized by

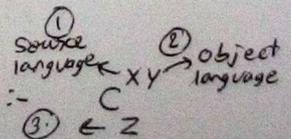
3 languages :-

1) Source language

2) Object

3) Language in which compiler itself is written.

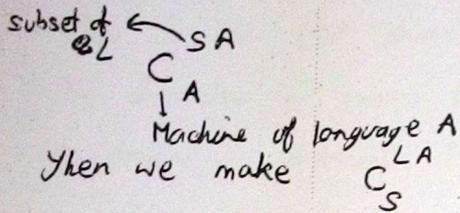
These all may be quite different. In this we use following notations :-



CROSS-COMPILER :- It is a compiler which runs on 1 machine & produces object code for some other machines.

Ex :- To implement the concept of bootstrapping, we take eg. of new language L that we want to make available on 2 machines A & B. We assume both A & B are similar to each other from hardware point of view.

For First for machine A, we write a small compiler,



After that  $C_S^{LA}$  is made to run through the compiler already made.

Equation is :-

$$C_S^{LA} \rightarrow \boxed{C_A^{SA}} \rightarrow C_A^{LA}$$

BOOTSTRAPPING

Now suppose we want to produce another compiler for L to run on machine B & produce object code for B. ( $C_B^{LB}$ ) For this first of all we write  $C_L^{LB}$  using the concept of bootstrapping. Then we use following equation :-

$$C_L^{LB} \rightarrow \boxed{C_A^{LA}} \rightarrow C_A^{LB} \quad (L-C)$$

$$C_L^{LB} \rightarrow \boxed{C_A^{LB}} \rightarrow C_A^{LB}$$

CHAPTER 3 : LEXICAL ANALYSIS :-

- Lexical analyser also known as Scanner.
- INPUT :- Programming Language i.e. source program.
- OUTPUT :- Tokens

This is the first phase of compiler. It reads the source program character by character & groups them into tokens. In this chapter we will study how to design & implement lexical analysis/analyser.

There are 2 issues in this phase :-

- 1) We need a method of describing all the possible tokens that can appear in the input program. For this we have a notation called as REGULAR EXPRESSION that describes all the possible tokens.
- 2) We need a mechanism to recognize these tokens. For this transition diagrams & finite automata are used.

The advantage of using regular expression is that from a regular expression we can automatically construct the token recognizers.

NEED FOR LEXICAL ANALYSIS AS A SEPARATE PHASE :-

The analysis of source program is split into the first 2 phases LA & SA (syntax analysis) - This is done :-

- to simplify the overall design of the compiler.
- It is easier to specify the structure of token than their syntax.

CONCEPT OF INPUT BUFFER :-

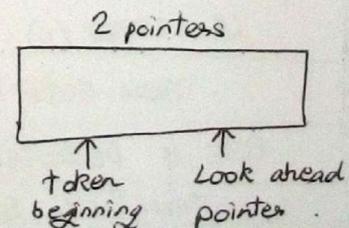
It is assumed that L.A. reads its input from an input buffer. This is needed for the following reason :-

- Buffer is a data structure.

Sometimes it is required by LA to examine many characters beyond the next token, even to determine the next token. Many schemes are used to design the buffers. For ex :-

The following fig. shows an example of it.

The first pointer marks the beginning of token while the second pointer scans the character until the next token is discovered.



Ex :- In PL/I program (language) there is a statement :-

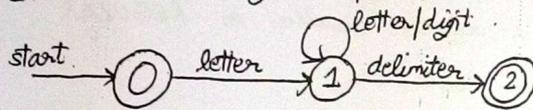
```
DECLARE (avg1, avg2, ... avgn)
```

In this language declare is a keyword & it is also used as an array name. So this is identified by L.A. as an array after scanning it till the right parenthesis. from the buffer

### DESIGN OF L.A. :

To design a software the simplest method used is a flowchart. To design we use special kind of flow chart used as TRANSITION DIAGRAMS. In these circles are used to show states & states are connected by arrows called Edges. The labels on the edges indicate the input character that appear in the source program.

Ex :- Transition diagram for an identifying identifier is as follows :-



→ delimiter is a character i.e. neither a digit nor letter & it is defined by programmer.

- An identifier is defined as a letter followed by any number of letters & digits.

19/1

P.T.O.

