

$$\text{Test for } s_4 \text{ is } - [y_{11}, y_{12}] = s_4 \times w \\ = [0 \ 0 \ 1 \ 0] \begin{bmatrix} 1 & 0 \\ 2 & 0 \\ 0 & 2 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 2 \end{bmatrix}$$

apply the activation function

$$\therefore [y_1 \ y_2] = [0 \ 1]$$

here the prediction is ~~not~~ same as t_4 .
So it recognise the s_4 .

o The art of Language Design :-

- Today there are thousands of HL programming languages, and new ones continue to emerge.
- Humans being use Assembly language only for special purpose applications.
- Why there are so many languages?
- There are several possible answers to it →

① Evolution :- The late 1960's & early 1970's saw a revolution in "structured programming" [extensively used → goto based] control flow languages like Fortran, Cobol, & Basic.

In late 1980's the nested block structure of language like

↓
Algol, Pascal, & Ada began to give way to the object oriented structures of C++, C#, Java, etc.

Gave a way
(To use while loops & switch cases.)

② special purpose :- Lisp → used for manipulating symbolic data and complex data structures.

③ Personal Preference :-

- Preference of programming is a matter of taste.
- Some people love the terseness of C, some hate it.
- So people uses recursion to solve problems so use iterations.

- Some people likes to work with pointer, whereas some people hate it and prefers lists & pointer.
 - ∴ The strength & variety of personal preference make it impossible that anyone will ever develop a universally acceptable programming language.
 - ∴ Of course some programming L. are more successful than others. like. [Python is more successful than BASIC]
- Expressive Power :-
- We can hear the arguments that one PL is more "powerful" than other,
LWLT is the meaning of, powerful?
∴ Powerful means → Language features.
- ∴ Language features clearly have huge impact on the programmer's ability to write clear, concise, and maintainable code, especially for large systems.
- Expressive power → abstraction facility
[use C/C++ & python as prog. example].
- Ease of use for the Novices -
- In Elementary school → Logo PL is used → even 5 yrs old can understand it. (further is the simplicity of logo PL).
 - C/C++ was used for object oriented prog. now Java has emerged to play the same role and Java is more like English L.

- Ease of Implementation :-

→ The PL should be easily implementable on tiny machines with limited resources.

For Ex:- Basic prog. can run on almost any system with any configuration.

However

If you want to run MatLab 2020 R2a

→ its min requirement is $\frac{8 \text{ GB RAM}}{16 \text{ GB RAM}}$

So, The language can be such that, it can be implemented on any machine.

- Standardization :- of both the Language and a broad set of Libraries.

- open source :- Most programming languages today have at least one open-source compiler or interpreter.

→ But some language — C in particular — are much more closely associated than others with ~~freely~~ freely distributed, peer-reviewed, community-supported computing.

Linux → written in 'C'.

- Excellent Compilers :-

Economics, Patronage, and Inertia

PL/I → depends or owes its life to IBM.

Cobol/Ada → U.S. Department of defence.

L/A Time

- ① What is the difference between ML & AL?
- ② In what way(s) are

o Von Neumann Architecture :- Is based on the stored-program computer concept, where instruction data and program data are stored in the same memory.

o Simula 67 :-
→ Simula is two simulation programming language.
Simula I & Simula 67. → 1960's

Simula 67
↓
introduced object, classes,
inheritance, subclasses, etc.

Norwegian
Computing Center
↓
Oslo

Translators - one language program → Another Language Program.
→ Error detection
→ Rule violation.

- (i) Compiler
(ii) Interpreter
(iii) Assembler

- :- Question / Answer :-

Q) What is the difference between machine & and Assembly L?

Ans:-

| ML | AL |
|----------------------------------------------------------------------------|-------------------------------------------------------------------------------|
| ① ML is only understand by humans. | ① AL is only understand by humans. |
| ② In ML data is represented using binary format %. | ② In AL data is represented using mnemonics such as. mov, Add, Sub, End, etc. |
| ③ ML is very difficult to understand by Humans beings. | ③ its easy to understand as compare to ML. |
| ④ Modification & Error fixing is nearly <u>impossible</u> to can't be done | ④ Modification & Error fixing can be done in AL. |
| ⑤ ML is very difficult to memorize, so its very difficult to learn ML. | ⑤ Easy to memorize & learn as compare to ML. |
| ⑥ Execution is fast in ML as all are in binary. | ⑥ Execution is slow. |
| ⑦ There is no need of translator for Execution. | ⑦ Assembler is used for conversion of AL \rightarrow ML. |
| ⑧ ML is hardware dependent | ⑧ AL is machine dependent & Not portable. |

- ② In what ways is the HLL an improvement on AL? In what circumstances does it still make sense to program in assembler?
- ③ Why are there so many programming languages?
- ④ What makes a programming L. successful?
- ⑤ What distinguishes declarative languages from imperative languages
- ⑥ Who spearheaded the development of Ada?
- ⑦ What is generally considered the first high-level prog.
L?
↳ Fortran → 1958/1960
- ⑧ What was the first functional L?
↳ LISP → late 1950's for IBM
700/7000 series of scientific computers by John McCarthy.
at MIT.

define circlearc(r) $(3.1415 * \pi * r^2)$.

$$\frac{\text{circlearc}(\text{Argument})}{\uparrow \text{replace}} \\ \frac{(3.1415 * \pi * r^2)}{r=5}$$

$$\frac{3.1415 * 5 * 25}{}$$

| DATE - 1.
| TIME - 5.

define PI 3.1415

include <stdio.h>

define PI 3.1415

define circlearc(r) (PI * r * r)

int main()

{

float radius, area;

printf ("Enter the radius: ");

scanf ("%f", &radius);

area = circlearc(radius);

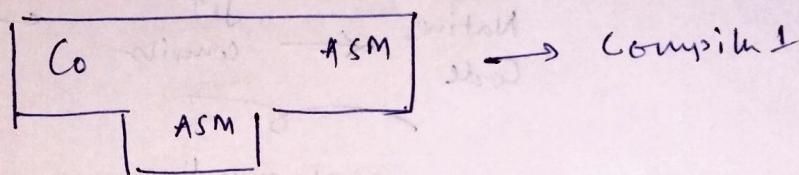
printf ("Area = %.2f", area);

return 0;

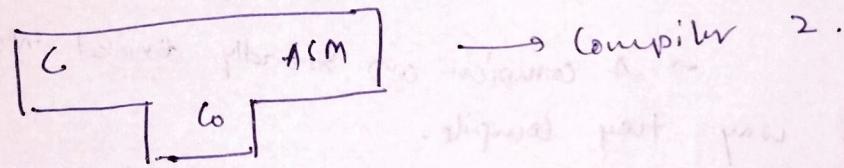
}

→ Bootstrapping is a process in which simple L. is used to translate more complex language. complicated prog. → b that can be used ↓ to translate even more complicated program ↓ and so on.

