

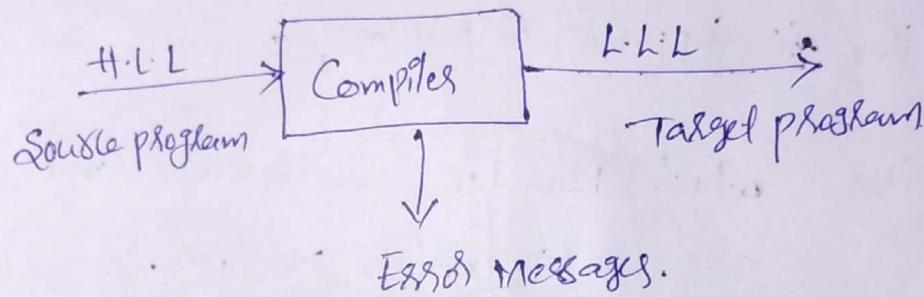
Introduction to Compiler Design!

①

Overview:-

Compiler: → compiler is a program that can read a program in one language - the source language - and translate it into an equivalent program in another language - the target language.

* An important role of the compiler is to report any errors in the source program that occurs during the translation process.



Why should we study Compiler Design?

Compilers are everywhere, Many applications of compiler technology.

1. Machine code generation for higher level languages
2. Software testing
3. Program optimization
4. Malicious code detection

→ Compiler is possibly the most complex system software.

System Software: Set of programs are to be used to perform a particular task.

Ex: operating system.

Application software: Set of applications which can be used to perform required task.

Ex: MS Excel, MS Word.

List of compilers:

- C compilers, • C++, • ALGOL, • BASIC, • COBOL,
- Java.

Depending on how they have been constructed (or) on what function they are supposed to perform, compilers are classified as.

Cross-compilers: which runs on one machine and generates a code for another machine.

Incremental Compilers: that allows a modified portion of a program to be recompiled.

Ideal compilers:

- * take less time for compilation
- * be smaller in size (if it is not suitable for large code)
- * be written in high-level languages.

Programming Languages:

(3)

* P.L; it is used to communicate between the instructions.

P.L classified into 3 types.

1, low level (or) Machine Level

2, high level

3, Assembly language (mnemonic code)

Ex: ADD, SUB, DIV, MUL.

Note: Machine language and Assembly language are example of low level language.

Translators: Translator is, convert one form of program to another form.

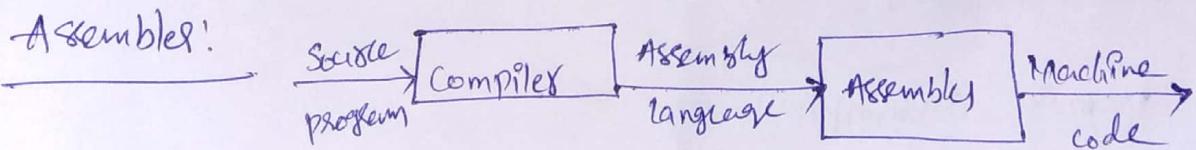
* Compilers

* Assemblers

* Interpreters.

* Assembler converts Assembly language program to Machine level language.

Assembler:



* An assembly code is a mnemonic version of machine code.

The typical assembly instructions are as given below.

Mov a, R,

MUL #15, R,

ADD R1, b.

The Assembler converts these instructions in the binary language which can be understood by the machine. Such a binary code is often called as machine code. This machine code is a Relocatable machine code that can be passed directly to the loader/linker for execution purpose.

The Assembler converts the assembly program to low level machine language using Two Passes.

* * Note: A pass means one complete scan of the input program.
The end of second pass is the Relocatable machine code.

Interpreters: An interpreter is a kind of translator.

The source program gets interpreted every time it is to be executed, and every time the source program is analysed. Hence interpretation is less efficient than compilers.

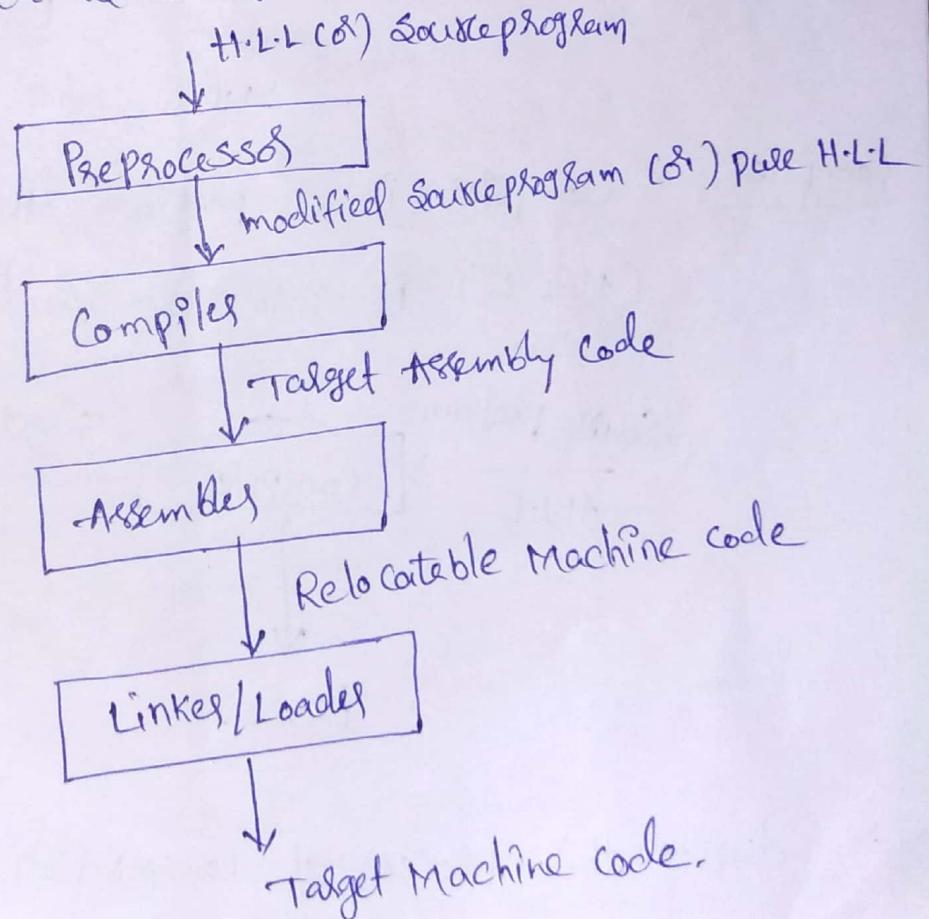
* The interpreters do not produce object code.

* Memory consumption is more.

Language Processing system:

(5)

- we have learn that any computers system is made of hardware and software. Hardware is just a piece of mechanical devices and it's functions are being controlled by a compatible software.
- Hardware understands instructions in the form of electronic charge, which is the counterpart of binary language in software programming.
- Hardware understands a language, which humans can not understand, so we write programs in high-level languages which is easier for us to understand and Remember.



Preprocessor: It is a part of a compiler, that produces input for a compiler.

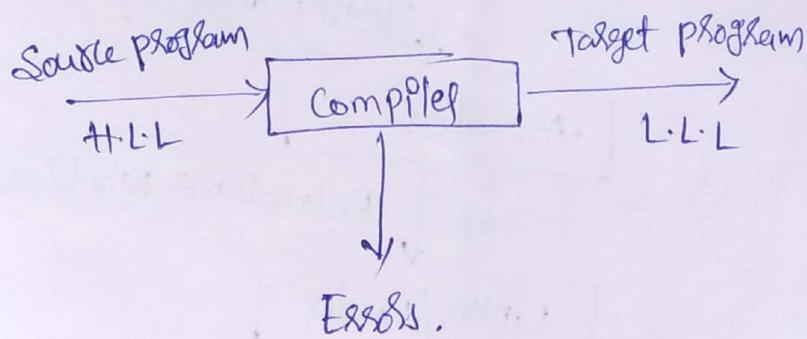
→ It deals with 2 types of functions File Inclusion Macro processing.

File Inclusion: A preprocessor may include header files into the program text.

Macro-processing: Macro is a small set of instructions, while we are running a program, whenever macro name is identified then that ~~the~~ name is replaced by "macro definition".

Ex: `#define PI 3.14` ↓ → Macro definition.
 Macro.

Compiler: Compiler is a program, that translates one language (H.L.L) into equivalent target language (L.L.L)



- * during this process of conversion, it returns errors.
→ those are handled by the error handler.

Assembles: To every platform (hardware to.s) we are having 2
an Assembler. These are not universal, since for
each platform we have one.

- It converts an Assembly language into machine code
i.e., binary representation.
- The o/p of an assembler is called an object file,
which contains a combination of machine instructions
as well as the data required to place these instructions
in memory.

Loader: Loader is a part of operating system and it is
responsible for loading executable files into memory
and execute them. It calculates size of the
program and creates memory space for it.

Linker: Linking, it is the process of putting together others
programs files and functions that are required by
the program.

→ Linker Resolves external memory addresses.

Ex: If the program is using exp() function, then the object
code of this function should be brought from the
"math library" of the system and linked to the main
program.

Goals of Compiler:

- It maintains the meaning of the program that is compiled
- Compiles must improve the performance of I/P programs i.e., speed of the programs.
- Reduce the time required for compilation so, as to support rapid development and debugging cycles.
- Give accurate error message.

Differences between Compiler and Interpreter:

Compiler

- Compiler checks (or) scans the entire h.l. program at once.
- The compilers produce object code
- The compiler is a complex program and it requires large amount of memory.
- Processing time is less
- Some compilation languages are C, C++, FORTRAN, PASCAL, Ada etc. (it is also called SW Translation)

Interpreter

- Interpreters translates one statement at a time.
- The interpreters do not produce object code.
- Interpreters are simple and give us improved debugging environment.
- Processing time is more.
- It is also called as software simulation. Some interpretation languages are ML, LISP, PYTHON, PROLOG etc.

Phases of compilation:-

⑦

There are Two parts of compilation

1, Analysis part

2, Synthesis part.

Analysis part consists of 3 phases.

1, Lexical analysis phase

2, syntax analysis phase

3, semantic analysis phase.

The synthesis part consists of 3 phases.

