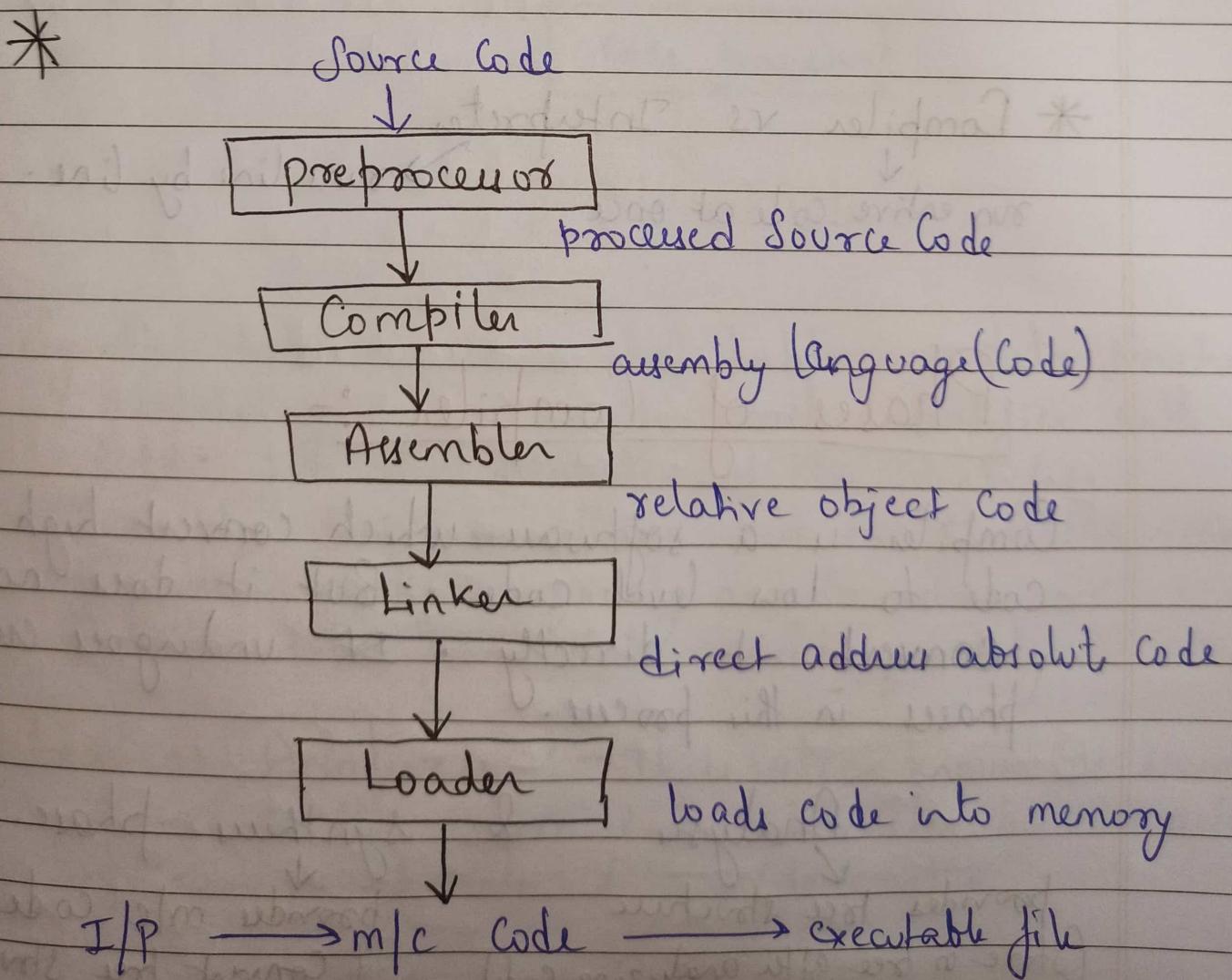
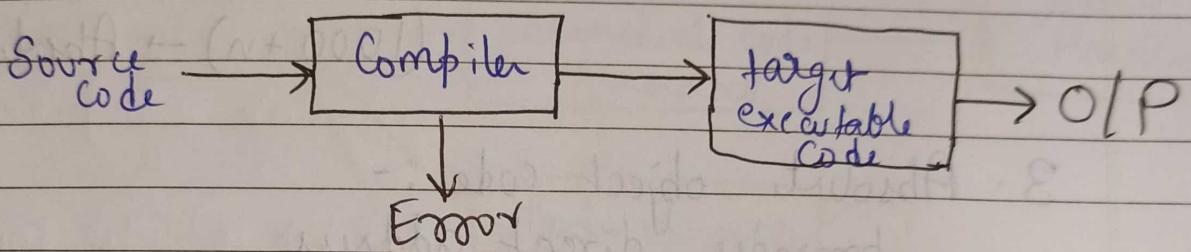


# Compiler Design

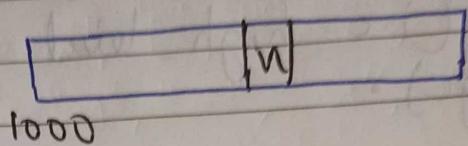
## Unit-1

Compiler :- Compiler is a execution environment that converts the source code to a target code (i.e high level language to machine level language).



1. Assembly Code → in between high & low level language

2. Relocatable object code → provide relative address.



$n \rightarrow$  Relative address

$(1000 + n) \rightarrow$  Absolute address

3. Absolute object code :-  
provides direct address

\* Compiler vs Interpreter  
↓  
run entire code at once      → run line by line

## Phases of Compiler :-

Compiler is a software which convert high level code to low level code. But it does not convert it directly. If we do you can see phases in this process.

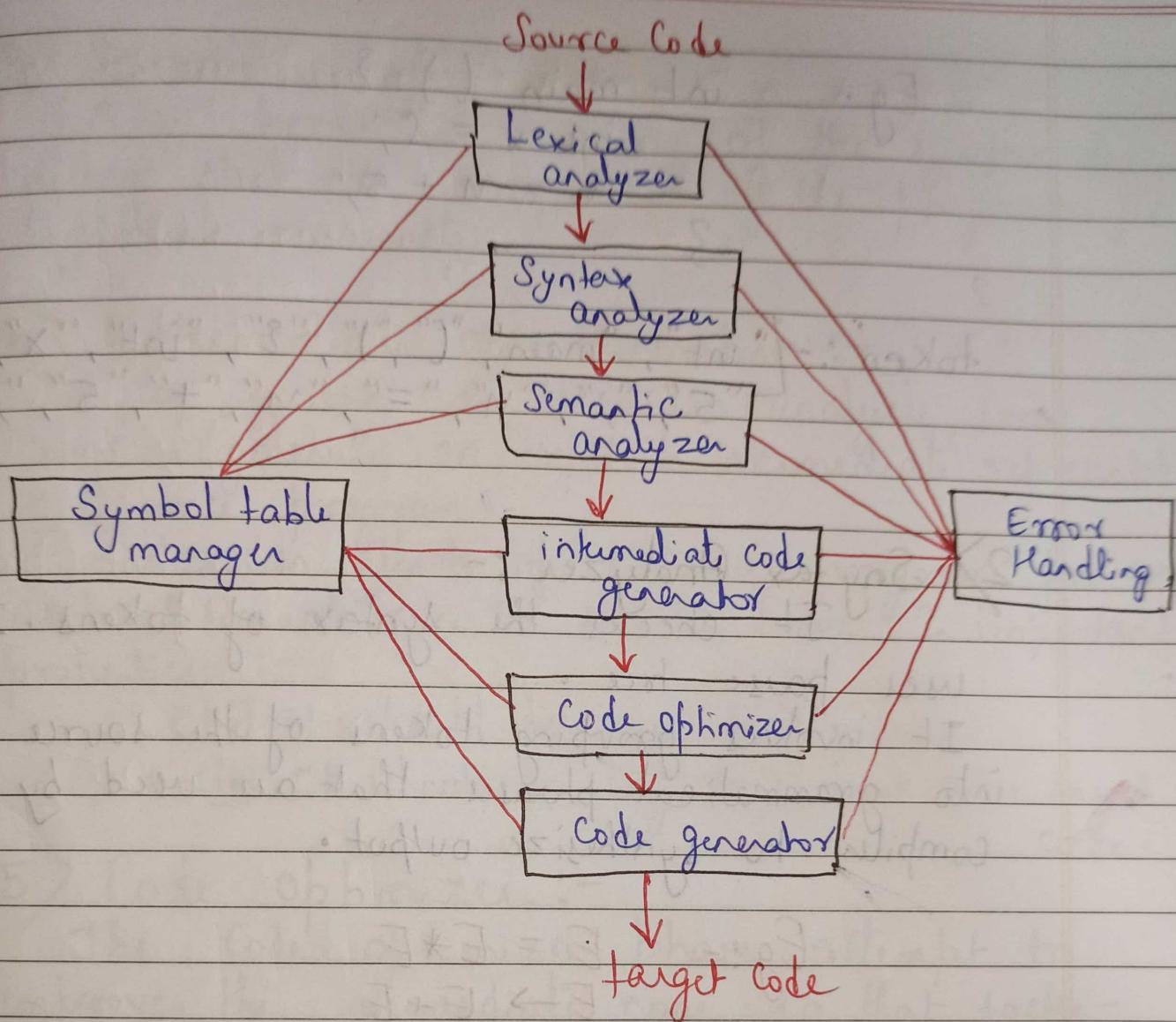
Analysis & Synthesis phase.