To convert transition diagram into a program, we write code for each state.

```
State 0 : c = GETCHAR ();  
IF LETTER (c) THEN GOTO  
STATE 1;  
ELSE  
FAIL ();
```

```
State 1 : c = GETCHAR ();  
IF LETTER (c) OR DIGIT (c)  
THEN GOTO STATE 1  
ELSE IF DELIMITER (c)  
THEN GOTO STATE 2
```

ELSE

FAIL();

STATE 2:- RETRACT();

return(id, INSTALL());

These are some functions used here :-

1) GETCHAR() :- This function returns the next character in the source program and advances the lookahead pointer.

2) LETTER(c) :- This function returns true if c is a letter

3) DIGIT(c) :- This function returns true if c is a digit

4) DELIMITER(c) :- " " " " " " " " " " delimiter

5) FAIL() :- This function retracts the lookahead pointer & start the next transition program or calls error routine.

6) RETRACT() :- This function retracts the lookahead pointer.

7) INSTALL() :- This function installs the identifier in the symbol table, if it is not present there.

In state 2, Lexical analyser returns to the parser a pair consisting of integer code for that particular token and a value that is the pointer to the particular symbol table returned by ~~the~~ install function.

TOKEN	CODE	VALUE
begin	1	-
end	2	-
if	3	-
then	4	-
else	5	-
identifier	6	pointer to symbol table
constant	7	pointer to symbol table
<	8	1

$\leq$	8	2
$=$	8	3
$\langle \rangle$	8	4
$\rangle$	8	5
$\geq$	8	6

REGULAR EXPRESSION :- It is a notation used by lexical analyser for describing tokens in the source program.

SOME IMPORTANT TABLE TERMS :-

1) Alphabet / Character class :- Finite set of symbols / characters.

Ex:  $\{0, 1\}$

2) String (Word / sentence) :- Finite set of symbols.

Ex: 001

3) Length of a string (x) is denoted by  $|x|$ . It is the no. of symbols present in the string.

4) Empty string  $\epsilon$  - string of length 0.

5) Operations :-

If  $x$  &  $y$  are strings then their concatenation is denoted by  $xy$  or  $x \cdot y$ . It is a string made up of symbols of  $x$  followed by symbols of  $y$ .

Ex:  $abc \cdot de = abcde$

$x \cdot \epsilon = x$ ,  $x^2 = x \cdot x$ ,  $x^i = x \cdot x \dots i \text{ times}$

6) Prefix of a string  $x$  :- It is a string formed by discarding 0 or more trailing symbols of  $x$ .

Ex:  $abc$  is a prefix of  $abcde$ .

→ Similarly suffix of string  $x$  is  $cde$  of  $abcde$ .

## 23/1 LANGUAGE :-

It is a set of strings form from a specific alphabet. Ex:-

if  $\{0,1\}$  is an alphabet then its language will be :-

$$\{0, 10, 0110, \dots\}$$

We use  $\phi$  to denote empty language or it is a set containing empty string.  $\phi = \{\epsilon\}$

## OPERATIONS ON LANGUAGE :-

1) CONCATENATION :- If  $L$  &  $M$  are 2 languages then their concatenation is shown as :-

$$L \cdot M = \{xy \mid x \text{ is in } L \text{ \& } y \text{ is in } M\}$$

$L \cdot M$  concatenation is all set of  $xy$  such that  $x \in L$  &  $y \in M$ .

Ex:-  $L = \{0, 01\}$

$$M = \{1\}$$

$$L \cdot M = \{011, 0111\}$$

Concatenation of  $L$  to the power  $i$   $L^i = L \cdot L \dots$   $i$  times

Gen  $L^0 = \phi = \{\epsilon\}$

2) UNION :- Union of 2 languages  $L$  &  $M$  is shown as follows :-

$$L \cup M = \{x \mid x \text{ is in } L \text{ or } x \text{ is in } M\}$$

3) CLOSURE :- Closure operation means concatenation with itself any number of times (from 0 to  $\infty$ ).

Suppose there is a language  $L$ , its closure is denoted by  $L^*$ . It means concatenation of  $L$  with itself any number of times.

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

$L^* \rightarrow$  Kleene Closure

$$L \cdot L^* = L \cdot \left( \bigcup_{i=0}^{\infty} L^i \right) \rightarrow \text{POSITIVE CLOSURE.}$$

$$L \cdot L^* = \bigcup_{i=0}^{\infty} L^{i+1}$$

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

$L^+ \rightarrow$  Positive Closure.