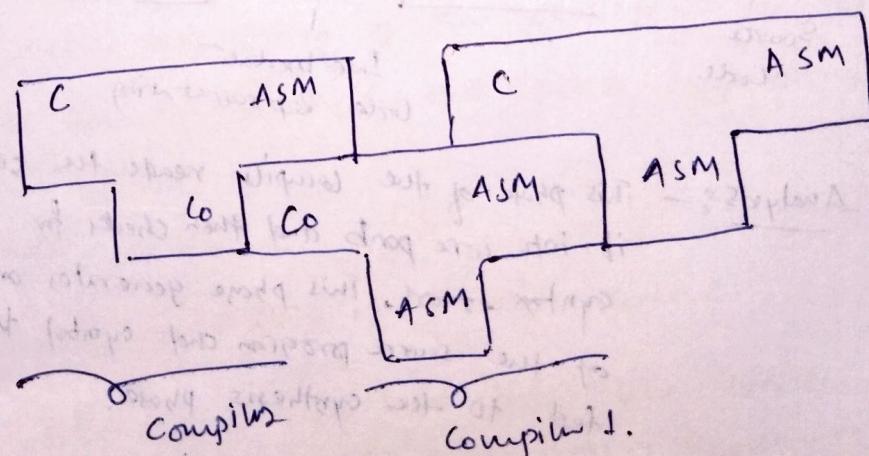
Step① → First we write a compiler for small C in AL.



Step② → Then using with small subset of C i.e. Co, for the Source Language C the compiler is written.



Step③ → Compile the second compiler using compiler 1, the compiler 2 is ~~compiler~~ compiled.

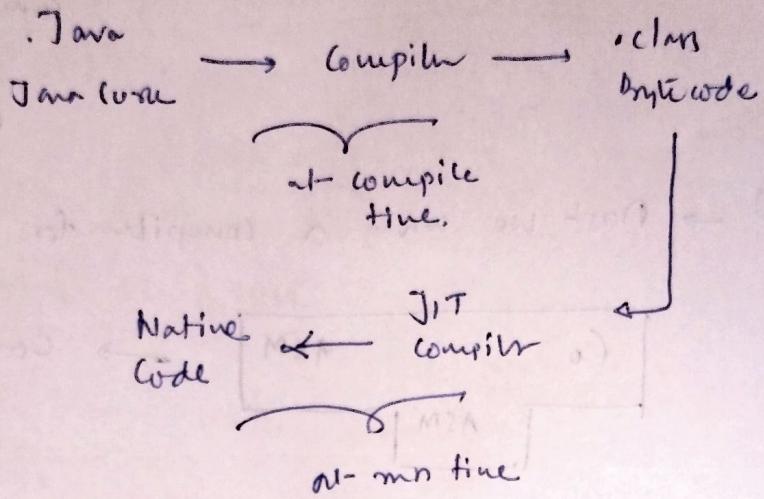


Step④ → We get a compiler written in ASM, which compiles C and generates code in ASM.

Just in Time Compiler

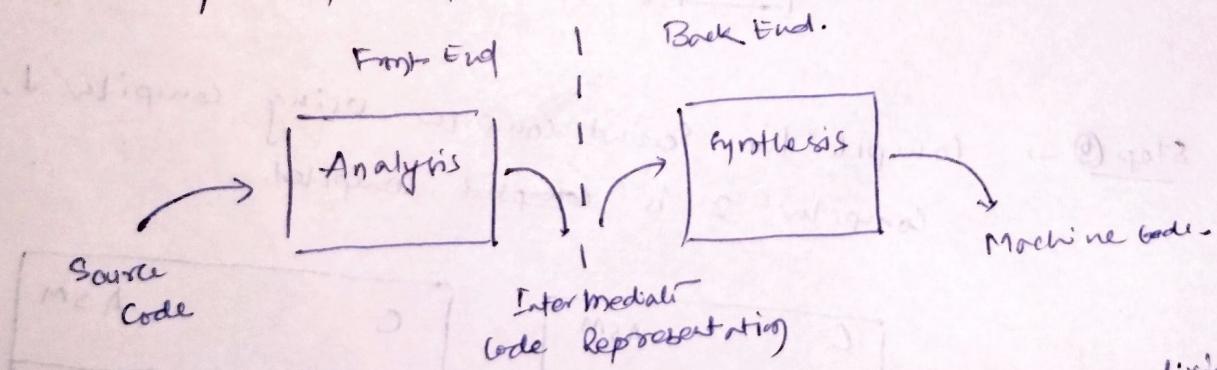
Just in Time compiler

At compile time



Structure of Compiler

→ A compiler can broadly divided into two phases, on the way they compile.



Analysis: → This phase of the compiler reads the source program, divides it into two parts and then checks for lexical, grammar, & syntax errors. This phase generates an intermediate representation of the source program and symbol table, which should be fed to the synthesis phase.

Synthesis: → It generates the target program with the help of intermediate code & symbol table.

7067

tokens.

Lexical Unit, Lexical Element, or Tokens are the smallest unit of program i.e. commands, variables, constants, operators, and many more symbols called operators & punctuators.

Phases of Compiler

Compiler operates in various phases, each phase transforms the source program from one representation to another. Every phase takes input from its previous stage and feeds its output to the next phase of the compiler.

- * There are six phases in compiler. Each phase help in converting the high-level language to the machine code. The phases of compiler are as follows:-

① Lexical Analysis :- (scanner)

- It's first phase of the compiler where compiler scans the code.
- This scanning process can be left to right.
- Character by character, and group these characters into tokens.

What happens:-

- ① Here the character stream from the source program is grouped into meaningful sequences by identifying tokens.
- ② It makes the entry of the corresponding tokens into the symbol table and passes the tokens to next phase.

Primary functions :-

- ① Identify lexical units in a source code.
- ② classify lexical units into classes like. constants, reserve words, and enter them in different tables.

- ⑥ Ignore and remove the comments.
- ⑦ Identify the tokens which are not part of the language.

Ex:- $x = y + 10$

Tokens

| | |
|----|---------------------|
| x | Identifier 1 |
| = | Assignment operator |
| y | Identifier 2 |
| + | Addition operator |
| 10 | Number (constant) |

- ⑧ Syntax Analysis :- Syntax Analysis is all about discovering structure in a code.

- It determines whether or not a text follows the expected/standard format or not.
- The main aim of this phase is to make sure that the source code was written by the programmer is correct or not.

- o Syntax Analysis is based on the specific programming language by constructing parse tree with the help of tokens.

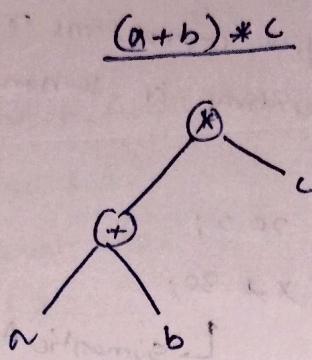
Main tasks :-

- ① Obtain tokens from the lexical Analyzer.
- ② Check if the expression is syntactically correct or not.

- ③ Report all syntax errors.

- ④ Construct a hierarchical structure which is known as parse tree.