↓  
provides tree structure

(Store in tree after analyzing code)

↓  
provides m/c code

(Convert tree structure to machine code)



**Symbol table :-** A symbol table is a data structure containing a record for each identifier with fields for the attributes of the identifier. The data structure allows us to find the record for each identifier quickly and to store or retrieve data from that record quickly.

\* **Lexical Analyzer :-** It is also known as scanner. It reads the entire code character by character and converts it into tokens (lexeme).

Eg. int main ()  
 int x = 5;  
 n = n + 5;

3

tokens :- ["int", "main", "(", ")", "?", "int", "x", "=", "  
 "5", ";", "n", "=", "n", "+", "5", ";", "3"]

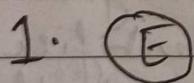
## 2) Syntax Analyzer :-

It checks the syntax of tokens. It uses parse tree.

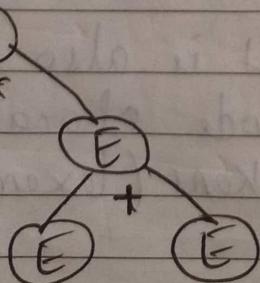
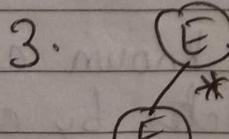
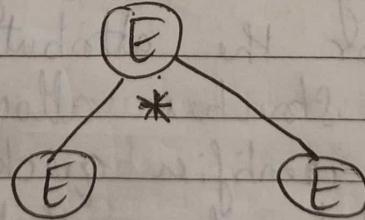
It involves grouping tokens of the source program into grammatical phrase that are used by compiler to synthesize output.

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

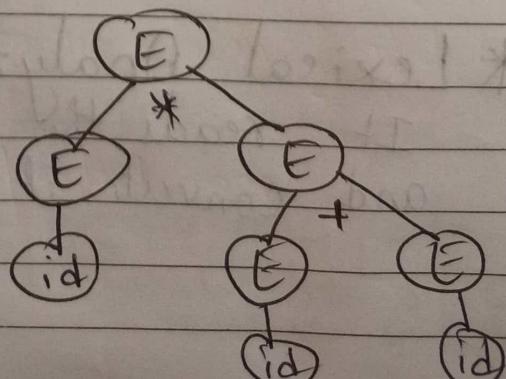
Generate:  $id * id + id$ .



2.



4.



### 3) Semantic analyzer :-

It checks the meaning of the syntax to ensure that the component of the program fit together manually.

### 4) Intermediate code Generation :-

After syntax and semantic analysis, some compiler generate an explicit intermediate representation of source program.

We can think of it as a program for an abstract machine. It should have 2 important properties :-

(i) It should be easy to produce

(ii) It should be easy to convert into target code.

### 5) Code optimizer :-

The code optimizer phase attempt to improve the intermediate code, so that faster-running machine code will result.

It removes the useless symbols without changing the meaning.

### 6) Code Generation :-

The final phase of the compiler is the generation of target code which is machine dependent. It is unique for all the core system (i.e. it is different for 32 bit & 64 bit and so on).

## Errors in various Analytic phase :-

1) Lexical Analyzer :- Misspelling  
Eg. main()

2) Syntax Analyzer :- Not in proper syntax  
Eg. int x = 5  
printf("Hello");

3) Semantics :- Syntax is valid or not (providing meaning or not).

Eg.  
int x = 5.5;  
char ch = 104.5;

## In Short Summary :-

machine independent code

- ① → info about tokens
- ② → data types, scope
- ③ → verify already stored information
- ④ → generate intermediate code based on info stored in symbol table
- ⑤ → optimize code based on info in symbol table
- ⑥ → machine dependent code

\* Front-end phase :- Machine independent  
 ↓  
(Analysis) Depend on programming language

Back-end phase :- Machine dependent  
 ↓  
(Synthesis) Don't Depend on programming language

\* Passes in Compiler :-

- 1) One-pass :- Single pass me convert kar dete h.  
 → faster but more space
- 2) Two-pass :- first convert to intermediate, then  
 in second pass convert to m/c code.  
 → slower but less space.

### Cousin's of Compiler :-

1) Preprocessor :- It produce input to compiler.  
 It may perform following function :-

- (a) Macro-processing :- #define PI 3.14
- (b) File inclusion :- #include <stdio.h>
- (c) Rational preprocessor :- It augments the old language to the new language by the use of built in macros.
- (d) Language extensions :-

2) Assemblers :- Some compilers produce assembly code that is passed to assembler which converts the assembly language to machine level code.

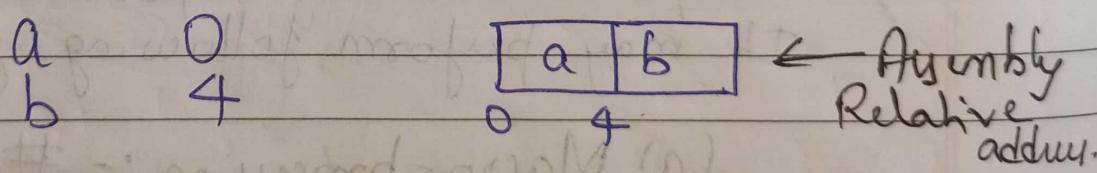
Assembly code is a mnemonic version of machine code, in which names are used instead of binary codes for operations and names are also given to memory address.

Eg:-

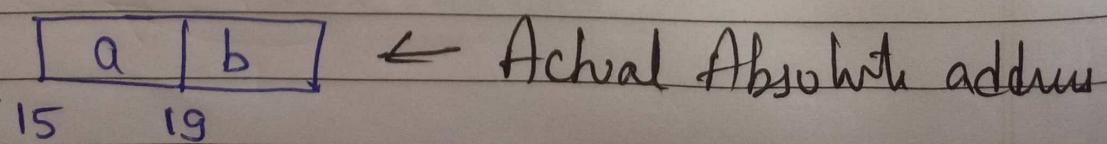
```
MOV a, B1
ADD #2, B1
MOV B1, b
```

3) Two Pass - Assembly :-

In first pass, all the identifiers that denote storage location are found and stored in symbol table. Basically it checks the code & assign the memory to variable according to requirements.



In Second pass, the assembler scans the input again. This time it translates each operation code into the sequence of bits representing that operation in machine language. The output of this pass is usually relocatable machine code.



MOV a, R1                    0001 01 00 00000000\*  
 ADD #2, R1 →              0011 01 10 00000010  
 MOV R1, b                    0010 01 00 00000100\*

#### 4) Loader & Link-Editors :-

Usually a program called Loader performs the two functions of loading and link-editing. The process of loading consists of taking relocatable machine code, altering the relocatable address and placing the altered instructions and data in memory at proper locations.

The Link-editor allows us to make a single program from several files from of relocatable machine code.

## Automata Theory :-

DFA :- 5 tuples :-

$Q \rightarrow$  Set of states

$\Sigma \rightarrow$  Set of input symbols

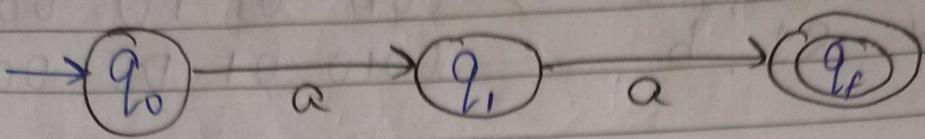
$\delta \rightarrow$  transition functions