The positive closure of language  $L$  is concatenation of  $L$  with itself from 1 to  $\infty$ .

REGULAR EXPRESSION NOTATION for token identifiers is :-

$$\text{identifier} = \text{letter}(\text{letter} \mid \text{digit})^*$$

UNION

→ Each regular expression denotes a LANGUAGE.

→ There are some rules :-

1) If  $E$  is a regular expression then it denotes a language  $\phi$ .

2) If  $a$  is a regular expression then it denotes a language  $\{a\}$ .

3) If  $R$  &  $S$  are 2 regular expressions denoting languages  $L_R$  &  $L_S$  respectively.

i) If  $R$  &  $S$  in  $R|S$  is a regular exp. then  $L_R \cup L_S$ .

ii) If  $RS$  is a regular exp then  $L_R \cdot L_S$ .

iii) If  $R^*$  is a regular exp then  $L_R^*$ .

→ There are some algebraic laws followed by Regular exp :-

1)  $R|S = S|R$  ( $|$  is commutative).

2)  $R|(S|T) = (R|S)|T$  ( $|$  is associative).

3)  $R \cdot (S \cdot T) = (R \cdot S) \cdot T$  ( $\cdot$  is associative).

4)  $R \cdot (S|T) = (R \cdot S) | R \cdot T$  (distributive).

### FINITE AUTOMATA

→ A finite automata works like a token recognizer.

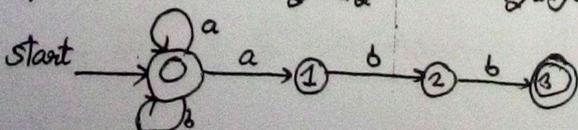
→ A token recognizer for a language  $L$  is a program that takes as input a string  $s$  & answers yes if  $s \in L$  & no otherwise.

It is of 2 types :-

- DFA - ~~NFA~~ NFA or NDFA.

NFA :- It is a labelled directed graph. Nodes are called 'states' & labelled edges are called TRANSITIONS. Edges can be labelled by  $\epsilon$  or any other character & some character can label 2 or more transitions out of a single state.

Ex:- A NFA recognizing the language  $(a/b)^* abb$



TRANSITION TABLE :-

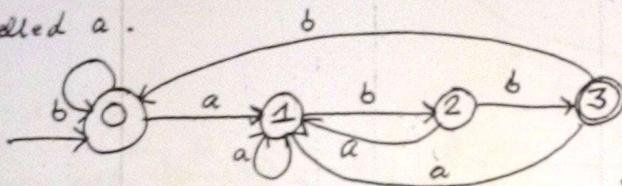
It is a table used to represent F.A. It has a row for each state & column for each input symbol.

States	a	b
0	{0,1}	{0}
1	-	{2}
2	-	{3}

A NFA accepts an input string  $x$  iff there is a path from start state to some accepting state such that the labels along that path spell out small  $x$ .

DFA :- It has 2 properties :-

- 1) There are no transitions on input  $\epsilon$ .
- 2) For each state  $s$  on input symbol  $a$  there is at most 1 edge labelled  $a$ .



It is easier to simulate DFA in programming -

How TO CONVERT REGULAR EXPRESSION TO NFA (ALGORITHM)

∴ Given :- A regular expression  $R$ .

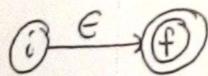
∴ Output :- A NFA accepting the language denoted by  $R$ .

∴ Method :-

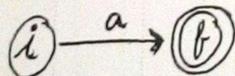
Firstly we decompose  $R$  into its basic components. For each component we draw a NFA.

