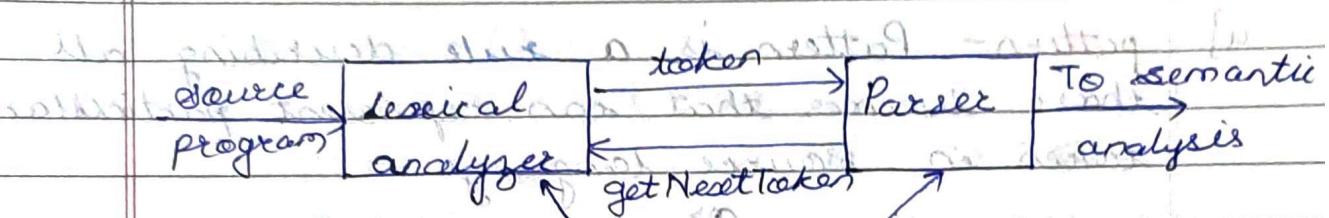


## Unit 2 - Compiler course -

### Lexical analysis

#### \* 1] The role of "lexical analysis" -



why to separate lexical analysis & parsing :-

- ① simplicity of design (tni) from tni
- ② improving compiler efficiency
- ③ enhancing compiler portability tni

2] Tokens - A token is a sequence of characters that can be treated as a unit with single logical entity.

Ex - keywords (for, if, while)

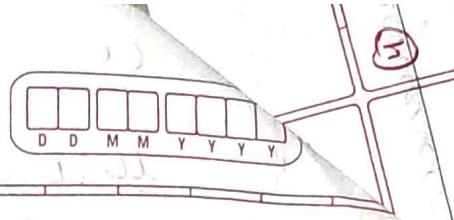
identifiers (variable name, function name)

operators (+, -, \*, /, ++, etc.)

separators (, ; ) etc.

3] lexeme - It is a sequence of characters in source program that is matched by pattern for a token.

or



A sequence of input characters that comprise a single token is called lexeme.  
ex - "float", "-", "723", ";"

- 4) pattern - Pattern is a rule describing all those lexemes that can represent particular token in source language.

or

There are some predefined rules for every lexeme to be identified as a valid token.  
These rules are defined by means of pattern.

- 5) count number of tokens.

① int max(int i);

② int main();

11/2 variables declared b. = 5.555555555555555

a = 10;

return 0;

{ (11+11+1)/2 ) / 2 = 18

= 18 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 = 28  
= 18 comments + not counted as tokens.

③ printf("Never give up");

string = 5 + 1 = 6 string is considered as one token

D	D	M	M	Y	Y	Y

④ printf("1.d Hello", &x);

⑤ int main()

```
{ int a=10, b=20;
    printf("sum is = %.d", a+b);
    return 0; }
```

### 6) Attributes for tokens.

<id, pointer to symbol table entry for E>

<assign-op> assignment operator

<id, pointer to symbol table entry for M>

<mult-op>

<id, pointer to symbol table entry for C>

<exp>, dimension from user input

<number, int value 22> literal example

### 7) Lexical errors.

Some errors are out of power of lexical analyzer to recognize:

fi o (a == f(x))...

problems with string precision or truncation

However it may be able to recognize errors like: d=2x.

Such errors are recognized when no pattern

D	D	M	M	Y	Y	Y
---	---	---	---	---	---	---

for tokens matches to a character sequence.

Types of errors detected by lexical analyzer

- ① long identifiers

Ex - ~~int sum\_of\_tw6 numbers entered by user; { int i;~~

- ② too long numerical literals

Ex - ~~int x = 987654321;~~

- ③ badly formed numerical literals

Ex - ~~char x = Vanita@123~~

- ④ input characters that are not present in source language.

Ex - ~~extra blank before or after , bi>~~

- ⑤ spelling mistakes

Ex - ~~extra punctuation before or after , bi>~~

~~<apn-Hum>~~

~~(extra & unnecessary) drops of remaining , bi>~~

- ① Panic mode: successive characters are ignored until we reach to a well defined token.