$q_0 \rightarrow$  Initial State

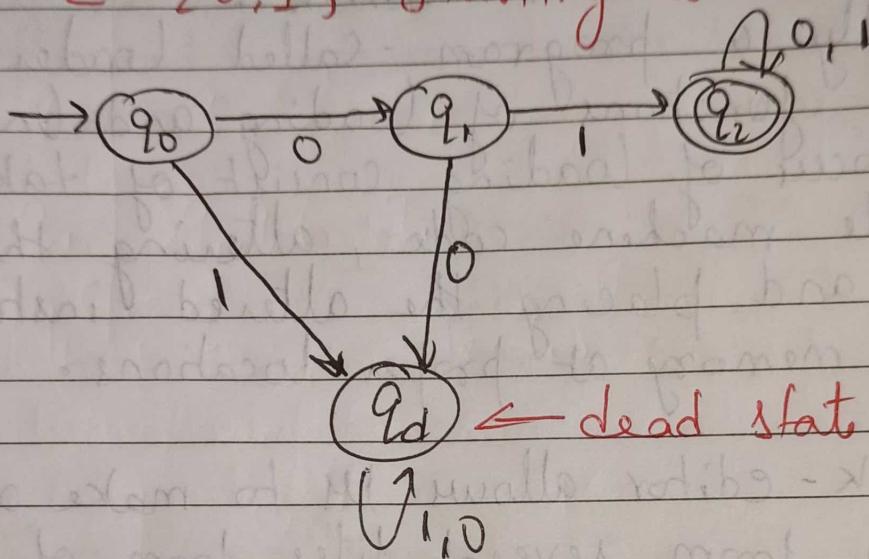
$q_f \rightarrow$  Final state

$$Q \times \Sigma \rightarrow Q$$

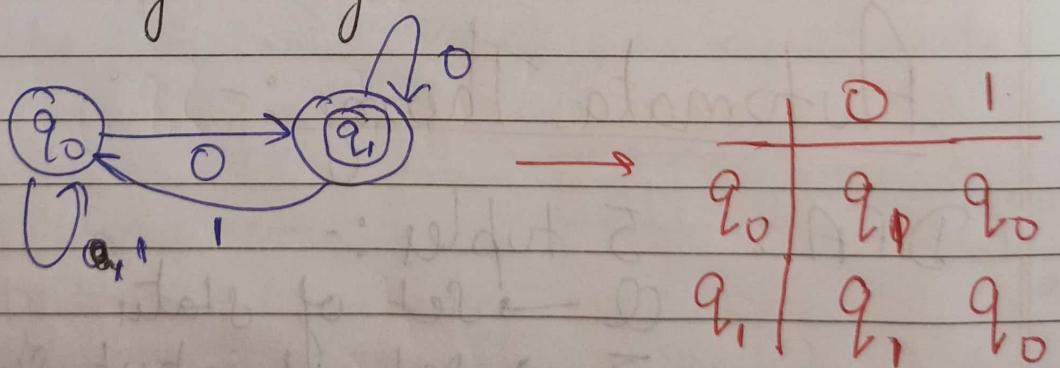
Q aa



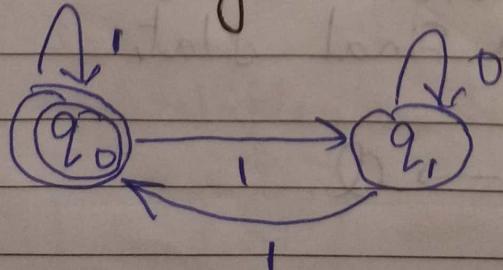
Q  $\Sigma = \{0, 1\}$  starting with 01.



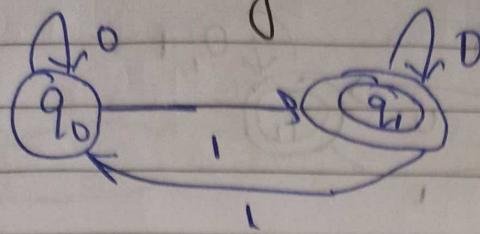
Q All string ending with 0.



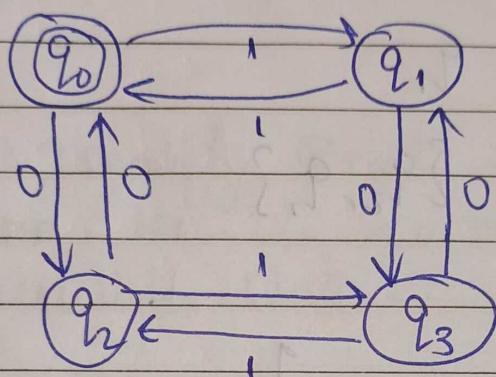
Q Even No. of 1's :-



Q. Odd no. of 1's :-



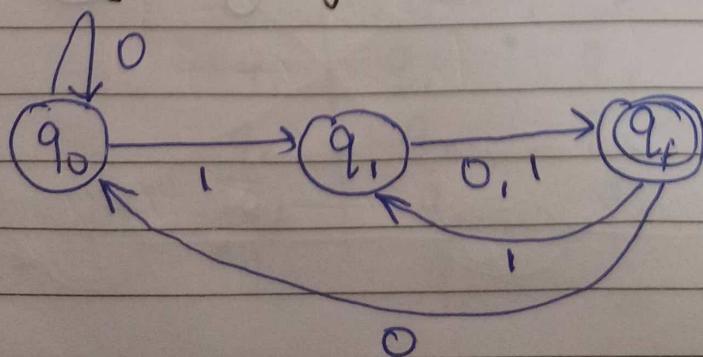
Q Even no. of 1's & even no. of 0's :-



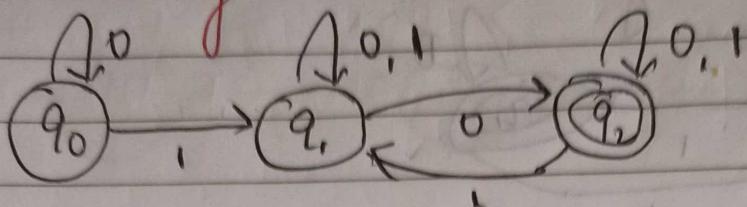
Practice Problem :-

- 1) Even no. of 1's or even No. of 1's.
- 2) odd no. of 1's and even no. of 0's.
- 3) odd no. of 1's or even no of 0.

Q. DFA for string with 2<sup>nd</sup> last bit as 1.

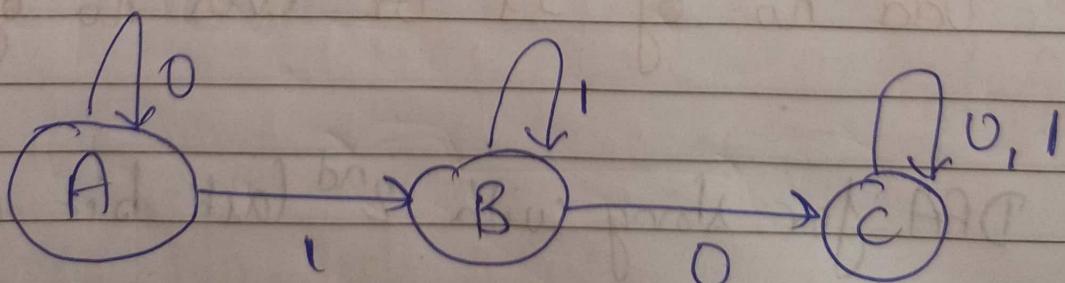


Converting NFA to DFA :-

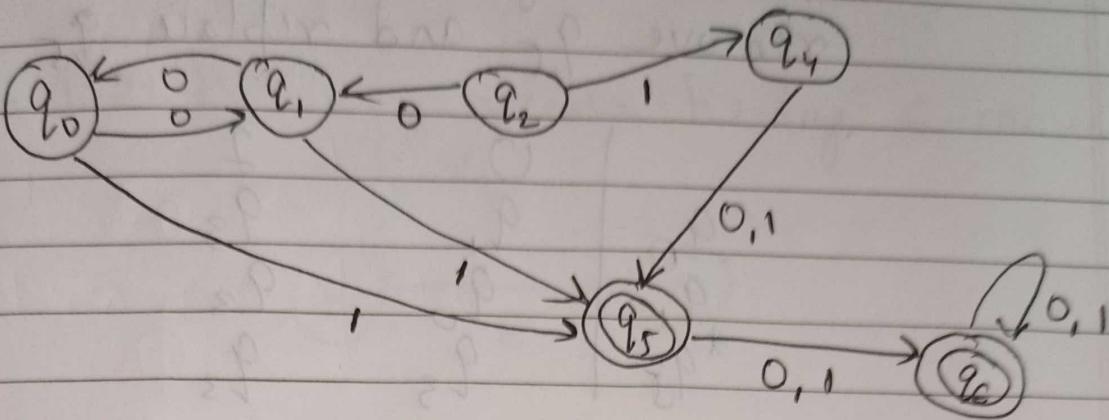


	0	1
$q_0$	$q_0$	$q_1$
$q_1$	$\{q_1, q_2\}$	$q_1$
$q_2$	$q_2$	$\{q_2, q_1\}$

- (A)  $q_0 \mid \begin{matrix} 0 \\ q_0 \\ \{q_1, q_2\} \end{matrix} \quad \begin{matrix} 1 \\ q_1 \\ q_1 \end{matrix}$
- (B)  $q_1 \mid \begin{matrix} 0 \\ \{q_1, q_2\} \\ q_1 \end{matrix} \quad \begin{matrix} 1 \\ q_1 \\ q_1 \end{matrix}$
- (C)  $q_1, q_2 \mid \begin{matrix} 0 \\ \{q_1, q_2\} \\ \{q_1, q_2\} \end{matrix} \quad \begin{matrix} 1 \\ \{q_1, q_2\} \\ \{q_1, q_2\} \end{matrix}$



# Minimization of DFA :-



Step 1 → Identify Dead state & Inaccessible state.  
 Here, there is no dead state but there is  
 2 inaccessible state.

Step 2 :- Draw Transition table will all the states  
 except dead & inaccessible one.

	0	1
→ q0	q1 q0	q3 q3
q1	q0 q1	q3 q1
* q3	q5 q5	q5 q5
* q5	q5 q5	q5 q5

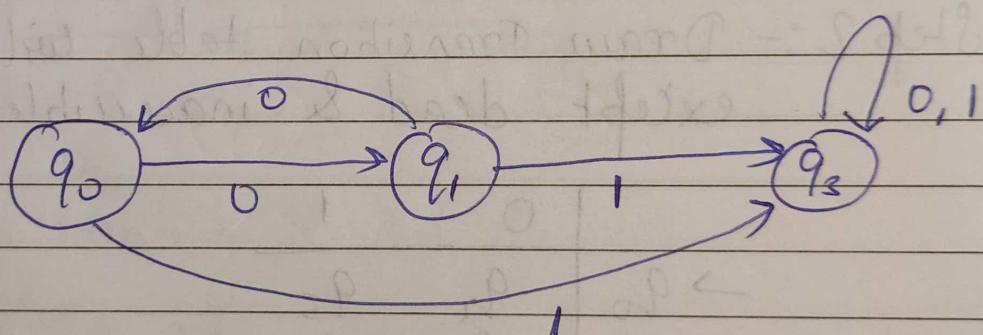
Step 3 :- Create partition between final stat and  
 all other state.