Ex :-

i) I<sup>st</sup> CASE :- If the component is ' $\epsilon$ ' then its NFA will be



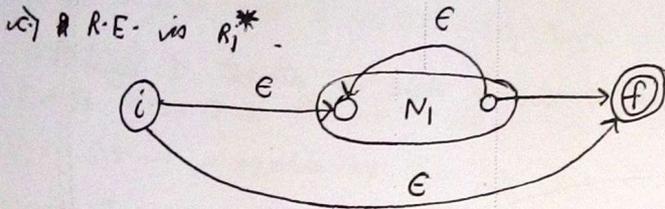
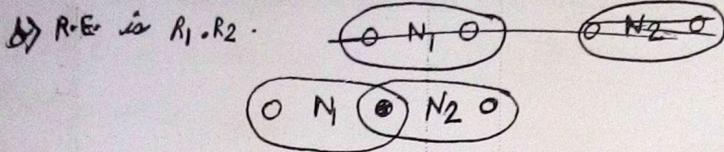
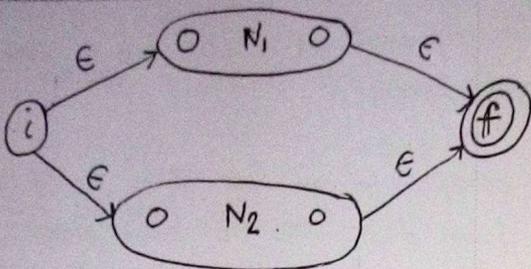
ii) II CASE :- If the component is ' $a$ ' then its NFA will be :-



iii) III CASE :- After this we combine the NFAs of basic components to produce NFA for complete regular expression.

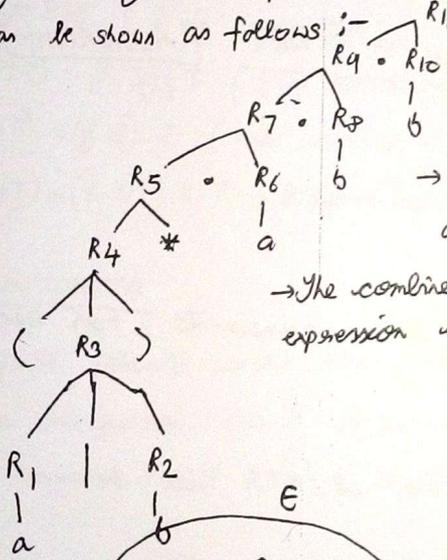
Now there are 3 possibilities :-

a) If the regular expression is  $R_1$  or  $R_2$  / ' $R_1 \cup R_2$ ' / ' $R_1 R_2$ ' -  
 The NFA for  $R_1$  is  $N_1$  & for  $R_2$  is  $N_2$  respectively, then their combined NFA will be drawn as follows :-



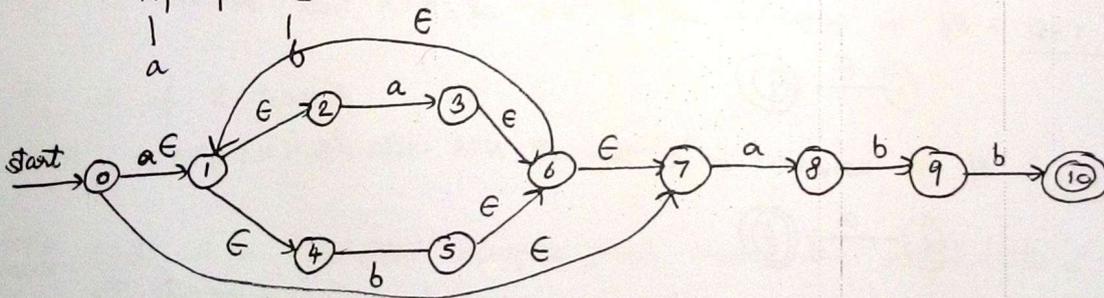
Ex: If we have R.E.  $(a|b)^*abb$ .

The given regular expression can be broken up into its basic components that can be shown as follows :-

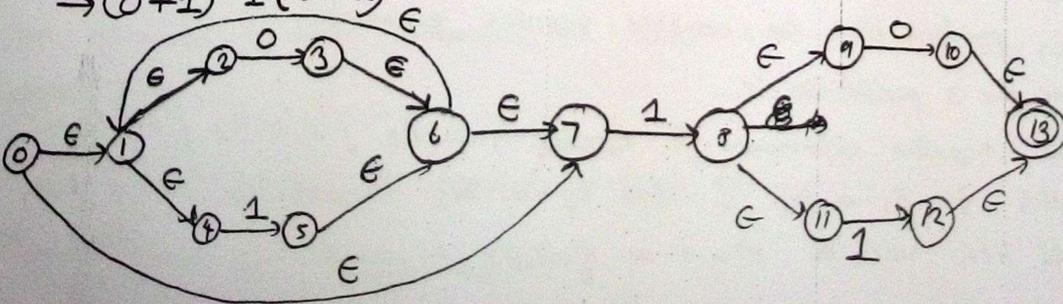


→ For each basic component (R) we draw its NFA.