~~comes back to~~

- ② delete one character from the remaining input

~~ignores all symbols~~

~~(left = right)~~

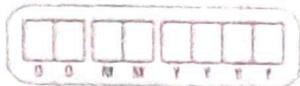
- ③ insert a missing character into remaining input

~~from the number~~

~~as . b . 0 . 1 . 2 . 3 . 4 . 5 . 6 . 7 . 8 . 9 .~~

- ④ replace a character by another character

- ⑤ transpose two adjacent characters



8)

Input buffering and token recognition

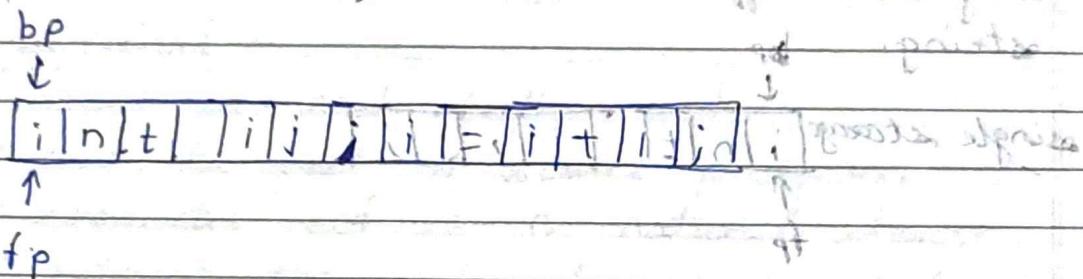
The lexical analyzer scans the input string from left to right one character at a time.

It uses two pointers -

- ① begin pointer (bp) - from where read/write head starts scanning or to read tokens
- ② forward pointer (fp) - which moves a step ahead.

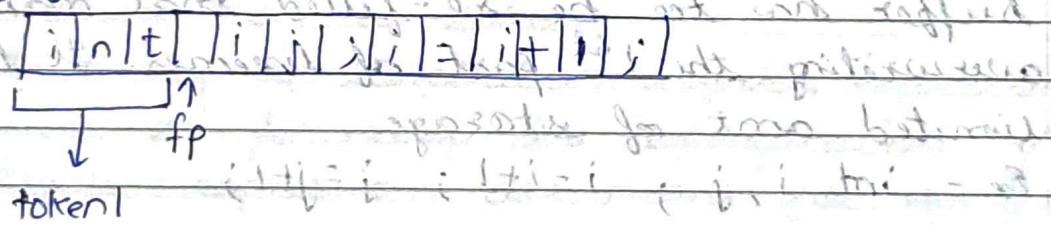
Ex - `int i, j;` in buffer

initial state: `i = int j;` in buffer



Initially both the pointer points the first character of input string.

bp



A forward pointer moves ahead to search the end of lexeme as soon as the blank space is encountered. It indicates the end of lexeme.

The input character thus read from secondary storage but reading in this way from secondary

storage is costly, hence buffering scheme is used.

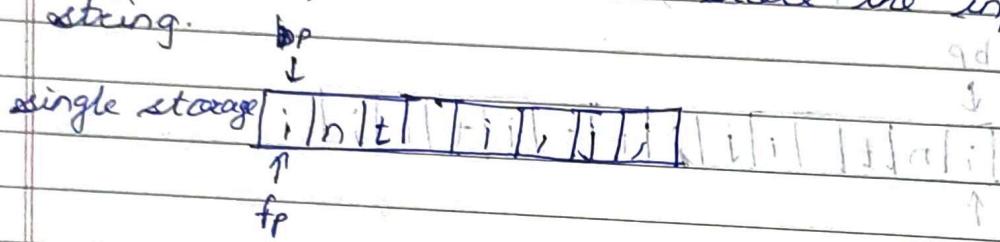
A block of data is first read into a buffer & then scanned by lexical analyzer.

Buffer means secondary storage.

Two methods are used:

- ① one buffer scheme
- ② two buffer scheme.

One buffer scheme :- In one buffer scheme, only one buffer is used to store the input string.



But the problem with this scheme is that if lexemes are very long then it processes the buffer boundary to scan rest of lexemes. The buffer has to be re-filled that makes overwriting the 1st part of lexemes. It has limited amt of storage.

Ex - int i, j; i = i + 1; j = j + 1;

Input

Two buffer scheme :- 2 buffers are used to store the input string. 1<sup>st</sup> & 2<sup>nd</sup> buffer are scanned alternatively. When the end of 1<sup>st</sup> buffer is used, 2<sup>nd</sup> buffer is filled.

D	D	M	M	T	T	T
---	---	---	---	---	---	---

Ex - int i, j; i = i + 1; j = j + 1;

[i | n | t] [i = i | + | 1] [buffer 1]  
nbp

[j | j = j | + | 1; | Eof] [buffer 2]  
↑  
fp

2<sup>nd</sup> buffer is used to improve execution speed of buffer.

token worth says hi → token  
token says token worth says hi

sentinel -

- added at each buffer end
- can't be part of source program
- character eof is a natural choice.
- retains the role of entire input end.
- when appears other than at end of buffer it means that input is not an end.

