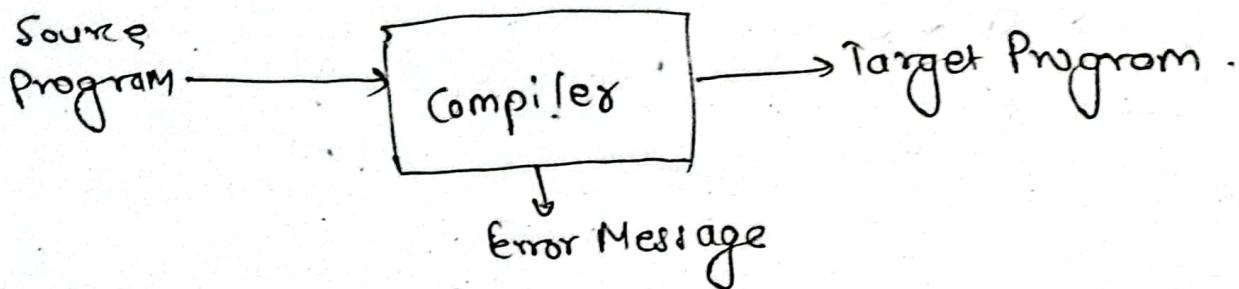


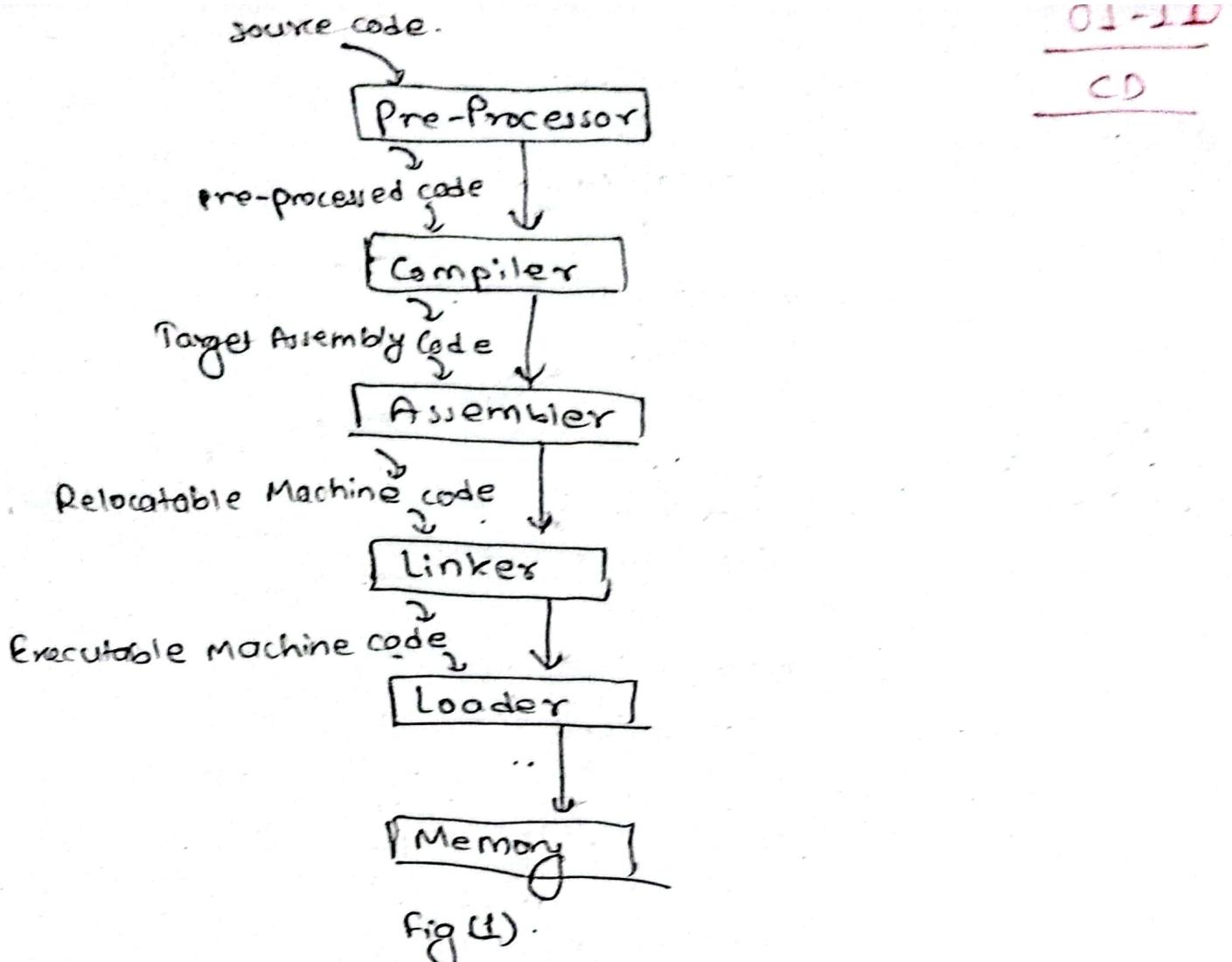
Unit-1: Overview of Compilation* Compiler:

- ↳ It is a program that reads a program written in source language and translates it into target language.
- ↳ It also shows the errors in the source program.



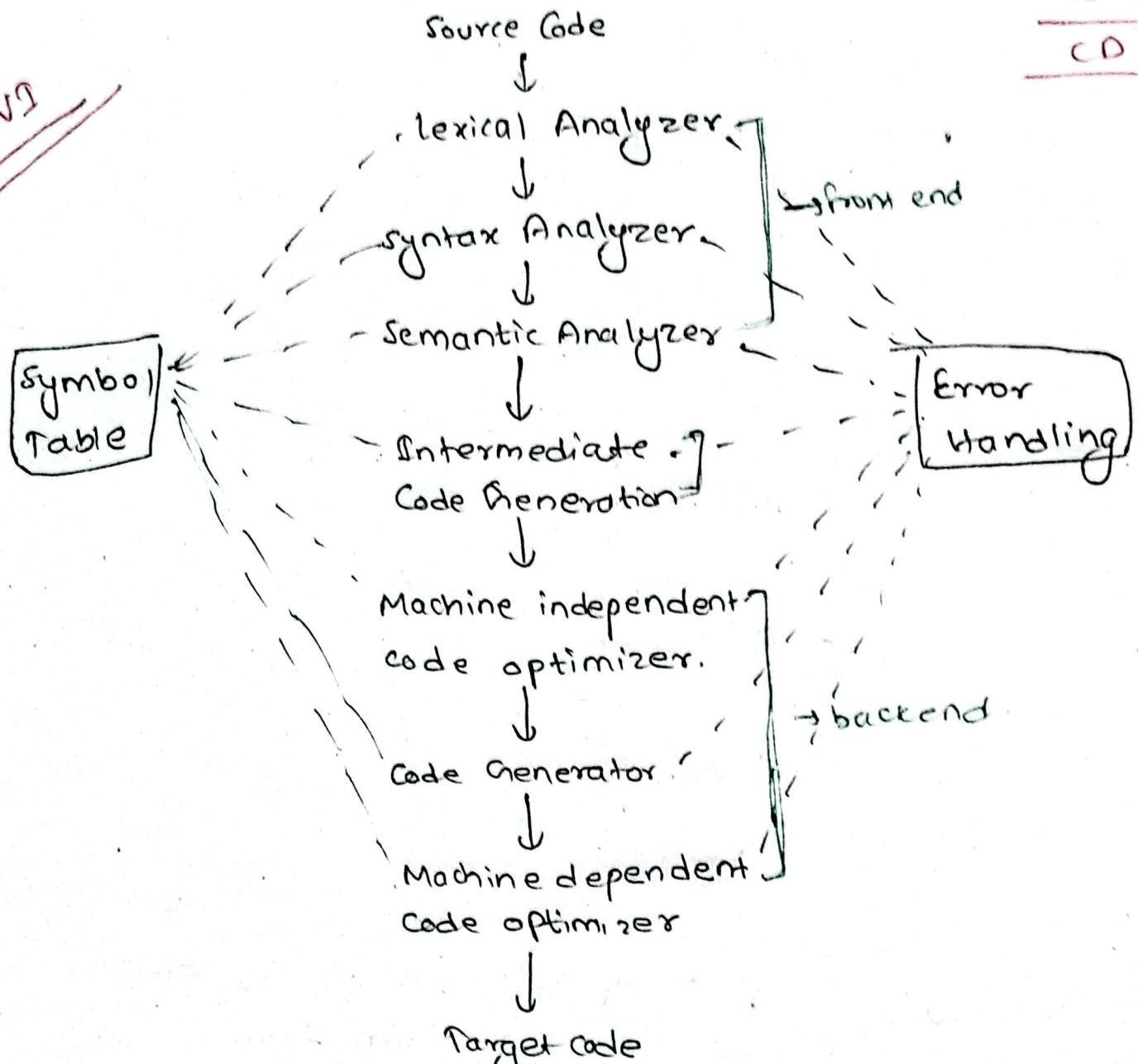
- ↳ There are two parts in compilation
 - i) Analysis
 - ii) Synthesis
- ↳ In Analysis, source program is broken down into constituent pieces and creates an intermediate representation of src program.
- ↳ In synthesis, desired target program is constructed from intermediate representation.
- ↳ To create an executable program, along with main compiler, other programs are also required as shown in fig(i).

CCD



→ Conceptually, a compiler operates in phases, each of which transforms the source program from one representation to another as in fig(2).

Code



fig(2).

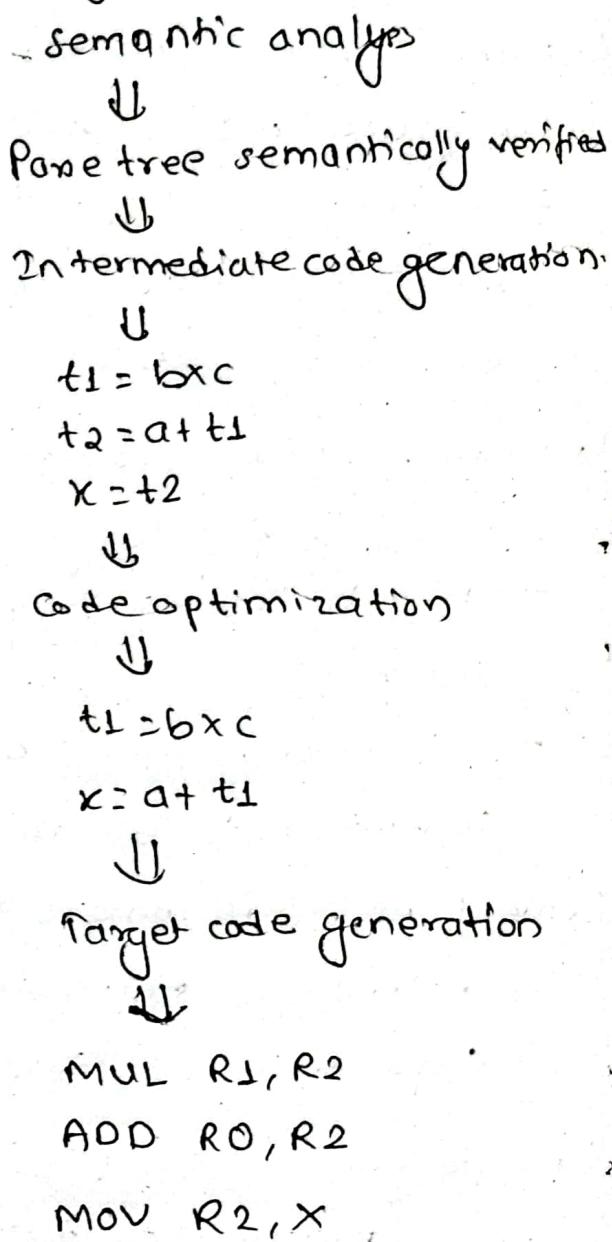
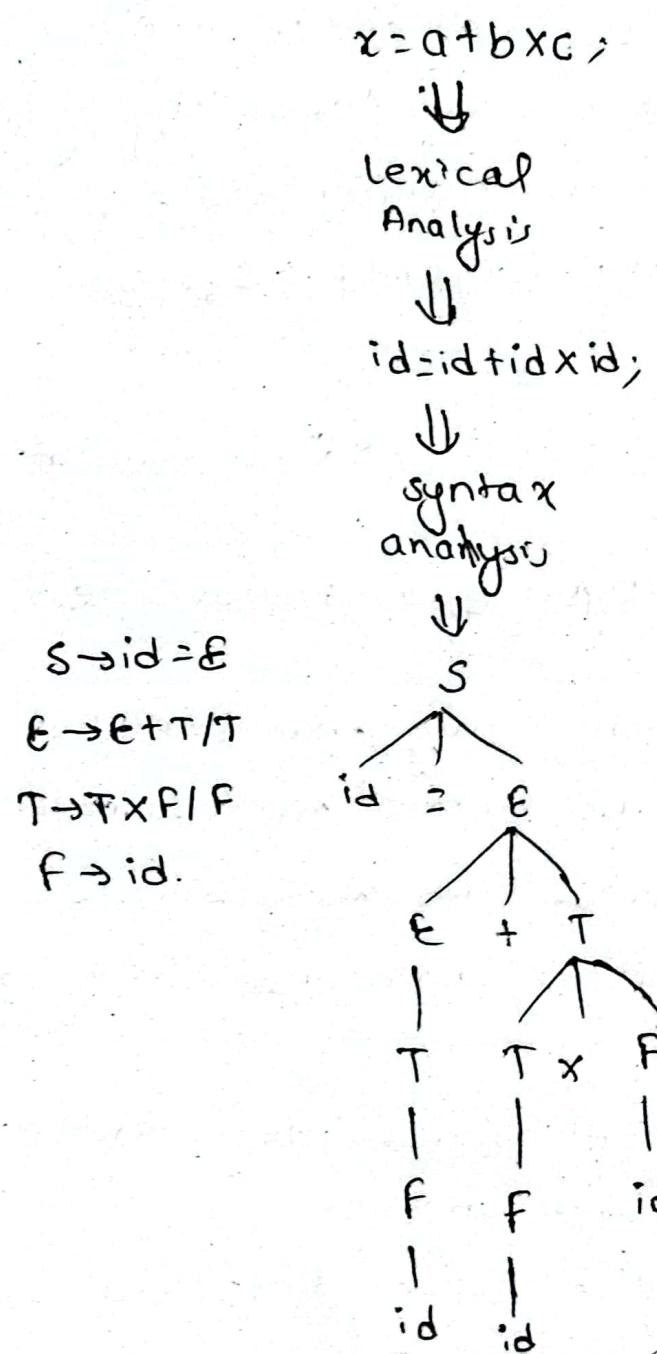
- ↳ first phase of a compiler is called lexical analyzer or scanning.
- ↳ lexical analyzer reads the stream of characters of source code and groups them into meaningful sequence called lexemes.
- ↳ for each lexeme, lexical analyzer produce as o/p a token of form <token name, attribute value>.

- ↳ Second phase of a compiler is called syntax analyzer or parsing. Q3-10
CD
- It takes the token produced by lexical analysis as i/p and generates a parse tree.
 - Here, token arrangements are checked against source code grammar.
- ↳ Next phase is semantic analysis which checks whether the parse tree constructed follows the rules of language.
- It keeps track of identifiers, their types and expression, whether identifiers are declared before use or not, etc.
- ↳ In next phase, compiler generates intermediate code of source code for target machine.
- It is between high level language and machine language.
 - It should be generated in such a way that it makes it easier to be translated into target machine code.
- ↳ Next phase is code optimization of intermediate code.
- It includes removing unnecessary code lines, arranging sequence of statements, etc.
- ↳ Code generation phase takes optimized intermediate code and map to target machine language.
- Code generator translates the intermediate code to sequence of relocatable machine code.

↳ symbol table store information about keyword 01-10
 and tokens found during lexical analysis. It is CD
 computed in almost all phase of compiler.

↳ Error handler handles all errors encountered.

Eg:



Unit-2:Scanner*Lexical Analysis

- ↳ initial part of reading and analyzing the program text.
- ↳ the text is read and then divided into tokens, each of which corresponds to a symbol in programming language. e.g.: variable name, keyword, or number, etc.
- ↳ lexical analyzer or lexer will discard comments and whitespaces.
- ↳ A lexical analyzer also called as scanner has the following functionality and characteristics.
 - 1) Its primary function is to convert from a sequence of characters into sequence of tokens.
 - 2) it must identify and categorize specific character sequence into tokens.
 - 3) handle lexical errors (illegal characters) by reporting them intelligently to users.
 - 4) token categories can be specified using regular expression.
 - 5) lexical analyzers can be written by hand or implemented automatically using finite automata.

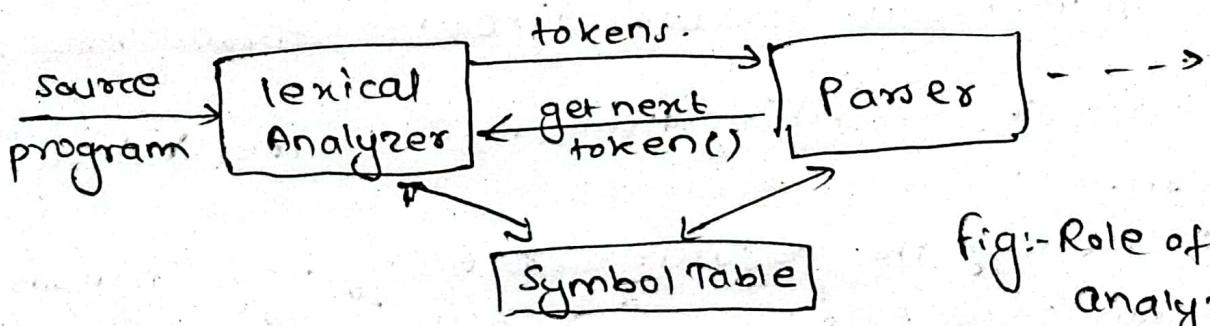


fig:-Role of lexical analyzer.