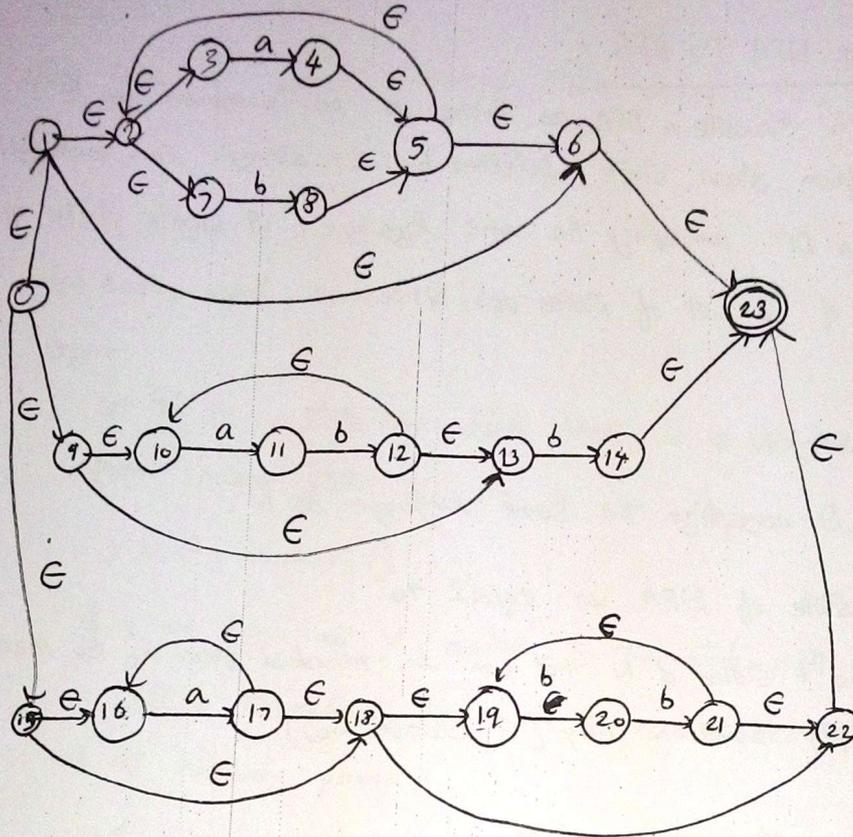
→ The combined NFA for the complete regular expression is drawn as follows :-



→  $(0+1)^* 1 (0+1)$



$$(ab)^* \mid (ab)^*b/a^*(bb)^*$$



DFA

States	a	b	Token accepted
<del>0 1 9 15</del>			
0 1 2 3 7 9 10 15 16 4 11 17 8			No
4 11 17	-	12 20	a
12 20	-	14 21	No
14 21			

## ALGO TO CONVERT NFA TO DFA :-

It is easier to simulate a DFA as there is no transition on  $\epsilon$ . A DFA has at most 1 path from start state labelled by 1 string. For each NFA, we have to draw a DFA accepting the same language. A single state of DFA will consist of a set of states of NFA.

Given :- A NFA,  $N$

Output :- A DFA,  $D$  accepting the same language as  $N$ .

Method :- Initial state of DFA is equal to

$= S_0 + \text{states of } N \text{ that can be reached from } S_0 \text{ by means of empty transitions only } (\epsilon \text{ closure}(S_0))$ .

where  $S_0$  is initial state of  $N$ .

### $\epsilon$ -CLOSURE( $s$ ) :-

This function gives us a set of states that can be reached from  $S$  by  $\epsilon$ -transitions (empty transitions).

1)  $S$  itself comes in  $\epsilon$ -CLOSURE( $s$ )

2) if  $t$  is in  $\epsilon$ -closure( $s$ ) & there is an edge labelled  $\epsilon$  from  $t$  to  $u$  then  $u$  is also added to  $\epsilon$ -closure( $s$ ).

This step is repeated till no more states can be added.

→ In the beginning, we assume that each state of DFA,  $D$  is unmarked.