9)

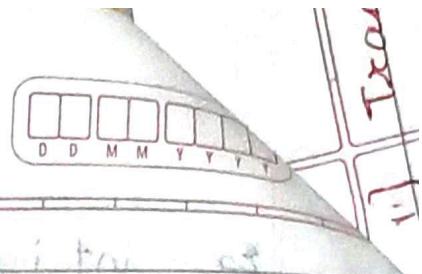
specification of tokens. [D-O] ← token

In theory of compilation regular expressions are used to formalize the specification of tokens. Regular expressions are means for specifying regular lang. \*(digit|letter)\* letter ← hi

Ex - letter -(letter|digit)\*

Each reg exp is a pattern specifying the form of strings.

letter ← hi



## 10] Recognition of tokens in relation to

token type	token value	state
------------	-------------	-------

tells the category of token - ex- identifier, constant, etc.

gives info about token.

also known as token value.

category of token here is suff.

stmt → if expr then stmt

| if expr then stmt else stmt  
| E

- Identifier

has suff. class to identify.

expr → term relop term term at start.

term → n or identifier

has suff. class to identify.

term → id next suff. suggests number

has no number. suffix looks same to

digit → [0-9]

digits → digits digit suff. for word of

number → digit (digit)\* | E[+/-] digit |

letter → [A-Z a-z]