	0	1
→ q0	q1 q0	q3 q3
q1	q0 q1	q3 q1
* q3	q5 q5	q5 q5
* q5	q5 q5	q5 q5

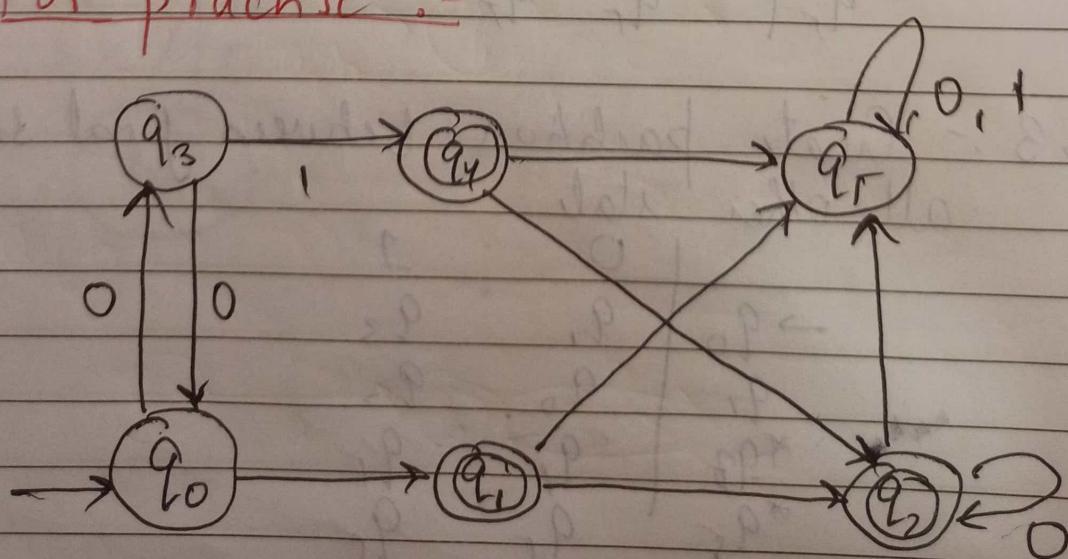
Step 4:- Since  $q_3$  and  $q_5$  both have same transition for 0 and 1. So, we can remove  $q_5$  and replace  $q_5$  with  $q_3$ .

	0	1
$q_0$	$q_1$	$q_3$
$q_1$	$q_0$	$q_3$
* $q_5$	$q_3$	$q_3$

Step 5:- Now, since there is no other pair of states with same transition, so, we stop and draw DFA for current transition table.

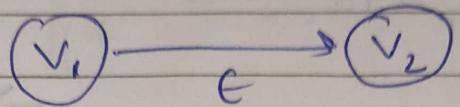


For practice :-



## Converting $\epsilon$ -NFA to NFA :-

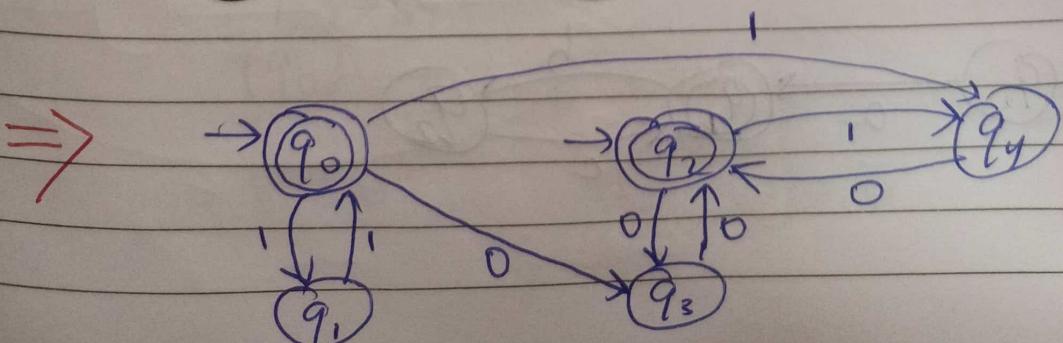
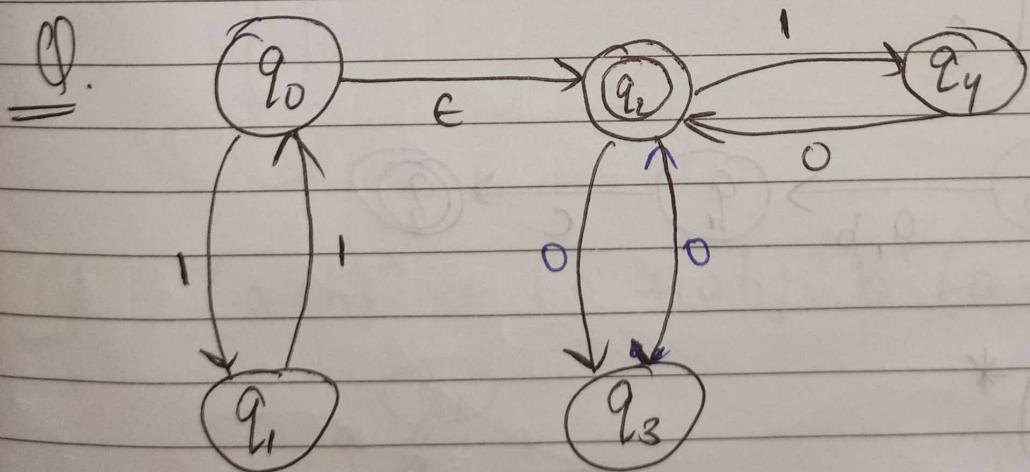
Step 1 :- Identify vertex having  $\epsilon$ -move as  $v_1$  and  $v_2$ .



Step 2 :- Copy all the moves of  $v_2$  and to  $v_1$ .  
Delete  $\epsilon$ -edge.

Step 3 :- If  $v_1$  is initial stat make  $v_2$  initial stati

Step 4 :- If  $v_2$  is final stat make  $v_1$  final stat.



## Regular Expression :

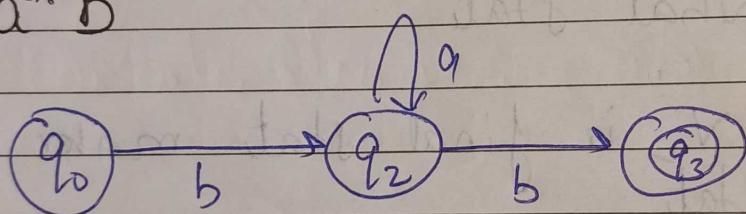
$$1. L = \{ \epsilon, a, aa, \dots \} \Rightarrow a^*$$