### ALGO:-

While there is ~~an~~ unmarked state  $x = \{s_1, s_2, \dots, s_n\}$  of  $D$  do

begin

Mask  $x$ ;

for each input symbol  $a$  do

begin

let  $T$  be a state to which there is a transition ~~from~~ <sup>on</sup> 'a' from source state  $s_i$  in  $x$ .

$$y = \epsilon\text{-closure}(T)$$

if  $y$  ~~is~~ has not  $\bullet$  added to the ~~set~~ set of states of  $D$ .

add a transition from  $x$  to  $y$  labelled  $a$ .

if not already present

end;

end;

Ex:- We take the same NFA,  $(a/b)^*abb$ .

$$\begin{aligned} \text{Initial state of } D &= 0 + \epsilon\text{-closure}(0) \\ &= \{0, 1, 2, 4, 7\} = A \end{aligned}$$

Transition diagram of  $D$  :-

	a	b
A	B	C
B	B	$\emptyset$
C	B	C
D	B	E
E	B	C

→ On input symbol  $a$  & state A following states have transitions & these next states are  $\{3, 8\}$ .

$$\text{Now } \epsilon\text{-closure}(\{3,8\}) \\ = \{1,2,3,4,6,7,8\} = B$$

Similarly on input symbol of  $b$  & state  $A = \{5\}$

$$\text{Now } \epsilon\text{-closure}(\{5\}) \\ = \{1,2,4,5,6,7\} = C$$

- On input symbol 'a' & state B,  $\{3,8\}$

$$\epsilon\text{-closure}(\{3,8\}) \\ = \{1,2,3,4,6,7,8\} = B$$

- On ~~an~~ input symbol  $b$  & state B,  $\{5,9\}$

$$\epsilon\text{-closure}(\{5,9\}) \\ = \{1,2,4,5,6,7,9\} = D$$

- On input symbol 'a' & state C  $\{3,8\}$

$$\epsilon\text{-closure}(\{3,8\}) \\ = \{1,2,3,4,6,7,8\} = B$$

- On input symbol 'b' & state C  $\{5\}$

$$\epsilon\text{-closure}(\{5\}) \\ = \{1,2,4,5,6,7\} = C$$

- On input symbol  $a$  &  $\emptyset$   $\{3,8\}$

$$\epsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = B$$

- On input symbol  $b$  &  $\emptyset$   $\{5,10\}$

$$\epsilon\text{-closure}(\{5,10\}) = \{1,2,4,5,6,7,10\} = E$$

- On input symbol  $a$  & state E  $\{3,8\}$

$$\epsilon\text{-closure}(\{3,8\}) = \{1,2,3,4,6,7,8\} = B$$

- On input symbol  $b$  & state E  $\{5\}$

$$\epsilon\text{-closure}(\{5\}) = \{1,2,4,5,6,7\} = C$$

## MINIMIZING THE NO. OF STATES OF DFA (if possible)

INPUT  $\rightarrow$  A DFA,  $\mathcal{D}$  with no. of states

OUTPUT  $\rightarrow$  A DFA,  $\mathcal{D}'$  accepting the same language as  $\mathcal{D}$  & having as few states as possible.

1. We construct a partition  $\pi$  of set of states. Initially  $\pi$  consists of 2 groups i.e. final state  $F$  & a non-final state  $S-F$ .

2. Then we construct a new partition  $\pi_{\text{new}}$  by the following methods:-

$\rightarrow$  for each group  $G$  of  $\pi$  do

begin

partition  $G$  into subgroups such that 2 states  $s$  &  $t$  of  $G$  are in the same subgroup if & only if for all input symbol, a state  $s$  &  $t$  have transitions to the states in the same group of  $\pi$ . Place all subgroups formed in  $\pi_{\text{new}}$ .

end

Ex:- For some DFA there were total 5 states:-

$(A, B, C, D, E)$ .

$$\pi_1 = (ABCD)(E)$$

$$\pi_2 = (ABC)(D)(E)$$

$$\pi_3 = (AC)(B)(D)(E)$$

Since  $A$  &  $C$  both gives same transitions. In such a case we select a representative out of the 2 states.

Suppose we select  $A$ .

$$\therefore \pi_4 = (A)(B)(D)(E).$$

Thus transition table of reduced DFA.