id → letter (letter | digit)\*

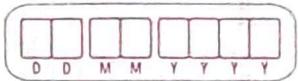
if → if

then → then

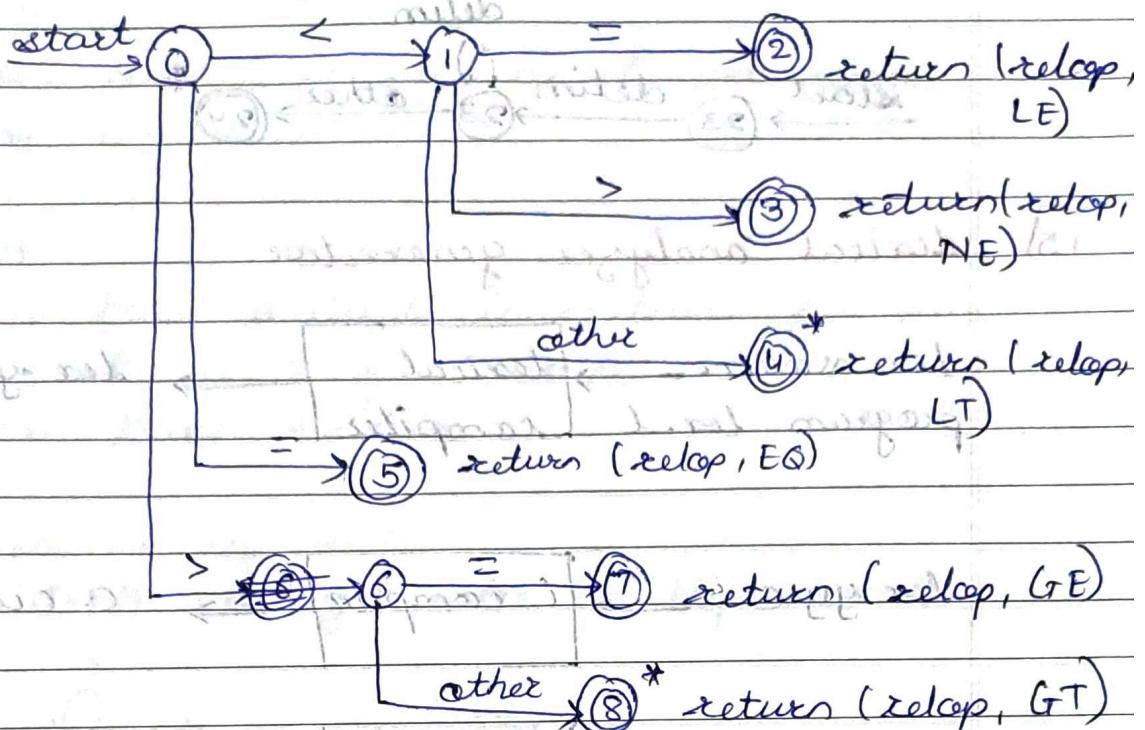
else → else

relop → < | > | <= | >= | = | <= | >=

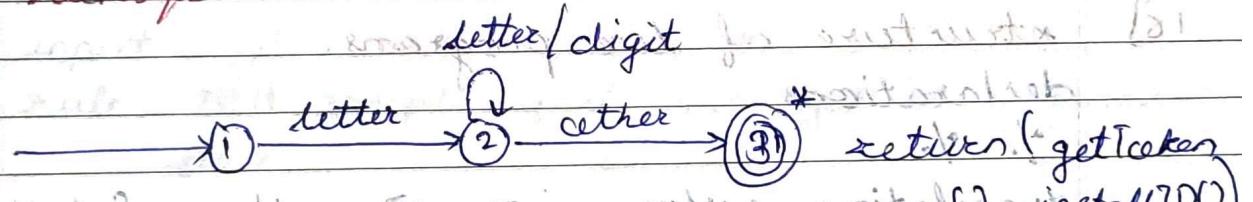
ws → (blank | tabl newline) \*



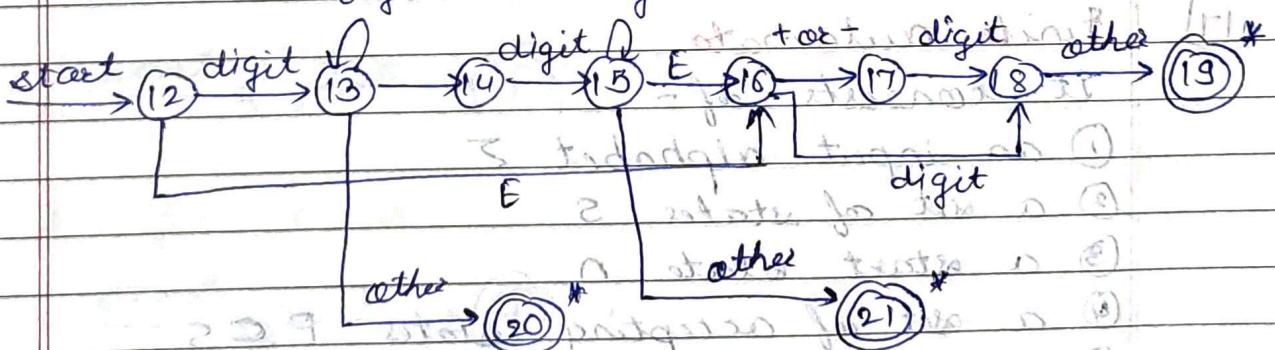
## 1] Transition diagram for relax

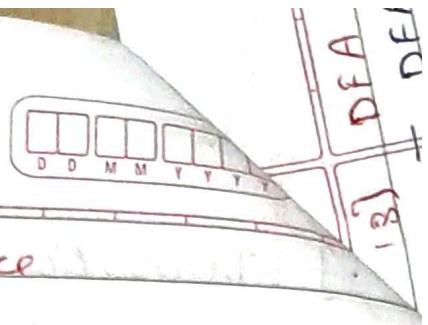


12] Transition diagram for reserved words & identifiers.



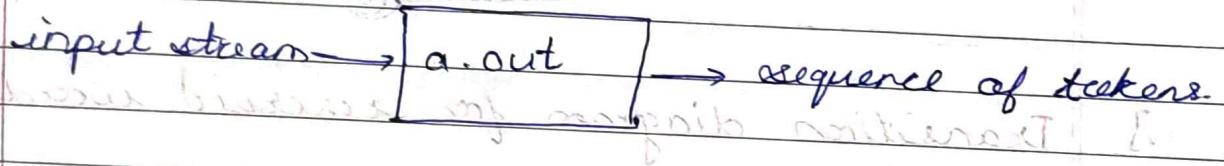
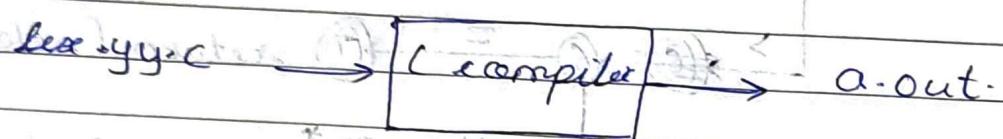
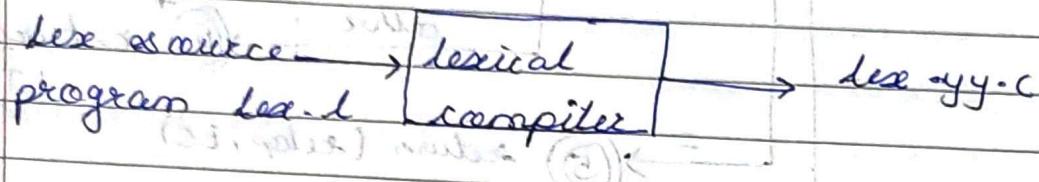
13) Transition diagram for unsigned numbers  
digit      digit





14) Transition diag for whitespace  
delim

15) lexical analyzer generator.



16) structure of lex programs.

declarations

% . % .

translation rules

% . % .