*Roles:

- ↳ lexical analyzer read the input stream one character at a time and translate it into valid tokens.
- ↳ the lexical analyzer works in lock step with the parser
- ↳ the parser request lexical analyzer for the next token whenever it requires one using getnexttoken().
- ↳ lexical analyzer may perform other operations like removing redundant white spaces, removing token separators, removing comments, providing line number to parser for error reporting.

*Tokens, patterns, lexemes

- ↳ A token is a single word of source code input.
- ↳ When a string representing a program is broken into sequence of substrings, such that each substring represents a constant, identifier, operator, keyword, etc of the language, these substrings are called tokens.
- ↳ they are building block of programming language.
eg:- if, else, identifier, etc.
- ↳ lexemes are the actual string matched as token.
—they are the specific characters that make up of a token. eg:- abc, 123.
- ↳ A token can represent more than one lexemes ie, token intnum can represent lexemes like 123, 4460, 78129, etc.
- ↳ Patterns are rules describing the set of lexemes belonging to a token.
—this is usually the regular expression . eg:- intnum token can be defined as $[0-9][0-9]^*$

<u>Token</u>	<u>Informal Definition</u>	<u>Sample lexemes</u>	01-15/ CD
if	characters i, f	if	
else	characters e, l, s, e	else	
comparison	< or > or <= or >= or == or !=	<=, !=	
id	letters followed by letters & digits.	pi, score, D2	
number	any numeric constant	3.1415, 0, 6.02e ²³	
literal	anything but surrounded by " "	"Hello"	

fig:- Examples of tokens.

Eg:

printf("Total = %d", score);

Here,

- ↳ printf and score are lexemes matching pattern of for token id.
- ↳ "Total = %d", is a lexeme matching taken literal.

↳ In many programming languages, the following classes are cover most or all of the tokens.

- i) one token for each keyword. here pattern for a keyword is same as keyword itself.
- ii) tokens for the operators either individually or in a classes such as token comparison is maintained.

iv) one or more tokens representing constants,
such that numbers ~~are~~ and literal strings.

03-15
CD

v) token for each punctuation symbol such as* left and
right parenthesis, comma, etc .

* Attributes of tokens

↳ when more than one lexeme can match a pattern,
the lexical analyzer must provide additional information about particular lexeme that matched.

↳ i.e., if there is more than one lexeme for a token,
we need to put extra information about token which
is called attribute of token.

↳ for eg:- token number matches both 0 and 1 but it is
extremely important for code generator to know
which lexeme was found in source program. Hence,
in this case, lexical analyzer must be able to represent
0 and 1 differently.

↳ for token id, information about an identifier is kept
in the symbol table so appropriate attribute value
for an identifier is a pointer to symbol table
entry for that identifier.

↳ eg:- ~~sum = val1 + val2.~~

Here,

token and associated attribute values are written as
<id, pointer to symbol table entry for ~~sum~~>
<assign-op>

<id, pointer to symbol table entry for val1>:

<add-op>

<id, & pointer to symbol table entry for val2>

01-25
CD

082-01-16

* Lexical errors

↳ Lexical analyzer, without help of other component cannot easily say there is error.

Let in C, we have

'fi (a==b)....'

Here,

lexical analyzer can't tell whether fi is misspelling of keyword if or an undeclared function identifier.

- For fi, the lexical analyzer return token id.

↳ Error in this phase is found when there is no matching string found as given by the pattern.

- When error occurs, the lexical analyzer must not halt the process so it can print the error message and continue.

* Specification of tokens

↳ Regular expressions are an important notation for specifying lexeme patterns.

↳ Each pattern matches a set of strings, so regular expression will serve as name for set of strings.

Some defⁿ:

* Alphabets

↳ An alphabet A is a set of symbols.

* strings:

↳ A string is a finite sequence of characters from alphabet A.

↳ The length of string 'w' is denoted by $|w|$ i.e. number of characters in w.

* Kleene Closure:

↳ Kleene closure over an alphabet A denoted by A^* is set of all strings of any length (0 also) possible from A.

$$A^* = A^0 \cup A^1 \cup A^2 \cup \dots$$

* language:

↳ Language L is set of strings such that $L \subseteq A^*$.

Operation

Union of L & M.

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

Concatenation of L & M.

$$LM = \{ st \mid s \text{ is in } L \text{ and } t \text{ is in } M \}$$

Kleene closure of L

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

Positive closure of L

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

fig:- Defⁿ of operations on language.

* Regular Expression (RE):

↳ is used to describe the tokens of programming language.

↳ has the capability to express finite languages by defining a pattern for finite strings of symbols.

- ↳ The grammar defined by RE is known as Regular Grammar.
- ↳ The language defined by Regular Grammar is called Regular Language.
- ↳ If r and s are regular expression denoting languages $L(r)$ and $L(s)$ respectively, then
 - $r+s$ is a regular expression denoting $L(r) \cup L(s)$.
 - rs is a regular expression denoting $L(r) \cdot L(s)$.
 - r^* is a regular expression denoting $(L(r))^*$.

Properties of regular expressions

- $r+s = s+r$; $+$ is commutative.
- $r+(s+t) = (r+s)+t$; $+$ is associative.
- $r(st) = (rs)t$; concatenation is associative.
- $r(st) = (rs)+rt$; concatenation distributes over $+$.
 $(s+t)r = sr+tr$
- $\epsilon r = r\epsilon = r$; ϵ is identity element for concatenation.
- $r^* = (r+\epsilon)^*$; ϵ is guaranteed a closure.
- $r^{**} = r^*$; $*$ is idempotent.

*Recognition of tokens:

- ↳ A recognizer for a language is a program that takes string ' w ' and answer 'yes' if w is a string of that language, otherwise 'No'.
- ↳ The tokens that are specified by RE are recognized by using FA.

- ↳ Here, transitions are defined and we start from Q1-Q6 start state of FA using transition diagram.
- if the transition leads to the accepting state, then the \$ token is matched and lexeme is returned.
- Otherwise, other transition diagram are tried out until failure is detected.

- ↳ Recognizer of tokens take the language L and string S as input and try to verify whether $s \in L$ or not.

*Recognizing words:

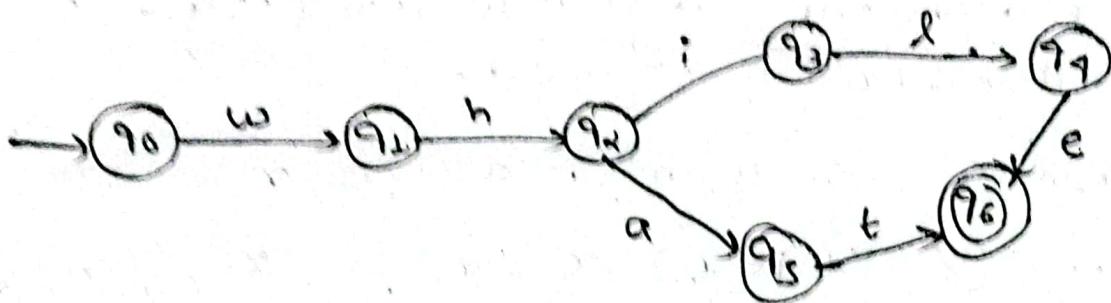
- ↳ is done by character by character formulation.
- ↳ Assume Nextchar() returns a character, sample code for recognizing the keyword, 'while' is like below.

```

C ← Nextchar();
if (C = 'w')
    then begin;
        C ← Nextchar();
        if (C = 'h')
            then begin;
                C ← Nextchar();
                if (C = 'i')
                    then begin;
                        C ← Nextchar();
                        if (C = 'l')
                            then begin;
                                C ← Nextchar();
                                if (C = 'e')
                                    then success;
                                    else failure;
                                end;
                                else failure;
                            end;
                            else failure;
                        end;
                        else failure;
                    end;
                    else failure;
                end;
                else failure;
            end;
            else failure;
        end;
        else failure;
    end;
    else failure;
end;
else failure;
    
```

↳ To recognize multiple words, we can create multiple edges, that leave a given state.

↳ e.g. one recognizer for both 'while' and 'what' be



Formalism for recognizer

↳ is through finite automata where FA is a tuple $(Q, \Sigma, \delta, q_0, f)$

↳ For above e.g. of 'while' or 'what' recognizer FA is like.

$$Q = \{q_0, q_1, q_2, q_3, q_4, q_5, q_6\}$$

$$\Sigma = \{w, h, i, l, e, a, t\}$$

$$q_0 = \{q_0\}$$

$$f = \{q_6\}$$

δ :

$$\delta(q_0, w) = q_1$$

$$\delta(q_1, h) = q_2$$

$$\delta(q_2, i) = q_3$$

$$\delta(q_2, a) = q_5$$

$$\delta(q_3, l) = q_4$$

$$\delta(q_4, e) = q_6$$

$$\delta(q_5, t) = q_6$$

Transition table

Q/Σ	w	h	i	l	e	a	t
$\rightarrow q_0$	q_1	-	-	-	-	-	-
q_1	-	q_2	-	-	-	-	-
q_2	-	-	q_3	-	-	q_5	-
q_3	-	-	-	q_4	-	-	-
q_4	-	-	-	-	q_6	-	-
q_5	-	-	-	-	-	-	q_6
$\star q_6$	-	-	-	-	-	-	-

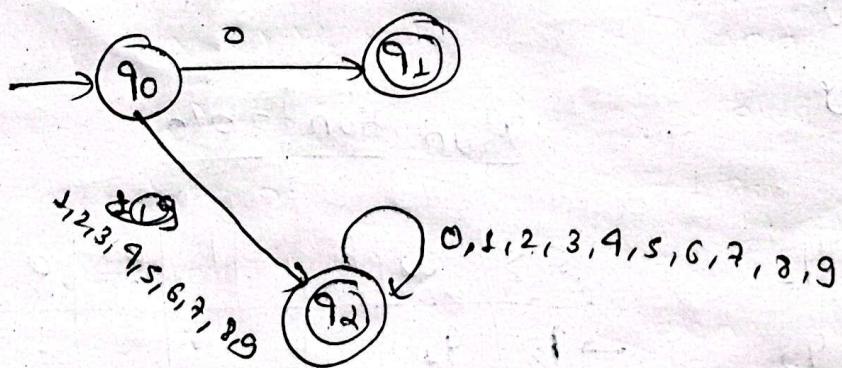
⑧ How does FA accept strings?

- The symbols of any strings are passed through starting state of FA individually. If at the end of string, it reaches accepting state then is accepted otherwise, rejected.
- If x is $x_1 x_2 x_3 \dots x_n$ then FA accepts x if and only if $\delta(\delta(\dots \delta(\delta(\delta(q_0, x_1), x_2), x_3), \dots x_{n-1}), x_n) \in q_f$.

082-01-202

*Recognizing more complex words:

- for unsigned integer(either 0 or it is a series of one or more digits where first digit is from 1 to 9 and subsequent digits are from 0 to 9).



*Closure properties of regular expression:

- Regular expressions are closed under

- union
- concatenation
- kleene closure

i) Union

↳ if r and s are two regular expressions, then
 $r \cup s$ is also regular expression.

ii) Concatenation

↳ if r and s are two regular expressions, then $r \cdot s$
is also regular expression.

iii) Kleene closure

↳ if r is regular expression, then r^* is also regular
expression.

→ There are two types of finite Automata.

i) Deterministic Finite Automata (DFA):

↳ DFA is deterministic, if there is exactly one transition
for each (state, input) pair.

↳ DFA is five tuple $(Q, \Sigma, \delta, q_0, F)$ where,

$Q \rightarrow$ finite set of states

$\Sigma \rightarrow$ finite set of inputs

$q_0 \rightarrow$ a start state

$F \rightarrow$ set of final states, $F \subseteq Q$.

$\delta \rightarrow$ transition fcn

$\delta: Q \times \Sigma \rightarrow Q$.

Implementing DFA:

↳ following is the algorithm for simulating DFA
to recognizing given string.

↳ for a given string w , in DFA D , with start state
 q_0 , the o/p is 'Yes' if D accepts w , else 'No'.

DFA sim(D, q₀)

```

q = q0;
c = getChar();
while(c != eof) {
    q = move(q, c);
    c = getChar();
}
if(q is in F)
    return 'Yes';
else
    return 'No';

```

Q

ii) Non-Deterministic Finite Automata (NFA):

- ↳ FA is non-deterministic if there is more than one transition for each (state, input) pair.
- ↳ NFA is five tuple $(Q, \Sigma, \delta, q_0, F)$ where,
 - $Q \rightarrow$ finite set of states
 - $\Sigma \rightarrow$ finite set of inputs
 - $q_0 \rightarrow$ a start state
 - $F \rightarrow$ set of final states; $F \subseteq Q$.
 - $\delta \rightarrow$ transition fn
 - $\delta: Q \times \Sigma \rightarrow 2^Q$

Algo:

$S = \epsilon\text{-closure}(\{s_0\})$

$c = \text{getChar};$

```

while(c != eof) {
    s = E-closure(move(s, c));
    c = getChar();
}
if (sNF != ∅) then
    return 'Yes';
else
    return 'No';

```

Q) For alphabet $\Sigma = \{0, 1, 3\}$, write regular expression for strings having 100 or 010 as substrings.

$$\boxed{R.E. = (0^*)^* 100 1^* 0^* + 0^* 1^* 010 1^* 0^*}$$

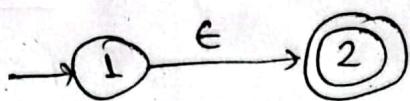
$$\begin{aligned}
R.E. &= (0+1)^* 100 (0+1)^* + (0+1)^* 010 (0+1)^* \\
&\text{or} \\
&= (0+1)^* (100 + 010) (0+1)^*
\end{aligned}$$

* Regular Expression to NFA (Thompson's Construction)

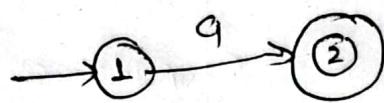
↳ Here,

- i/p is regular expression, σ over alphabet A.
- o/p is E-NFA accepting $L(\sigma)$.
- procedure is process in bottom up manner by creating E-NFA for each symbol in A including ϵ .
- then, recursively create for other operations.
- we will use the rules which define a regular expression as a basis for the construction.

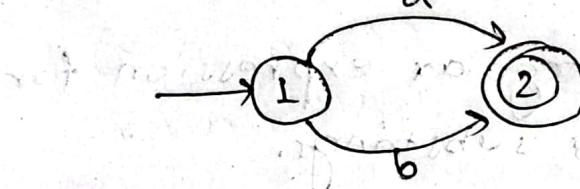
a) the NFA representing empty string.



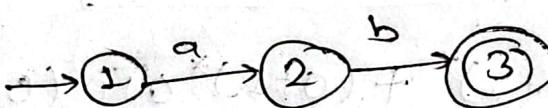
b) the NFA representing $a \in A$.



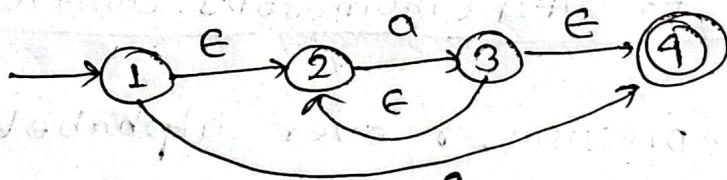
c) the NFA representing $a+b$ or $a|b$ or $a \cup b$.



d) the NFA representing $a \cdot b$



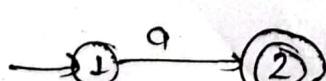
e) the NFA representing a^*



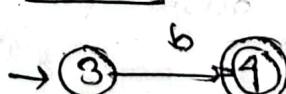
④ Construct NFA of given R.E using thompson's construction.

$$(a \cdot b)^* \cap C$$

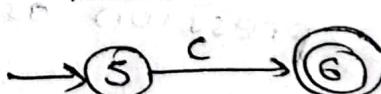
for a



for b



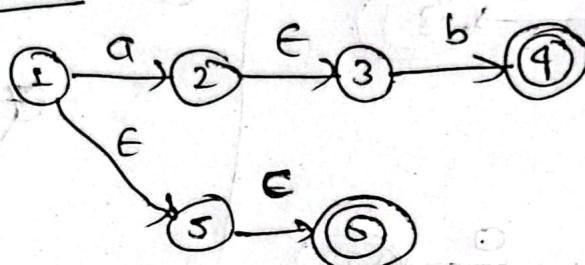
for c



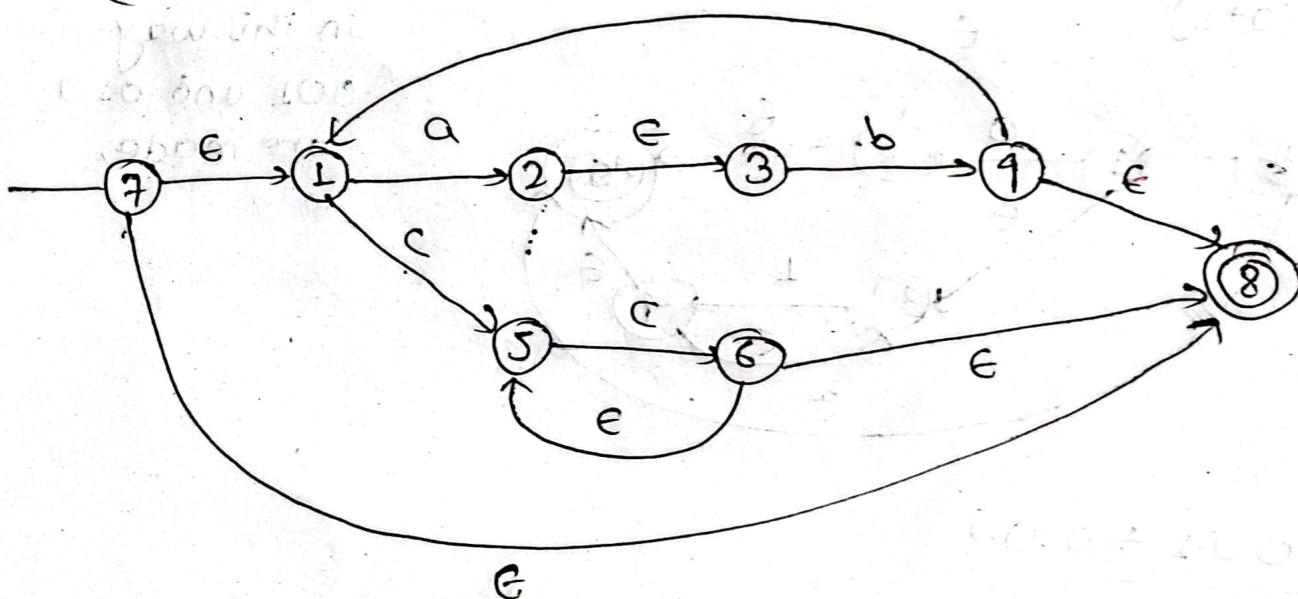
for (a.b);



for ((a.b)|c)