4, Intermediate code Generation

5, code optimization

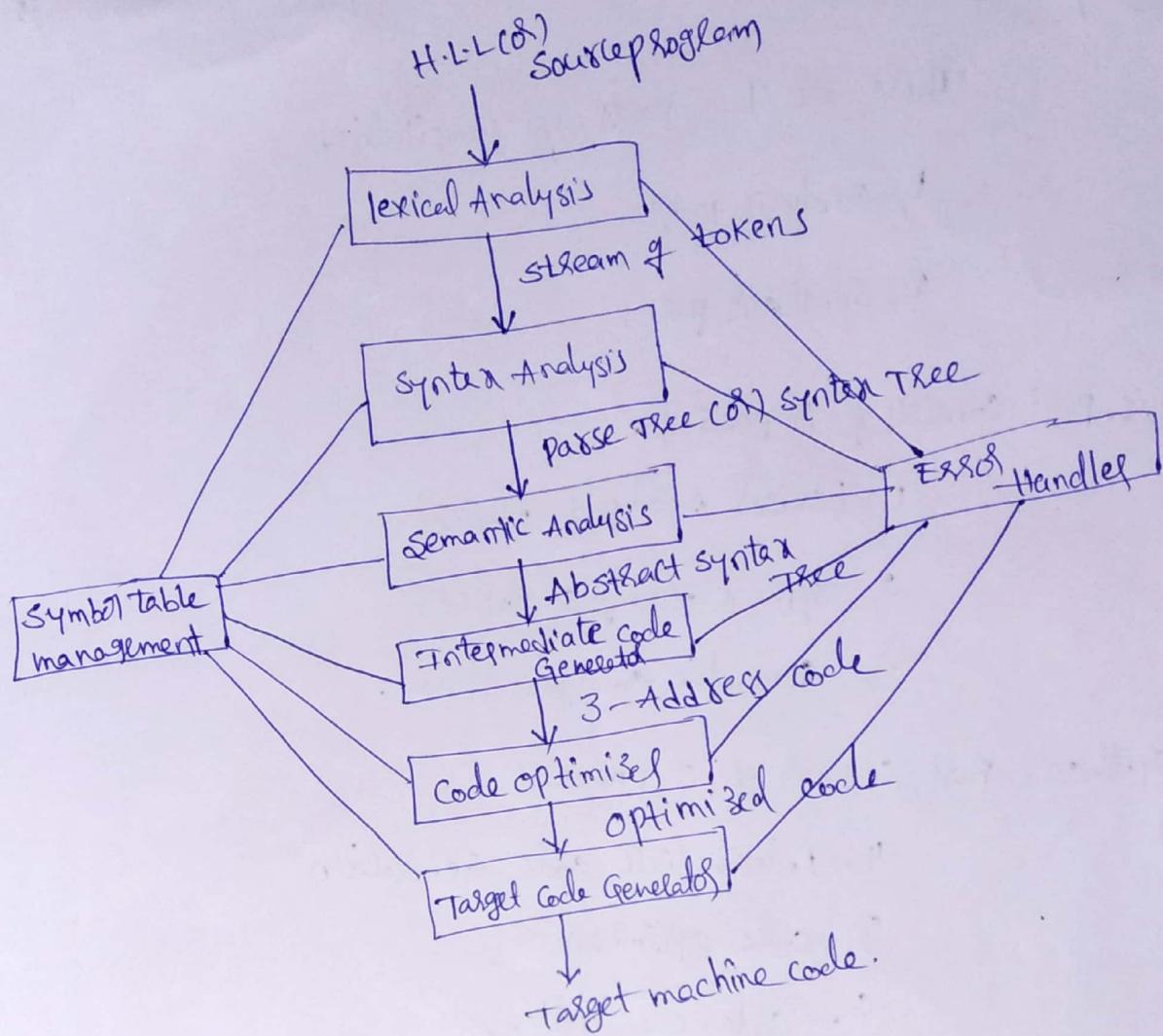
6, Code Generation.

* Analysis part, it is also called as Front end of a compiler.

* synthesis part, it is also called as Back end of a compiler.

→ Analysis Phase breaks up the source program into constituted pieces and imposes a grammatical structure on them. It uses this structure to create an Intermediate Representation of the source program.

→ synthesis phase constructs the desired target program from the intermediate representation and the information in the symbol table.



Lexical Analysis Phase:

- * Lexical Analysis phase is also called scanner.
- * This phase Ready the characters in the source program and groups them into a stream of tokens.
- * token is a sequence of character (or) string having a collective meaning.
- * the character sequence forming a token is called the lexeme for that token.

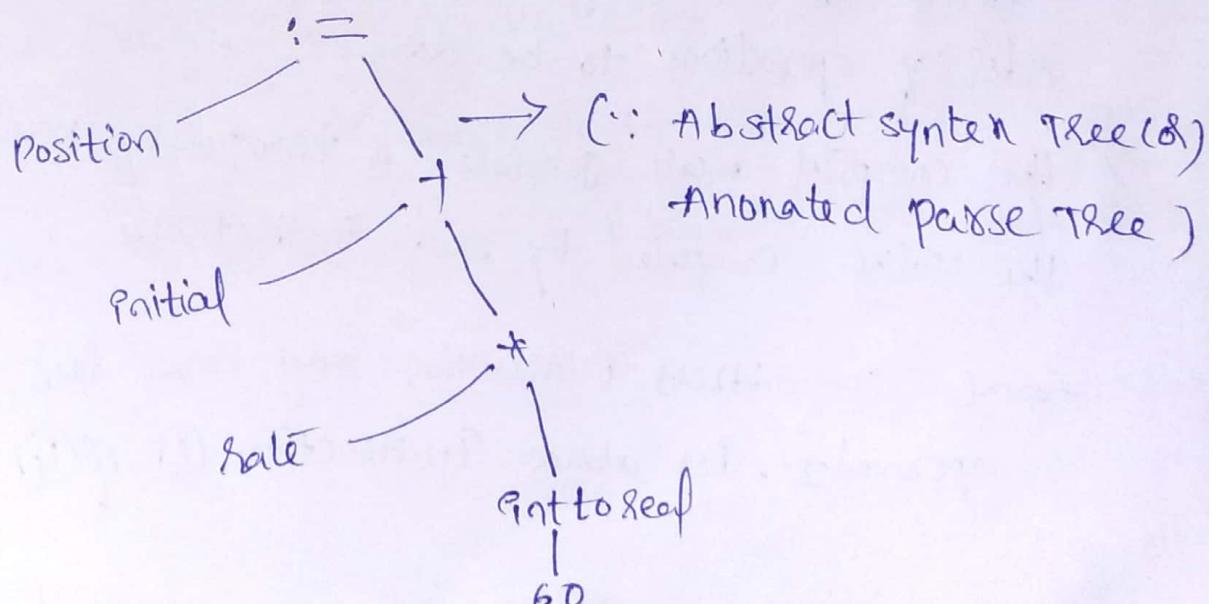
Semantic Analysis phase:

(P)

- * Semantic analysis checks the source program for semantic errors. (matching if-else stmt, missing parenthesis, typechecking).
- * It uses the hierarchical structure determined by the syntax analysis phase to identify the operators and operands of expressions and statements.
- * Semantic analysis performs Type checking. i.e., it checks that each operator has operands that are permitted by the source language specifications.

Ex: If a real number is used to index an array i.e., A[1.5] then the compiler will report an error. This error is handled during Semantic Analysis.

Suppose it is assumed that all identifiers (id₁, id₂, id₃) have been declared to be real. But 60 assumes itself to be an integer. When * is applied to a real rate and integer 60, integer is converted into real during the type checking process. This is achieved from the operator Int to Real, which converts an integer into real.



Intermediate Code Generation:

- * After syntax and semantic analysis compilers generate an intermediate representation of the source program. This representation should be easy to produce and easy to translate into the target program.
- * Thus the Intermediate code Generation phase transforms the parse tree into an intermediate-language representation of the source program.
One popular type of intermediate language is called as "Three-address" code which is like the assembly language. Three-address code for the statement
$$\boxed{\text{Position} := \text{Initial} + \text{Rate} \times \text{Goto}}$$
 is given as.

$\text{temp}_1 := \text{Pint to Real (60)}$

$\text{temp}_2 := \text{Id}_3 * \text{temp}_1$

$\text{temp}_3 := \text{Id}_2 + \text{temp}_2$

$\text{Id}_1 := \text{temp}_3$.

3-Address code has several properties. They are,

- Each 3-Address instruction has almost one operator in addition to the assignment. compiler has to decide the order of operations to be done.
- The compiler must generate a temporary name to hold the value computed by each instruction.
- Some 3-Address instruction may have less than 3 ~~operator~~ 3-operands. In above instruction ① of ④.

Consider an Assignment statement

(B) (A)

$$\boxed{\text{Position} := \text{Initial} + \text{Rate} * 60}$$

This would be grouped into the following tokens.

- 1, the Identifier position
- 2, the Assignment symbol :=
- 3, the Identifier Initial
- 4, the plus sign
- 5, the Identifier Rate
- 6, the multiplication sign
- 7, the number.

The blanks separating the characters will be eliminated during lexical analysis. When the lexical Analyser finds the Identifier position, it generates a token, say id. The Identifier position, rate and initial are the lexemes. The statement after lexical analysis is given by,

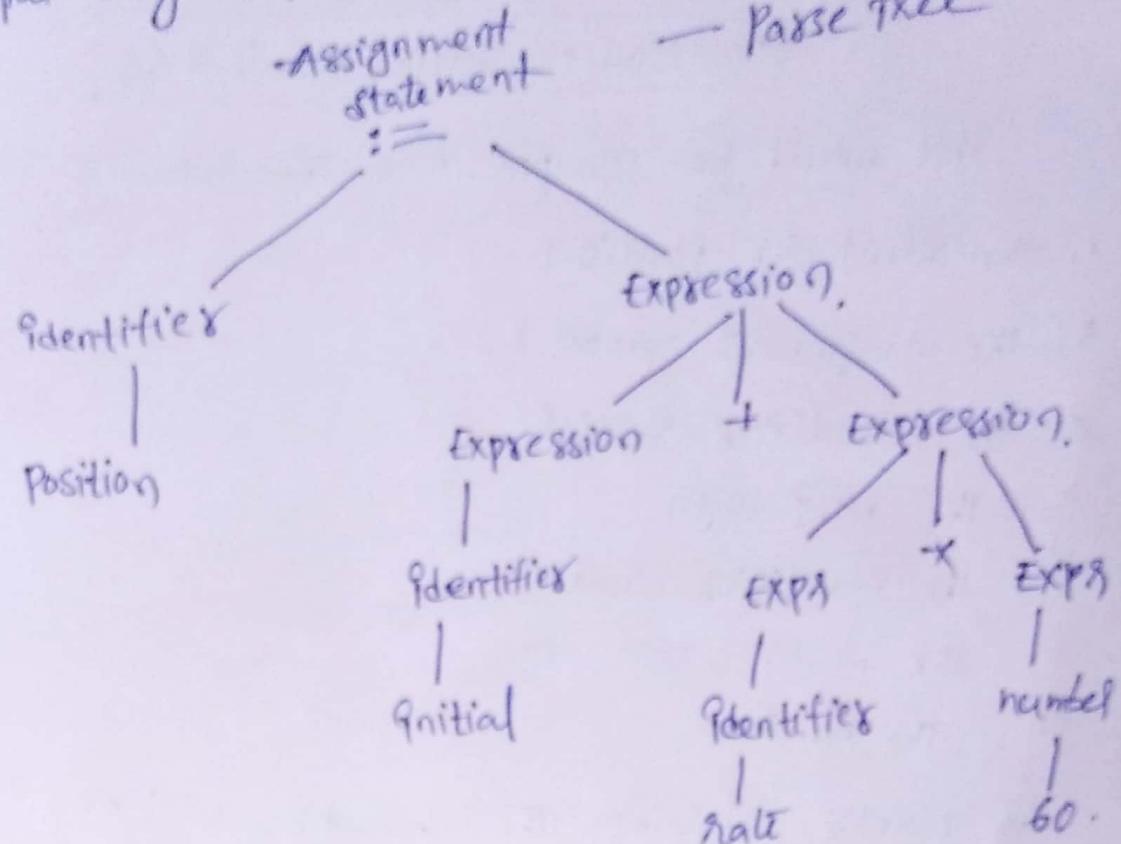
$$\boxed{\text{id}_1 := \text{id}_2 + \text{id}_3 * 60}$$

($\because \text{id}_1, \text{id}_2, \text{id}_3$ are tokens for initial, position, rate).