→ pattern { actions }.

auxiliary functions of memory management for

17) finite automata

It consists of -

- ① an input alphabet  $\Sigma$
- ② a set of states  $S$
- ③ a start state  $n$
- ④ a set of accepting states  $F \subseteq S$
- ⑤ a set of transitions state  $\rightarrow$  input state

D	D	M	M	Y	Y	Y
---	---	---	---	---	---	---

18) DFA & NFA.

DFA.

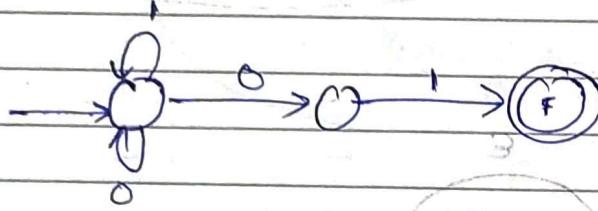
- one transition per input per state
- no  $\epsilon$ -moves.

NFA

- can have multiple transitions for one input in a given state
- can have  $\epsilon$ -moves.

19) acceptance of NFA.

An NFA can get into multiple states.



input : 101

rule: NFA accepts if it can get in a final state.

20) Regular exp to NFA (Thompson's construction)

① for  $\sigma = \epsilon$



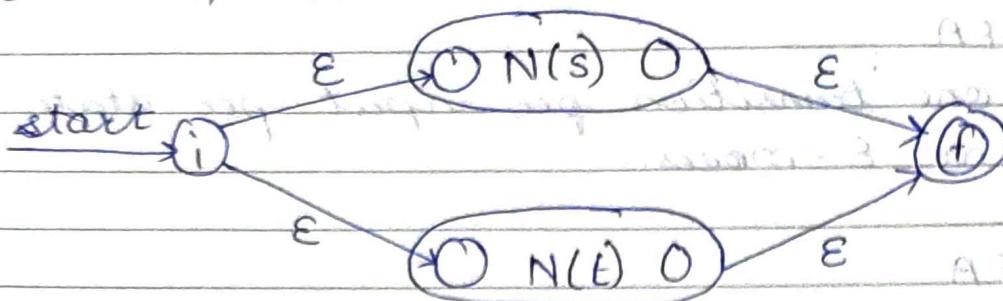
② for  $\sigma$  reg exp.



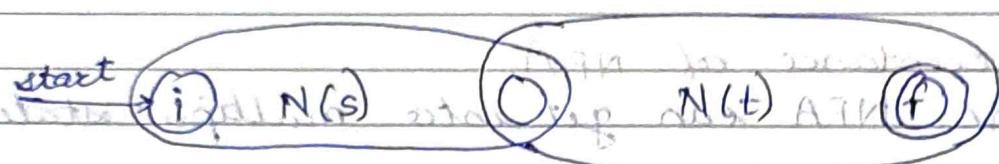
③ for  $\sigma = a$



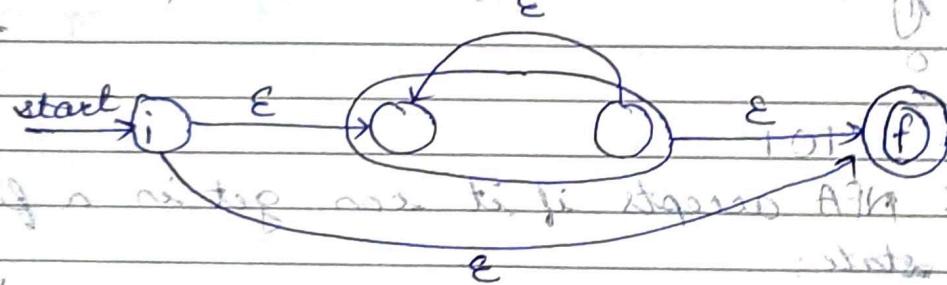
③  $s|t$ ,  $N(s|t)$ .



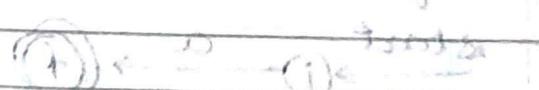
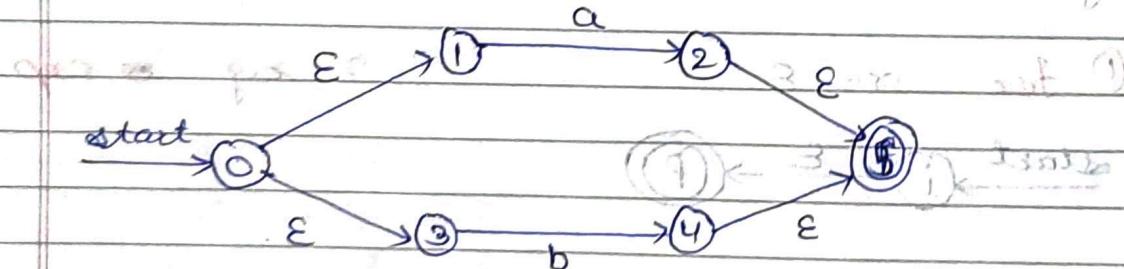
4)  $N(st)$