Ex:- $(a+b)*c$



In parse tree :-

- ① Interior nodes :-
- ② Leaf node.

③ Semantic Analysis :-

- Semantic Analysis checks the semantic consistency of the code.
- It uses the syntax tree of the previous phase along with the symbol table to verify that the source code is semantically consistent.
- It also checks whether the code is conveying an appropriate meaning or not.
 - o It will check for → ① Type mismatch.
 - ② incompatible operands.
 - ③ a function called with improper arguments.
 - ④ undleared variables, etc.

o Main Task :-

- ① Helps you to store type information gathered and save it in symbol table or syntax tree.
- ② Allows you to perform type checking.
- ③ collects type information and checks for type compatibility.

- ⑭ checks if the source language permits the operand or not.
- ⑮ In the case of type mismatch, where there are no exact type corrections unless which satisfy the desired operation, a semantic error is shown.

Ex:-

float x = 20.2;

float y = x * 30;

↳ Semantic Analyzer will type convert

integer 30 to float 30.0.

④ Intermediate Code Generation :-

Once the semantic analysis phase is over, the compiler generates intermediate code for the target machine.

→ Intermediate code is between the high-level and machine level language.

→ This intermediate code needs to be generated in such a manner that makes it easy to translate it into the target machine code.

Main Tasks :-

- ① It should be generated from the semantic representation of the source program.
- ② Hold the value computed during the process of translation.
- ③ Helps you to translate the intermediate code into target machine code.
- ④ Allows you to maintain precedence ordering of the source language.
- ⑤ It holds the correct number of operands of the instruction.

Ex:-

total = count + rate * 5;

[Intermediate code with the help of address code method.

t₁ := int-to-float(5);

t₂ := rate * 5;

t₃ := count + t₂;

total = t₃;

④ Code optimization :-

- This phase removes unnecessary code lines and arrange the sequence of statements to speed up the execution of program without wasting resources.
- The main goal of this phase is to improve on the intermediate code to generate a code that runs faster and occupies less space.

Primary functions :-

- ① It helps you to establish a trade-off between execution and compilation speed.
- ② Improves the running time of the program.
- ③ Generates fast running code with same representation.
- ④ Remove unreachable code and unused variables.
- ⑤ Remove the statement which are not altered from a loop.

Ex:-

a = int-to-float(10);

b = c * a;

d = e * b;

f = d;

b = c * a;

f = e + b;

- ⑥ Code Generation :- It is a final phase of compiler.
- It gets input from the code optimization phase and produces the page code or object code as a result.
 - The objective of this phase is to allocate storage and generate relocatable machine code.
 - It also allocates memory locations for the variable.
 - The instructions in the intermediate code are converted into machine instructions.
 - This phase converts the optimized or intermediate code to the target language.
 - In this phase, all the memory locations and registers are also selected and allotted during this phase.

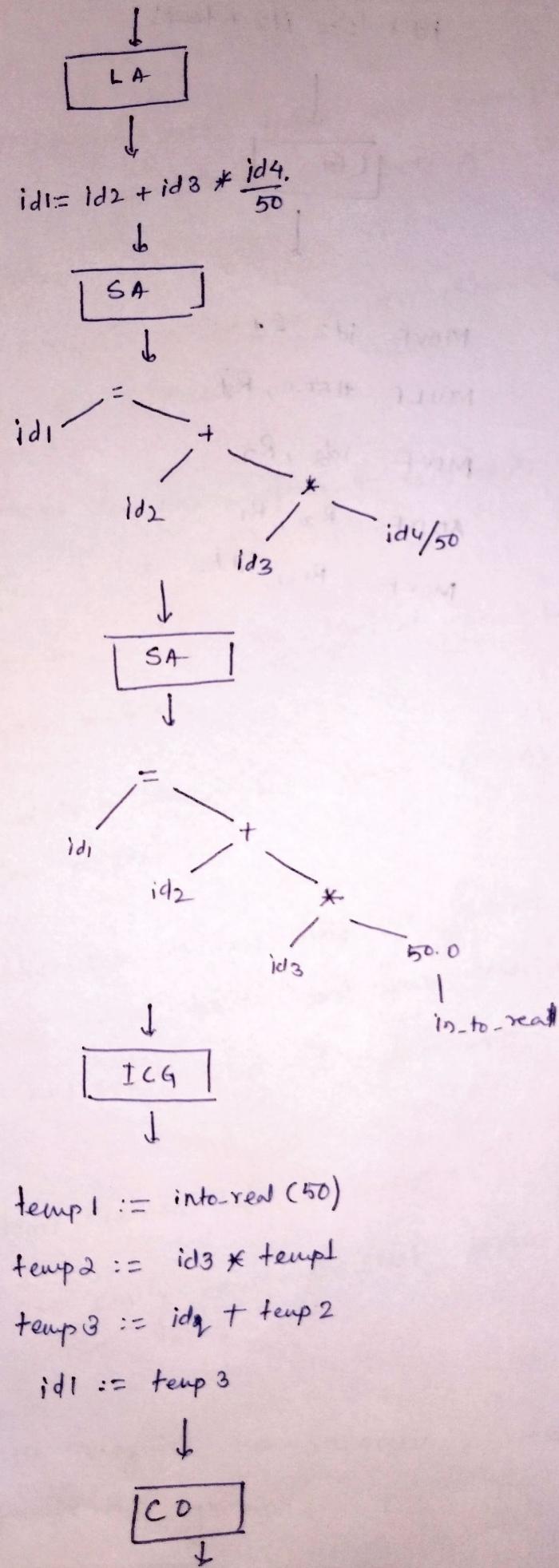
Ex :- $a = b * 60.0$

Would be possibly translated to registers.

MOV F B, R1
MUL F #60.0, R1
MOV F R1, A



Sum := old Sum + Rate * 50



temp1 := int-to-real(50)

temp2 := id3 * temp1

temp3 := id2 + temp2

id1 := temp3



CO

temp1 := id3 * 50.0

id1 := id2 + temp1

[CG]

MOVF id3, R1

MULF #50.0, R2

MOVF id2, R2

ADDF R2, R1

MOVF R1, id1.

[AD]

[out]

(id) var star == 1 quest

temp * cbi == 1 quest

quest + phi == 1 quest

1 quest == 1 bi

[03]

Lexical Analysis

- Its very first phase is CD.
- A Lexical analyzer (Lexer) takes the modified source code which is written in the form of sentences.
- o In other words →
 - "Lexer converts a sequence of characters into a sequence of tokens."
- o It also removes the extra space & comments written in the source code.
- So, programs that performs lexical analysis are called "Lexical analyzer or Lexer".
- ∵ If lexer detects that the token is invalid, it generates an error.
- It reads character streams from the source code, checks for the legal tokens, and pass the data to the syntax analyzer when it demands.
- Example:- "How Pleasant Is The Weather?"
 - ∴ Here we can easily recognize that there are fine words.
 - Humans can recognize the separators, blanks, and the punctuation symbols.