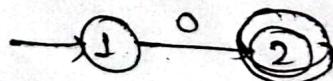


for (((a.b)|c))*

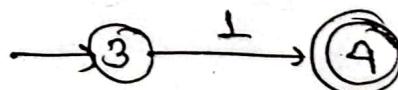


⑨ $(0+1)^*(001 + 010)(0+1)^*$

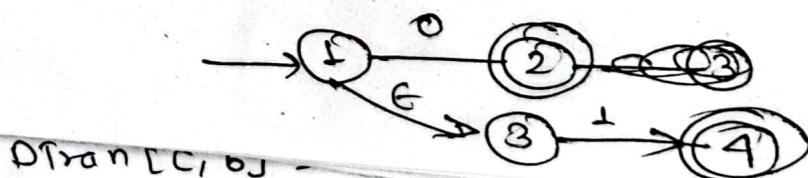
for 0,



for 1,



for 0+1,

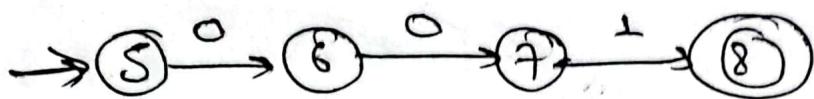
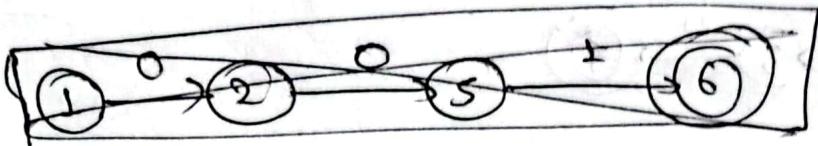


DTrans [C, b]

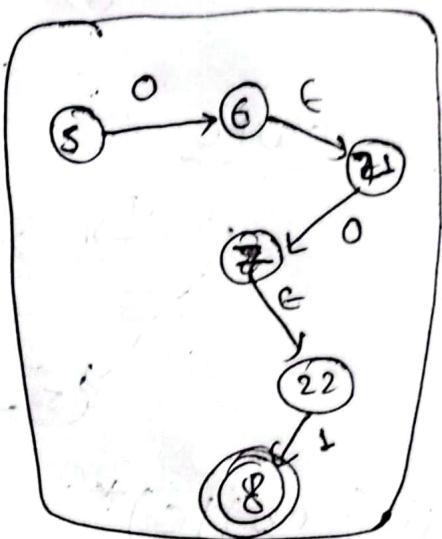
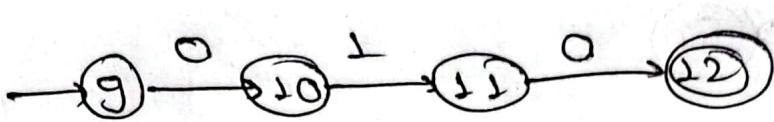
01-22

CD

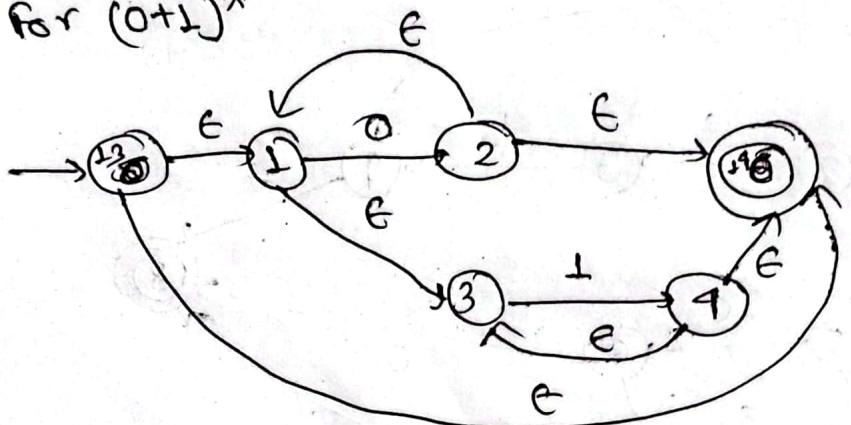
for 001



for 010

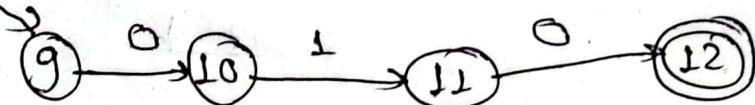
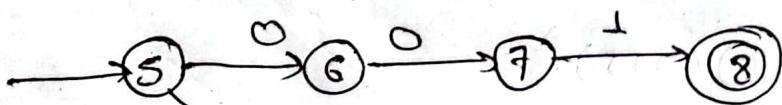


for $(0+1)^*$



In this way,
001 and 010
are made,

for $(001 + 010)$

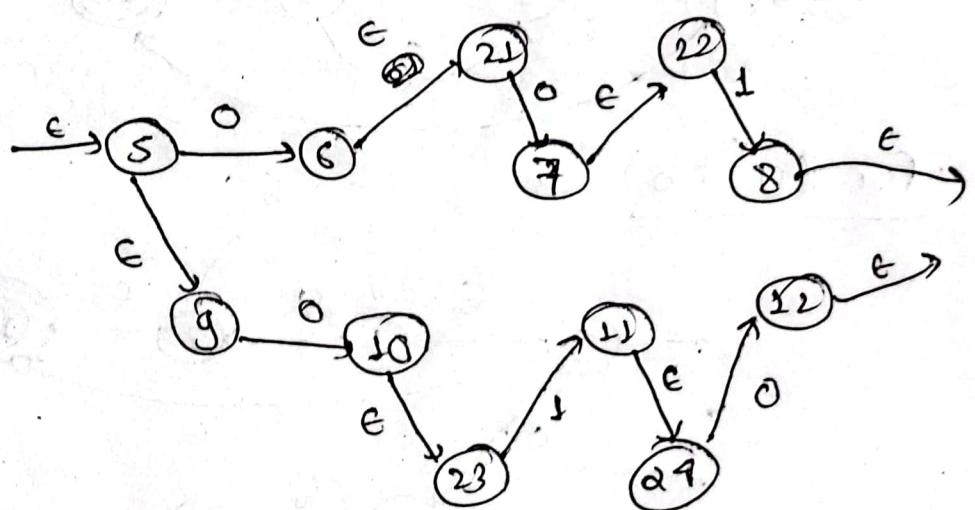
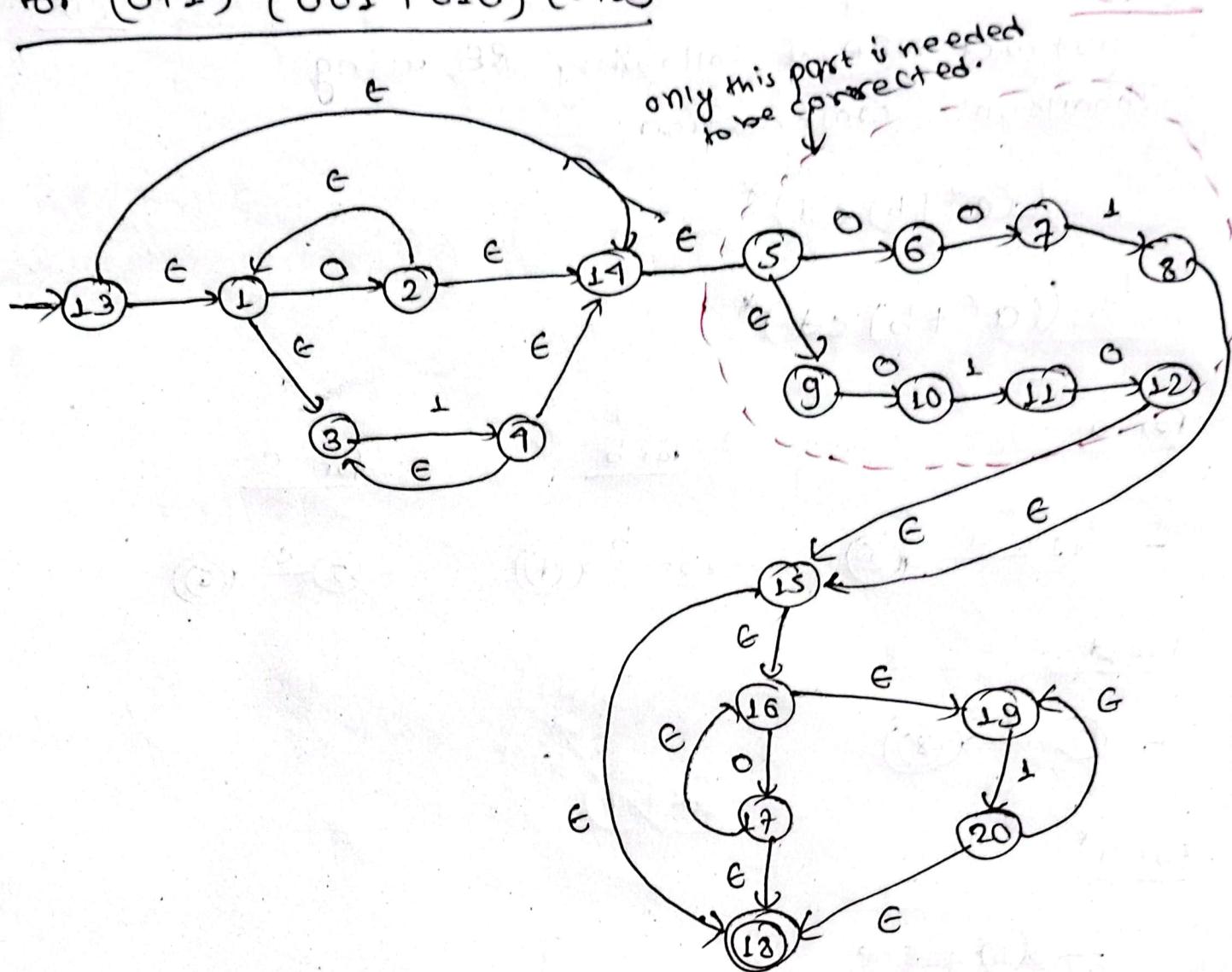


for $(0+1)^*(001+010)(0+1)^*$

02-22

0

3



081-01-23

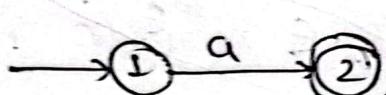
01-23

Q) Construct NFA of following RE using Thompson's construction:

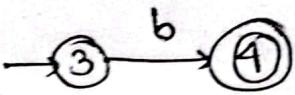
$$((a^* 1-b) cd)^*$$

$$\hookrightarrow ((a^* + b) cd)^*$$

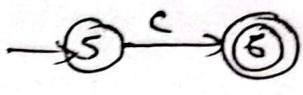
for a



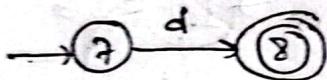
for b



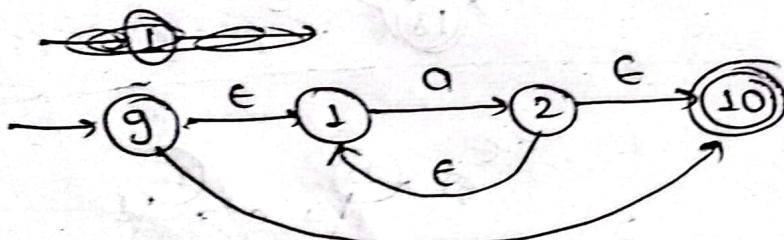
for c



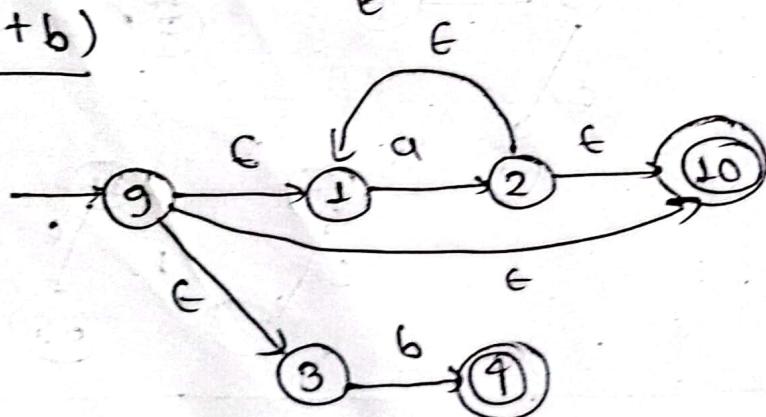
for d



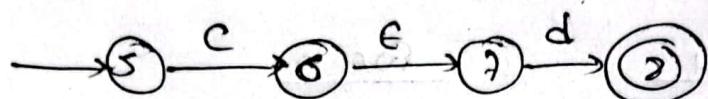
for a^*



for $(a^* + b)$



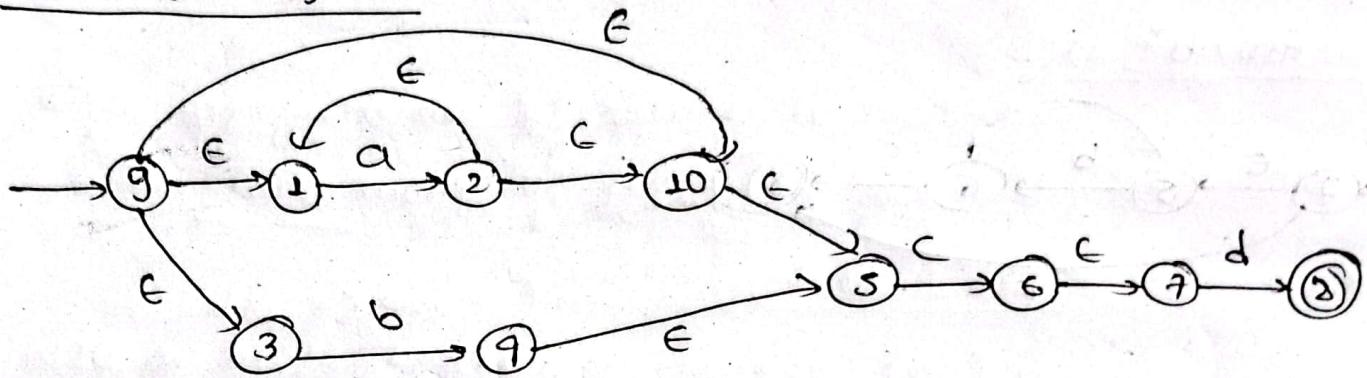
for cd



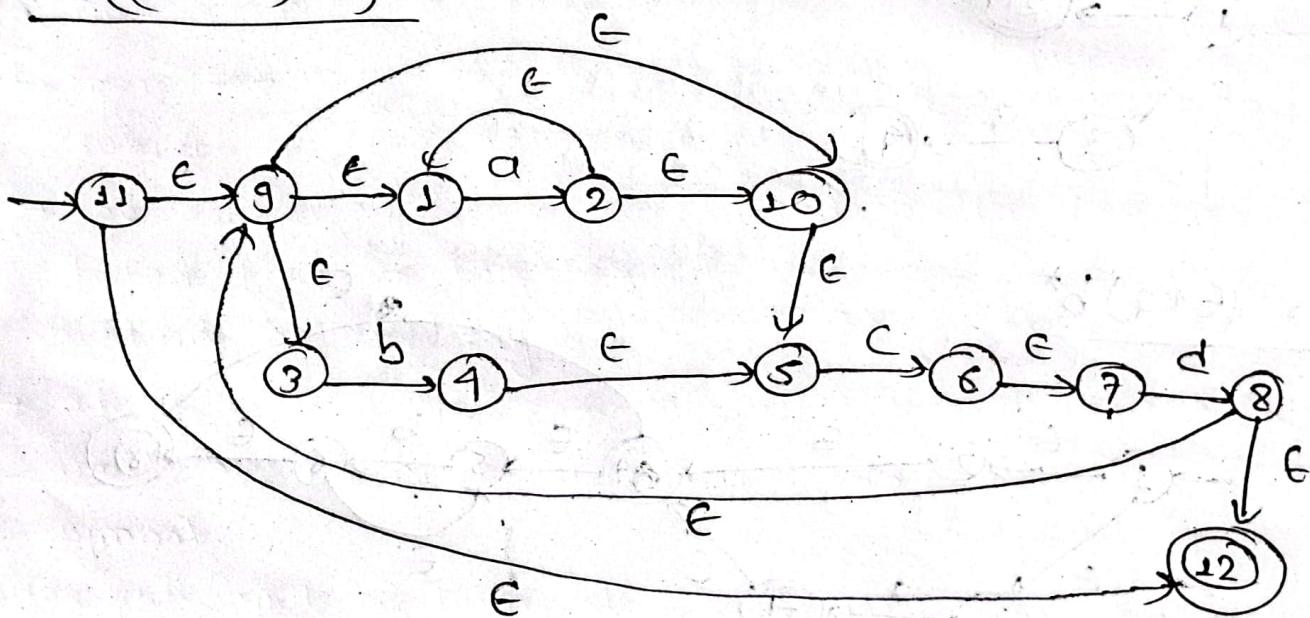
Q1-23

CD

for $(a^* + b)cd$



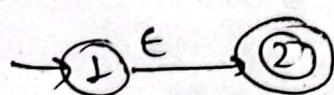
for $((a^* + b)cd)^*$



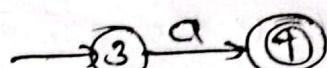
Q1-23

$$\textcircled{9} : RE = ((\epsilon + a) b^*)^*$$

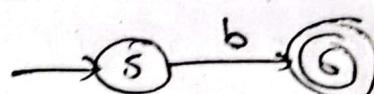
for ϵ



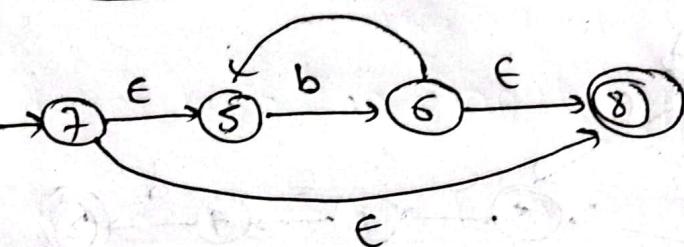
for a



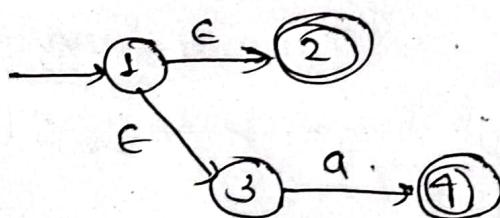
for b



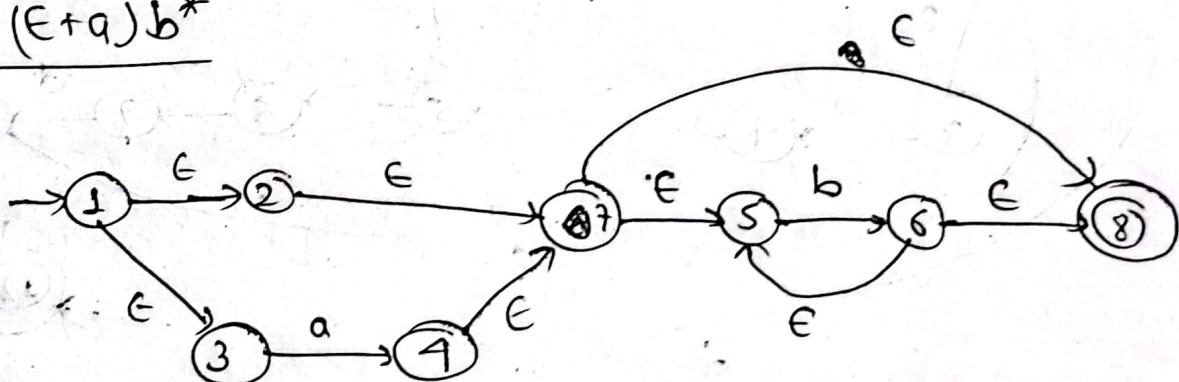
for b^*



for $(\epsilon + a)$



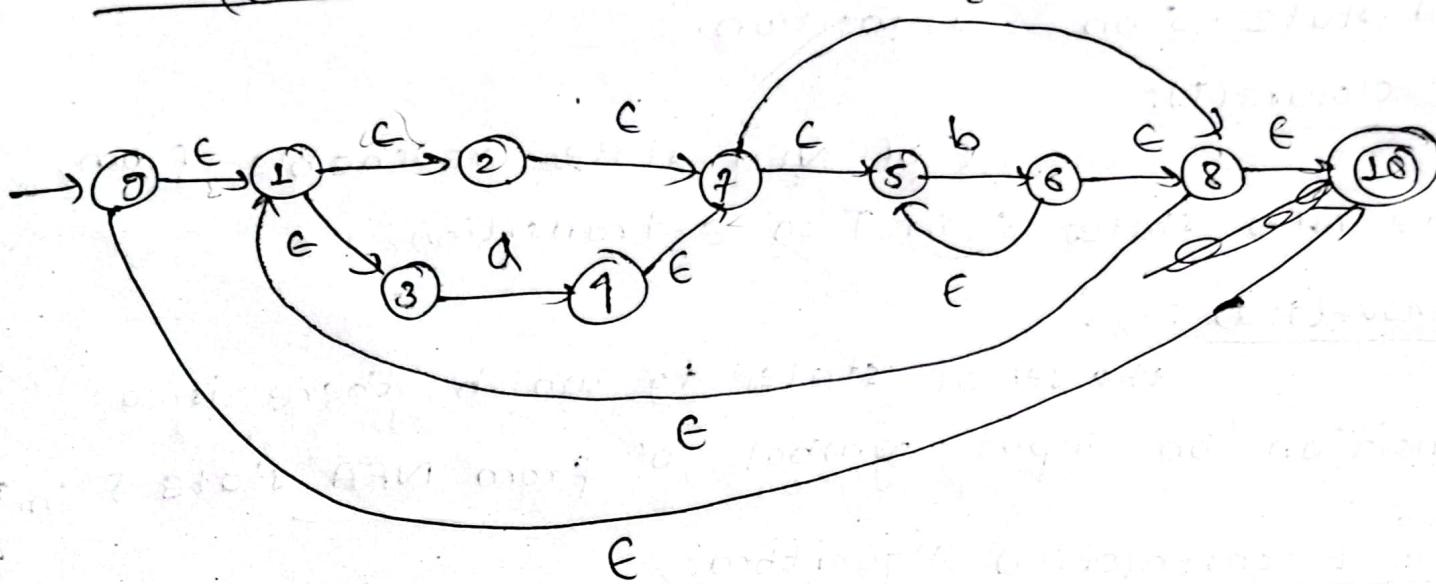
for $(\epsilon + a) b^*$



Aug. 01-23

CD

for $((\epsilon + a)b^*)^*$



* Conversion from NFA to DFA (subset construction):

- *Conversion from NFA

 - ↳ Subset construction is an approach that constructs DFA from NFA, that recognize the same language.
 - ↳ There may be several accepting states in a given subset of Non deterministic states.
 - ↳ The accepting state corresponding to the pattern listed first in the lexical analyzer generator specification has priority.
 - ↳ State transitions are made until a state is reached which has no next state for the current input symbol.
 - ↳ The last input position at which the DFA entered an accepting state gives the lexemes.
 - ↳ Following actions are needed:

1) ϵ -closure(S):

The set of NFA states reachable from NFA state S on ϵ -transition.

2) ϵ -closure(T):

The set of NFA states reachable from some NFA states S in T on ϵ -transition.

3) Move(T, a):

The set of states to which there is a transition on input symbol 'a' from NFA state S in T .

*Subset Construction Algorithm:

Put ϵ -closure (S_0) as an unmarked state in Dstates.

while there is an unmarked state T in Dstates do

Mark T .

for each input symbol $a \in A$ do

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

if U is not in Dstates then

Add U as an unmarked state to Dstates.

end if;

$RTrans[T, a] = U;$

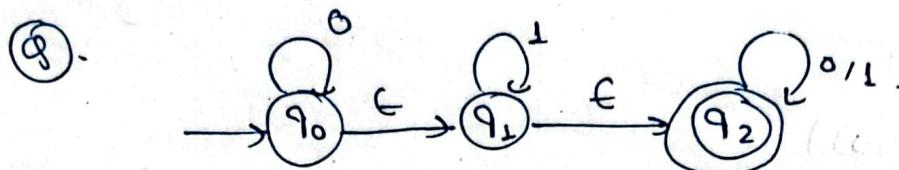
end do;

end do;

↳ Here,

- Dstates is the set of states of the new DFA consisting the set of states of NFA.

- DTran is the transition table of the new DFA. Q2-23
CD
- A set of Dstates is an accepting state of DFA if it is a set of NFA states containing at least one accepting state of NFA.
- The start state of DFA is ϵ -closure(S_0) where, S_0 is start state of NFA.



Ans;

Initially,

ϵ -closure(q_0) = $\{q_0, q_1, q_2\} = A$ as unmarked state in Dstate

Mark A,

for i/p symbol '0',

$$U = \epsilon\text{-closure}(\text{move}(A, 0))$$

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

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

$$= \{q_0, q_1, q_2\} \cup \{q_2\}$$

$$= \{q_0, q_1, q_2\} = A, \text{ already in Dstates.}$$

for i/p symbol '1',

$$U = \epsilon\text{-closure}(\text{move}(A, 1))$$

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

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

$$= \{q_1, q_2\} \cup \{q_2\}$$

DTran[0, 0]

$$= \{q_1, q_2\} = B.$$

add B as unmarked state in Dstates.

$$DTran[A, 0] = A$$

$$DTran[A, 1] = B$$

Mark B.

for i/p symbol '0'

$$U = \epsilon\text{-closure}(\text{move}(B, 0))$$

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

$$= \{q_2\} = C.$$

add C as unmarked state in Dstates.

for i/p symbol '1'

$$U = \epsilon\text{-closure}(\text{move}(B, 1))$$

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

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

$$= \{q_1, q_2\} \cup \{q_2\}$$

$$= \{q_1, q_2\} = B, \text{ already in Dstates.}$$

$$DTran[B, 0] = C$$

$$DTran[B, 1] = B.$$

Mark C

for i/p symbol '0'

$$U = \epsilon\text{-closure}(\text{move}(C, 0))$$

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

$\Sigma \{ q_2 \} = C$, already in Dstates.

01-23

CD

for i/p symbol '1',

$U = \epsilon\text{-closure}(\text{move}(C, 1))$

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

$= \{ q_2 \} = C$, already in Dstates.

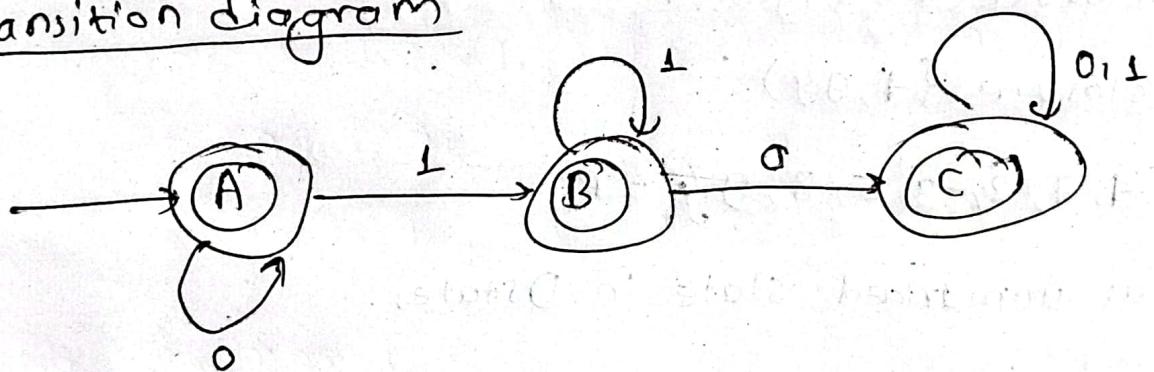
$D\text{tran}[C, 0] = C$

$D\text{tran}[C, 1] = C$

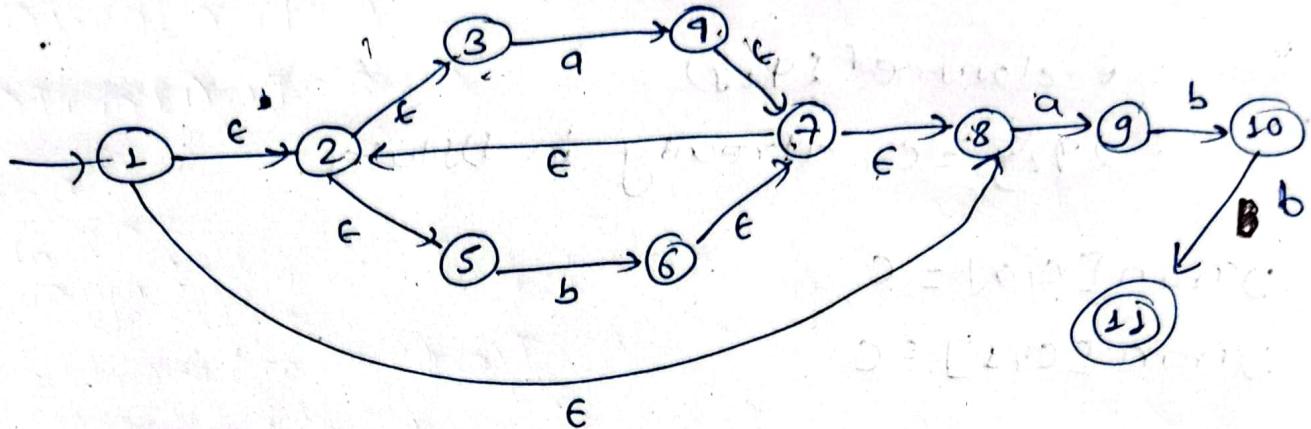
finally, transition table of DFA

Q/Σ	0	1
$\rightarrow *A$	A	B
*B	C	B
*C	C	C

transition diagram



Q) find equivalent DFA of following NFA using subset construction algorithm.



To solve:

Initially,

ϵ -closure({1}) = {1, 2, 3, 5, 8} = A as unmarked state in Dstates.

Mark A

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(A, a))$$

$$= \epsilon\text{-closure}(\{4, 9\})$$

$$= \{4, 7, 2, 3, 5, 8, 9\} = B,$$

add B as unmarked state in Dstates.

for i/p symbol 'b'

$$U = \epsilon\text{-closure}(\text{move}(A, b))$$

$$= \epsilon\text{-closure}(\{6\})$$

$$= \{6, 7, 2, 3, 5, 8\} = C,$$

add C as unmarked state in Dstates.

$$\text{DTrans } [A, a] = B$$

$$\text{DTrans } [A, b] = C$$

Mark B

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(B, a))$$

$$= \epsilon\text{-closure}(\{4, 9\})$$

$= \{\emptyset, 4, 7, 2, 3, 5, 8, 9\} = B$, already in Dstates.

for i/p symbol 'b'

$$U = \epsilon\text{-closure}(\text{move}(B, b))$$

$$= \epsilon\text{-closure}(\{6, 10\})$$

$$= \{6, 7, 2, 3, 5, 8, 10\} = D,$$

~~add D as unmarked states in Dstates.~~

$$D\text{Tran}[B, a] = B$$

$$D\text{Tran}[B, b] = D$$

Mark C

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(C, a))$$

$$= \epsilon\text{-closure}(\{4, 9\})$$

$= \{4, 7, 2, 3, 5, 8, 9\} = B$, already in Dstates

for i/p symbol 'b'

$$U = \epsilon\text{-closure}(\text{move}(C, b))$$

$$= \epsilon\text{-closure}(\{6\})$$

$= \{6, 7, 2, 3, 5, 8\} = C$, already in Dstates

$$D\text{Tran}[C, a] = B,$$

$$D\text{Tran}[C, b] = C$$

Mark D.

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(D, a))$$

$$= \epsilon\text{-closure}(\{4, 9\})$$

$$= \{4, 7, 2, 3, 5, 8, 9\} = B, \text{ already in D states.}$$

for i/p symbol 'b'

$$U = \epsilon\text{-closure}(\text{move}(D, b))$$

$$= \epsilon\text{-closure}(\{6, 11\})$$

$$= \{6, 7, 2, 3, 5, 8, 11\} = E_\alpha$$

add E as unmarked state in D states.

$$\text{Diran}[D, a] = B$$

$$\text{Diran}[D, b] = E$$

Mark E

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(E, a))$$

$$= \epsilon\text{-closure}(\{4, 9\})$$

$$= \{4, 7, 2, 3, 5, 8, 9\} = B, \text{ already in D states}$$

for i/p symbol 'b'

$$U = \epsilon\text{-closure}(\text{move}(E, b))$$

$$= \epsilon\text{-closure}(\{6\})$$

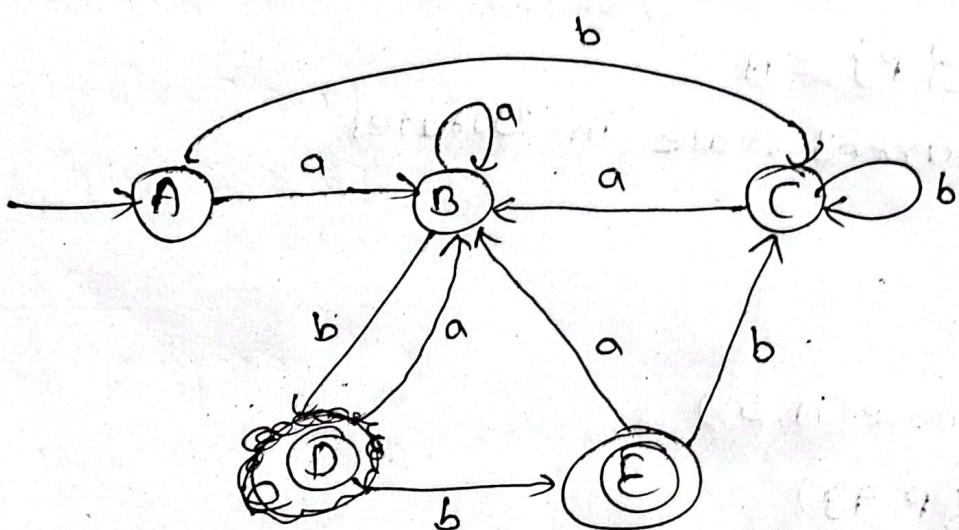
$$= \{6, 7, 2, 3, 5, 8\} = C, \text{ already in D states.}$$

$$D\text{ran}[E, a] = B$$

$$D\text{ran}[E, b] = C$$

Transition table

Q/Σ	a	b
A	B	C
B	B	D
C	B	C
D	B	E
*E	B	C



Q) Consider the following transition table

Q/Σ	0	1
$\rightarrow P$	$\{P, Q\}$	$\{P\}$
Q	$\{R\}$	$\{R\}$
R	$\{S\}$	—
$\ast S$	$\{S\}$	$\{S\}$

Now, find equivalent DFA using subset construction algorithm.