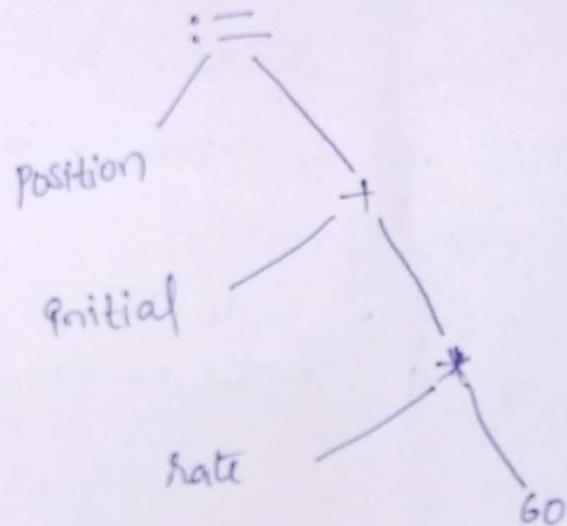
Syntax Analysis Phase:

- * The syntax analysis phase is also called Parser.
- * In this phase, the tokens generated by the lexical analyser are grouped together to form a hierarchical structure called "Parse Tree".

For the input string $q_{d_1} := q_{d_2} + q_{d_3} \times 60$



Syntax Tree:



* A syntax tree is a compressed representation of the parse tree. In which, operators appear as the internal nodes and operands of an operator are the children of node to that operator.

Code Optimization: code optimization phase improves the intermediate code i.e., it reduces the code by removing the repeated (or) unwanted instructions from the intermediate code. The above given intermediate code can be optimized to the following code.

$$\begin{aligned} \text{temp}_1 &:= f_{d_3} + 60.0 \\ f_{d_1} &:= f_{d_2} + \text{temp}_1 \end{aligned}$$

Target code Generation: code generation phase converts the intermediate code into a target code consisting of sequenced machine code (or) assembly code that performs the same task. The above optimized code can be written using Registers R₁ and R₂ which is as given below.

```
LDF R2, fd3
MULF R2, R2, #60.0
LDF R1, fd2
ADDF R1, R1, R2
STF Rd1, R1.
```

Here # signifies that 60.0 is treated as constant.
f is a floating point number.

Symbol Table management: It collects the information about "data objects" that appear in source program
Ex: whether the variable represents int, float, etc

attribute may provide information about the storage allocated for a name, its type, its scope.

- * The symbol table is a data structure containing a record for each variable name, with fields for the attribute of the name.

Error Handler: In compilation, each phase detects errors, these errors are must be reported to error handles to handle the errors so that the compilation can proceed.

LEX: LEX is a unix utility which generates the lexical analyzer. A LEX lexer is very much faster in finding the tokens as compared to the handwritten LEX program part. The LEX program consists of 3 parts
1. Declaration section 2. Rule section 3. Procedure section.

1. {
Declaration section → declaration of variable constants can be done.

% }

% . %.
Rule section → consists of Regular Expression with associated actions.

% . %.

Auxiliary procedure section. → in which all the required procedures are defined.

- * Sometimes these procedures are required by the actions in the rule section.

Describe the output for the various phases of compiler with respect to the following statements.

(7)

$$\text{position} := \text{Pinitial} + \text{Rate} * 60$$

Sol: The O/P at each phase during compilation is as given below.

$$\text{position} := \text{Pinitial} + \text{Rate} * 60$$



Lexical Analysis



$$Pd_1 := Pd_2 + Pd_3 * 60.$$

Syntax Analysis



position
:=

Pinitial

rate

60

Semantic Analysis



Position
:=

Pinitial

rate

int to float

60.

Intermediate code generated



$$\text{temp}_1 := \text{Pinitial}$$

$$\text{temp}_2 := Pd_3 * \text{temp}_1$$

$$\text{temp}_3 := Pd_2 + \text{temp}_2$$

$$Pd_1 := \text{temp}_3$$



Code Optimized



$$\text{temp}_1 := Pd_3 * 60.0$$

$$Pd_1 := Pd_2 + \text{temp}_1$$



Code Generator



LDF R2, Pd_3

MULF R2, R2, #60.0

LDF R1, Pd_2

ADDF R1, R1, R2

STF Pd_1, R1

(Machine code)

1	position	...
2	Pinitial	...
3	Rate	...

symbol table



$$\text{temp}_1 := Pd_3 * 60.0$$

$$Pd_1 := Pd_2 + \text{temp}_1$$



Code Generator



LDF R2, Pd_3

MULF R2, R2, #60.0

LDF R1, Pd_2

ADDF R1, R1, R2

STF Pd_1, R1

(Machine code)

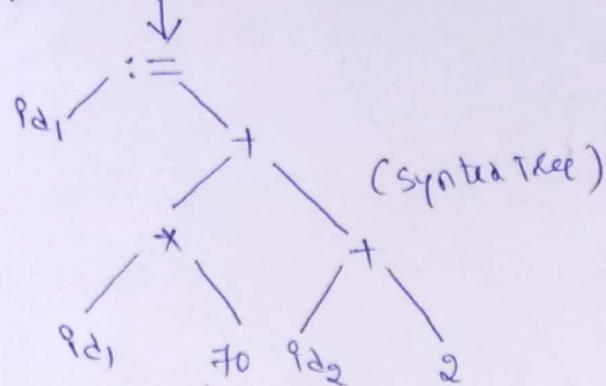
2. write off all phases of the compiler for the below statement.

$$i = i * 70 + i + 2.$$

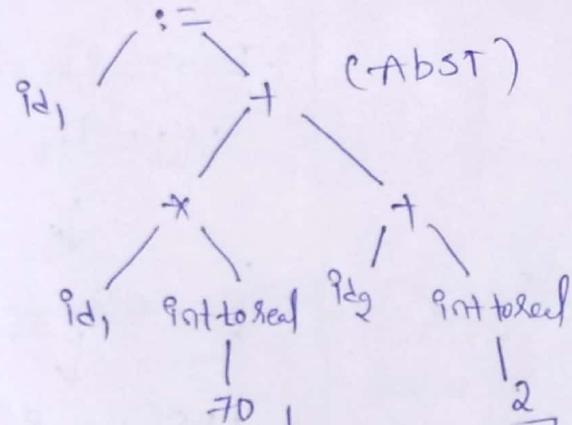
↓
lexical analysis

$i := i * 70 + i + 2$ (tokens)

↓
syntax analysis



↓
semantic analysis



↓
Intermediate Code Generation

$\text{temp}_1 := \text{int to real}(70)$

$\text{temp}_2 := \text{id}_1 * \text{temp}_1$ (3-Addres code)

$\text{temp}_3 := \text{int to real}(2)$

$\text{temp}_4 := \text{id}_2 + \text{temp}_3$

$\text{temp}_5 := \text{temp}_2 + \text{temp}_4$

$\text{id}_1 := \text{temp}_5$

↓
code optimization

$\text{temp}_1 := \text{id}_1 * 60.0$

$\text{temp}_2 := \text{id}_2 + 2.0$

$\text{temp}_3 := \text{temp}_1 + \text{temp}_2$

$\text{id}_1 := \text{temp}_3$ (optimized code)

↓
code generation

MOVF id_1, R_1

MULF #60.0, R_1

MOVF id_2, R_2

ADDF #2.0, R_2

ADDF R_2, R_1

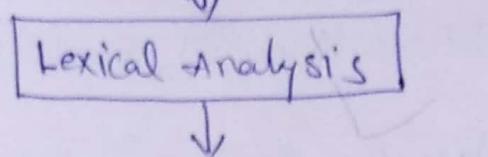
MOVF R_1, id_1

(Machine code)

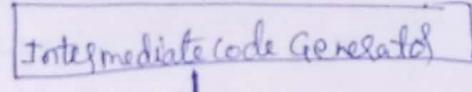
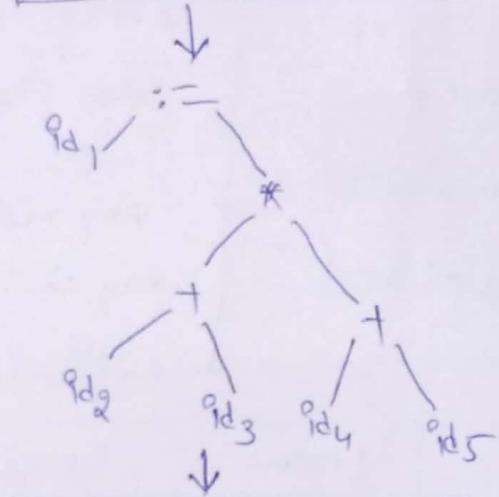
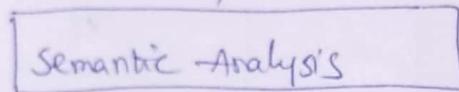
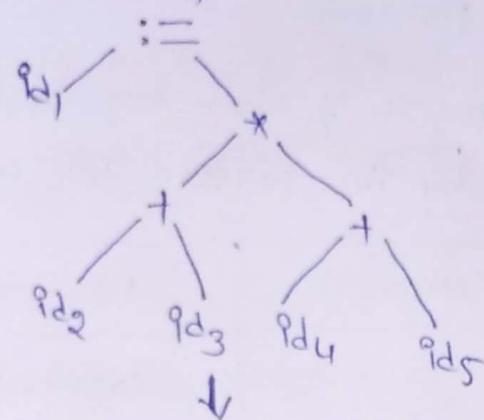
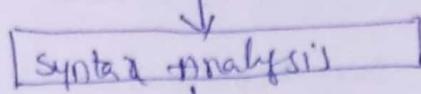
Show the o/p generated by each phase for the following expression.

$$x := (a+b) * (c+d),$$

19



$$q_{\mathcal{D}_1} := (\text{id}_2 + q_{\mathcal{D}_3}) * (\text{id}_4 + q_{\mathcal{D}_5})$$

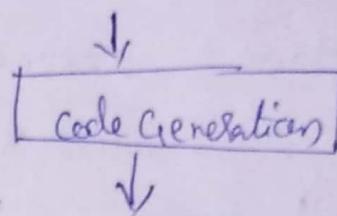


$$\text{temp}_1 := \vartheta_{d_2} + \vartheta_{d_3}$$

$$\text{temp}_2 = q_{d4} + i_{d5}$$

$$\text{temp}_3 := \text{temp}_1 \neq \text{temp}_2$$

$$P_{d_1} = \text{temp}_3$$



Mov %d2, R6

ADD QD₃, R₀

MOV R4, R1

APPENDIX R

MUL R₁, R₀

Mov R0, id₃.

(Machine code)

TOKEN, pattern, Lexeme.

TOKEN: it describes the stream of characters having a collective meaning. Ex: identifiers, keywords, constants, operators.

Pattern: set of rules that describe the token. which ^{the} same token is produced as output & it is a rule describing the set of lexemes that can represent a particular token in the source program.

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

Ex:

TOKEN	Sample Lexemes	Description of the pattern
const	const	constant
if	if	if
relation	<, <=, =, >, >=	< (0+) <= 08 = 08 > 08 >=
id	Pi, count, D2	letter followed by letters (0+) digits.
num	3.1416, 0, 6.02E23	Any numeric constant
literal	"language processor"	any characters enclosed b/w " and "

Ex: 2 int MAX (int a, int b)

{
if (a > b)

return a;

else

return b;

}

(*) The blank and new line characters can be ignored.