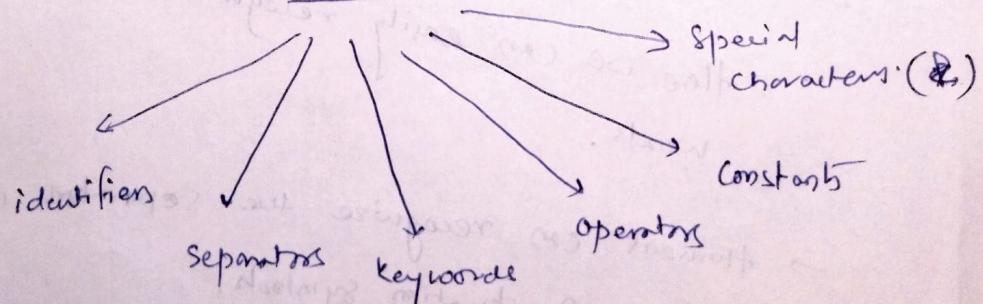
lets try again →

How PI agent is The else other?
→ now it is difficult to read but still we
can manage to read it.
→ Lexer will check for all the valid tokens.

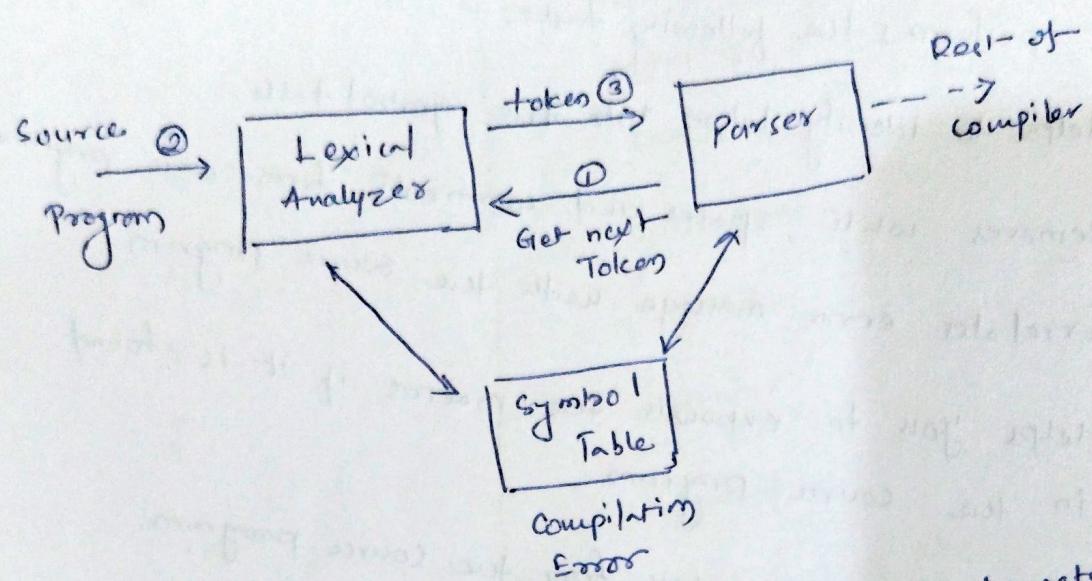
o Basic Terminologies :-

- ① Lexeme :- A lexeme is a sequence of characters that are included in the source program according to the matching pattern of a token.
• It's nothing but an instance of token.
- ② Token :- The token is a sequence of characters which represents a unit of information in the source program.
- ③ Pattern :- A pattern is a description which is used by the token.
→ In case of keyword which uses as a token, the pattern is a sequence of characters.

Tokens



0 Lexical Analyzer Architecture :-



- The main task of lexer is to read input characters in the code and produce tokens.
- Lexer scans the entire source code of the program.
- It identifies each token one by one.
- scanners are usually implemented to produce tokens only when requested by a parser.

Working :-

1. "Get the next token" command which is sent from the parser to the lexical analyzer.
 2. On receiving this command, the lexer scans the input until it finds the next token.
 3. It returns the token to the parser.
- Lexer will skip whitespace and comments while creating tokens.
- If any error is present, then lexer will correlate that error with the source file and line number.

o Roles of Lexical analyzer :-

it performs the following tasks:-

- ① Helps to identify tokens into the symbol table.
- ② Removes white spaces and comments from source program.
- ③ Correlates error message with the source program.
- ④ Helps you to expands the macros if it is found in the source program.
- ⑤ Read input character from the source program.

Example 8 -

```
#include <stdio.h>
int maximum(int x, int y)
{
    // This will compare 2 numbers.
    if (x > y)
        return x;
    else
        return y;
}
```

An:-

| Lexeme | Tokens |
|---------|------------|
| int | keyword |
| maximum | identifier |
| (| operator |

| | |
|-----|------------|
| int | keyword |
| x | identifier |
| , | operator |
| int | keyword |
| y | identifier |
|) | operator |
| { | operator |
| if | keyword |
| (| operator |
| x | identifier |
| > | operator |
| y | identifier |

Non-tokens:

| Type | Examples |
|------------------------|---------------------------------|
| comment | // This will compare 2 numbers. |
| Pre-procesor directive | # include <stdio.h> |
| Pre-procesor directive | # define NUMS 8,9 |
| Macro | NUMS |
| whitespace | /n, /b, /t |

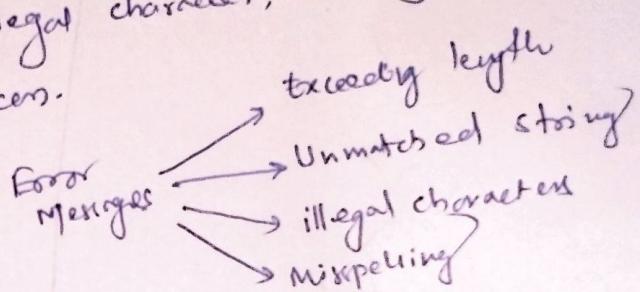
① Lexical Errors :-

→ A character sequence which is not possible to scan into any valid tokens is a lexical error.

- Lexical errors are not very common but it should be managed by a scanner.

- Mis-spelling of identifiers, operators, keywords are considered as lexical errors.

- Generally, a lexical error is caused by the appearance of some illegal character, mostly at the beginning of a tokens.



② Error Recovery in Lexer :-

so few most common error recovery techniques:-

- ① Removes one character from the remaining input.
- ② In the panic mode, the successive characters are always ignored until we reach a well-formed tokens.
- ③ By inserting the missing character into the remaining input.
- ④ Replace a character with another character.