5)  $s^*$ ,  $N(s^*)$

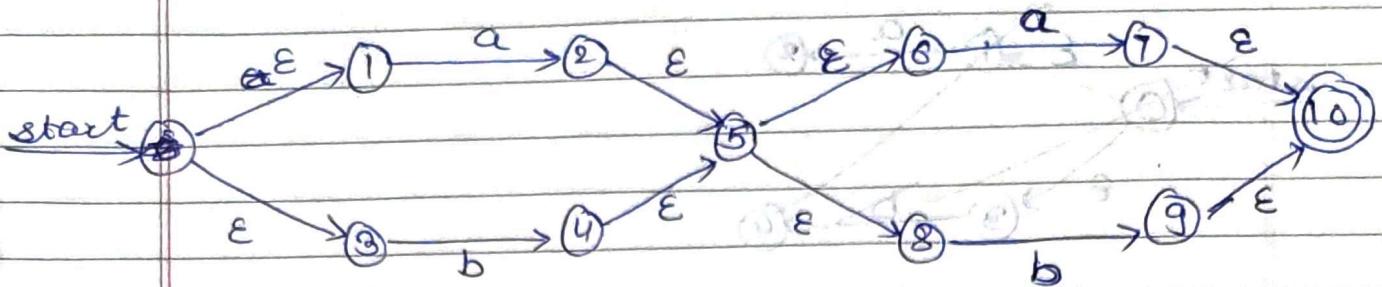


6)  $(a|b)^*$ ,  $N((a|b)^*)$

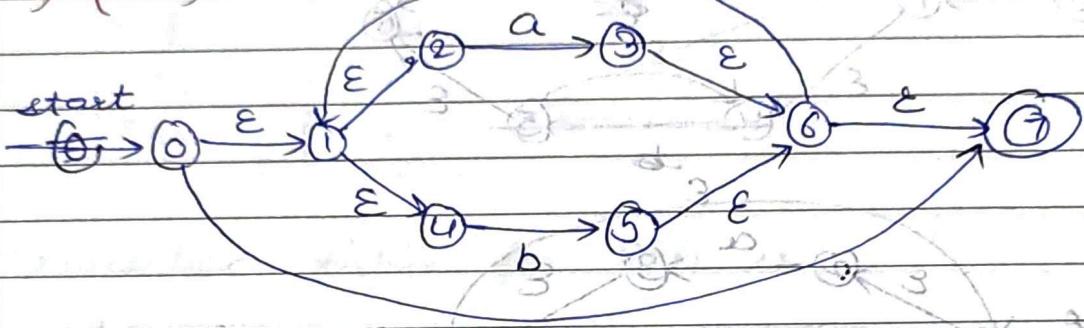


DD	MM	YY	YY
----	----	----	----

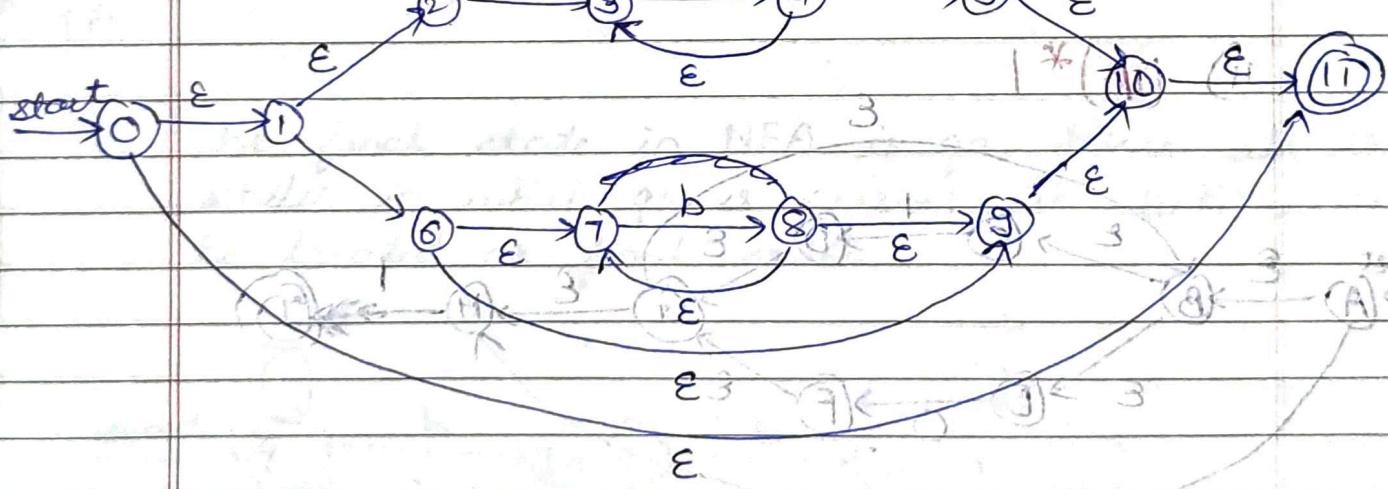
7)  $(a|b)(a|b)$ .