$$2. L = \{ a, aa, aaa, \dots \} \Rightarrow a^+$$

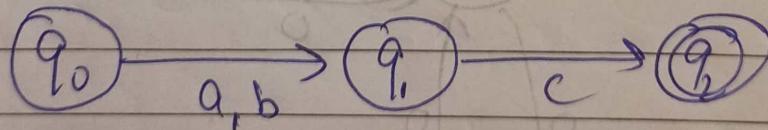
$$3. L = \{ \} \Rightarrow \emptyset$$

$$4. L = \{ \epsilon \} \Rightarrow \epsilon$$

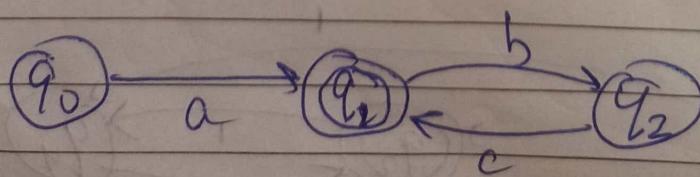
Q.  $ba^*b$



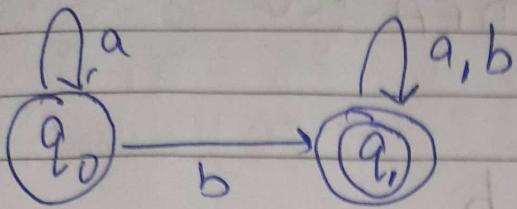
Q.  $(a+b)c$



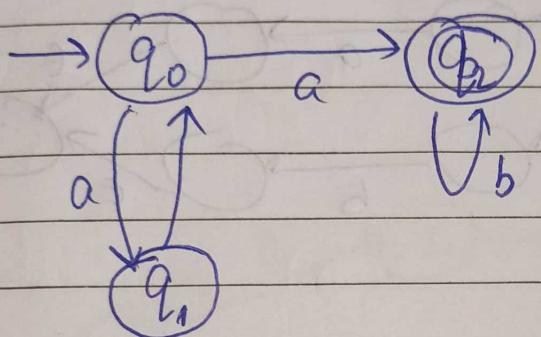
Q.  $a(bc)^*$



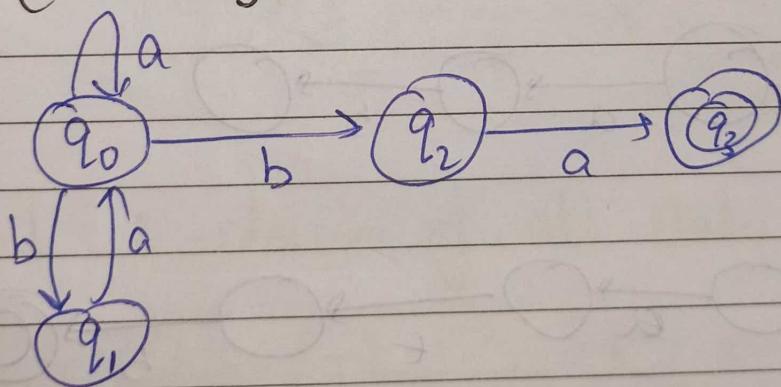
Q.  $a^* b (a+b)^*$



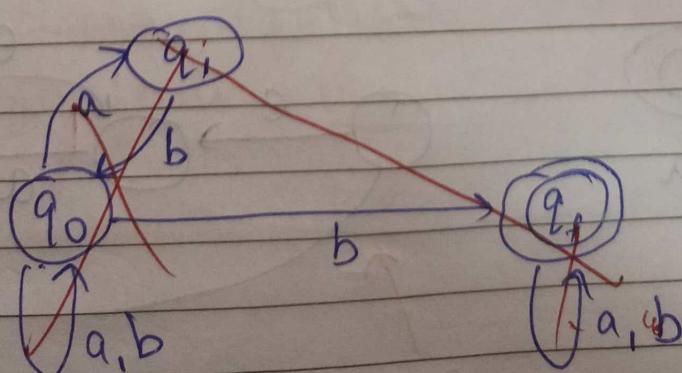
Q.  $(ab)^* a b^*$



Q.  $(a+ba)^* ba$

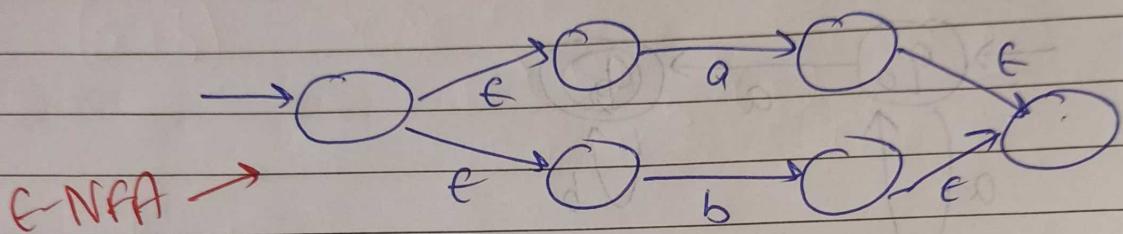
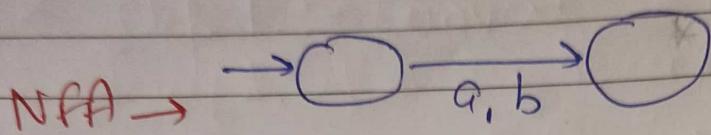


Q.  $(a+b)^* + (a+ab)^* b^+ (a+b)^*$

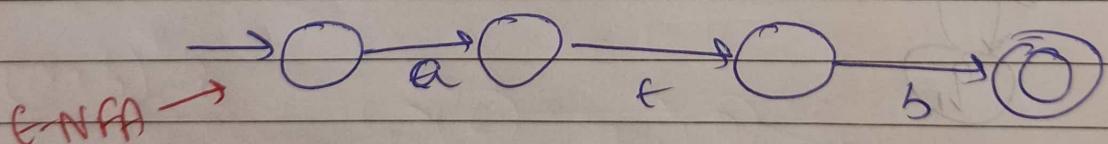
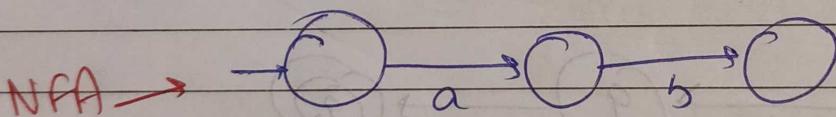


## Thompson's Algorithm for Converting RE to DFA :-

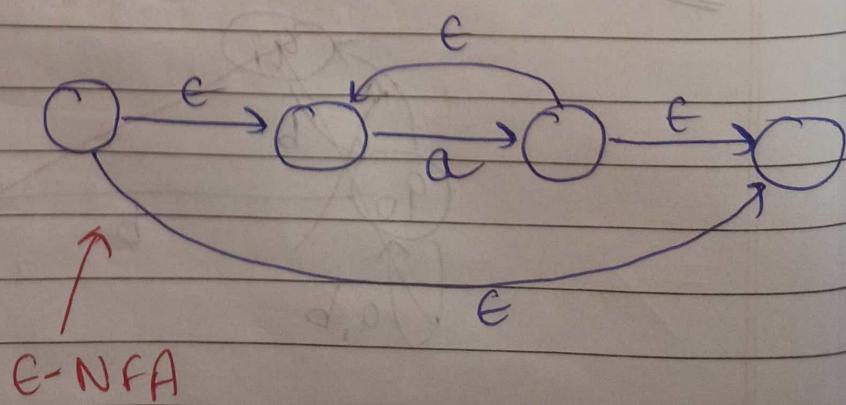
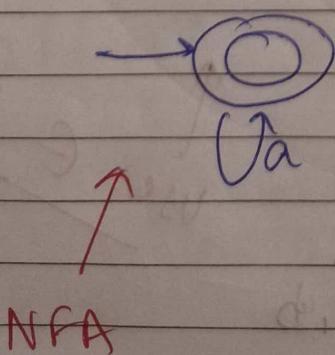
1. Union :-  $a + b$