Ex :- `printf("i=%d, & i=%x", i, &i);`

Q Lexical Analyzer & Parser :-

| LA | Parser |
|------------------------------------|--------------------------------------------------|
| ① Scans Input program. | ① Performs syntax analysis |
| ② Identify tokens. | ② Create an abstract representation of the code. |
| ③ Insert tokens into symbol Table. | ③ Update symbol Table entries. |
| ④ It generates lexical errors. | ④ It generates syntax errors. |

Q why separate Lexical & Parser :-

- ① The simplicity of the design: It eases the process of lexical analysis and the syntax analysis by eliminating unwanted tokens.
- ② To improve compiler efficiency: Helps you to improve compiler efficiency.
- ③ Specialization: specialized techniques can be applied to improve the lexical analysis process.
- ④ Portability: only the scanner requires to communicate with the outside world.
- ⑤ higher Portability:

② Adv. of Lexer :-

- ① Lexer can be used by the programs like compilers which helps to generate compiled binary executable code.
- ② It is used by the web browsers to format and display a web page with the help of parsed data from JavaScript, HTML, CSS.
- ③ A separate lexical analyzer helps you to construct - a specialized and potentially more efficient processor for the task.

③ Dis. Adv. of Lexer :-

- ① Time consuming:- You need significant time to read source code & form tokens.
- ② Some regular expressions are quite difficult to understand.
- ③ More effort is needed to develop and debug the lexer and its token descriptions.
- ④ Additional runtime overhead is required to generate the lexer tables and construct the tokens.

Ex:- Find the lexemes and tokens from the following
'C' code?

int get (int a, int b)

{
 while (a != b)

{
 if (a > b)

 a = a - b ;

 else
 b = b - a ;

}
return a ;

{

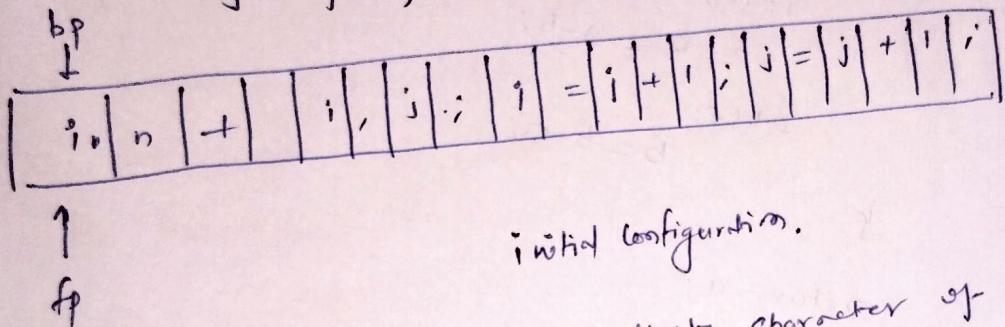
Input Buffering in Compiler Design

- The lexer scans the input from L to R & char by char at a time.
- It uses two pointers begin ptr (bp) and forward ptr (fp) to keep track of the pointer of input scanned.

∴ int i, j;

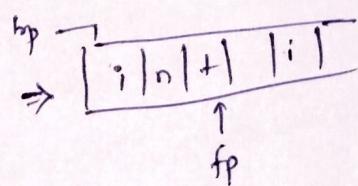
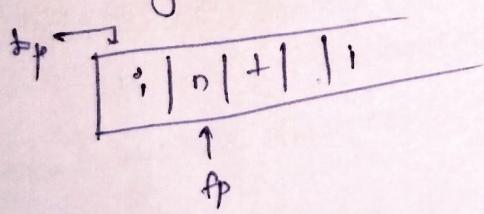
i = j + 1;

j = j + 1;



initial configuration.

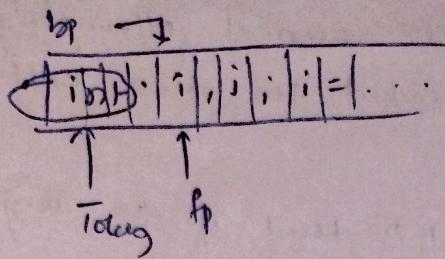
- initially both pointers point to the first character of the input string as shown below.



- The fp will move ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates the end of lexeme.

↳ space occurs i.e. \rightarrow (int) lexeme is identified.

- The fp will move ahead at white space, when fp encounters white space, it ignore and move ahead.
- Then both pointer [bp+fp] are set at next tokens.



→ The input character is read from secondary storage. but reading in this way from secondary storage is costly.

→ ∴ Buffering technique is required / used.

→ A block of data is first read into a buffer, and then second by lexical analyzer.

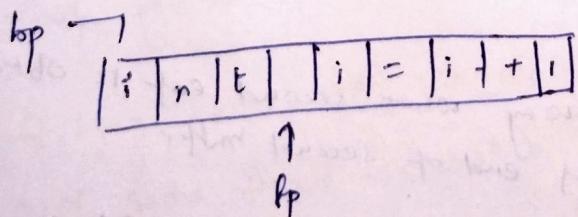
→ There are two methods used in this context:

- 1) One Buffer Scheme
- 2) Two Buffer Scheme.

① 1 Buffer Scheme: In this scheme, only one buffer is used to store the input string.

→ Problem with this scheme → Lexem Length.

→ Problem with this scheme → Lexeme crosses the buffer boundary. ∵ to scan rest of the lexeme the buffer has to be refilled, that makes overwriting of first lexeme.

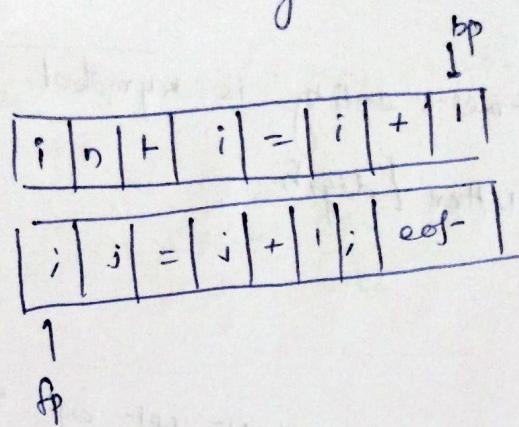


② 2 Buffer Scheme:

→ To overcome the problem of one buffer scheme, two buffer scheme is used to store the input string.

- The first buffer and second buffer are scanned alternately.
- When end of current buffer is reached the other buffer is filled.
- The only problem with this method is that if the length of buffer the lexeme is longer than the length of buffer then scanning symbol cannot be scanned completely.
- Initially both bp and fp will point to the first char of the first buffer.
- Then fp moves forward to search for the end of lexeme
- As soon as the blank character is recognized, the string between bp and fp is identified as Token.
- To identify the boundary boundary of first buffer, end of buffer char is placed at the end of first buffer.
- Similarly end of second buffer is also recognized by the end of buffer mark present at the end of second buffer.
- When fp encounters first eof, then one can recognize end of first buffer and hence filling up second buffer is started.
- In the same way when second eof is obtained then it indicates of end of second buffer.
- Alternatively both the buffers can be filled up until end of input program and stream of tokens is identified.

- This self-character introduced at the end is called "Sentinel". which is used to identify the end of buffer.



FSR

Specifications of Tokens

There are three specifications of Tokens

- 1) Strings
- 2) Language
- 3) Regular Expression.

* Strings and Languages:-

- An alphabet or character class is a finite set of symbols.
- A string over an alphabet is a finite sequence of symbols drawn from that alphabet.
- A language is any countable set of strings over some fixed alphabet.

∴ The length of string is generally written as $|s|$.