Soln;

Initially,

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

add A as unmarked state in Dstates.

Mark A

for i/p symbol 0,

$$V = \epsilon\text{-closure}(\text{move}(A, 0))$$

$$= \epsilon\text{-closure}(\{P, Q\})$$

$$= \{P, Q\} = B$$

add B as unmarked state in Dstates.

for i/p symbol 1;

$$U = \epsilon\text{-closure}(\text{move}(A, 1))$$

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

= $\{P\} = A$, already in Dstates.

$$DTran[A, 0] = B$$

$$DTran[A, 1] = A$$

Q1-24

CD

Mark B

for i/p symbol 0

$$U = \epsilon\text{-closure}(\text{move}(B, 0))$$

$$= \epsilon\text{-closure}(\{P, q_0, \gamma\})$$

$$= \{P, q_0, \gamma\} = C$$

add C as unmarked state in Dstates.

for i/p symbol 1

$$U = \epsilon\text{-closure}(\text{move}(B, 1))$$

$$= \epsilon\text{-closure}(\{P, \gamma\})$$

$$= \{P, \gamma\} = D$$

add D as unmarked state in Dstates.

~~Mark~~ $DTran[B, 0] = C$

$DTran[B, 1] = D$

Mark C

for o/i/p symbol 0

$$U = \epsilon\text{-closure}(\text{move}(C, 0))$$

$$= \epsilon\text{-closure}(\{P, q, r, s\})$$

$$= \{P, q, r, s\} = E$$

add E as unmarked state in Dstates.

for i/p symbol 1

$$U = \epsilon\text{-closure}(\text{move}(C, 1))$$

$$= \epsilon\text{-closure}(\{P, \gamma\})$$

$\{P, \pi\} = D$, already in Dstates.

DTran [C, 0] = E

DTran [C, 1] = D

Mark D @

for i/p symbol '0'

$U = \epsilon\text{-closure}(\text{move}(D, 0))$

$= \epsilon\text{-closure}(\{P, q, s\})$

$= \{P, q, s\} = F$

add f as unmarked state in Dstates.

for i/p symbol '1'

$U = \epsilon\text{-closure}(\text{move}(D, 1))$

$= \epsilon\text{-closure}(\emptyset \{P\})$ ~~(A, already)~~

$= \{P\} = A$, already in Dstates

DTran [D, 0] = F

DTran [D, 1] = A

Mark E

for i/p symbol '0'

$U = \epsilon\text{-closure}(\text{move}(E, 0))$

$= \epsilon\text{-closure}(\{P, q, r, s\})$

$= \{P, q, r, s\} = E$, already in Dstates.

for i/p symbol '1'

$U = \epsilon\text{-closure}(\text{move}(E, 1))$

$= \epsilon\text{-closure}(\{P, r, s\})$

$$\{P, \delta, S\} = G_1$$

Add η as unmarked state in Dstates.

$$DTran[F, 0] = \emptyset$$

$$DTran[\emptyset, 1] = \eta$$

Mark F

for i/p symbol '0'

$$U = \epsilon\text{-closure}(\text{move}(F, 0))$$

$$= \epsilon\text{-closure}(\{P, Q, R, S\}) \quad \text{~~already~~ states}$$

$$= \{P, Q, R, S\} = E, \text{ already in Dstates}$$

for i/p symbol '1'

$$U = \epsilon\text{-closure}(\text{move}(F, 1))$$

$$= \epsilon\text{-closure}(\{P, R, S\}) \quad \text{~~already~~ states}$$

$$= \{P, R, S\} = G_1, \text{ already in Dstates}$$

~~∴~~ $DTran[F, 0] = \emptyset$

$$DTran[F, 1] = G_1$$

Mark G

for i/p symbol '0'

$$U = \epsilon\text{-closure}(\text{move}(G_1, 0))$$

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

$$= \{P, Q, S\} = F_1 \text{, already in Dstates.}$$

for i/p symbol '1'

$$U = \epsilon\text{-closure}(\text{move}(G_1, 1))$$

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

$$\{P, S\} = H$$

add H as unmarked state in Dstates,

$$DTran[H, 0] = F$$

$$\underline{\text{Mark H}} \quad DTran[H, 1] = H$$

for :/p symbol 0.

$$U = \epsilon\text{-closure}(\text{move}(H, 0))$$

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

$$= \{P, Q, S\} = F, \text{ already in Dstates}$$

for :/p symbol 1

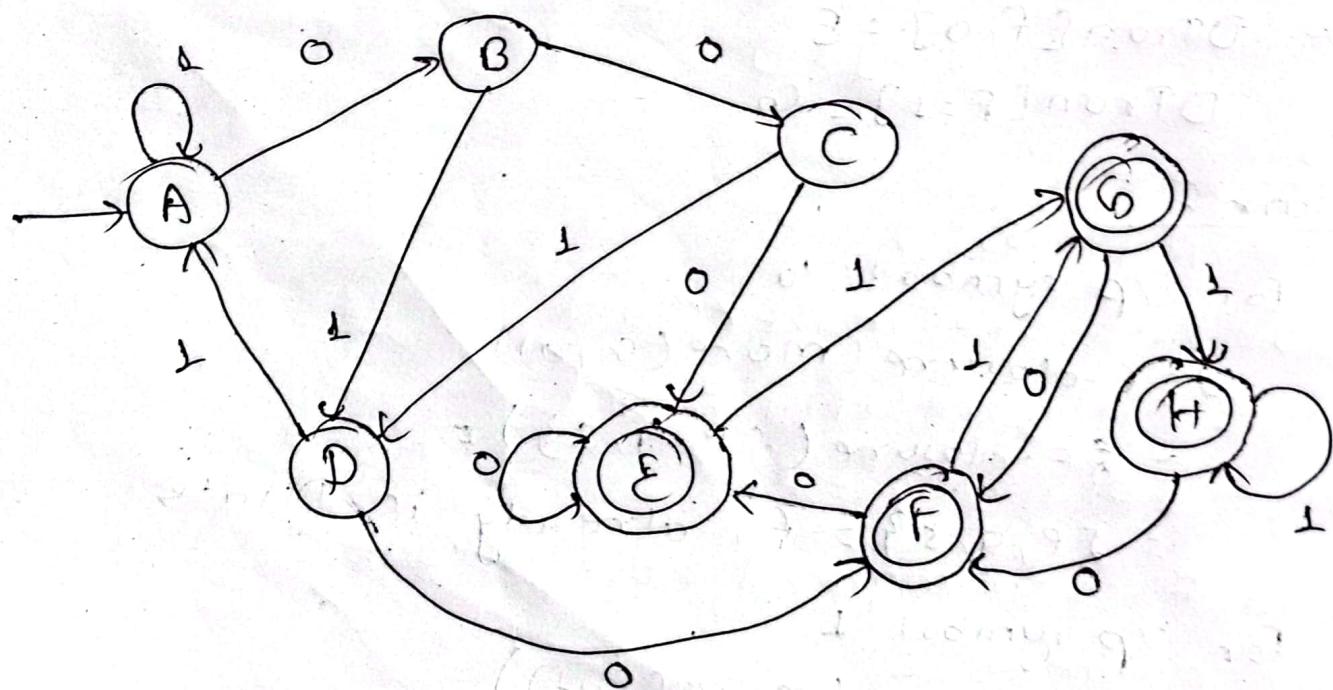
$$U = \epsilon\text{-closure}(\text{move}(H, 1))$$

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

$$= \{P, S\} = H, \text{ already in Dstates}$$

$$DTran[H, 0] = F$$

$$DTran[H, 1] = H$$



* State Minimization in DFA :

- ↳ It refers to process of transforming a given DFA into an equivalent DFA with minimum no. of states.
- ↳ Two states p and q are called equivalent, if for all i/p string w , $\delta(p, w)$ is an accepting state iff $\delta(q, w)$ is an accepting state.
 - otherwise distinguishable states.
- ↳ We say that, string w distinguishes state p from state q if by starting with DFA M in state p and feeding it input w , we end up in an accepting state.
- ↳ But starting in state q and feeding it with same i/p w , we end up in a non accepting state.
OR, vice-versa.
- ↳ Each group of states that can't be distinguished is then merged into single state.

* Procedure:

- 1) Partition the set of states into two partitions.
 - a) set of accepting states.
 - b) set of non-accepting states.
- 2) Split the partition on the basis of distinguishable state and keep equivalent states in group.
- 3) To split, we process the transition from state in a group with all i/p symbols.
 - If the transitions on any i/p from the state, in a group is on different group of states, then they are distinguishable so, remove those states from

the current partition and create groups. 02-07
CD

4) Process until all the partitions contain equivalent states only or have single states.

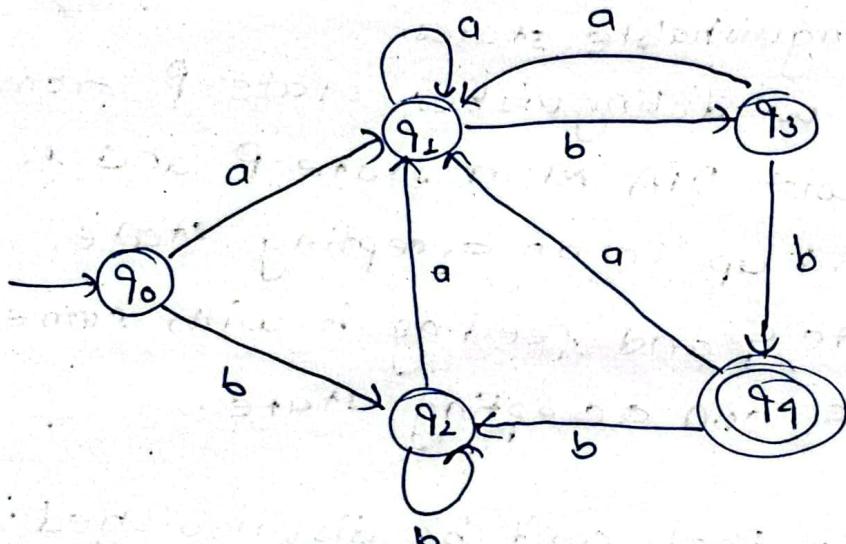
→ state p and q are equivalent if

$$\delta(p, w) \in F \& \delta(q, w) \in F$$

or $\delta(p, w) \notin F \& \delta(q, w) \notin F$

$$\delta(p, w) \notin F \& \delta(q, w) \in F$$

Eg:



Transition table:

Q / Σ	a	b
→ q0	q1	q2
q1	q1	q3
q2	q1	q2
q3	q1	q4
* q4	q1	q2

Now,

$$0 \text{ equivalent set} = [q_0, q_1, q_2, q_3] \ [q_4]$$

For 1 equivalent,

for q_0, q_1 : they are equivalent

\because for a and b, they are in same group.

for q_0, q_2 : they are equivalent.

for q_0, q_3 : they are not equivalent.

\because for a they are in same group but for b they are in diff. group.

$$\therefore 1 \text{ equivalent sets} = [q_0, q_1, q_2] [q_3] [q_4]$$

for 2 equivalent,

for q_0, q_1 : they are not equivalent.

\because for a they are in same group but for b they are not in same group.

for q_0, q_2 : they are equivalent.

\because for both a and b they are in same group.

$$\therefore 2 \text{ equivalent sets} = [q_0, q_2] [q_1] [q_3] [q_4]$$

for 3 equivalent,

for q_0, q_2 : they are equivalent.

$$\therefore 3 \text{ equivalent sets} = [q_0, q_2] [q_1] [q_3] [q_4]$$

↳ Here, 2 and 3 equivalent sets are equal and further on we will get same result.

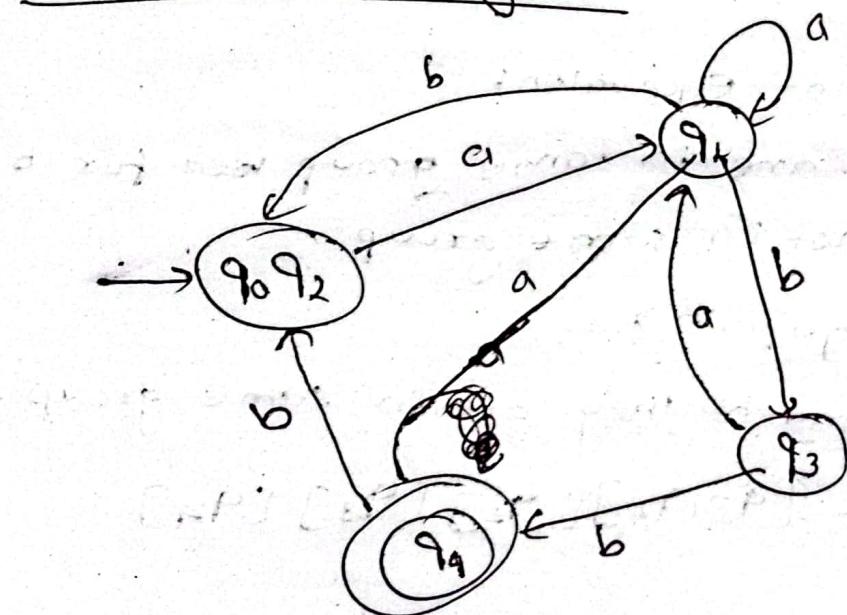
\therefore we have four states $\{q_0, q_2\}, \{q_1\}, \{q_3\}, \{q_4\}$

Transition tables

02-07
CD

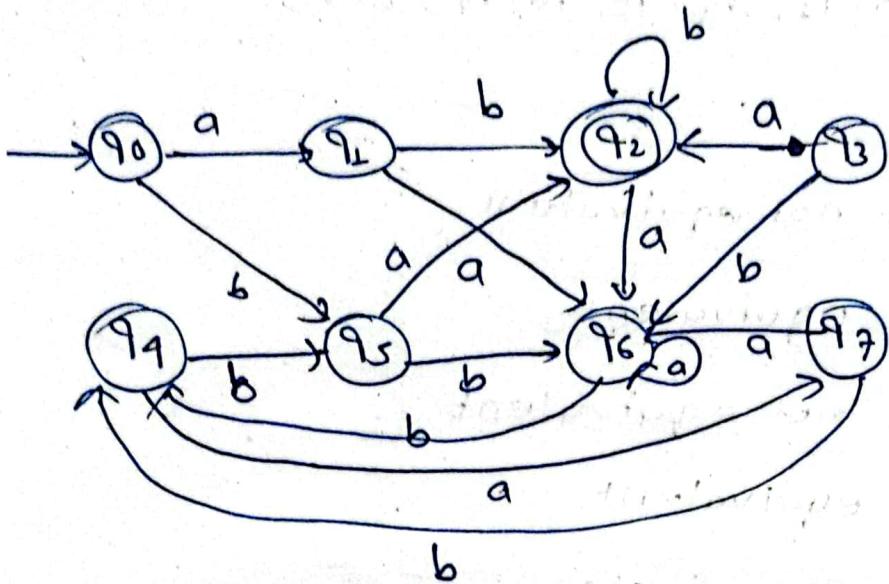
Q/Σ	a	b
$\rightarrow q_0, q_2$	q_1	q_0, q_2
q_1	q_2	q_3
q_3	q_2	q_4
q_4	q_2	q, q_0, q_2

Minimized state diagram



CD

⑧ Minimize the following DFA with equivalence methods.



Transition table:

Q/Σ	a	b
$\rightarrow q_0$	q_1 q_5	
q_1	q_6 q_2	
* q_2	q_6 q_2	
q_3	q_2 q_6	
q_4	q_7 q_5	
q_5	q_2 q_6	
q_6	q_6 q_4	
q_7	q_6 q_4	

Now, Remove 9₃

Expt

02-07

CD

0 equivalent sets = [9₀, 9₁, 9₉, 9₅, 9₆, 9₇] [9₂]

for 1 equivalent

for 9₀, 9₁ : they are not equivalent

for 9₀, 9₉ : they are equivalent

for 9₀, 9₅ : they are not equivalent

for 9₀, 9₆ : they are equivalent

for 9₀, 9₇ : they are equivalent.

∴ 1 equivalent sets = [9₀, 9₉, 9₆, 9₇] [9₁] [9₅] [9₂]

for 2 equivalent

for 9₀, 9₉ : they are not equivalent

for 9₀, 9₆ : they are not equivalent

for 9₀, 9₇ : they are not equivalent

for 9₉, 9₆ : they are not equivalent

for 9₉, 9₇ : they are not equivalent.

for 9₆, 9₇ : they are equivalent

∴ 2 equivalent sets = [9₆, 9₇] [9₀] [9₉] [9₁] [9₅] [9₂]

for 3 equivalent

for 9₆, 9₇ : they are equivalent.