Lexeme	Token
int	keyword
MAX	identifier
(operator
int	keyword
a	identifier
)	operator
:	:

Role of lexical Analyser (or) Need for lexical analyser:

②

- * lexical Analyser it is first phase of compiler, the lexical Analyser reads the character from source program from left to right one character at a time and generates the sequence of tokens.
- * Each token is a single logical cohesive unit such as identifiers, keyword, operators and punctuation marks.
- * and it also count, to tell the how many tokens are generated from given source program.

Ex:- `int max(x,y)`

`int x,y;`

`/* find max of x and y */`

`{`

`return (x>y ? x : y);`

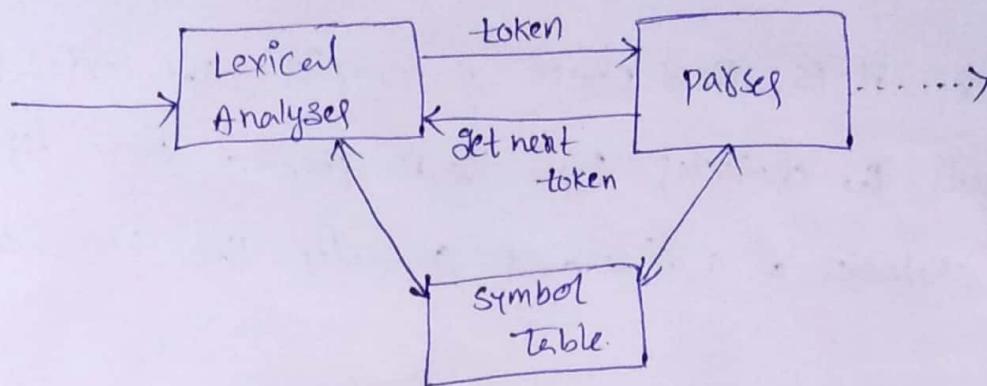
`}`

`total 25 tokens`

Ex:- `printf ("%d\n", 8);`

`8 tokens`

- * LEX, which is an automated tool it is used to generate lexical Analyser.



- * Lexical Analyser Reads the characters from source program, line by line.
- * Lexical Analyser use the symbol table to produce tokens into parser.
- * upon Receiving a "get next token" command from the parser, the lexical analyser reads I/p characters until it can identify the next token.

Functions of lexical Analyser:

- 1, it produces stream of tokens
- 2, it eliminates comments, and white space which are in the form of blank, tab and new line characters from the source program.
- 3, it generates a symbol table which stores the information about identifiers, constants encountered in the I/p.
- 4, it keep track of line numbers.

5, It Reports the errors encountered while generating
the tokens.

(23)

Lexical Errors: These type of errors can be detected during lexical analysis phase. Typical lexical phase errors are

1, Exceeding length of identifiers (or) numeric constants

2, Appearance of illegal characters (Ex: `Print("HelloWorld");$`)

3, unmatched string. (Ex: comment /*, and if the comment is present, but beginning of comment is not)

* if the string `#` is encountered in 'C' program i.e,

`fi(a == f(n))...`

Lexical Analyzer can not tell whether `fi` is a misspelling of the keyword `if` (or) an undeclared function identifier.

because, None of the patterns for tokens matches a prefix of the remaining in input. So, This type of errors are solved by Error Recovery action (or) Error Recovery Solution.

Panic mode Recovery:

* Deleting successive character — `Pointf f`

* Deleting missing character — `Pointf → Pointf`

* Replacing incorrect characters by correct characters — `Pointf → Printf`

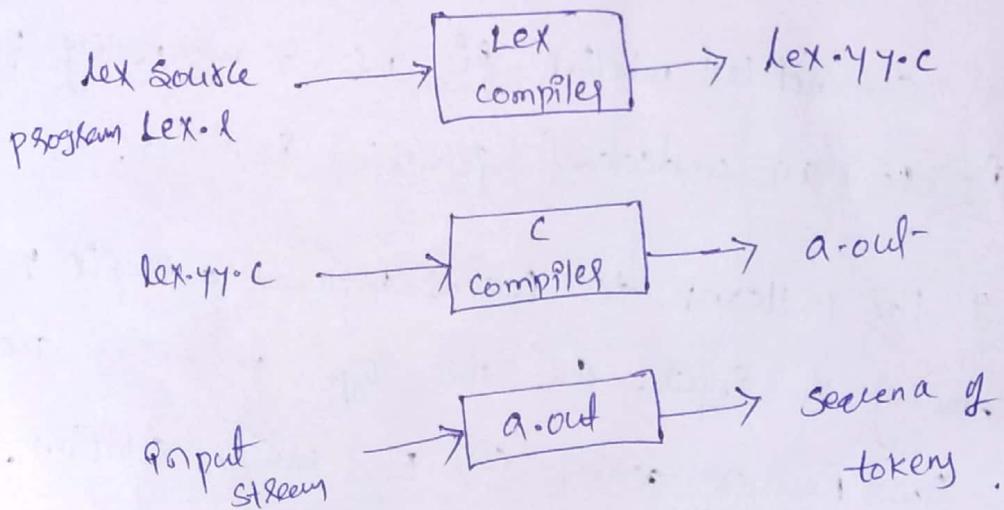
* Transposing Two adjacent characters — `Printf ↔ Pointf`.

Attribute FD Tokens:

The tokens and associated attribute values for the statement
 $E = M * C ** 2.$

<id, pointer to symbol-table entry for E>
 <assign-op>
 <id, pointer to symbol-table entry FD M>
 <mult-op>
 <~~num~~ id, pointer to symbol-table entry for C>
 <exp-op>
 <num, integer value 2>.

Lexical Analyzer Generation:



Input Buffering:

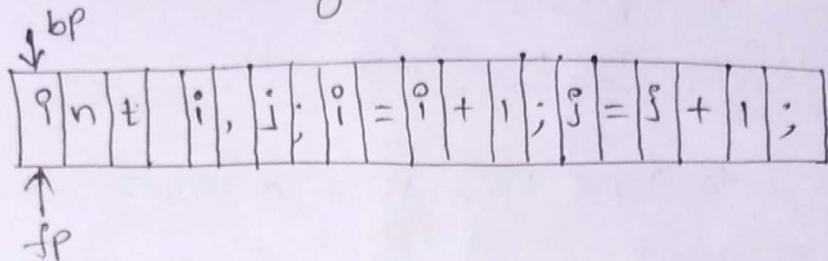
18

- * The lexical analyzer scans the I/p string from left to right one character at a time. It uses two pointers.

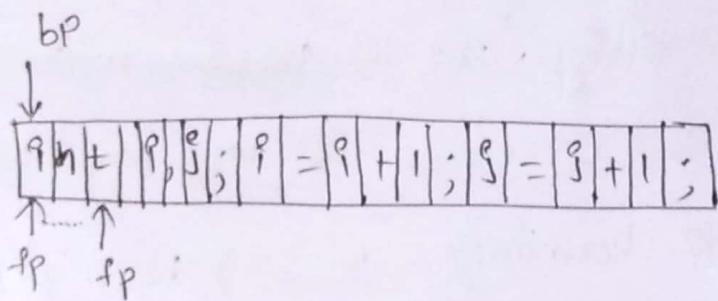
- ## 1. Lexeme begin -ptx

2. T8ward-Pt8.

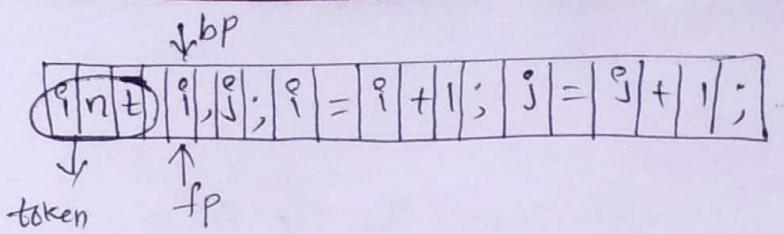
- * These are the pointers to keep track of the position of the IIP scanned. Initially both the pointers point to the first character of the IIP string as shown below.



- * The forward_ptr moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme. In above as soon as forward_ptr (fp) encounters a blank space the lexeme "int" is identified.



- * the fp will be moved a head at white space. when fp encounters white space, it ignore and moves a head. Then both the begin_Pt8(bp) and forward_Pt8(fp) are set next token f.



* The I/O character is thus read from secondary storage. but Reading in this way from secondary storage is costly. (System calls). i.e, if the program size is large & one pt needs more system calls. to read the tokens from memory we can overcome this problem with the help of buffering technique

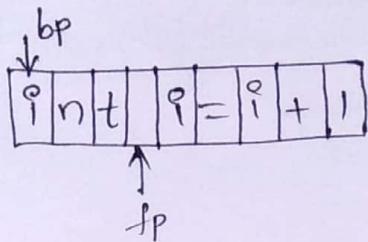
buffering: A block of characters to be read in to the buffer using only one system call.

→ A block of Data first Read Into a Buffer, and then scanned by lexical analyzer. There are Two methods used in this context

1, one buffer scheme

2, Two buffer scheme.

one buffer scheme: In this, only one buffer is used to store the I/O string. But the problem with this scheme is that if lexeme is very long then it crosses the buffer boundary, to scan Rest of the lexeme the buffer has to be refilled, that makes overwriting the first part of lexeme.



Two Buffer scheme: To overcome the problem of one buffer scheme, in this method two buffers are used to store the I/p string. The first buffer and second buffer are scanned alternately. When end of current buffer is reached the other buffer is filled. To identify the boundary of first buffer end of file (eof) character should be placed at the end of ~~next~~ present at the ~~end~~ first buffer. Similarly in second buffer also.



* This eof character introduced at the ~~end~~ ^{fp} end is called sentinel which is used to identify the end of buffer.

Code:

```

if (fp == eof(buff1)) /* encounters end of first buffer */
{
    /* Refill buffer2 */
    fp++
}

else if (fp == eof(buff2)) /* encounters end of second buffer */
{
    /* Refill buffer1 */
    fp++
}

else if (fp == eof(iInput))
    return; /* terminate scanning */
}

```

else
 $fp++;$
 /* Still Remaining I/p has
 to be scanned */

(8)

(4)

switch (*forward++) {

case eof:

If (forward is at end of first buffer)

{

reload second buffer;

forward = beginning of second buffer;

}

else if (forward is at end of second buffer)

{

reload first buffer;

forward = beginning of first buffer;

}

else /* eof within a buffer marks the end of input */
break; terminate lexical analysis;

cases for the other characters

}

Specification of Tokens: (Regulated Grammars and Regulated Expressions) (21)

(String, language and operations on language)

String:- (Sequence of characters)

* A string over an alphabet set Σ is a finite sequence of symbols from Σ , denoted by w :

Example: If $\Sigma = \{0, 1\}$ then
 $0, 1, 01, 10, 01, \dots$
 $a, b, c, abc, bca, cab, \dots$ are strings.

Operations on Strings:

The length of string s , usually written $|s|$, is the number of occurrences of symbols in s .

Ex: RamaKrishna
 $|w| = 10$.

The empty string can be denoted by ' ϵ ' & ' λ '. (A string that consists of zero symbols).

Prefix: Prefix of string ' s ' is any string obtained by removing zero or more symbols from the end of ' s '. i.e., removing leading symbols of that string.

Ex: banana \rightarrow ban and e are prefixes of banana.