String and Language:-

→ Symbol :- An abstract entity is symbol.

Ex:- letter / digit

String and finite

→ Char class / Alphabet class :- Its finite set of symbols.
 $\{0, 1\}$ $\{A, B, C\}$

→ String :- A string over an alphabet is a finite set of symbols drawn from that alphabet.

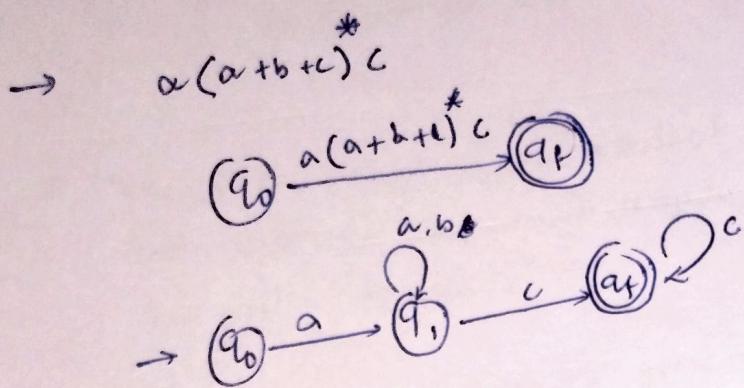
Ex:- ① $abcb \rightarrow$ one string over $a, b, + c$.
 cab**o**
② $|w|$ is the length of string
③ ϵ is an empty string of length 0.

→ Language :- is any countable set of strings over some fixed alphabet.

- ① \emptyset & $\{\epsilon\}$ are language.
- ② The set of palindromes over $\{0, 1\}$ is infinite.
L → The set of palindromes of 5 char over $\{0, 1\}$ is a f. l.
- ③ The set of strings $\{01, 10, 11, 00\}$ over $\{0, 1\}$ is f. l.

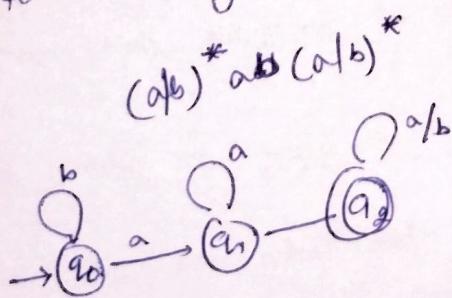
∴ If Σ is an alphabet
 \downarrow
 Σ^* is the set of all strings over Σ .

- Each subset of Σ^* is a language.
- RE (T3), CFG (T2 or CFL), NL (T1 or CSL), or
to L. or REL.
- REC CFL CCSL C REL

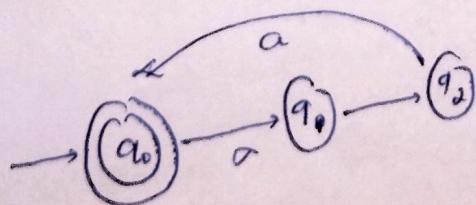


Ex → string ending with

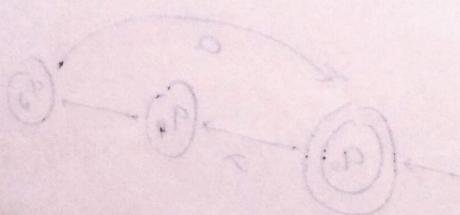
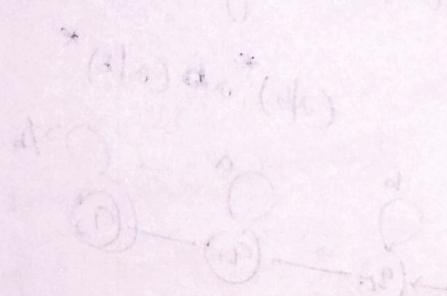
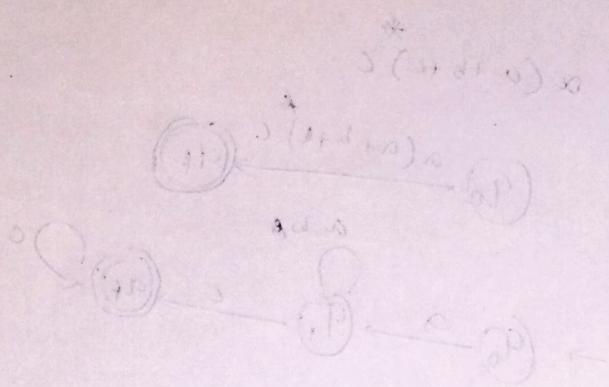
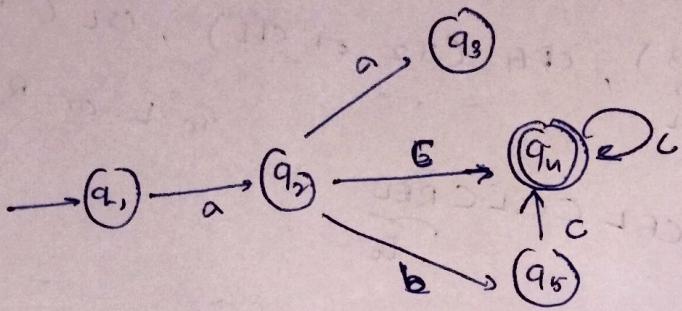
Ex: all strings with sub-string 'ab'



Ex: strings in which no of a is divisible by 3. $L = \{a^{3n} | n \geq 0\}$



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



-: Specification of Tokens :-

① ∵ Regular definitions, a mechanism based on regular expressions are very popular for specification of Tokens.

① It has been implemented in the lexical analyzer generator too → LEX

② We study regular expressions first, and then specify of tokens using LEX.

③ Transition diagrams, a variant of finite automata, are used to implement regular definitions and to recognize tokens.

① LEX automatically generate FSA from regular definition.

② So we will study FSA and their generation from regular expressions in order to understand transition diagrams in LEX.

But first we need to understand → The following terms.

1. Symbols → Its an abstract entity, not defined.
Ex:- digits, letters, special symbols.

2. char class / Alphabet class → A finite set of symbols.
Ex:- $\{0, 1\}$, $\Sigma = \{A, B, C\}$, etc.

3. String or String over an set of alphabet Σ in the finite sequence of symbol drawn from Σ .

Ex. 101010...
 $\Sigma = \{0, 1\}$ with alternati 1's and 0's

3. Prefix of string :- string obtained by removing zero or more symbols from end of string.

Ex:- ban, bana, banana, e, are prefix of "banana".

4. Suffix of string :- string obtained by removing zero or more symbols from beginning of string.

Ex:- ana, na, banana, e, are suffix of "banana".

5. Substring :- string obtained by deleting any prefix and any suffix from string.

Ex:- banana, na, ana, e are substring of "banana"

6. Proper Prefix :- Prefix which is not equal to ϵ or string it self.

Ex:- ban, banan

7. Proper Suffix :- Suffix which is not equal to ϵ or string it self.

Ex:- na, anana, etc.