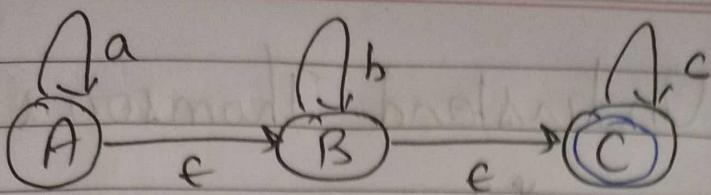


2. Concatenation :-  $ab$



3. Closure :-  $a^*$

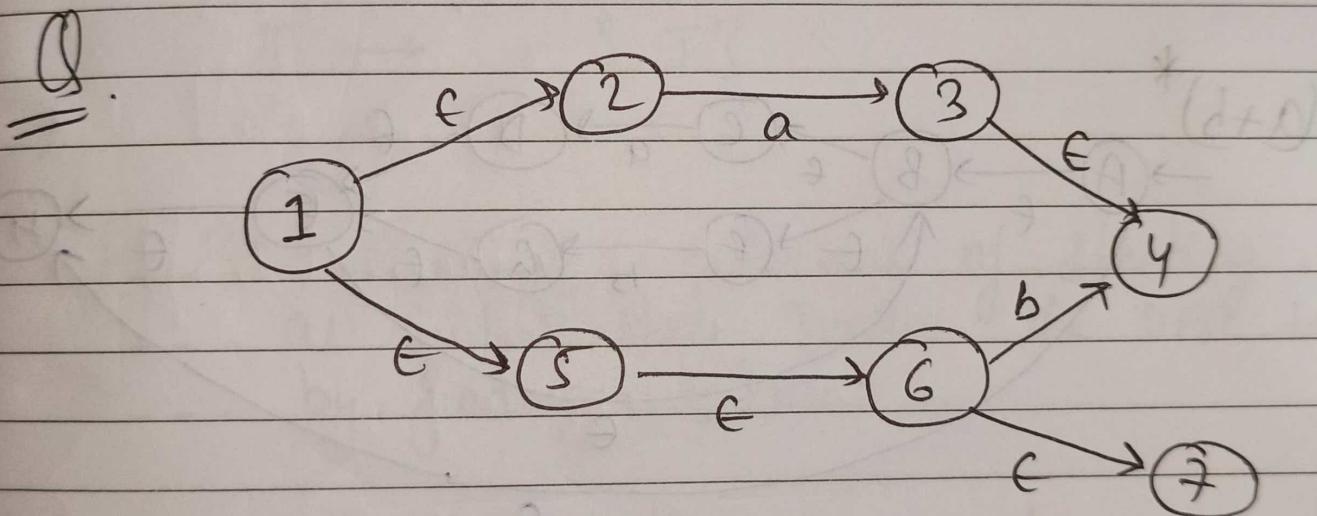




$\epsilon$ -closure(A) :-  $\{A, B, C\}$

$\epsilon$ -closure(B) :-  $\{B, C\}$

$\epsilon$ -closure(C) :-  $\{C\}$



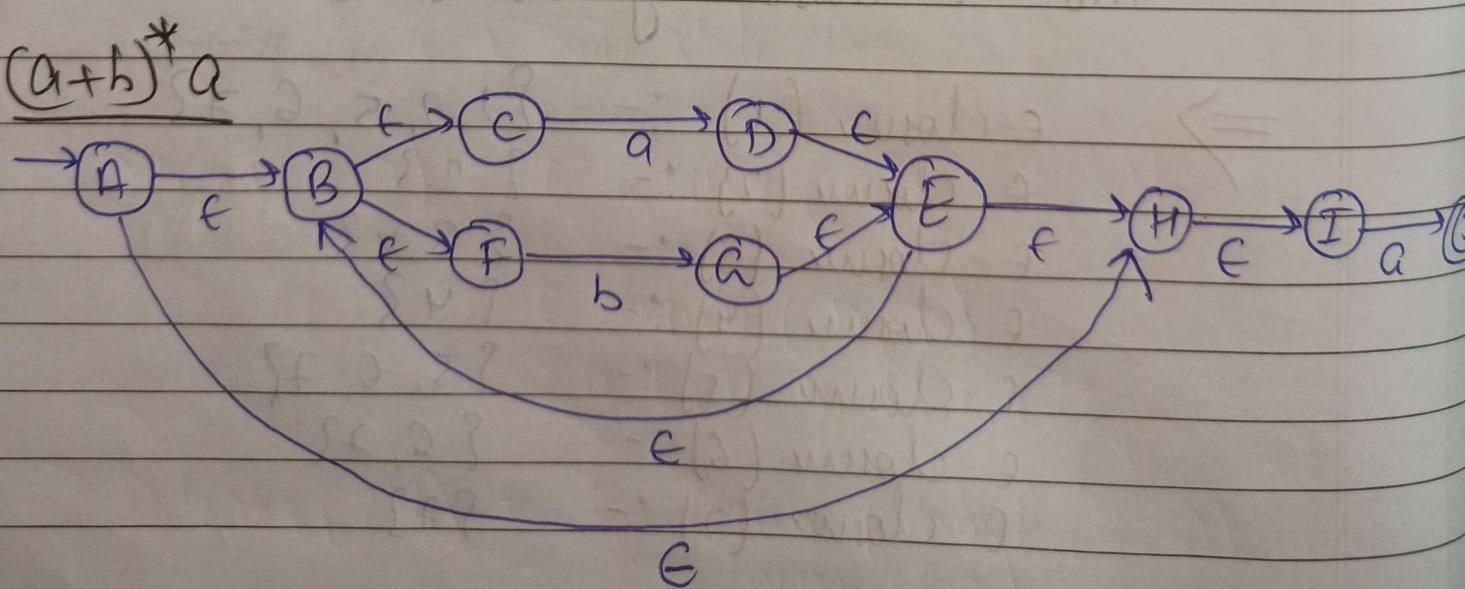
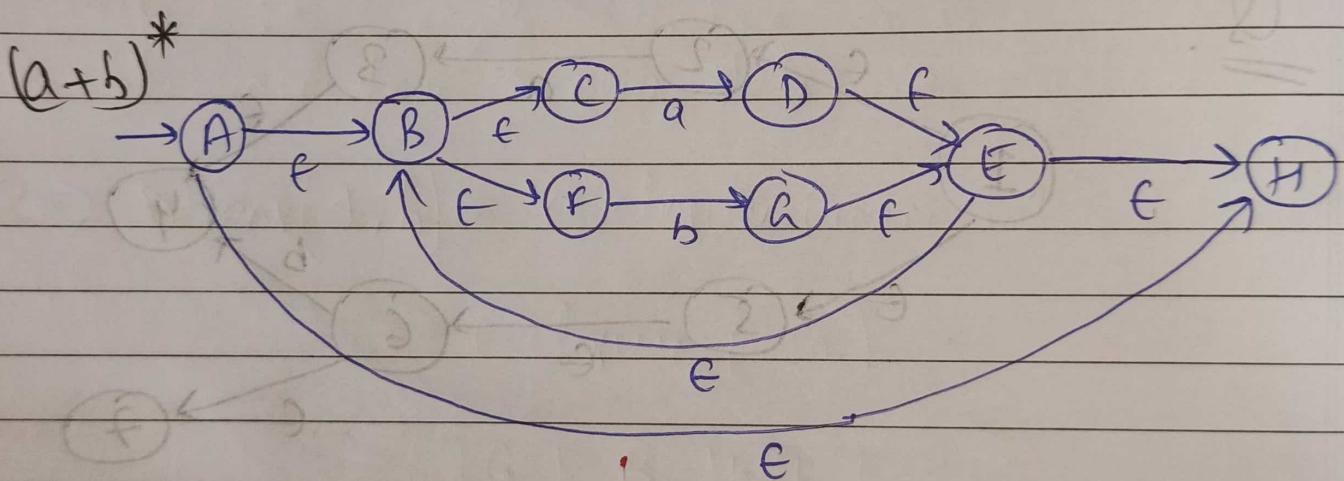
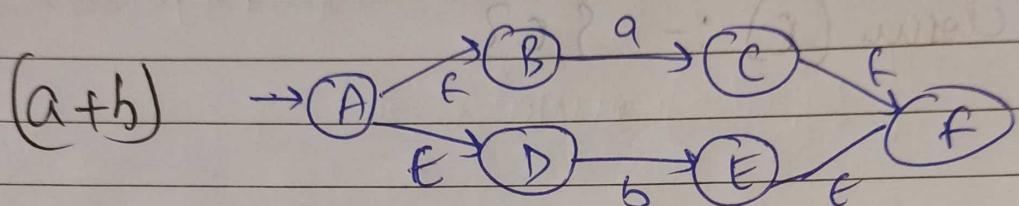
find  $\epsilon$ -closure of all states

- $\Rightarrow$
- $\epsilon$ -closure(1) :-  $\{1, 2, 5, 6, 7\}$
  - $\epsilon$ -closure(2) :-  $\{2, 3\}$
  - $\epsilon$ -closure(3) :-  $\{3, 4\}$
  - $\epsilon$ -closure(4) :-  $\{4\}$
  - $\epsilon$ -closure(5) :-  $\{5, 6, 7\}$
  - $\epsilon$ -closure(6) :-  $\{6, 7\}$
  - $\epsilon$ -closure(7) :-  $\{7\}$

Problem to Understand Thomson's Method :-

Q Convert  $(a+b)^*a$  to DFA.

⇒ Step 1 :- Convert RE to  $\epsilon$ -NFA



**Step 2 :-** Find  $\epsilon$ -closure of all states :-

$$\begin{aligned}
 A &\rightarrow \{A, B, C, F, H, I\} \\
 B &\rightarrow \{B, C, F\} \\
 C &\rightarrow \{C\} \\
 D &\rightarrow \{D, E, B, C, F, H, I\} \\
 E &\rightarrow \{E, H, I, B, C, F\} \\
 F &\rightarrow \{F\} \\
 G &\rightarrow \{G, E, B, C, F, H, I\} \\
 H &\rightarrow \{H, I\} \\
 I &\rightarrow \{I\} \\
 J &\rightarrow \{J\}
 \end{aligned}$$

**Step 3 :-** Taking  $\epsilon$ -closure of starting stat as initial stat Draw the DFA table by finding DTrans .

$\epsilon$ (NFA) stat	DFA stat	a	b.
$\rightarrow \{A, B, C, F, H, I\}$	1	2	3
$\{D, E, H, B, I, C, F, J\}$	2	2	3
$\{G, E, B, H, I, F, J\}$	3	2	3

final state  
it contains  
J

Formula :-

$$D\text{Trans}(\text{State}, \text{input-symbol}) = \epsilon\text{-closure}(\text{move}(\text{State}, \text{input-symbol}))$$

$$\Rightarrow D\text{Trans}(1, a) = \epsilon\text{-closure}(\text{move}(1, a))$$

$$= \epsilon\text{-closure}(D, S)$$

$$= \{D, E, H, B, I, C, F, S\}$$

Since it is a new state we call it ②

$$\Rightarrow D\text{Trans}(1, b) = \epsilon\text{-closure}(\text{move}(1, b))$$

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

$$= \{G, E, B, H, I, E, F\}$$

Since it is a new state we call it ③.

$$\Rightarrow D\text{Trans}(2, a) = \epsilon\text{-closure}(\text{move}(2, a))$$

$$= \epsilon\text{-closure}(D, S)$$

$$= \{D, E, H, B, I, C, F, S\}$$

$$\Rightarrow D\text{Trans}(2, b) = \epsilon\text{-closure}(\text{move}(2, b))$$

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

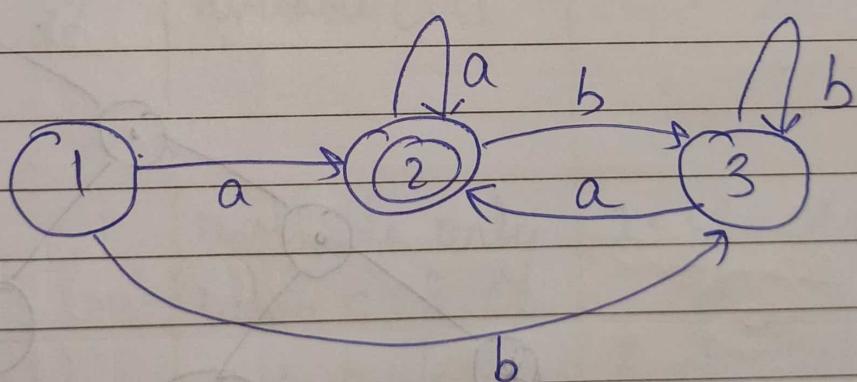
$$= \{G, E, B, H, I, E, F\}$$

$$\Rightarrow D_{trans}(3, a) = \text{c-closure}(\text{move}(3, a)) \\ \Rightarrow \text{c-closure}(D, J) \\ \Rightarrow \{D, E, H, B, I, C, F, T\}$$

$$\Rightarrow D_{trans}(3, b) = \text{c-closure}(\text{move}(3, b)) \\ \Rightarrow \text{c-closure}(G) \\ \Rightarrow \{A, E, B, H, I, E, F\}$$

↳ since there is no new stat available now,  
so we stop here.