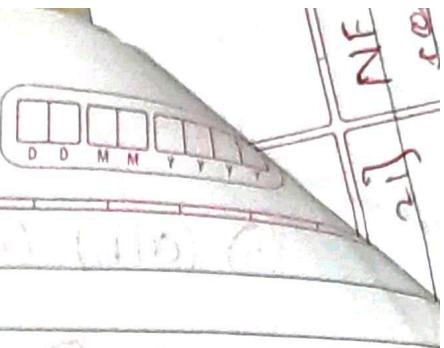


8)  $(a|b)^*$ .

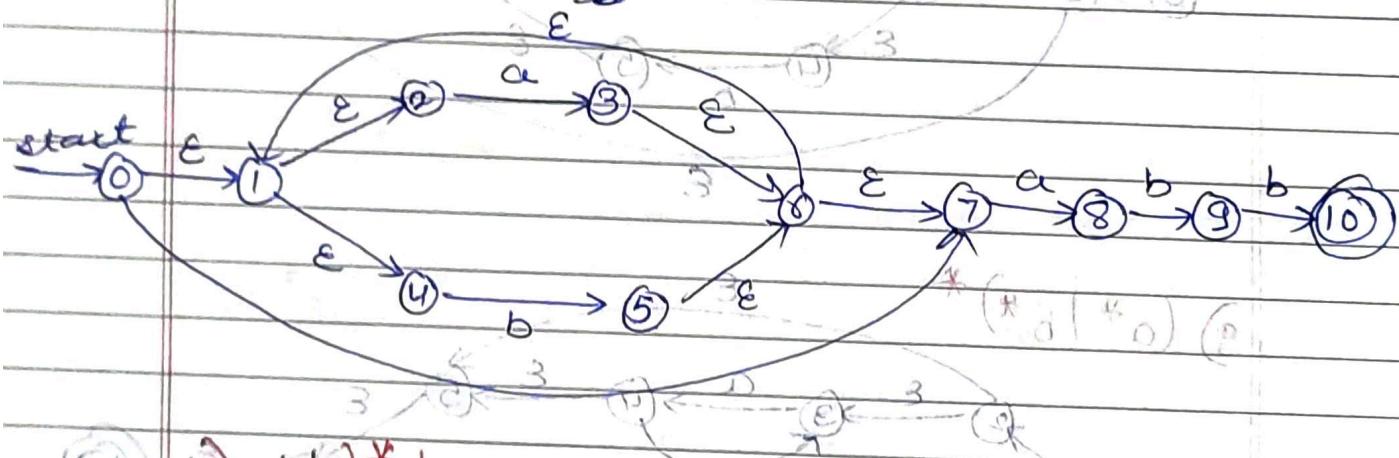
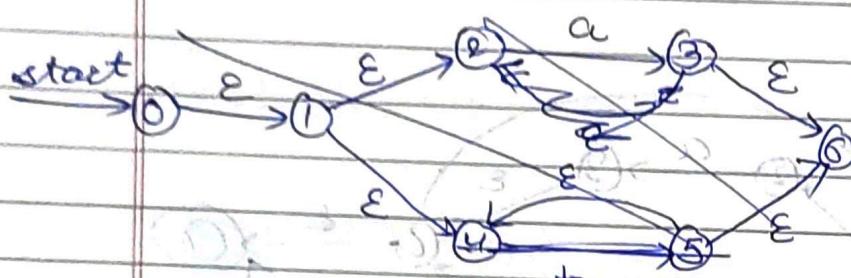
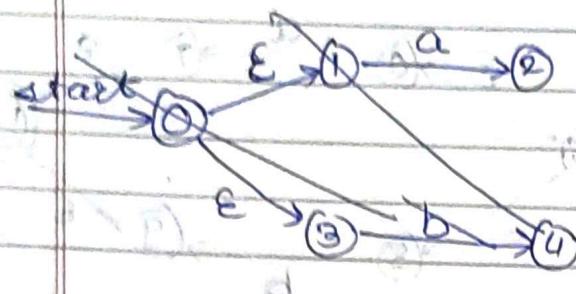


9)  $(a^*|b^*)^*$ .