8. Sub-sequence :- string formed by deleting zero or more not necessarily consecutive positions of string.

Ex:- baan, anaa, ...

Operations on Language:-

| operations | definition & notation |
|--------------------------|-------------------------------------------------------------------------------------------------------------------|
| ① Union of L & M. | $L \cup M = \{s s \text{ is in } L \text{ or } s \text{ is in } M\}$ $= \{s s \in L \text{ or } s \in M\}$ |
| ② Concatenation of L & M | $LM = \{st s \text{ is in } L \text{ and } t \text{ is in } M\}$ $= \{st s \in L \text{ and } t \in M\}$ |
| ③ Kleene closure of L | $L^* = \bigcup_{i=0}^{\infty} L^i$ |
| ④ Positive closure of L | $L^+ = \bigcup_{i=1}^{\infty} L^i$ |

⑤ Regular Definition :- (Used for token specification).

→ for some alphabet set Σ , sequence of regular

definitions :

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$\vdots$$

$$d_n \rightarrow r_n \quad \text{where,}$$

(1) each d_i is new symbol (not in Σ or other d_j)

(2) r_i is regular expression over $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

* Regular Expression :-

- ① The Language accepted by finite automata can be easily described by simple expressions are called Regular Expressions.
 - It is most effective way to represent any language.
- ② "A regular expression is an algebraic formula whose value is a pattern consisting of a set of strings, called the language of expression."
- ③ Regular Expressions can also be described as a sequence of patterns that defines a string.

∴ For instance →

→ In RE, x^* means zero or more occurrence of x .
{ i.e. $\epsilon, x, xx, xxx, xxxx, \dots$ }

→ In RE, x^+ means one or more occurrence of x
{ i.e. $xx, xxx, xxxx, \dots$ }

* R.E. can be recursively defined as follows:-

- ϵ is RE indicates the language containing an empty string $L(\epsilon) = \{\epsilon\}$
- \emptyset is a RE denoting an empty language $L(\emptyset) = \{\emptyset\}$.
- x is a RE where $L = \{x\}$
- If x is RE denoting the language $L(x)$ and y is a RE denoting the language $L(y)$, then
 - (1) $x+y$ is a R.E. corresponding to the language $L(x) \cup L(y)$ ∴ $L(x+y) = L(x) \cup L(y)$

(i) $X \cdot Y$ is a R.E. corresponding to the language $L(X) \cdot L(Y)$
 $\therefore L(X \cdot Y) = L(X) \cdot L(Y)$

(ii) r^* is an R.E. then $L(r^*) = R^* = \bigcup_{i=0}^{\infty} R^i$
 $\therefore R = L(r)$
 $\therefore R^* = (L(r))^*$

\therefore If we apply any of the above rule, they are RE.

Some RE examples :-

| RE | Regular set - |
|------------------------------|-----------------------------------------------------------------------------------------------------------------------|
| ① $(0+10)^*$ | $L = \{0, 1, 10, 100, 1000, 10000, \dots\}$ |
| ② $(0^* 1 0^*)$ | $L = \{1, 01, 10, 010, 001, 100, 0010, 00100, \dots\}$ |
| ③ $(0+\epsilon)(1+\epsilon)$ | $L = \{\epsilon, 0, 1, 01\}$ |
| ④ $(a+b)^*$ | $L \Rightarrow$ set of strings of a's and b's of any length including null string. |
| ⑤ $(a+b)^* ab$ | $L = \{a, b, ab, aa, bab, \dots\}$ $L \Rightarrow$ set of strings of a's and b's ending with a'b. |
| ⑥ $(11)^*$ | $L = \{11, 111, 1111, \dots\}$ $L \Rightarrow$ set consisting even number of 1's. |
| ⑦ $(aab)^* (bbb)^* b$ | $L \Rightarrow$ set of strings consisting even no of a's followed by odd no's of b's $L = \{b, aab, aabb, \dots\}$ |

$$④. (aa + ab + ba + bb)^*$$

$L \Rightarrow$ string of 'a's and 'b's with even length obtained by concatenating any combination of 'aa, bb, ab, ba.
 $L = \{aa, ab, bb, ba, aabb, abba, \dots\}$

o Regular Set → Any set that represents the value of the Regular Expression is called Regular set.

∴ Properties of Regular Set →

① The union of two regular set is regular.

② The intersection of two regular set is regular.

Proof →

$$RE_1 = a(aa)^*$$

$$RE_2 = (aa)^*$$

$\therefore L_1 = \{a, aaa, aaaa, \dots\} \rightarrow$ odd length of 'a'

$L_2 = \{\epsilon, aa, aaaa, \dots\} \rightarrow$ even length of 'a'

$\therefore L_1 \cup L_2 = \{\epsilon, a, aa, aaa, \dots\} \rightarrow$ strings of all possible lengths of 'a'

$$RE(L_1 \cup L_2) = a^*$$

Proof → $RE_1 = a^*$ $RE_2 = (aa)^*$

$\therefore L_1 = \{a, aa, aaa, \dots\}$ $L_2 = \{\epsilon, aa, aaaa, \dots\}$

$L_1 \cap L_2 = \{aa, aaaa, aaaaaa, \dots\}$

$$RE(L_1 \cap L_2) = a(aaa)^*$$

③ The complement of regular set is regular.

$$\therefore RE = (aa)^*$$

$$L = \{ \cancel{a}, aa, aaa, \dots \}$$

$$\therefore L' = \{ a, aaa, aaaaa, \dots \}$$

$$RE(L') = a(aa)^*$$

④ The difference of two regular set is regular.

$$\therefore RE_1 = a(aa)^*$$

$$RE_2 = (aa)^*$$

$$\therefore L = \{ a, aa, aaa, aaaa, \dots \}$$

$$L_2 = \{ \cancel{a}, aa, aaaa, \dots \}$$

$$L_1 - L_2 = \{ a, aaa, aaaaa, \dots \}$$

$$RE(L_1 - L_2) = a(aa)^*$$

⑤ The reversal of regular set is regular.

$$\therefore L = \{ 10, 10, 00, 11 \}$$

$$RE(L) = 01 + 10 + 00 + 11$$

$$L^R = \{ 10, 01, 00, 11 \}$$

$$RE(L^R) = 10 + 01 + 00 + 11 \quad \text{i.e. } RE.$$

Q. 4. →

- ① Write a RE for a language accepting all combinations of a's over $\Sigma = \{a\}$.

Ans. → RE = a^*

- ② Except null string

Ans: RE = $a(a)^* = a^+$

- ③ Any no of a's and b's

Ans: $(a+b)^*$

- ④ RE for the language accepting all strings which are starting with 1 and ending with 0. over $\Sigma = \{0, 1\}$.

Ans: RE = $1(0+1)^*0$

- ⑤ RE for the language starting and ending a and having any combinations of b's in between over $\Sigma = \{a, b\}$

Ans: RE = ab^*a .