∴ 3 equivalent sets = [9₆, 9₇] [9₀] [9₉] [9₁] [9₅] [9₂]

∴ The states we got are:

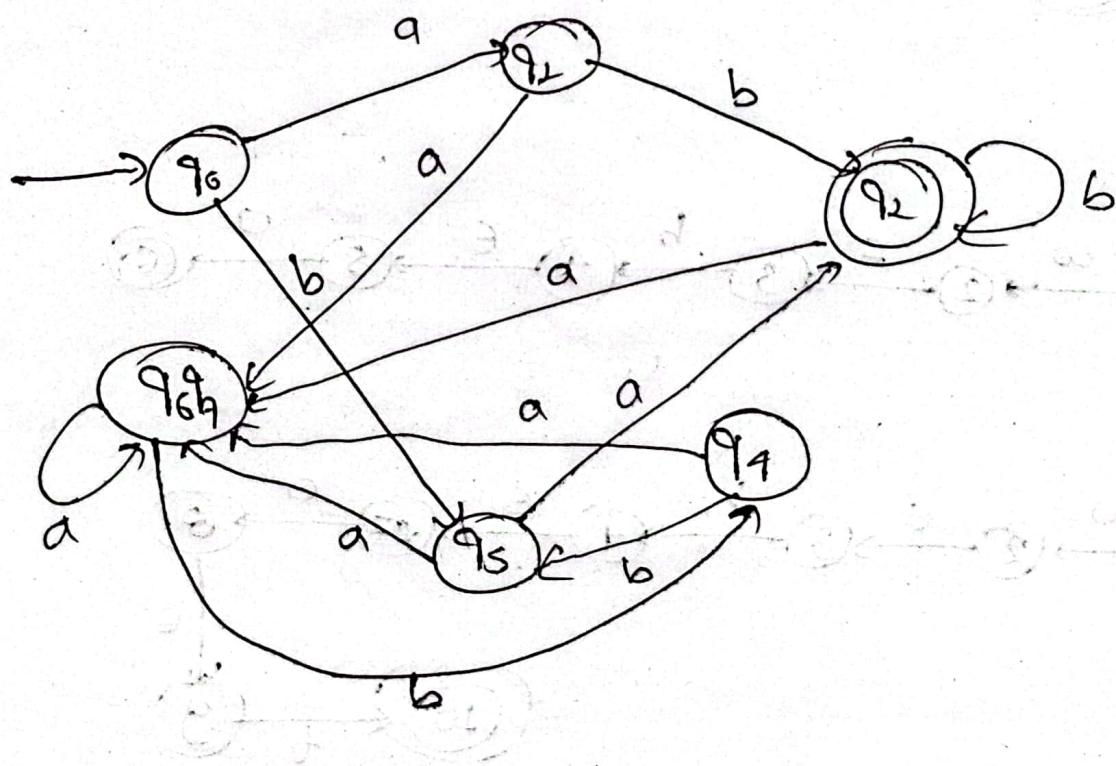
02-07
CP

{ q_6, q_7 }, { q_8 }, { q_1 }, { q_2 }, { q_3 }, { q_5 }

transition table

Q/Σ	a	b
$\rightarrow q_0$	q_1	q_5
q_1	q_6, q_7	q_2
$* q_2$	q_6, q_7	q_2
q_4	q_6, q_7	q_5
q_5	q_2	q_6, q_7
q_6, q_7	q_2	q_4

Minimized state diagram



Q) Construct the regular expression

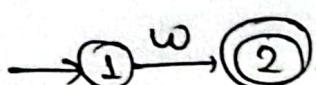
02-07
CD

who | what | where . Use thompson's construction
to build NFA from RE . Use subset construction
to build DFA from NFA . Then, minimize DFA.

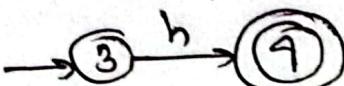
Ans:

R.E: who + what + where

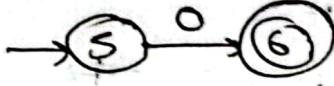
for ω ,



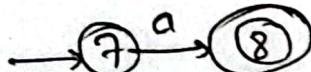
for b ,



for o ,



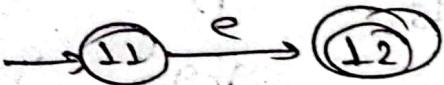
for a ,



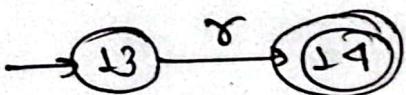
for t ,



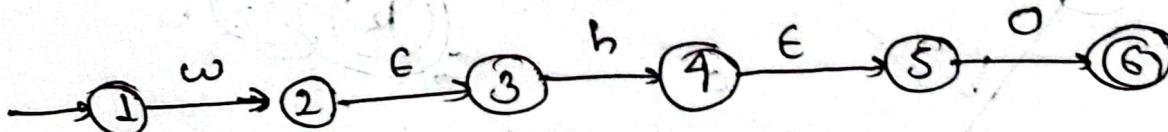
for e ,



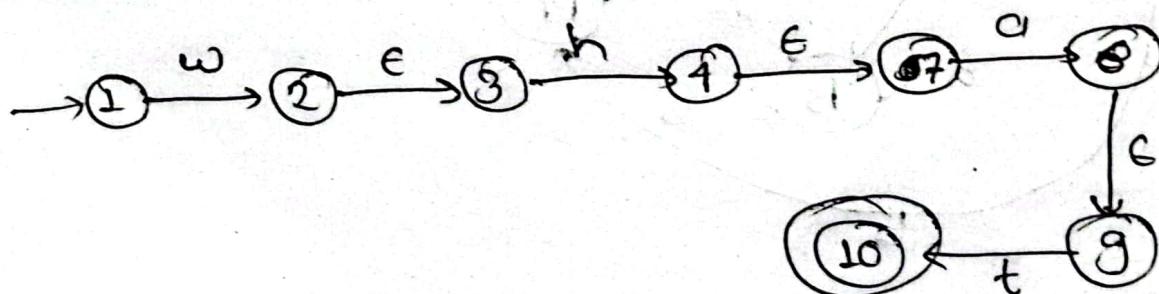
for r ,



for who ,



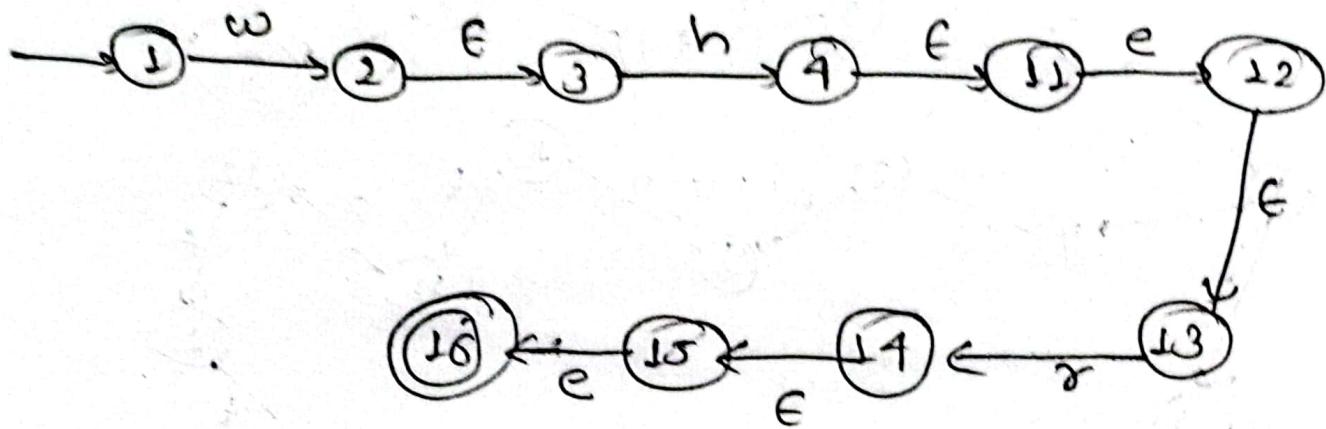
for what ,



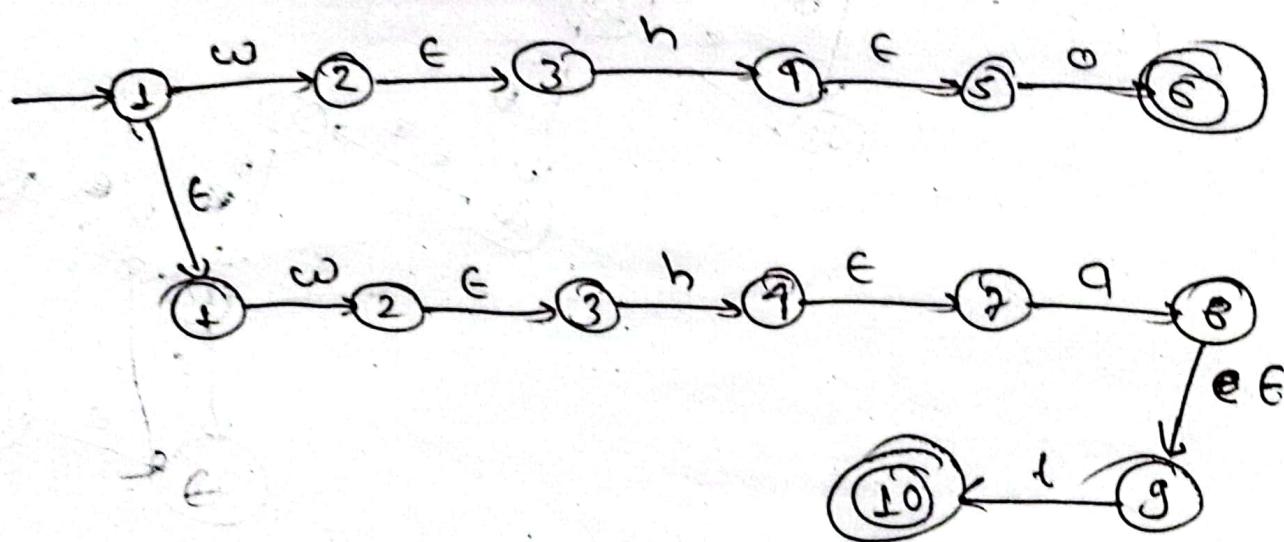
for where

02:02

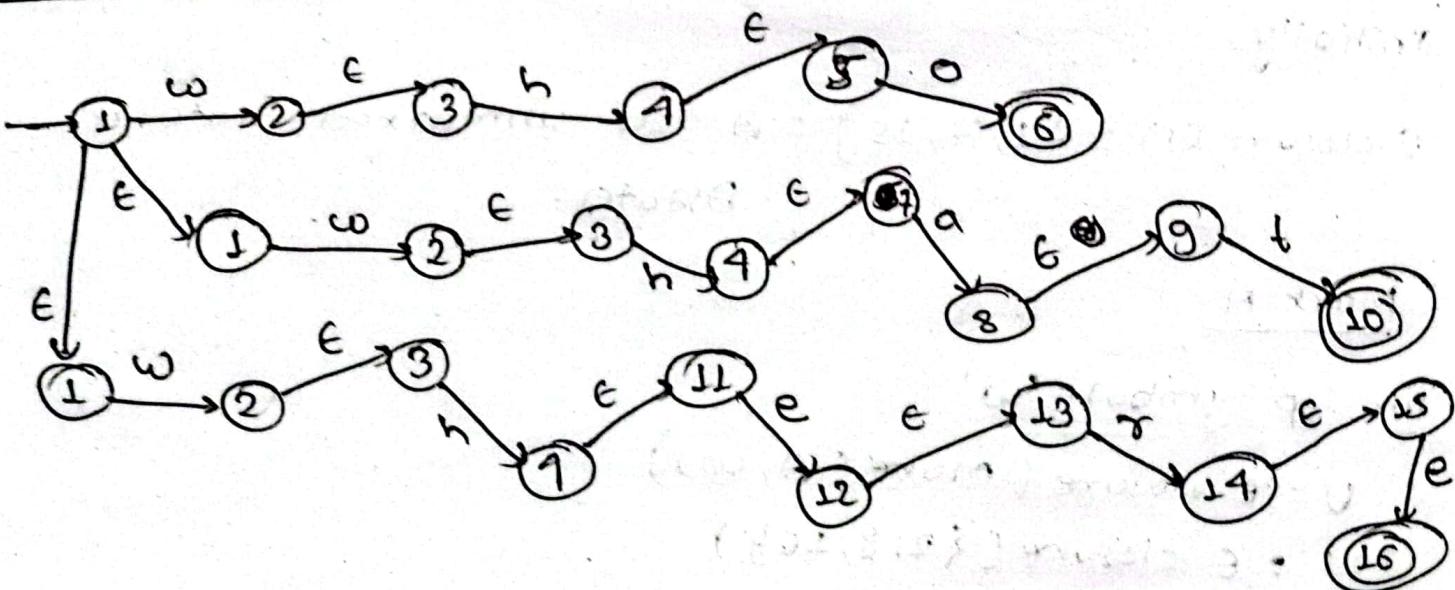
10



for who + what

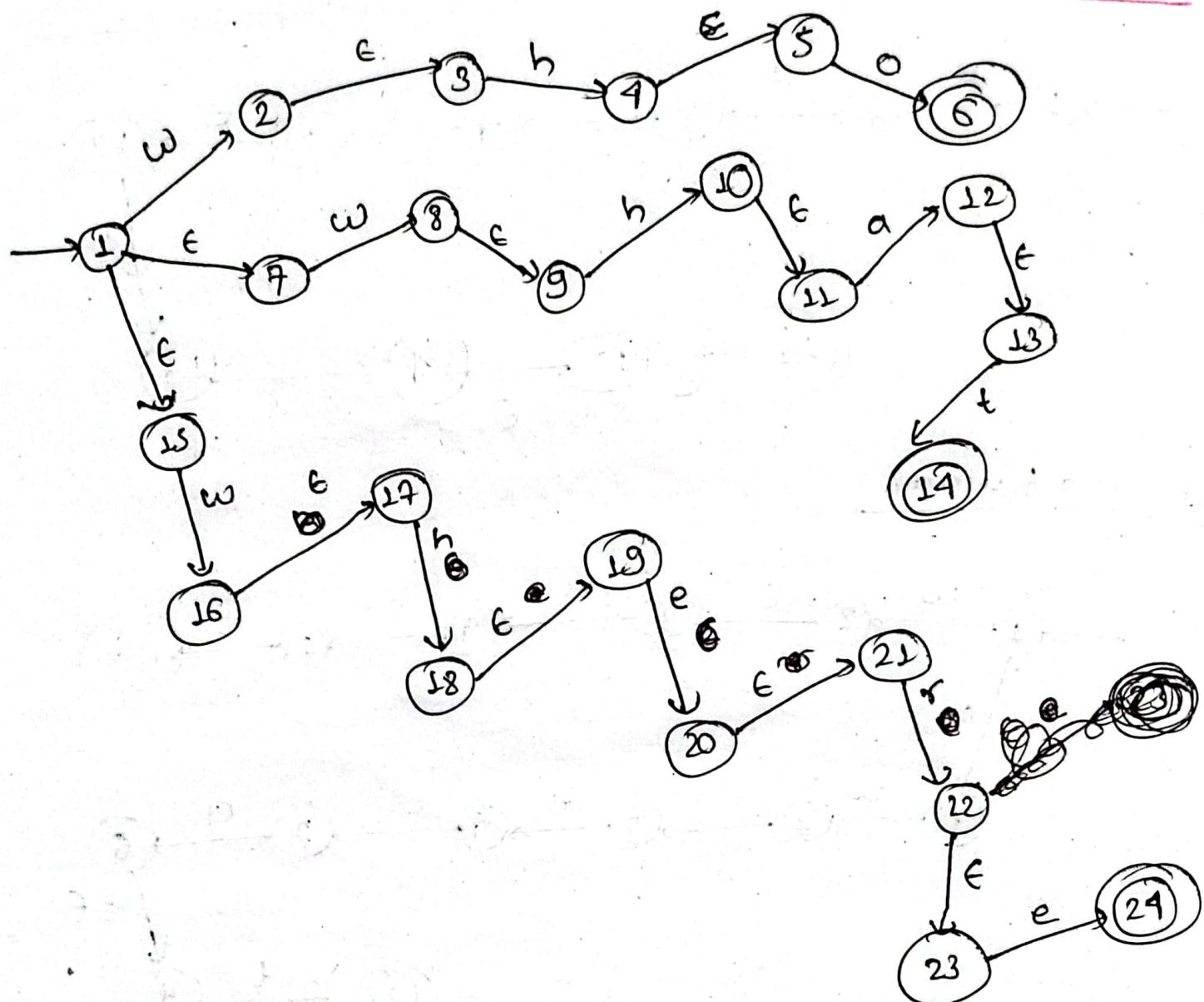


for what what + where



finally, changing start state.

Q2 ~ Q7
CD



Conversion of NFA to DFA using subset construction:

Initially,

ϵ -closure(1) = {1, 7, 15} = A as unmarked state in
Dstates.

Mark A

for i/p symbol 'w'

$$U = \epsilon\text{-closure}(\text{move}(A, w))$$

$$= \epsilon\text{-closure}(\{2, 8, 16\})$$

$$= \{2, 3, 8, 9, 16, 17\} = B$$

add B as unmarked state in Dstates.

for i/p symbol 'h'

$$U = \epsilon\text{-closure}(\text{move}(A, h))$$

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

$$= \emptyset$$

for i/p symbol '0'

$$U = \emptyset$$

for i/p symbol 'a'

$$U = \emptyset$$

for i/p symbol 't'

$$U = \emptyset$$

for i/p symbol 'e'

$$U = \emptyset$$

for i/p symbol 'x'

$$U = \emptyset$$

$$DTran[A, w] = B$$

$$DTran[A, h] = \text{Dead}$$

$$DTran[A, 0] = \text{Dead}$$

$$DTran[A, a] = \text{Dead}$$

$$DTran[A, t] = \text{Dead}$$

$$DTran[A, e] = \text{Dead}$$

$$DTran[A, x] = \text{Dead}$$

Mark B.