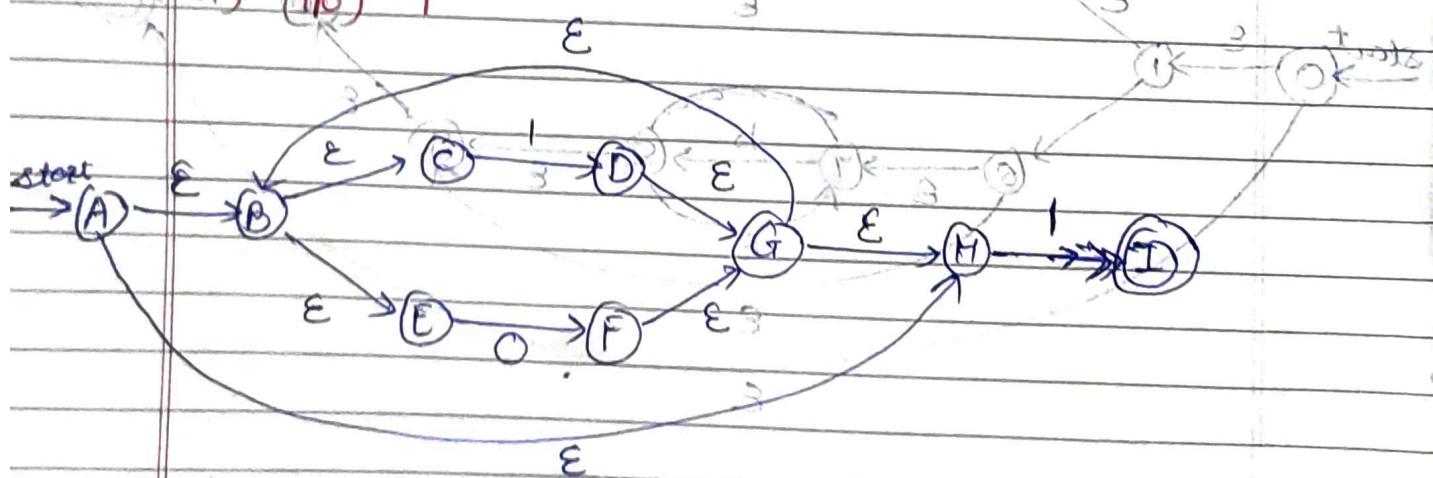




10)  $(a|b)^* \cdot abb$



11)  $(1|0)^* \cdot 1$



D	D	M	M	Y	Y	Y
---	---	---	---	---	---	---

no. of states

Max number of states in DFA = 2

Q. Explain CC tools with example.

\* → There are many types of software tools that have been developed to create one or more phases of compiler. These tools are called compiler construction tools. These tools use specialized language for defining & executing components & can use algorithms that are completely refined.

Following are examples of CC tools -

① Scanner generators -

This tool takes regular exp. as input & creates lexical analysis. The fundamental lexical analysis is produced by finite automata.

Ex - LEX

② Parser generators -

This software produces syntax analysis which takes input in the form of syntax of programming lang. depending on context free grammar.

③ Syntax-directed translation engines -

These software tools offer an intermediate code by using parse tree. It has goal of associating one or more translations with each node of parse tree.



#### ④ automatic code generators -

Takes intermediate code & converts them into machine language.

Q.

What is role of lexical analysis?

- \* → - The lexical analysis is first phase of compiler where a lexical analyzer operate as an interface between source code & rest of phases of compiler.
- It reads the input characters of the source program, groups them into lexemes and produces sequence of tokens for each lexeme.
- The tokens are sent to parser for syntax analysis.
- The lexical analyzer also interacts with the symbol table while passing tokens to parser.

Q. Write note on attributes of token.

→

- A token is pair of a token name & an optional token value.
- A pattern is a description of the form that lexeme of a token may take.
- A lexeme is a sequence of characters in the source program that matches the pattern for a token.
- When more than one lexeme matches a pattern, the lexical analyzer must provide additional information about the particular lexeme that matched to the subsequent phases of the compiler.
- The lexical analyzer collects information about tokens into their associated attributes.
- Example - The program statement  
 $E = M^* C^{**} 2$ .

yields the following token-attribute pairs :-

- ① <id, pointer to symbol table entry for E>
- ② <assign - op>
- ③ <id, pointer to symbol table entry for M>
- ④ <mult - op>
- ⑤ <id, pointer to symbol table entry for C>
- ⑥ <exp - op>
- ⑦ <number, integer value 2>.