Suffix: Suffix of a string ' s ' is any string obtained by removing zero or more symbols from beginning of ' s '. i.e., removing trailing symbols of that string.

Ex: banana, \rightarrow nana, & are suffixes of banana.

Substring: A substring of s obtained by deleting any prefix and any suffix from s . Removing prefix and suffix. ②

Ex: banana \rightarrow nam

Subsequence: A subsequence of s is any string formed by deleting zero or more not necessarily consecutive positions of s .

Ex:- Banana \rightarrow ban is a subsequence of banana.

Concatenation: If x and y are strings, then concatenation of x and y , denoted by xy .

Ex: $x = 010, y = 1$
 $z = 0101$

Identity Element: The empty string is the identity under concatenation. That is, for any string s , $es = se = s$.

Language: A language is a set of strings of symbols from some alphabet. Languages like \emptyset , an empty set ϵ are also examples.

Ex: $Z = \{a\}$

$Z^* = \{\epsilon, a, aa, aaa, \dots\}$

\hookrightarrow language.

Ex: $\{0, 1\} = Z$

$Z^* = \{\epsilon, 0, 1, 00, 01, 10, 11, \dots\}$

Operations on languages:-

L_1, L_2 are 2 languages. The union denoted $L_1 \cup L_2$ is a language containing all strings from both languages.

$L_1 \cup L_2 = \{w / w \in L_1 \text{ or } w \in L_2\}$

(31)

Example: Let $L = \{0, 01, 10, 11\}$

Ex-1

$$L_2 = \{0, 01, 001\}$$

$$L \cup L_2 = \{0, 01, 10, 11, 001\}$$

Ex-2 $L_1 = \{0, 00, 000, \dots\}$ and $L_2 = \{0, 00, 000, \dots\}$

$$L \cup L_2 = \{0\}^*$$

Concatenation: $L_1 L_2$ denoted by $L_1 L_2$ is a language containing the strings from L_1 followed by strings from L_2 .

E.g. $L_1 L_2 = \{w_1 w_2 / w_1 \in L_1 \text{ and } w_2 \in L_2\}$

Ex-1 $L_1 = \{aa, ab, bb, b\}$

~~given~~ $L_2 = \{b\}$

$$L_1 L_2 = \{aab, abb, bab, bb\}$$

Kleene closure: Language L is denoted by L^* ; it is a language containing all the strings obtained by concatenating zero or more strings from L .

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

(+)

Positive closure: it is denoted as L^+ , is the same as the Kleene closure, but without the term L^0 . that is $\underline{\epsilon}$ will not be in L^+ unless it is in L itself.

L^4 : it is the set of all 4-letter strings.

Operations

Definition:

(7)

union

$$L_1 \cup L_2 = \{s | s \in L_1 \text{ or } s \in L_2\}$$

Concatenation

$$L_1 L_2 = \{sp | s \in L_1 \text{ and } s \in L_2\}$$

Kleene closure

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

positive closure

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

Regular Expression: — the language accepted by FA are easily describe by simple expressions call Regular expression. R.E useful for representing certain sets of strings in an algebraic way. i.e, we are able to describe identities by giving names to sets of letters and digits and using the language operators union, concatenation, and closure.

Let Σ be an alphabet. The Regular expression over Σ and the sets they denote are defined as follows.

- ϕ is a regular expression and denotes empty set
- ϵ is a regular expression and denotes the set $\{\epsilon\}$, i.e., the set containing empty string.
- * The union of two R.E R_1 and R_2 . Return by $R_1 + R_2$ is also in R.E.
- * The concatenation of two R.E R_1 and R_2 return by $R_1 R_2$ is also in R.E

* If R is an RE then R^* is also RE.

(33) ⑤

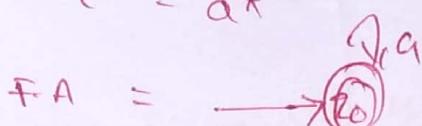
Regular Set:- Any set represented by RE is called as "Regular Set".

It is a set of strings from which \exists some FA which accepts the set it is represented by ' $\{ \cdot \}$ '.

Ex:

$$\text{Regular Set} = \{ \epsilon, a, aa, aaa, \dots \}$$

$$R.E = a^*$$



R.E

a

ab

ab

at | *pd. bbaabaa*

(ab)^{*}

Regular Set

{101}

{abba}

{01,10}

{ε,ab}

{abb,a,b,ba}

{ε,0,00,000...}

{1,11,111...}

{ε,(a,aa,ab,b,a,b,...)}

Regular Set

{a}

{a,b}

{a,b}

{ε,a,aa,aaa,...}

{a,b}^{*}

R.E

101

abba

01+10

ε+ab

abb+a+b+ba

0*

1+

(ab)^{*}

Problems on R.F:

- * all the strings of 0's and 1's - $\{0, 01, 10, 11, 01, 00, \dots\} = (0+1)^*$
- * ending with 00 - $(0+1)^* 00$
- * beginning with 0 ending with 1 - $0(0+1)^* 1$
- * set of all strings having even no of 1's - $\{0, 11, 1111, \dots\} = (11)^*$
- * odd no of 1's - $1(11)^*$
- * strings of 0's and 1's with at least one consecutive 0's
 - $= (1+0)^* 00 (1+0)^*$
 - set of all strings with 1100 as substring
 $(0+1)^* 1100 (0+1)^*$
 - set of all strings with one (8) more 0's followed by 1 / 0
 $0^8 1$

Algebraic laws for R.F:

- these laws are used to identify the base in which the expression can be modified or apply the following associative law

1) union operation on R.F are compatible.

$$\text{i.e., } R \cup S = S \cup R$$

2) union is associative

$$(R \cup S) \cup T = R \cup (S \cup T)$$

3) concatenation operation are associative

$$(R \cdot S) \cdot T = R \cdot (S \cdot T)$$

4, Concatenation is Right distributive over union

(35) ⑦

$$(R+S)t = Rt+St$$

5, Concatenation is left distributive over union

$$t(R+S) = tR+tS.$$

$$\emptyset^* = \epsilon.$$

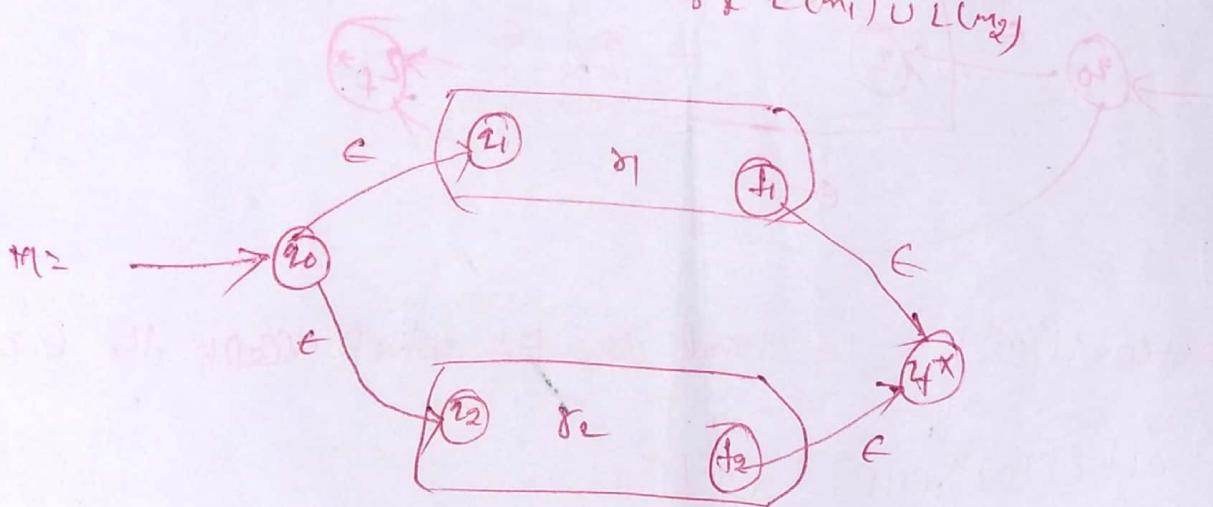
Operation on R.E: (Union)

Let δ_1, δ_2 be two R.E's with less than 3 operations that of F.A's
making after η

$$m_1 = \{Q_1, I_1, \delta_1, z_1, f_1\} \text{ and}$$

$$m_2 = \{Q_2, I_2, \delta_2, z_2, f_2\}$$

To show that R formed using union operation of $\delta_1 + \delta_2$ by
automate on M which accept the language $L(m_1) \cup L(m_2)$



Concatenation: Let δ_1, δ_2 be two R.E's with less than 3 operations that
of F.A's - making after η

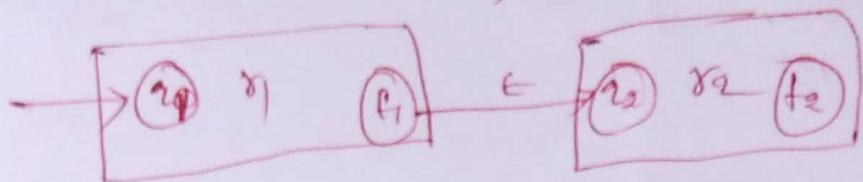
$$m_1 = \{Q_1, I_1, \delta_1, z_1, f_1\}$$

$$m_2 = \{Q_2, I_2, \delta_2, z_2, f_2\}$$

To show that γ is formed using concatenation operator as (8)

$\gamma = \gamma_1 \cdot \gamma_2$ as automata on M . which accept the language

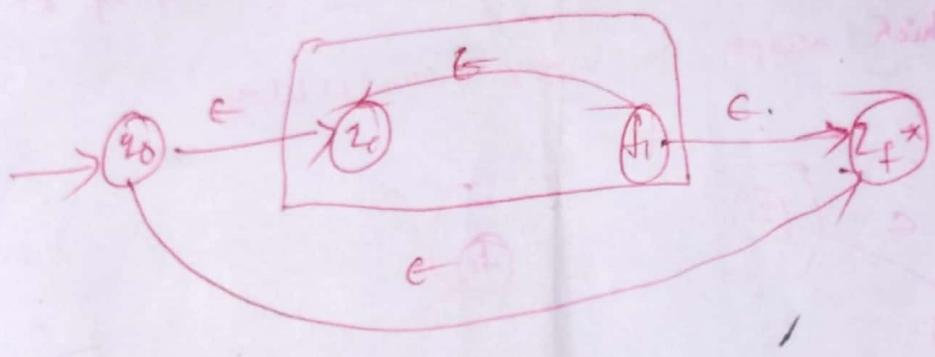
$$L(M_1) \cup L(M_2)$$



Closure operation: let γ_1 be the R.E but less than γ operations
out of the finite automate m_1 giving 13

$$m_1 = \{q_0, q_1, s_1, z_1, f_1\}$$

To show that γ is formed using closure operators as $\gamma_1 = \gamma_1 \cdot \gamma$ as
automate m which accept the language $L(m, *)$



Ex: 1, $10 + (0+11)^*$ — construct the FA which accept the R.E

$$\text{2, } \gamma = 01[(10)^* + 11]^* + 0J^*$$

$$3, (ab+ba)^* abb$$

& construct NFA with ϵ to DFA.

Components of Regular Expressions

B7 @

$x \rightarrow$ the character

$\cdot \rightarrow$ any character usually accept a new one.