02-07

(D)

for i/p symbol 'w'

$$U = \epsilon\text{-closure}(\text{move}(B, w))$$

$$= \emptyset$$

for i/p symbol 'h'

$$U = \epsilon\text{-closure}(\text{move}(B, h))$$

$$= \epsilon\text{-closure}\{2^q, 10, 18\}$$

$$= \{q, 5, 10, 12, 18, 19\} = C$$

add C as unmarked state in DStates.

for i/p symbol 'o'

$$U = \epsilon\text{-closure}(\text{move}(B, o))$$

$$= \emptyset$$

for i/p symbol 'a'

$$U = \epsilon\text{-closure}(\text{move}(B, a))$$

$$= \emptyset$$

for i/p symbol 't'

$$U = \epsilon\text{-closure}(\text{move}(B, t))$$

$$= \emptyset$$

for i/p symbol 'e'

$$U = \epsilon\text{-closure}(\text{move}(B, e))$$

$$= \emptyset$$

for i/p symbol 's'

$$U = \epsilon\text{-closure}(\text{move}(B, s))$$

$$= \emptyset$$

$DTran[B, w] = \text{Dead}$

Q2-Q7

CD

$DTran[B, h] = C$

$DTran[B, o] = \text{Dead}$

$DTran[B, a] = \text{Dead}$

$DTran[B, t] = \text{Dead}$

$DTran[B, e] = \text{Dead}$

$DTran[B, x] = \text{Dead}$

Mark C.

for i/p symbol 'w'

$U = \epsilon\text{-closure}(\text{move}(C, w))$

$= \emptyset$

for i/p symbol 'h'

$U = \epsilon\text{-closure}(\text{move}(C, h))$

$= \emptyset$

for i/p symbol 'o'

$U = \epsilon\text{-closure}(\text{move}(C, o))$

$= \epsilon\text{-closure}(\{G, Y\})$

$= \{G, Y\} = D$

add D as unmarked state in Dstates.

for i/p symbol 'a'

$U = \epsilon\text{-closure}(\text{move}(C, a))$

$= \epsilon\text{-closure}(\{L2\})$

~~$= \{L2, L3\} = E$~~

add E as unmarked state in Dstates.

end

for i/p symbol 't',

$V = \epsilon\text{-closure}(\text{move}(C, t))$

$= \text{enclosed } \emptyset$

for i/p symbol 'e',

$V = \epsilon\text{-closure}(\text{move}(C, e))$

$= \epsilon\text{-closure}(\{20\})$

$= \{20, 21\} = F$

add f as unmarked state in D states.

for i/p symbol 'g'

$V = \epsilon\text{-closure}(\text{move}(C, g))$

$= \emptyset$

$D\text{tran}[C, w] = \text{Dead}$

$D\text{tran}[C, h] = \text{Dead}$

$D\text{tran}[C, o] = D$

$D\text{tran}[C, a] = E$

$D\text{tran}[C, t] = \text{Dead}$

$D\text{tran}[C, e] = F$

$D\text{tran}[C, r] = \text{Dead}$

Mark D.

for i/p symbol 'w'

$V = \epsilon\text{-closure}(\text{move}(D, w))$

$= \emptyset$

Q3

for i/p symbol 'h' $U = \epsilon\text{-closure}(\text{move}(D, h))$ $= \emptyset$ for i/p symbol 'o' $U = \epsilon\text{-closure}(\text{move}(D, o))$ $= \emptyset$ for i/p symbol 'a' $U = \epsilon\text{-closure}(\text{move}(D, a))$ $= \emptyset$ for i/p symbol 't' $U = \epsilon\text{-closure}(\text{move}(D, t))$ $= \emptyset$ for i/p symbol 'e' $U = \epsilon\text{-closure}(\text{move}(D, e))$ $= \emptyset$ for i/p symbol 'g' $U = \epsilon\text{-closure}(\text{move}(D, g))$ $= \emptyset$ $D\text{Tran}[D, w] = \text{Dead}$ $D\text{Tran}[D, h] = \text{Dead}$ $D\text{Tran}[D, o] = \text{Dead}$ $D\text{Tran}[D, a] = \text{Dead}$ $D\text{Tran}[D, t] = \text{Dead}$ $D\text{Tran}[D, e] = \text{Dead}$ $D\text{Tran}[D, g] = \text{Dead}$

Mark E,

02-07
CD

for i/p symbol 'w'

$U = \epsilon\text{-closure}(\text{move}(E, w))$

$= \emptyset$

for i/p symbol 'h'

$U = \epsilon\text{-closure}(\text{move}(E, h))$

$= \emptyset$

for i/p symbol 'o'

$U = \epsilon\text{-closure}(\text{move}(E, o))$

$= \emptyset$

for i/p symbol 'a'

$U = \epsilon\text{-closure}(\text{move}(E, a))$

$= \emptyset$

for i/p symbol 't'

$U = \epsilon\text{-closure}(\text{move}(E, t))$

$= \epsilon\text{-closure}(\{14\})$

$= \{14\} = G$

add G as unmarked state in Ostate

for i/p symbol 'e'

$U = \epsilon\text{-closure}(\text{move}(E, e))$

$= \emptyset$

for i/p symbol 'g'

$U = \epsilon\text{-closure}(\text{move}(E, g))$

$= \emptyset$

$D\text{tran}[E, w] = \text{Dead}$

$D\text{tran}[E, h] = \text{Dead}$... Do it yourself - - -

Unit - 3Parser* Syntax Analysis:

- ↳ is the second phase of compiler process, which checks the syntactical structure of the given i/p ie, whether the given i/p is in correct syntax or not.
- ↳ also called parsing.
- ↳ It does so by building a data structure called a parse tree or syntax tree; which is constructed by using pre-defined grammar of language and the i/p string.
- ↳ Syntaxes are represented using CFG.
- ↳ Parsing is the act of performing syntax analysis to verify an i/p program's compliance with source language.
- ↳ it ensures we have valid sequence of tokens.

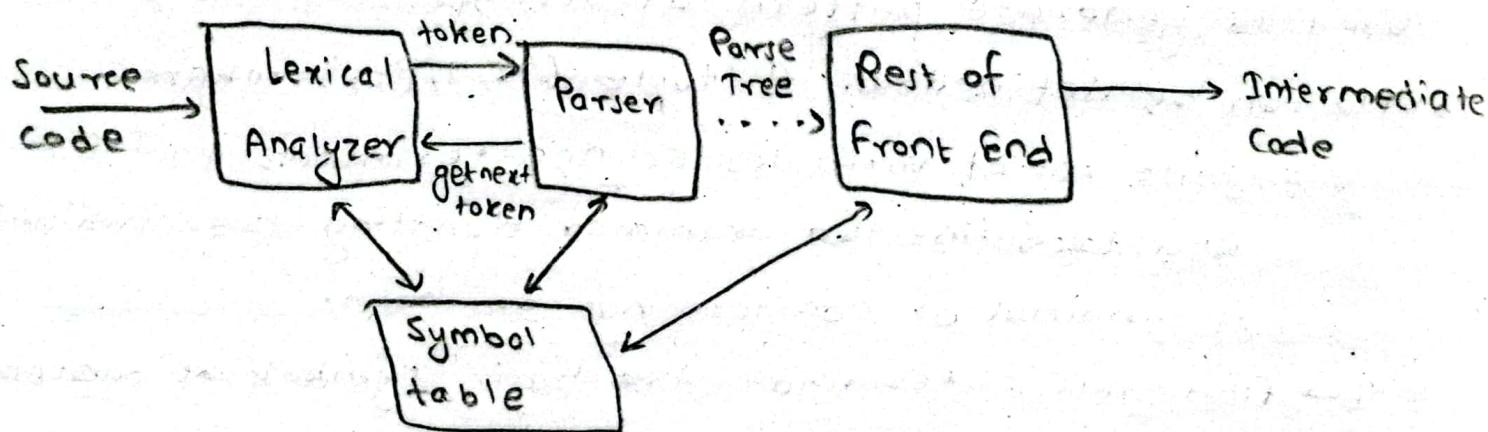
* Role of Parser:

Fig:- role of parser in compiler model.

- ↳ Parser analyze the source program based on the U2-13
CD definition of its syntax.
- ↳ it works in lock-step with lexical analyzer and is responsible for creating a parse tree of source code.
- ↳ it ~~uses~~ provides the mechanism for context sensitive analysis.
- ↳ it determines error and tries to handle them.

*Types of Parser

1) Top Down Parser

- ↳ builds parse tree from top to bottom.
- ↳ easy to build.
- ↳ works on subset of grammars like LL grammar.

2) Bottom Up Parser

- ↳ builds parse tree from bottom to top.
- ↳ difficult to build.
- ↳ works on subset of grammars like LR grammar.

*Content free Grammar (CFG):

- ↳ CFG is a set of recursive rewriting rule or production used to generate pattern of strings.
- ↳ CFG can be defined as A tuple (V, T, P, S) where
 - $V \rightarrow$ finite set of variables or non-terminals that are used to define the grammar denoting the combination of terminal or non-terminal or both.
 - $T \rightarrow$ finite set of terminals, the basic symbols of sentence.
 - $P \rightarrow$ set of productions, which are rules that defines the combination of terminal or non-terminal or both for particular non-terminal.

$S \rightarrow$ start symbol) which is a special non-terminal symbol that appears in the initial string generated by grammar.

02-23

CD

↳ lower case letters and symbols, bold string are used for denoting terminal. eg:- a, b, c, 0, 1 ET

↳ Uppercase letters, italicized strings are used for denoting nonterminals. eg:- A, S EV.

↳ Production rules are of the form $A \rightarrow \alpha$

↳ strings comprising of both terminals and non-terminals are denoted by greek letters like $\alpha, \beta, \gamma, \delta$.

* Derivations:

↳ Verification of the sentence defined by the grammar is done using the production rule to obtain the particular sentence by expansion of non terminal is called derivation.

↳ strings of non terminals and terminals is called sentential form, whereas strings of terminals is called sentence.

↳ During derivation, if we always expand the left most non terminal first then it is called left-most derivation.

↳ During derivation, if we always expand the right most non terminal first then it is called rightmost derivation.

↳ $\alpha \Rightarrow \beta$ is used to represent that β can be derived from α

↳ $\alpha \xrightarrow{*} \beta$ represents that β can be derived from α in zero or more steps.

↳ $\alpha \xrightarrow{+} \beta$ represents that β can be derived from α in one or more steps.

↳ $L(G)$ is the language of G which is set of sentences.

↳ A sentence of $L(G)$ is string of terminals of G .

b) i.e., if s is start symbol of G , then w is a sentence of $L(G)$ if and only if $s \Rightarrow^* w$, where w is string of terminals of G .

↳ If G is a CFG then, $L(G)$ is a CFL.

↳ Two grammars are equivalent if they produce the same language.

Eg:-

Consider a grammar $G = (V, T, P, S)$

where,

$$V = \{E\}$$

$$T = \{+, *, -, (,), id\}$$

$$P = \{ E \rightarrow E + E \}$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow E - E$$

$$E \rightarrow id$$

$$S = E$$

for $id * id + id$

L.M.D

$$\overline{E} \xrightarrow{m} \overline{E} + \overline{E}$$

$$\overline{E} \xrightarrow{m} \overline{E} * \overline{E} + \overline{E}$$

$$\overline{E} \xrightarrow{m} id * \overline{E} + \overline{E}$$

$$\overline{E} \xrightarrow{m} id * id + \overline{E}$$

$$\overline{E} \xrightarrow{m} id * id + id$$

R.M.D

$$E \xrightarrow{m} E * E$$

$$\xrightarrow{m} E * E + E$$

$$\xrightarrow{m} E * E + id$$

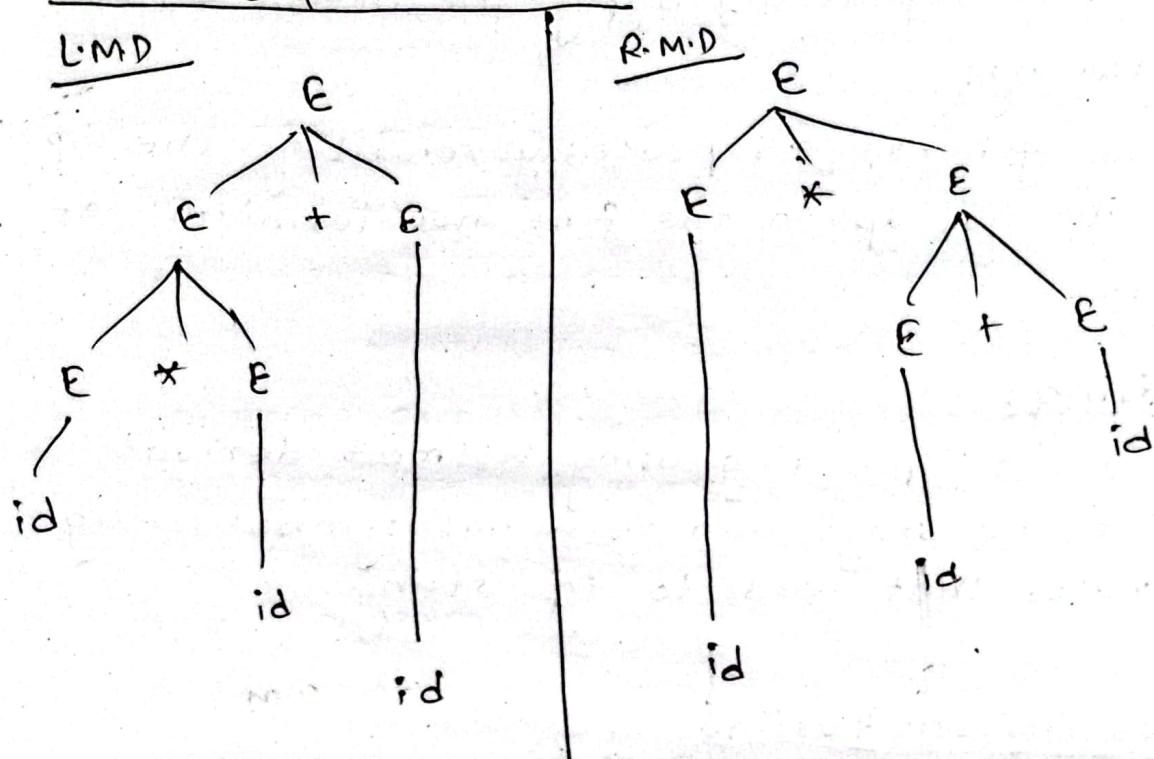
$$\xrightarrow{m} E * id + id$$

$$\xrightarrow{m} id * id + id$$

*Parse Tree:

↳ is a graphical representation of CFG derivation.

↳ here, root is the start symbol, internal nodes represent non terminals and leaf nodes represent terminals.

Parse tree for $id * id + id$ 

↳ if a grammar produce more than one unique parse tree for a sentence or derivation (left or right) then it is called ambiguous grammar.

Eg:

$$E \rightarrow E+E$$

$$E \rightarrow E \times E$$

$$E \rightarrow id$$

for $id + id * id$ ① LMD

$$\stackrel{1m}{\Rightarrow} E + E$$

$$\stackrel{2m}{\Rightarrow} id + E$$

$$\stackrel{3m}{\Rightarrow} id + E \times E$$

$$\stackrel{4m}{\Rightarrow} id + id \times E$$

$$\stackrel{5m}{\Rightarrow} id + id \times id$$

② LMD

$$E \stackrel{1m}{\Rightarrow} E \times E$$

$$\stackrel{2m}{\Rightarrow} (E + E) \times E$$

$$\stackrel{3m}{\Rightarrow} (id + E) \times E$$

$$\stackrel{4m}{\Rightarrow} id + id \times E$$

$$\stackrel{5m}{\Rightarrow} id + id \times id$$

It is ambiguous.

*Types of Parsing:1) Top Down Parsing:

↳ Parsing involves generating the string from the first non-terminal and repeatedly applying production rules.

2) Bottom Up Parsing:

↳ Parsing involves generating repeatedly rewriting the i/p string until it ends up in the first non-terminal of the grammar.

*Top Down Parsing:

↳ Tries to derive the i/p string using leftmost derivation i.e. starting at the start non terminal symbol using production rules that leads to i/p string.

↳ It includes:

a) Recursive Descent Parsing:

↳ Backtracking is needed i.e. if a choice of production rule does not work, we have to back track to check for next production rule.

↳ Not widely used, not efficient.

b) Non-recursive Predictive Parsing:

↳ No backtracking.

↳ Efficient.

↳ Use LL (left to right scanning, left to right derivation) method.

↳ needs special form of grammar (LL(1) grammar)

↳ Non-recursive (table driven) predictive parser is also known as LL(2) parser.

*Recursion:
 ↳ start symbol
 first mat
Rules:
 1) U
 2)
 3)
 4)

* Recursive Descent Parsing :

↳ Start with the starting non-terminal and use the first production, verify with input and if there is no match backtrack and apply another.

Rules:

- 1) Use two pointers: iptr for pointing i/p symbol to be read and optr for pointing the symbol of o/p string that is ~~use~~ initially start symbol s.
- 2) If the symbol pointed by optr is non terminal, use the first production ^{rule} for expansion.
- 3) While the symbol pointed by iptr and optr is same increment both pointers.
- 4) The loop at the above step terminates when
 - a) A non terminal at o/p (Case A)
 - b) the end of i/p (Case B).
 - c) Unmatching terminal symbol pointed by iptr and optr is seen (Case C).
- 5) If Case A is true, expand non terminal using first rule and goto step 2.
- 6) If Case B is true, terminate with success.
- 7) If Case C is true, decrement both pointers to the place of last non terminal expansion and use the next production rule for non terminal.
- 8) If there is no more production and B is not true, report error.