**Step 4 :-** Draw DFA from the calculated table.



**Question for practice :-**

$$\rightarrow (ab)^* + ab$$

$$\rightarrow (a + ab)^*$$

## Direct Method for Converting RE to DFA:

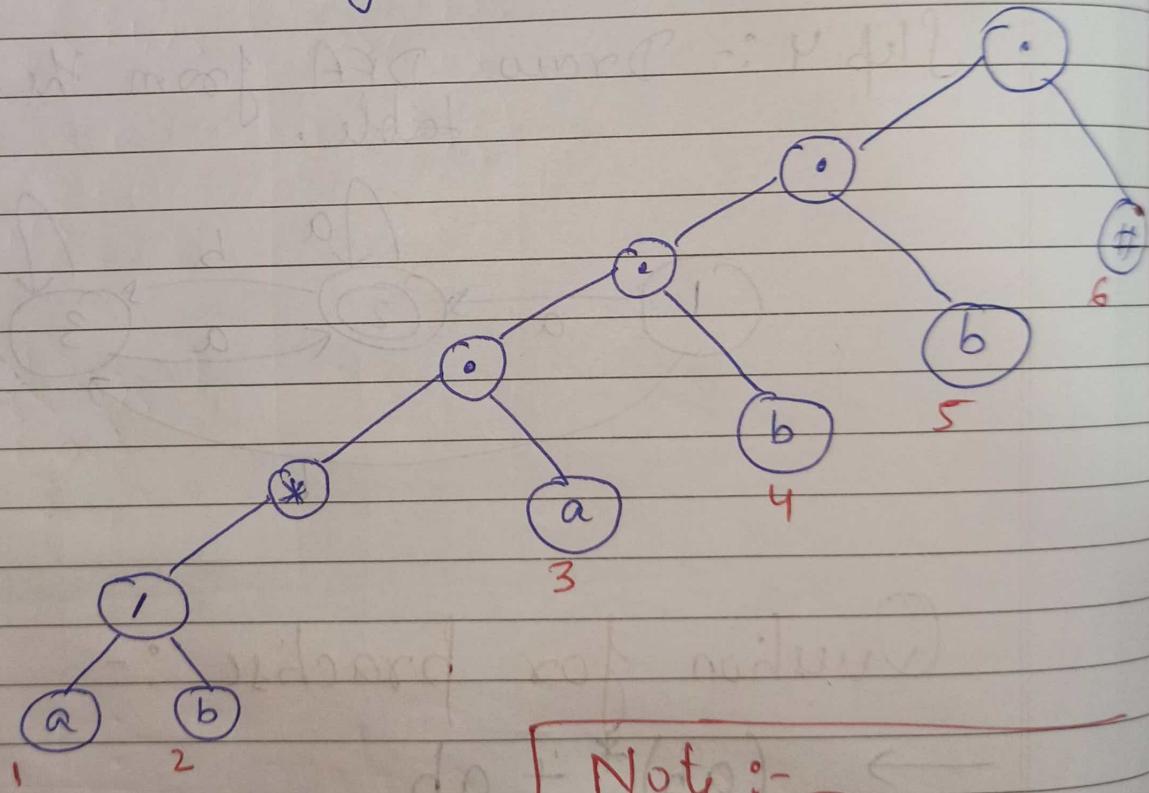
Q  $(a/b)^*abb$

Step 1 :- Augment  $\#$  at the end of regular expression.

$$\gamma' = (a/b)^*abb\#$$

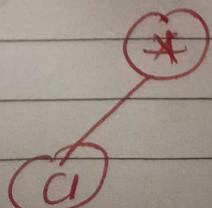
$$\Rightarrow (a/b)^* \cdot a \cdot b \cdot b \cdot \#$$

Step 2 :- Create Syntax tree for  $\gamma'$ .

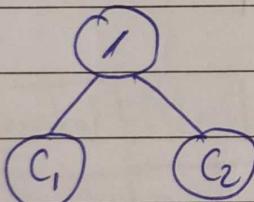
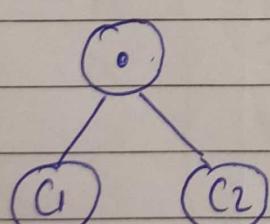
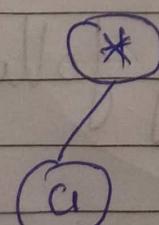


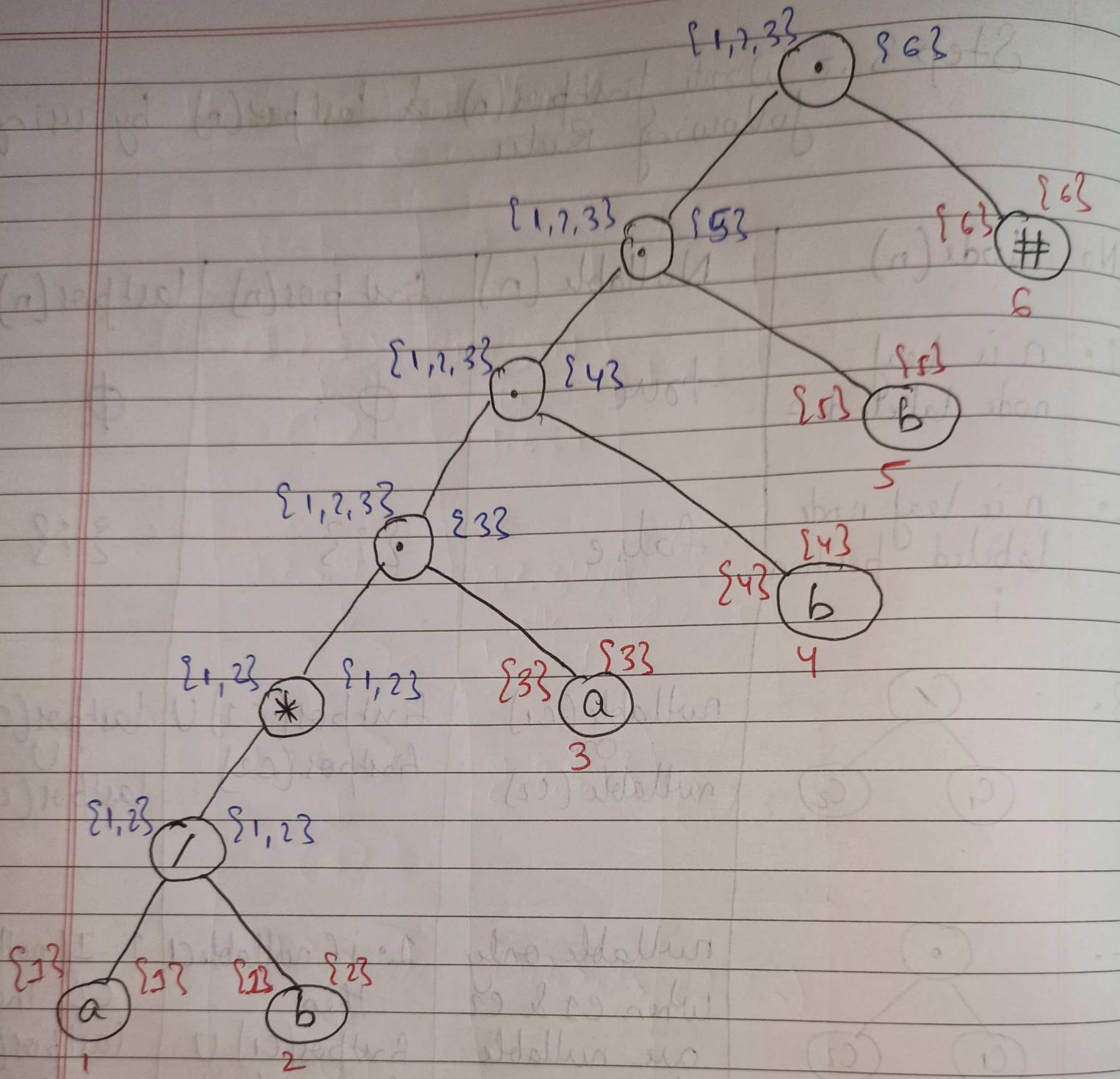
Note :-

ek li (\*) ka child hota h-



Step 3 :- Write  $\text{firstpos}(n)$  &  $\text{lastpos}(n)$  by using following Rules.

S.No	Node(n)	Nullable(n)	$\text{firstpos}(n)$	$\text{lastpos}(n)$
1.	$n$ is leaf node labeled $e$	true	$\emptyset$	$\emptyset$
2.	$n$ is leaf node labeled position ;	false	$\Sigma^3$	$\Sigma^3$
3.	 $\text{nullable}(c_1) \text{ or } \text{nullable}(c_2)$	$\text{nullable}(c_1)$ $\text{nullable}(c_2)$	$\text{firstpos}(c_1) \cup \text{lastpos}(c_1)$ $\text{firstpos}(c_2)$	$\text{lastpos}(c_1) \cup$ $\text{lastpos}(c_2)$
4.	 nullable only when $c_1$ & $c_2$ are nullable else not nullable	nullable only when $c_1$ & $c_2$ are nullable else not nullable	1. if nullable then $\text{firstpos}(c_1) \cup$ $\text{firstpos}(c_2)$ 2. Else $\text{firstpos}(c_1)$	1. If nullable $c_2$ then $\text{lastpos}(c_1)$ $\cup$ $\text{lastpos}(c_2)$ 2. Else $\text{lastpos}(c_2)$
5.		true	$\text{firstpos}(c_1)$	$\text{lastpos}(c_1)$



Step 4: Create a followpos(n) table by using Rule of computing followpos.