$R^* \rightarrow$ 0 or more occurrence

$R^+ \rightarrow$ one & more occurrence

$R_1 R_2 \rightarrow R_1$ followed by R_2

R_{2k}

$R^k \rightarrow$ the set of k occurrence of string.

Regulated definition: For notation convenience, we may wish to give names to certain Regular Expression and use those names, p_j Subsequent expression, as the names were symbols may.

If Z is an alphabet of basic symbol then a regulated definition is a sequence of the form.

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

:

$$d_n \rightarrow r_n$$

where

Each d_i is a new symbol, not in Z and not the same as any other of the d_j 's

Each r_i is a R.E over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{n-1}\}$

Ex: Σ

Identifiers are strings of letters, digits, and underscores.
Here is a regulated definition for the language of Σ Identifiers.

We shall conventionally use Postfix for the symbols defined in R.D. ⑩

letter $\rightarrow [A-Z] \dots [Z/a-zA-Z]$

digit $\rightarrow 0/1 \dots 9$

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

Ex: Unsigned numbers (integer or floating point) are stored such as
5280, 0.01234, 6.336E4 & 1.89E-4

digit $\rightarrow 0/1 \dots 9$

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

optional fraction $\rightarrow \cdot \text{digit} | e$

optional exponent $\rightarrow (E (+-|e) \text{digit}) | e$

number $\rightarrow \text{digit} \text{ optional fraction} \text{ optional exponent}$

using NFA, we can rewrite the regular definition

letter $\rightarrow [A-Za-zA-Z]$

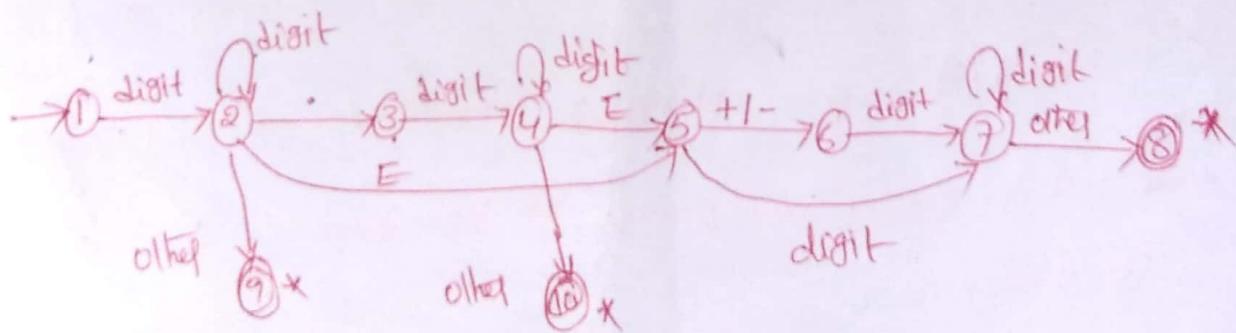
digit $\rightarrow [0-9]$

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

digit $\rightarrow \text{digit}^+$

number $\rightarrow \text{digit} (\cdot \text{digit})? (E (+-|e) \text{digit})?$

Transition Diagram:



Recognition of Tokens:

(38)

- * For a programming language there are various types of tokens such as Identifiers, keywords, constants and operators and so on. The token is usually represented by a token type and token value.

Token Type	Token Value.
------------	--------------

- * The token type tells us the category of token and token value gives us the information regarding token. The token value is also called token attribute. During lexical analysis process the symbol table is maintained. The token value can be a pointer to symbol table in case of identifiers and constants. The lexical analyzer reads the input program and generates a symbol table for tokens.

Ex! We will consider some encoding of tokens as follows.

Token	code	value
if	1	-
else	2	-
while	3	-
for	4	-
Identifier	5	ptr to symbol table
constant	6	ptr to symbol table
<	7	1
<=	7	2
>	7	3
>=	7	3
<>	7	5

e	8	1
)	8	2
+	9	1
-	9	2
=	10	-

Our lexical Analyser will generate following token stream.

If ($a < 10$) $p = p + 2;$ else $p = p - 2;$	1, (8,1), (5,100), (7,11), (6,105), (8,2), (5,107), 10, (5,107), (9,1), (6,110), 2, (5,107), (5,107), (9,2), (6,110).
---	---

The corresponding symbol table of identifiers and constants will be,

Location counter	Type	Value
100	Identifier	a
:	:	:
105	Constant	10
:	:	:
107	Identifier	i
:	:	:
110	Constant	2

Transition Diagram for Recognition of token:

Lexical Analyser to scan the input string, and search for where the lexeme end.

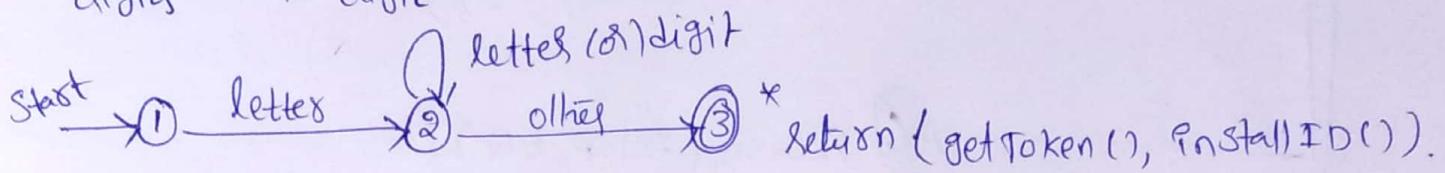
- * Recognition of Identifiers
- * Recognition of Delimiter
- * " " Keyword
- * " " Operator
- * a Number.

Recognition of Identifiers:

letter \rightarrow [A-Z a-z]

digit \rightarrow [0-9]

digits \rightarrow digit⁺



Recognition of Delimiters:

WS \rightarrow (blank | tab | newline)⁺

WS \rightarrow Delimiter⁺



Recognition of Keyword:

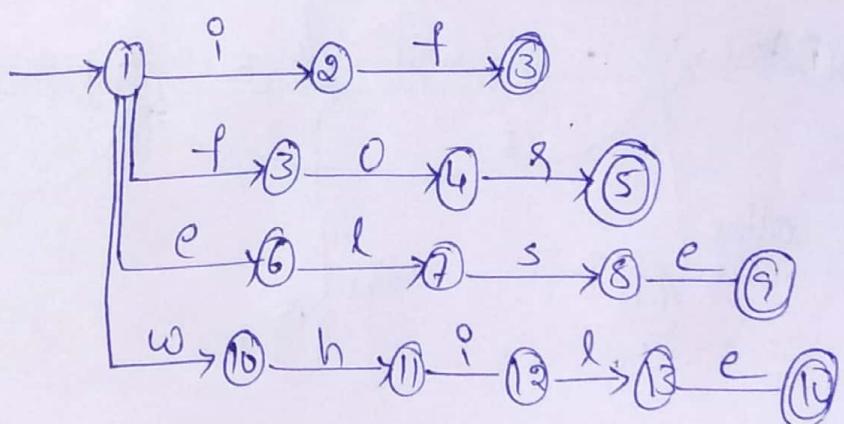
Keywords are Reserved words, These values can not change during compilation.

if \rightarrow if

for \rightarrow for

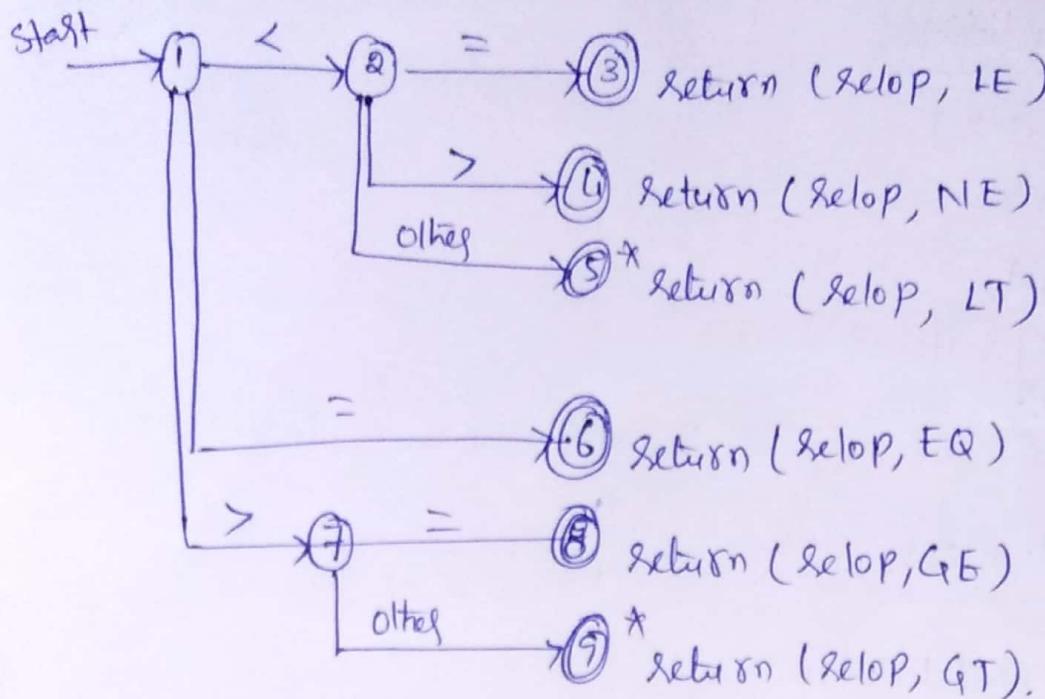
else \rightarrow else

while \rightarrow while.



Recognition of operators:

SetOp \rightarrow > | >= | < | <= | = | < >

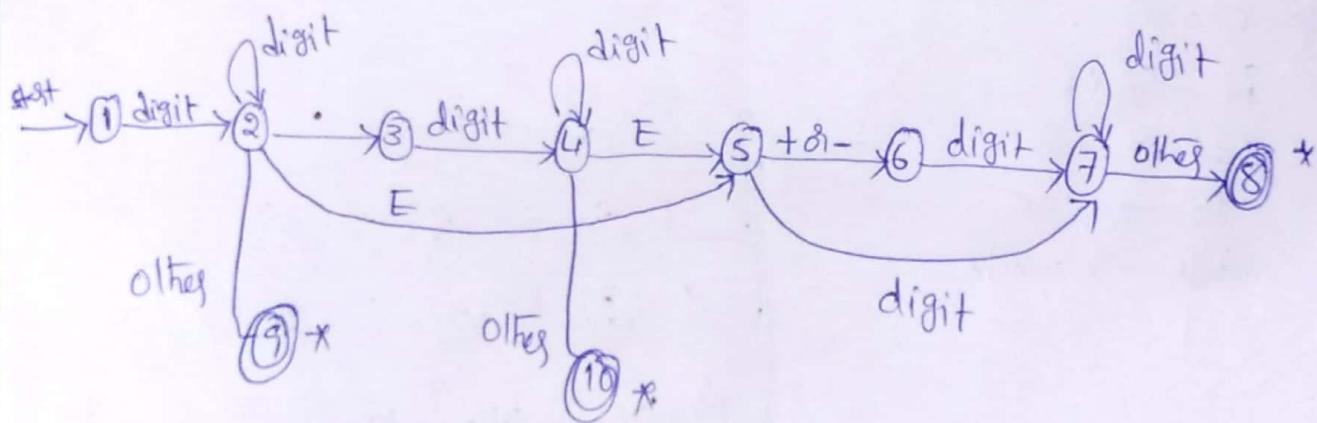


Recognition of numbers:

digit \rightarrow [0-9]

digit \rightarrow digit $^+$

Number \rightarrow digits (digits)? (E [+ -]? digits)?