Example:

let grammar G1 where,

$$S \rightarrow cAd$$

$$A \rightarrow ab|a$$

for i/p cad;

input	output	Rules
(iptr) cad	(optr) S	Rule 1
(iptr) cad	(optr) CAD	Rule 2
c(iptr)ad	c (optr) Ad	Rule 3, match c.
c(iptr)ad	C (optr) abd	Rule 5.
ca(iptr)d	ca (optr) bd	Rule 3, match a
c(iptr)ad	c (optr) Ad	Rule 7, back track
c(iptr)ad	c (optr) ad	Rule 5.
ca(iptr)d	ca (optr) d	Rule 3, match a
cad(iptr)	cad (optr)	Rule 3, match d.

Now, process terminates with success.

g) consider a grammar G1 where

$$A \rightarrow oAi$$

$$A \rightarrow B$$

$$B \rightarrow H$$

Perform recursive descent parsing for string OOHLL.

CS

for i/p 00H11,

02-14
CD

input	output	Rule
(iptr) 00H11	(optr) A	Rule 1
(iptr) 00H11	(optr) OA1	Rule 2
00(iptr) 0H11	0(iptr) A 1	Rule 3, match 0.
0(iptr) 0H11	0(iptr)'OA11	Rule 2 .
00(iptr) #11	00(iptr) A 11	Rule 3, match 0
00(iptr) #11	00(iptr) OA111	Rule 2
00(iptr) H11	00(iptr) A 11	Rule 7, backtrack
00(iptr) H11	00(iptr) B 11	Rule 2
00(iptr) H11	00(iptr) # 11	Rule 2
00#(iptr) 11	00#(optr) 11	Rule 3, match #
00#1(iptr) 1	00#1(iptr) 1	Rule 3, match 1
00#11(iptr)	00#11(iptr)	Rule 3, match 1

Now,
process terminates with success.

for i/p 0011,

→ Ctd

input	output	Rules
(iptr) 0 0 1 1	0 (ptr) A	Rule 1
(iptr) 0 0 1 1	0 (ptr) 0 A 1	Rule 2
0 (iptr) 0 1 1	0 0 (ptr) A 1 1	Rule 3, match 0
0 0 (iptr) 0 1 1	0 0 (ptr) 0 A 1 1	Rule 2
0 0 0 (iptr) 1 1	0 0 0 (ptr) A 1 1 1	Rule 3, match 0
0 0 0 (iptr) 1 1	0 0 0 (ptr) A 1 1	Rule 2
0 0 0 (iptr) 1 1	0 0 0 (ptr) B 1 1	Rule 7, backtrack
0 0 0 (iptr) 1 1	0 0 0 (ptr) # 1 1	Rule 2
0 0 0 (iptr) 1 1	0 0 0 (ptr) B 1 1	Rule 7, backtrack

Error! reported as ~~no~~ no production rules are there and ips are not ended yet.

*Left Recursion:

- ↳ the grammar with recursive non terminal at the left of production is called left recursive grammar.
- ↳ A grammar is left recursive if it has a non-terminal 'x' such that there is derivation
 $x \Rightarrow x\alpha$ for some string α .
- ↳ Left recursion causes recursive descent parser to go to infinite loop.
- ↳ Top down parsing techniques can't handle left recursive grammar.
- ↳ Hence, we have to convert left recursive grammar into an equivalent grammar which is not left recursive.

↳ the left recursion may appear in a single step of derivation called immediate left recursion or co may appear in more than one step of the derivation.

02-19/2L

CO

082-02-21

↳ example of non immediate left recursion

$$S \rightarrow Aablc$$

$$A \rightarrow Sclc$$

- this grammar is not immediate left recursive, but is still left recursion i.e,

$$S \rightarrow Aab \rightarrow \text{Scab}$$

or
 $A \rightarrow Sc \rightarrow Aabc$

*Removing left recursion

↳ if we have the grammar of the form $A \rightarrow A\alpha | \beta$ the rule for removing the left recursion is,

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' | \epsilon$$

↳ Any production of the form

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_n | \beta_1 | \beta_2 | \dots | \beta_n$$

we have,

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

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon$$

*left factoring:

↳ Backtracking is expensive operation so if we are able to predict the unique prodⁿ rule while parsing then we are saving a lot of time resources.

↳ left factoring is the approach of transforming the grammar that has two or more prodⁿ rules with same

initial substring to the one that has not, so that it will be useful for predictive parsing.

02-21
CD

↳ For the grammar of the form,

$$A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \dots | \alpha \beta_n$$

left factoring yield,

$$A \rightarrow \alpha A'$$

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

* Predictive Parsing:

- ↳ Recursive descent parsing is inefficient as it causes lot of backtracking and suffers from left recursion problem.
- ↳ Predictive parsing tries to predict which production produces the least chance of backtracking and infinite loop.
- ↳ Predictive parser can choose the production to apply solely on the basis of next input symbol and the current non-terminal being processed.
- ↳ To enable this, the grammar must take a particular form, such a grammar is called LL(1).
- ↳ first L means scan the i/p from left to right.
- ↳ second L means create leftmost derivation.
- ↳ 1 means one i/p symbol of lookahead.
- ↳ LL(1) has no left recursion prodn and has been left factored.

Eg:-

$$\text{stmt} \rightarrow \text{if } \dots \dots \mid$$

 • while $\dots \dots$)

 begin $\dots \dots$ |
 for $\dots \dots$

↳ Here, for non terminal 'stmt', if next token is 'begin', then we use third production rule instead of first production as in recursive descent parsing.

* Non Recursive Predictive Parsing :

↳ It is built by maintaining a stack explicitly rather than implicitly through recursive calls.

↳ Is a table driven parsing.

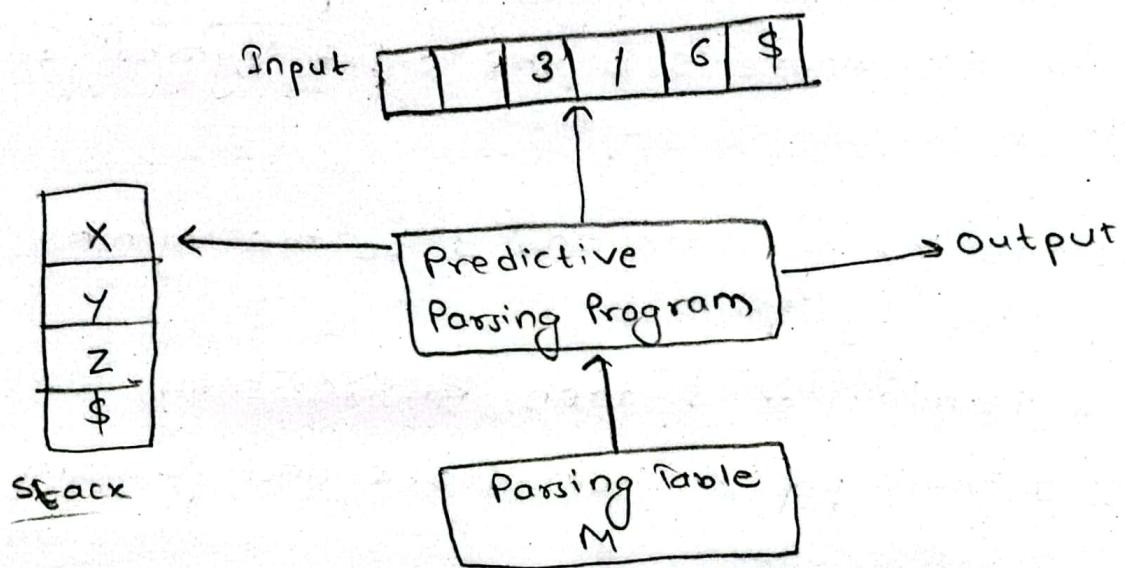


Fig:-Non-Recursive Predictive Parsing

↳ A table driven predictive parser has stack input buffer, parsing table and output stream.

↳ Input buffer contains sentence to be parsed and at the end \$ as end marker.

↳ Stack contains symbols of grammar, initially stack contains start symbol on top of \$.

↳ When stack is empty (only \$ left on stack) parsing is completed.

↳ Parsing table is two dimensional array containing the information about the production to be used on seeing terminal at input and non terminal at stack.

↳ each entry in parsing table holds prodⁿ rule.

↳ each row holds non terminal symbol.

↳ each column holds terminal or \$.

*Algorithm for Parsing

input: a string w.

output: if w is in L(G), a leftmost derivation of w,
else error.

ip → input pointer

→ 1. set ip to point the first symbol of input stream.

2. set stack to \$s where s is the start symbol of grammar.

3. do

let x be the top of stack and a be the symbol
~~pointed by ip~~ pointed by ip.

if x is a terminal or \$ then ~~pop x~~

if x=a then pop x from stack and forward p.

else

error()

else

if $M[x, a] = x \rightarrow y_1, y_2, \dots, y_n$, then

- pop x from stack

- push y_n, y_{n-1}, \dots, y_1 onto stack with y_1 on top

- output production $x \rightarrow y_1, y_2, \dots, y_n$

else

error()

4. until $x = \$$

*Constructing LL(1) Parsing Table:

02-21

CD

- ↳ Parsing table requires two functions FIRST and FOLLOW.
- ↳ for a string of grammar symbols α .
FIRST(α) is the set of terminals that begin all possible strings derived from α .
- ↳ if $\alpha \xrightarrow{*} \epsilon$, then ϵ is also in FIRST(α).
- ↳ FOLLOW(A) for non terminal A is the set of terminals that can appear immediately to the right of A in some sentential form
 - if A can be last symbol in a sentential form then \$ is also in FOLLOW(A).
 - a terminal 'a' is in FOLLOW(A) if $s \xrightarrow{*} \alpha A a \beta$.
 - \$ is in FOLLOW(A) if $s \xrightarrow{*} @ A @$.

Computation of FIRST

- ↳ FIRST(α) = {set of terminals that begins all strings derived from α }.
- if 'a' is a terminal symbol then
 $FIRST(a) = \{a\}$.
- if 'X' is a non terminal symbol and $X \rightarrow \epsilon$ is a prodⁿ rule
then, $FIRST(X) = FIRST(X) \cup \epsilon$
- if 'X' is a non terminal symbol and $X \rightarrow Y_1 Y_2 \dots Y_n$ is a prodⁿ rule, then
i) if a terminal 'a' is $FIRST(Y_1)$ then
 $FIRST(X) = FIRST(X) \cup FIRST(Y_1)$.

ctd

ii) if a terminal 'a' is in $\text{FIRST}(Y_i)$ and ϵ is in all $\text{FIRST}(Y_j)$ for $j=1 \dots i-1$ then
 $\text{FIRST}(X) = \text{FIRST}(X) \cup a$

iii) if ϵ is in all $\text{FIRST}(Y_j)$ for $j=1 \dots n$ then

$$\text{FIRST}(X) = \text{FIRST}(X) \cup \epsilon$$

↳ if X is ϵ then, $\text{FIRST}(X) = \{\epsilon\}$

↳ if X is $Y_1 Y_2 \dots Y_n$

i) if a terminal 'a' is in $\text{FIRST}(Y_i)$ and ϵ is in all $\text{FIRST}(Y_j)$ for $j=1 \dots i-1$ then

$$\text{FIRST}(X) = \text{FIRST}(X) \cup a$$

ii) if ϵ is in all $\text{FIRST}(Y_j)$ for $j=1 \dots n$ then

$$\text{FIRST}(X) = \text{FIRST}(X) \cup \epsilon$$

Computation of FOLLOW

↳ $\text{FOLLOW}(A) = \{\text{set of terminals that can immediately follow non-terminal } A\}$

↳ apply the following rules until nothing can be added to any FOLLOW set.

i) if S is the start symbol, then $\$$ is in $\text{FOLLOW}(S)$.

ii) if $A \rightarrow \alpha \beta \beta$ is a prodⁿ rule then everything in $\text{FIRST}(\beta)$ is placed in $\text{FOLLOW}(\beta)$ except ϵ .

iii) if ($A \rightarrow \alpha \beta$ is a prodⁿ rule) or ($A \rightarrow \alpha \beta \beta$) is a rule and ϵ is in $\text{FIRST}(\beta)$ then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(\beta)$.

Algorithm for Parsing Table

input: LL(1) Grammar

output: Parsing Table M

for each prodⁿ rule $A \rightarrow \alpha$ in grammar G do

 for each $a \in \text{FIRST}(\alpha)$ do

 Add $A \rightarrow \alpha$ to $M[A, a]$

 End do

 if $\epsilon \in \text{FIRST}(\alpha)$ then

 for each $b \in \text{FOLLOW}(A)$ do

 Add $A \rightarrow \alpha$ to $M[A, b]$

 End do

End if.

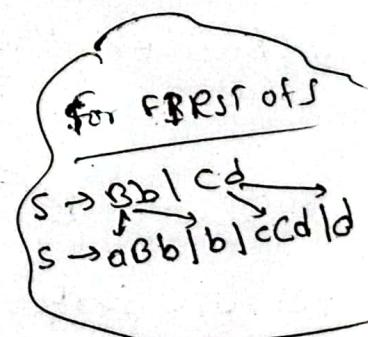
End do

Mark each undefined entry in M error

082-02-22

①	FIRST	FOLLOW
$S \rightarrow aABb$	{a}	{\$}
$A \rightarrow c \epsilon$	{c, ϵ }	{d, b}
$B \rightarrow d \epsilon$	{d, ϵ }	{b}

②	FIRST	FOLLOW
$S \rightarrow Bb Cd$	{a, b, c, d}	{\$}
$B \rightarrow ab \epsilon$	{a, ϵ }	{b}
$C \rightarrow cc \epsilon$	{c, ϵ }	{d}



③.

	FIRST	FOLLOW
$S \rightarrow ACB \mid CbB \mid Ba$	$\{d, g, h, b, a, \epsilon\}$	$\{\$, \}$
$A \rightarrow da \mid BC$	$\{d, g, h, \epsilon\}$	$\{h, g, \}, \$\}$
$B \rightarrow g \mid \epsilon$	$\{g, \epsilon\}$	$\{\$, a, h, g\}$
$C \rightarrow h \mid \epsilon$	$\{h, \epsilon\}$	$\{h, g, b, \}, \$\}$

In FOLLOW, if non-terminal in right, write the FIRST of that non-terminal, if ϵ is there replace it into main rule

④.

	FIRST	FOLLOW
$S \rightarrow ABCDE$	$\{a, b, c\}$	$\{\$\}$
$A \rightarrow a \mid \epsilon$	$\{a, \epsilon\}$	$\{b, c\}$
$B \rightarrow b \mid \epsilon$	$\{b, \epsilon\}$	$\{c\}$
$C \rightarrow c$	$\{\epsilon\}$	$\{d, e, \$\}$
$D \rightarrow d \mid \epsilon$	$\{d, \epsilon\}$	$\{e, \$\}$
$E \rightarrow e \mid \epsilon$	$\{e, \epsilon\}$	$\{\$\}$

Q. Consider the grammar

$$S \rightarrow aBa$$

$$B \rightarrow bB \mid \epsilon$$

Now, perform non recursive predictive parsing on string 'abba'

Soln:

Initially, find FIRST and FOLLOW

02-22

CD

	FIRST	FOLLOW
$S \rightarrow aB\alpha$	{a}	{\\$}
$B \rightarrow bB \mid \epsilon$	{b, \epsilon}	{a}

Now, construct parsing table.

	a	b	\$
S	$S \rightarrow aB\alpha$.	.
B	$B \rightarrow \epsilon$	$B \rightarrow bB$.

Stack	Input	Output
\$S	a b b a \$	$S \rightarrow aB\alpha$
\$aB\alpha	↑ a b b a \$	
\$aB	↑ a b b a \$	$B \rightarrow bB$
\$aBb	↑ a b b a \$	
\$aB	↑ a b b a \$	$B \rightarrow bB$
\$aBb	↑ a b b a \$	
\$aB	↑ a b b a \$	$B \rightarrow \epsilon$
\$a	↑ a b b a \$	
\$	↑ a b b a \$	

successfully complete i.e. accept.

CS

Output:

$$S \rightarrow aB_0, B \rightarrow bB, B \rightarrow bB, B \rightarrow \epsilon \quad \textcircled{2}$$

∴ derivation.

$$\begin{aligned} S &\Rightarrow aB_0 \\ &\Rightarrow abB_0 \\ &\Rightarrow abbB_0 \\ &\Rightarrow abba \end{aligned}$$

Q

$$S \rightarrow aABb$$

$$A \rightarrow c | \epsilon$$

$$B \rightarrow d | \epsilon$$

Perform non-recursive predictive parsing for string
 'acdb' ~~and check~~

Soln:

Initially

[Take reference from previously done i.e. FIRST & FOLLOW]

Now, construct parsing table

	a	b	c	d	\$
s	$s \rightarrow aABb$				
A		$A \rightarrow \epsilon$	$A \rightarrow c$	$A \rightarrow \epsilon$	
B		$B \rightarrow \epsilon$		$B \rightarrow d$	

Stack

\$S

\$bBAa

\$bBA

\$bBc

\$bB

\$bd

\$b

\$

input

~~acdb\$~~a ~~acdb\$~~~~acdb\$~~

acdb\$

acd b\$

acd b\$

acd b\$

acd b\$

acd b\$

O/p.

s → aABb

A → c

B → d

Successfully complete i.e accept

Output:

s → aABb, A → c, B → d

Derivation

s ⇒ aABb

⇒ acBb

⇒ acdb