Rules for Computing followpos :-

1. If ( $n$  is a star \* node)

for (each  $i$  in  $\text{lastpos}(n)$ )

$$\text{followpos}(i) = \text{followpos}(i) \cup \text{firstpos}(n);$$

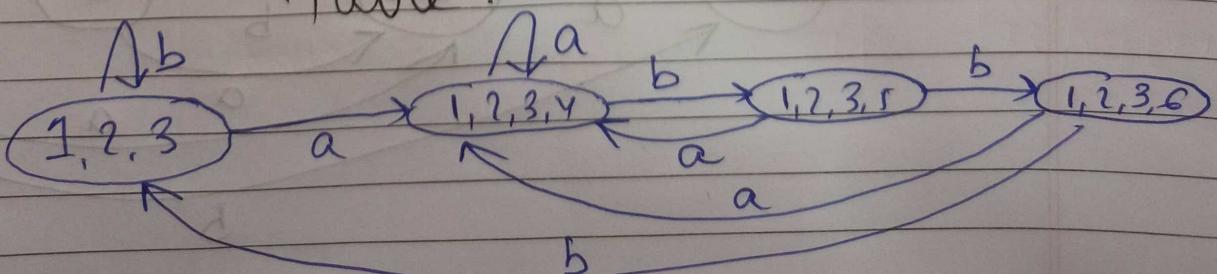
2. If ( $n$  is a cat (.) node)

for (each  $i$  in  $\text{lastpos}(c_1)$ )

$$\text{followpos}(i) = \text{followpos}(i) \cup \text{firstpos}(c_2)$$

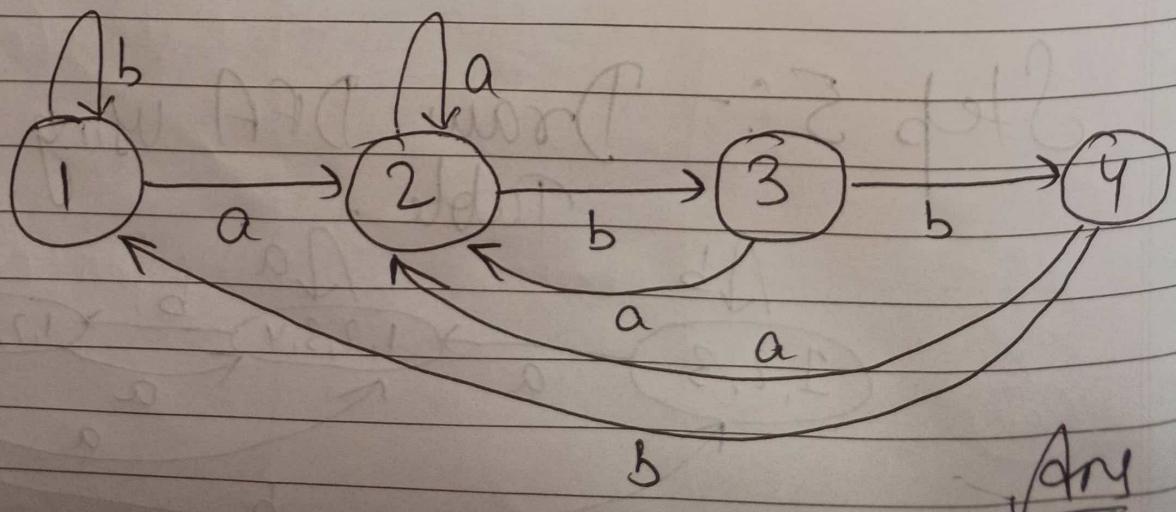
Symbol	Node	followpos
a	1	$\Sigma 1, 2, 3 \}$
b	2	$\Sigma 1, 2, 3 \}$
a	3	$\Sigma 4 \}$
b	4	$\Sigma 5 \}$
b	5	$\Sigma 6 \}$
#	6	—

Step 5:- Draw DFA using followpos Table.



How we got this?

State	a	b
1, 2, 3 ↓ 1	$\{1, 2, 3\} \cup \{g\}$ $= 1, 2, 3, 4$	1, 2, 3
1, 2, 3, 4 ↓ 2	$\{1, 2, 3, 4\} \cup \{f\}$ $= 1, 2, 3, 4$	1, 2, 3, 5
1, 2, 3, 5 ↓ 3	1, 2, 3, 4, 3	1, 2, 3, 6
1, 2, 3, 6 ↓ 4	1, 2, 3, 4, 3	1, 2, 3

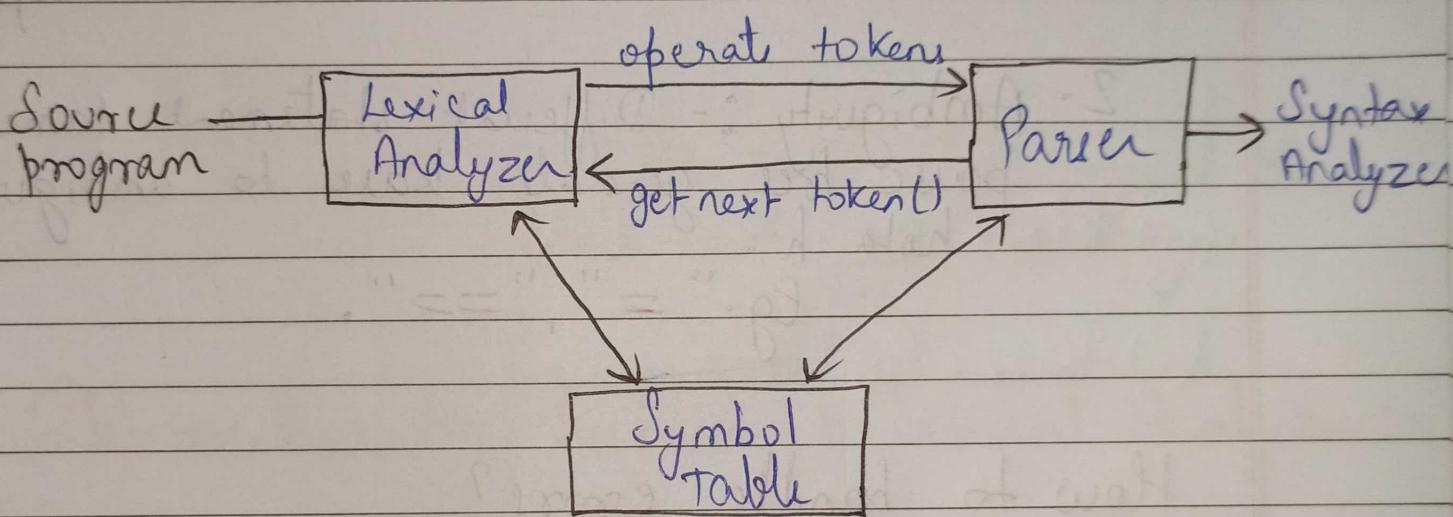


Problem for practice :-

$$\rightarrow (a+b)^* + (a \cdot c)^*$$

$$\rightarrow (a+b)^* \cdot (a \cdot b)^*$$

## Lexical Analyser



\* Lexical Analyzer does :-

(i) Scanning :-

(ii) Tokenization :-

a. Lexeme (charaku)

Eg. int a+b;

→ i, n, t, a, +, b, ;

b. Patterns (rule)

c. Tokens → meaningful words generate karega  
int, a, +, b.

\* Parser :- Lexical analyzer ko request karta h  
ki next line ke tokens generate kare.

Why we need lexical analyzer?  
→ generate tokens & identify typing error

### Issues in Lexical Analysis?

1. Lookahead :- Basically difficulty in deciding where to stop.

Eg.  $a + b + c = d + e$

in this it is not easy to identify when to stop.

2. Ambiguity :- Different statement same fa  
pse code generate have to ambiguity  
note h.

Eg. " = " , " == " .

### How to handle errors?

- ~~Delete~~ Input one char → int  $\xrightarrow{\text{int}}$  ~~del~~
- ~~Delete~~ Input one char → point  $\xrightarrow{\text{point}}$  Input
- Replace one char → int  $\xrightarrow{\text{int}}$  int
- transpose two characters → int  $\xrightarrow{\text{int}}$  int  
 $\uparrow$   
swap.

### Input Buffering :-

Buffer → temporary space where we store data.

i | n | t | a | = | b | + | c | ; | )

↑  
lexeme Begin  
forward

2

- 2 pointers work on this.
- Lexeme Begin which identify 1<sup>st</sup> character of the word
  - forward which moves forward in search of whitespace.

inF Eoc

↑  
Buffer space which match the word with our directory.

## Approaches to Implement Lexical Analyzer:-

- 1) LEX TOOLS
- 2) Basic Programming Language
- 3) Assembly Language

## Specification of Tokens :-

- 1. String
- 2. Language
- 3. Regular Expression

1. String :- Set of characters

(i) Prefix :-

(ii) Suffix :-

(iii) Substring :-

2. Language :-  $L = \{0, 1\}$ ,  $\Sigma = \{a, b, c\}$

(i) Union  $\rightarrow L \cup \Sigma = \{0, 1, a, b, c\}$

(ii) Concatenation  $\rightarrow LS = \{0a, 0b, 0c, 1a, 1b, 1c\}$

(iii) Kleene Closure  $\rightarrow L^*$

(iv) ~~Kleene~~ Positive closure

3. Regular Expression :- used to represent regular language

Q. RE for second element as 0.  
 $\Rightarrow (0+1)0(0+1)^*$

Q. RE for first element as 1 & last element as 0  
 $1(0+1)^*0$

Q. RE for either start with 0 or end with 1.  
 $\Rightarrow 0(0+1)^* + (0+1)^*1$ .