Transition Diagram:-

(43) (11)

An Automation is self-operating machine, a system which obtains inputs, transmit and uses information to perform the function without direct human participation.

* An Automation is a model of computation consisting of set of states and discrete I/P's and O/P's.

* A start state in I/P alphabet and transition function that maps I/P symbol and current state to the next state.

* computation begins in the start state with an I/P string. It changes to new state depending on the transition function.

F.A can be represented a 5-tuple

$$M = (Q, \Sigma, S, q_0, F)$$

A70

A71

A72

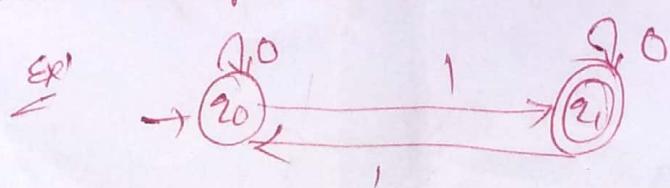
Representation of F.A: It can be done in 3 ways

1, Transition diagram

2, Transition table

3, Transition function

T-D: F.A is represented in the form of a finite directed labelled graph called as transition graph. In which each vertex represents as state and the directed edges indicate the transition of state.



T.T. It is a conventional tabular representation of T.D. (R)

It accepts 2 arguments and returns a value.

Rows - States

Column - I/p symbol.

	0	1
q0	q0	q1
q1	q1	q0

Transition Function:

This is a function which describes the change

of state during the transition. It maps a transition from one state to another state on some I/p symbol.

representation of DFA

$$S = Q \times \Sigma \rightarrow Q$$

Types of F.A:

1. DFA

2. NFA.

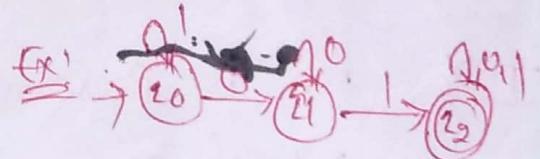
DFA:

have, for each state, and for each symbol of its input alphabet exactly one edge with that symbol leaving that state.

* It doesn't accept ϵ^* or P/P

* Atmost one transition

mathematical function of DFA



It is in D.F.A.

T.T

$$M_D = (Q, \Sigma, S, q_0, F)$$

$$S = Q \times \Sigma \rightarrow Q$$

N.F.A! have no restriction on the labels of their edges.

A symbol can label several edges out of the same state, and ϵ , the empty string is possible label.

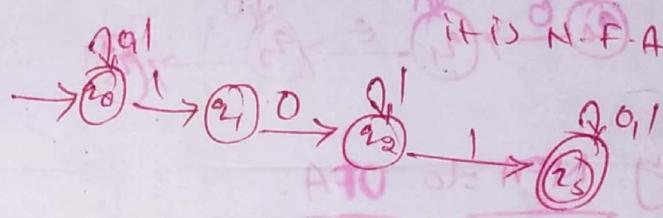
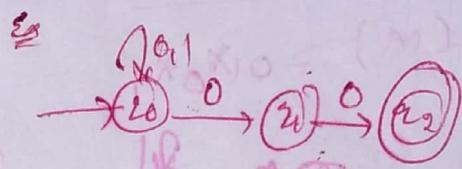
* at least one transition

it accepts Σ^* as p/p symbol

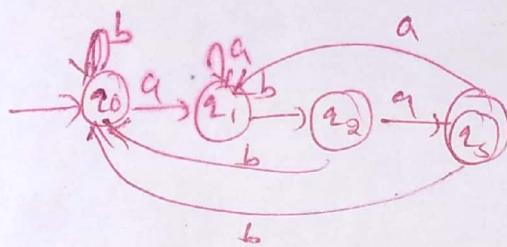
Mathematical Def of NFA

$$M_N = \{Q, I, S, \Sigma_0, F\}$$

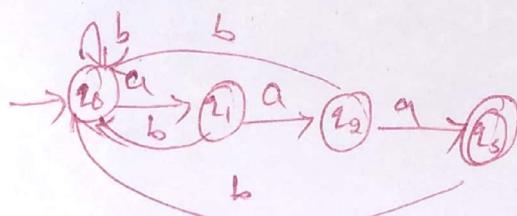
$$S = Q \times I \rightarrow 2^Q \text{ states.}$$



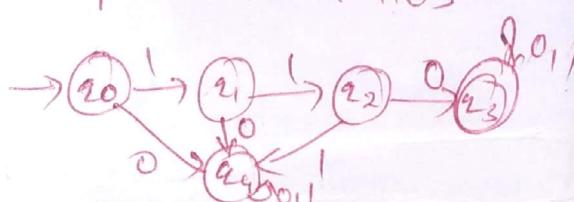
Constructing DFA accepting all string over $\{a, b\}$ ending in $a^k b^l$



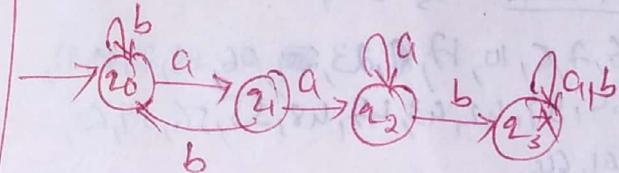
ending with $a^k b^l$



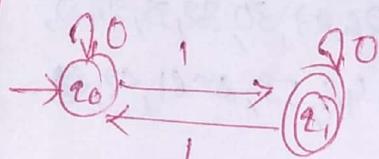
$L = \{w | w \text{ starts with } 110\}$



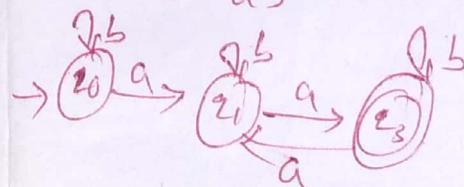
substring



odd no of 1's



even no a's

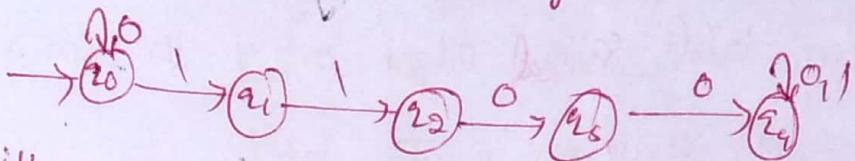


Example
aabaae

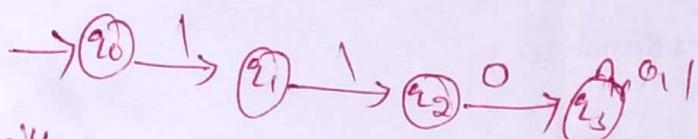
NFA:

containing 1100 as substring

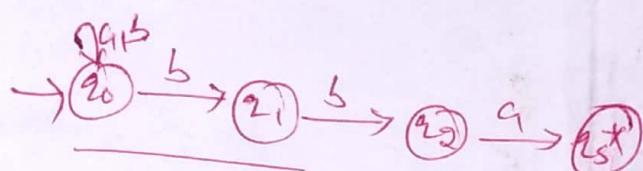
(19)



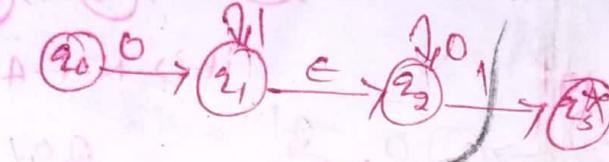
starts with 110



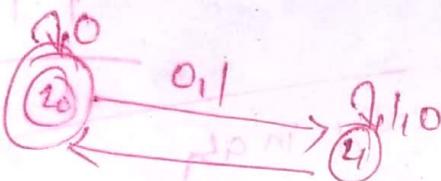
ending with bba



$$L(M) = 0^* 1^* 0^* 1$$



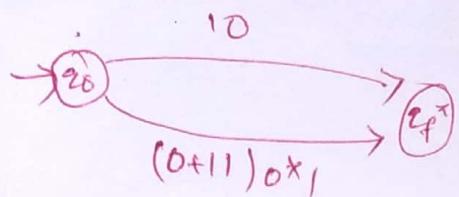
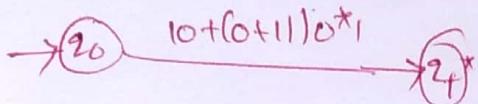
Converting NFA to DFA:



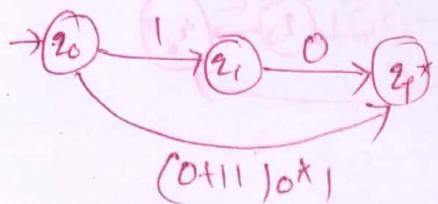
Construct the FA which accept the R.F

$$10 + (0+11)0^*$$

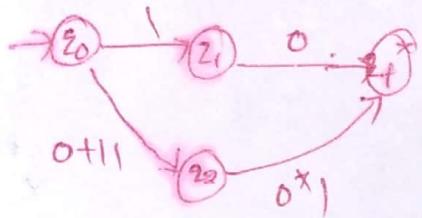
$$\delta = 10 + (0+11)0^*$$



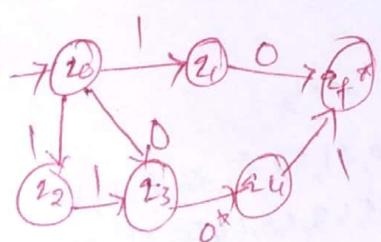
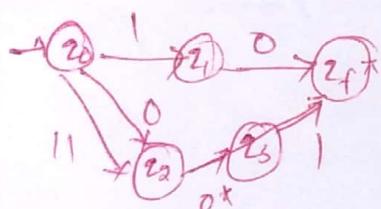
on eliminating union



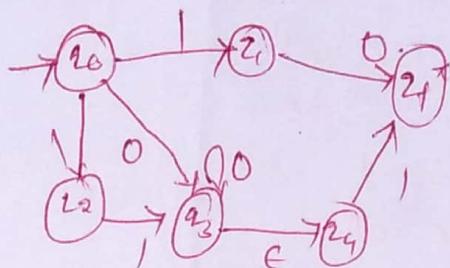
on eliminating concatenation



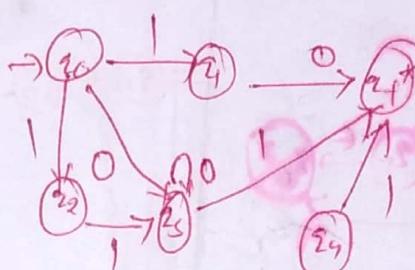
on eliminating concatenation



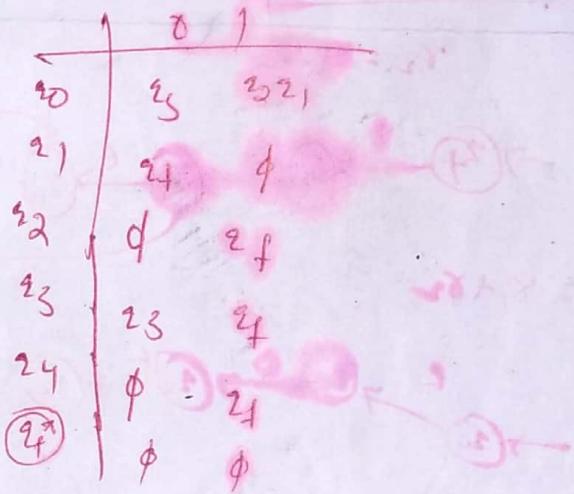
elimination



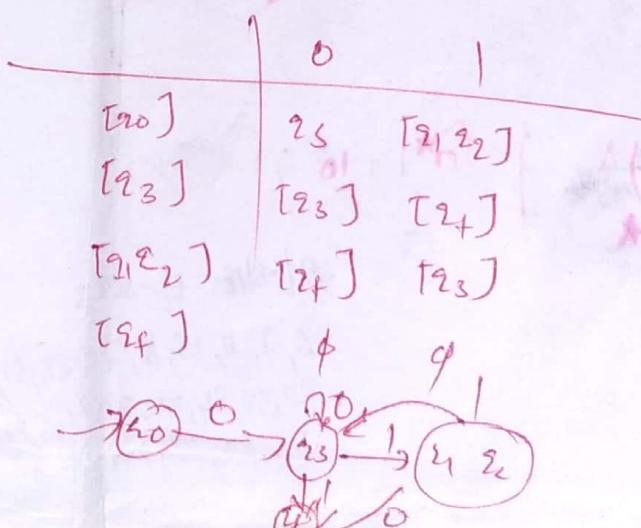
elimination E



construction of DFA



$$Q = \{[q_0], [q_3], [q_1, q_2], [q_f]\}$$



Construct NFA with E-moves for the following RE

$$01^* + 10^*$$

$$\gamma = 01^* + 10^*$$

$\delta_1 \quad \delta_2$

$$\delta_3 = 01^*$$

$\delta_3 \delta_4$

$$\delta_3 = 0$$

$$\delta_4 = 1^*$$

$$\delta_2 = 10^*$$

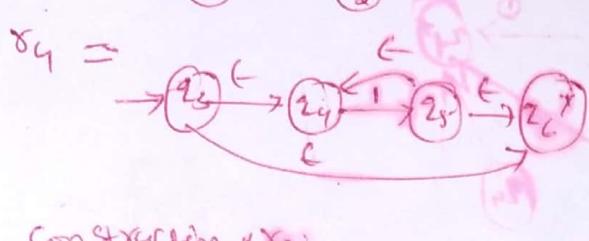
$\delta_5 \delta_6$

$$\delta_5 = 1$$

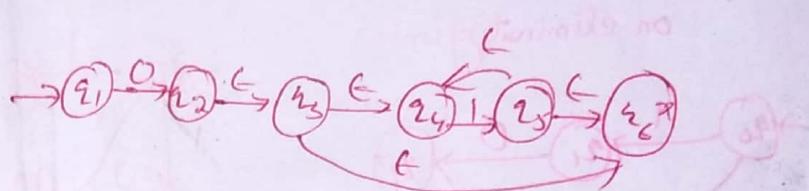
$$\delta_6 = 0^*$$

Construction of δ_1 :

$$\delta_3 = \rightarrow (2_1) \xrightarrow{0} (2_2)$$

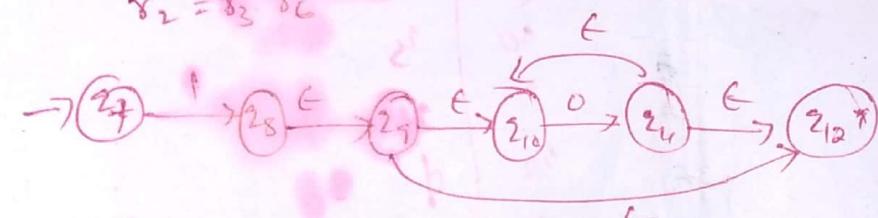


$$\delta_1 = \delta_3 \cdot \delta_4$$

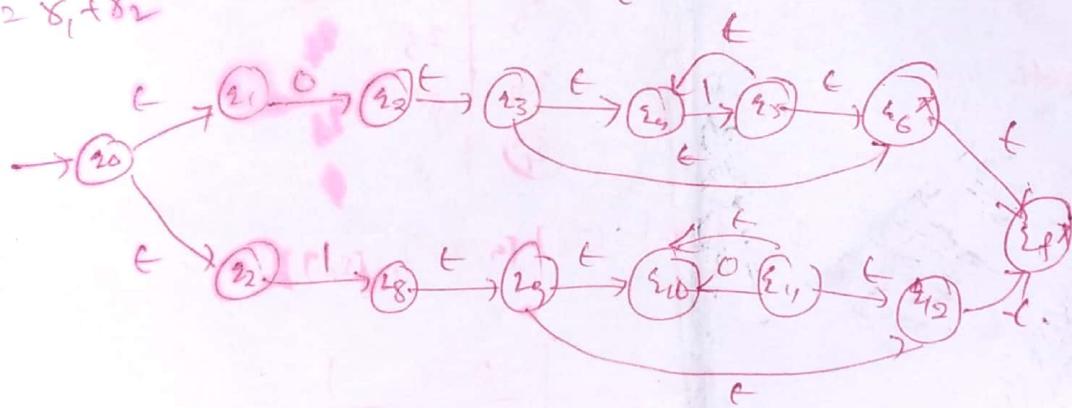


Construction of δ_2 :

$$\delta_2 = \delta_3 \delta_6$$



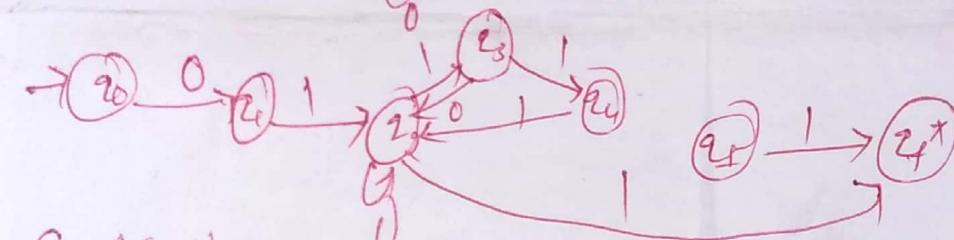
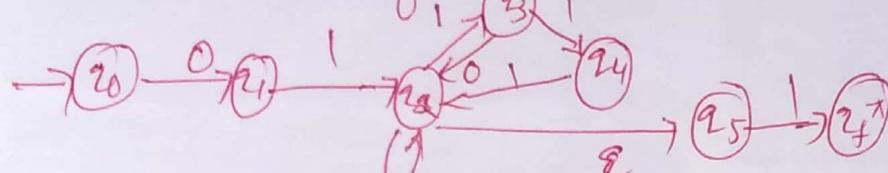
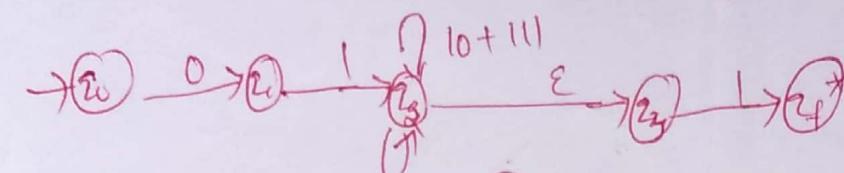
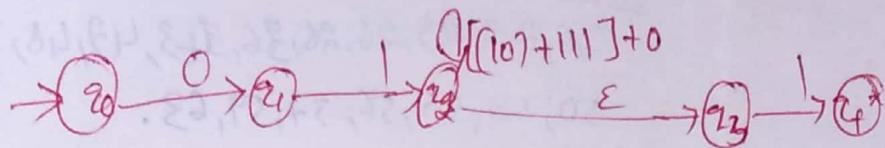
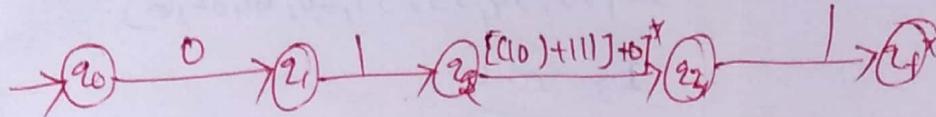
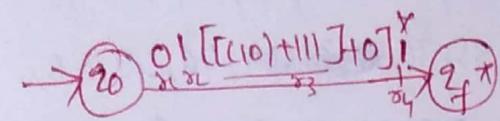
$$= \delta_1 + \delta_2$$



(49)

Construct FA equivalence to the following R.F.

$$R = 01[(10) + 111] + 0^*$$

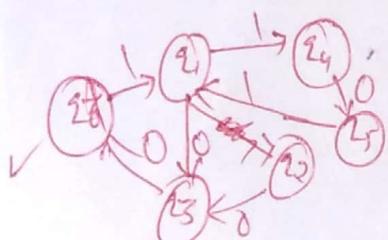
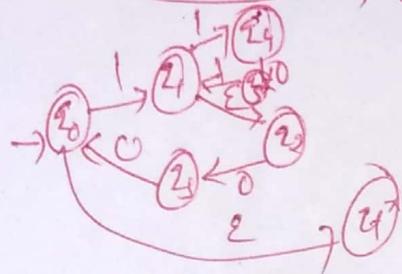
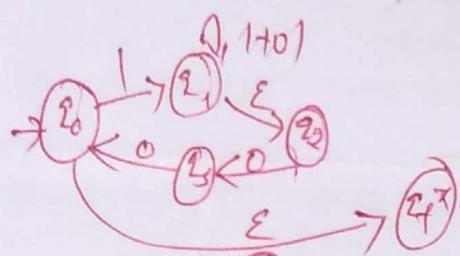
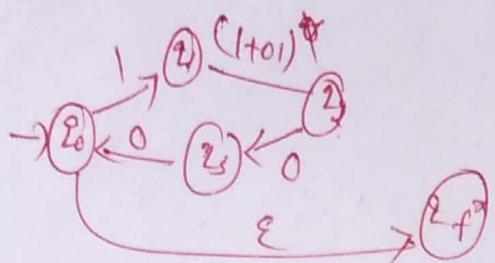
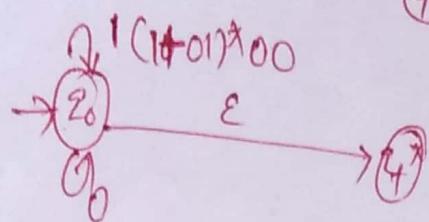
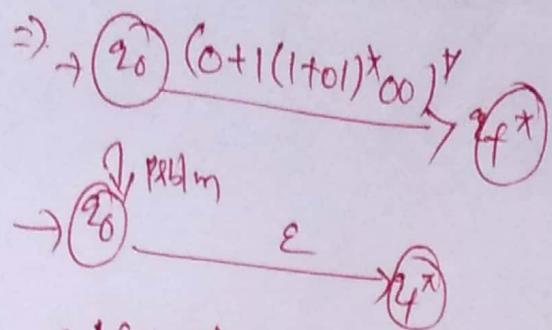


Construction of DFA

	0	1
q0	q1	∅
q1	∅	q2
q2	q2	{q3, q4}
q3	q2	q4
q4	q2	q2
q5	∅	q4*
q4*	∅	∅

	0	1
q0	{q3}	q
q1	∅	{q2, q3}
q2	{q2, q3}	{q3, q4}
q3	{q2, q3}	{q2}
q4*	∅	{q2, q3}

$$(0 + 1((1+01)^* 00)^*)^*$$



(4_f^*)