- ⑥ RE for → starting with a but not having consecutive b's

Ans: $(a+ab)^*$

⑦ Write a RE for the language having a string which should have at least one 0 and at least one 1.

$$\text{Ans:-- } \text{RE} = [(0+1)^* 0 (0+1)^*, (0+1)^*] + [(0+1)^* 1 (0+1)^* 0 (0+1)^*]$$

⑧ $L = \{0, 1\}$ "RE for string do not contain substring '01'."

$$L = \{\epsilon, 0, 1, 100, 1000, 1100, 1110, \dots\}$$

$$\text{Ans:-- } \text{RE} = (1^* 0^*) = (1^* 0^*).$$

⑨ RE for L containing the string in which every 0 is immediately followed by 1.

$$\text{Ans:-- } \text{RE} = (011 + 1)^*$$

⑩ Regular definition:-

A regular definition is a sequence of definitions in the form:-

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$d_3 \rightarrow r_3$$

:

$$d_n \rightarrow r_n$$

where:-

(i) each d_i is a distinct symbol

(ii) each r_i is a regular expression over the symbols

$$\Sigma \cup \{d_1, d_2, d_3, \dots, d_n\}$$

Ex: ① C identifiers are strings of letters, digits, and underscores. Here is a regular definition of the language of C identifiers. We shall conventionally use italics for the symbols defined in regular definitions.

letters $\rightarrow A | a | \dots | Z | z | b | \dots | z | -$

digit $\rightarrow 0 | 1 | \dots | 9$

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

Ex: ② Unsigned numbers (integer or floating point) are strings such as 5280, 0.01234, 63.36E4, or 1.89E-4. The regular definition is,

digit $\rightarrow 0 | 1 | \dots | 9$

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

optional fraction $\rightarrow \cdot \text{digit}^* | \epsilon$

optional exponent $\rightarrow (E (+|-|-\epsilon)) \text{digit}^* | \epsilon$

number $\rightarrow \text{digits} \cdot \text{optional fraction} \cdot \text{optional exponent}$.

* Extensions of Regular Expressions \rightarrow

① One or More Instances \rightarrow The unary, postfix '+' operator represents the positive closure of a regular expression. And its language. That is if r is a RE, then $(r)^+$ denotes the language $(L(r))^+$.

① Zero or one instance → The unary postfix operator '?' means 'zero or one occurrence'. That is $\sigma ?$ is equivalent to $\sigma \in$ or we can say $L(\sigma?) = L(\sigma) \cap \{ \sigma \in \}$

② Character class → A regular expression $a_1 | a_2 | \dots | a_n$ where a_i 's are each symbols of the alphabet, can be replaced by $[a_1, a_2, \dots, a_n]$. More importantly, when a_1, a_2, \dots, a_n form a logical sequence e.g. (consecutive upper letters, lowercase letters, digits) we can replace them by $a_1 - a_n$.

Thus. $|a|b|c|\dots|z \rightarrow [a, b, c, \dots, z] \rightarrow [a-z]$

Let's visit Ex ① & ② again →

Ex :- ① → letter $\rightarrow [A - Z a - z -]$
digit $\rightarrow [0 - 9]$
id $\rightarrow \text{letter} - (\text{letter} | \text{digit})^*$

Ex :- ② → digit $\rightarrow [0 - 9]$
digits $\rightarrow \text{digit}^+$
number $\rightarrow \text{digits} (\cdot \text{digits})? (E [+ -] ? \text{digits})?$

o Regd. Recognition of Tokens.

stmt → if expr then stmt
| if expr then stmt else stmt
| ε

expr → term relop term
| term

term → id

| number

↑ "A grammar for branching statements"

↓ Pattern of Tokens (Regular Definition)

number → [0 - 9]

digit → digit⁺

number → digits^{*} (. digits)? (E [+ -]? digits)?

letter → [A - Z a - z -]

id → letter (letter | digits)*

if → if

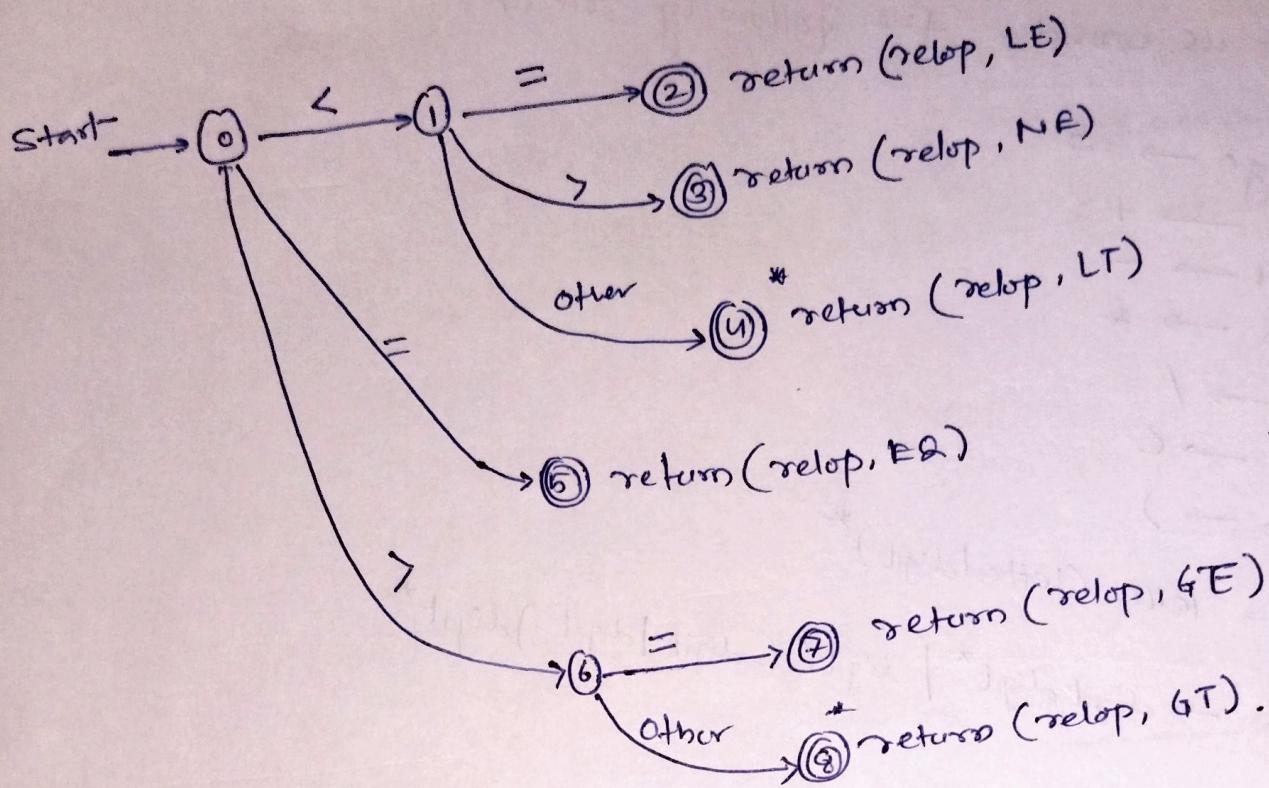
then → then

else → else

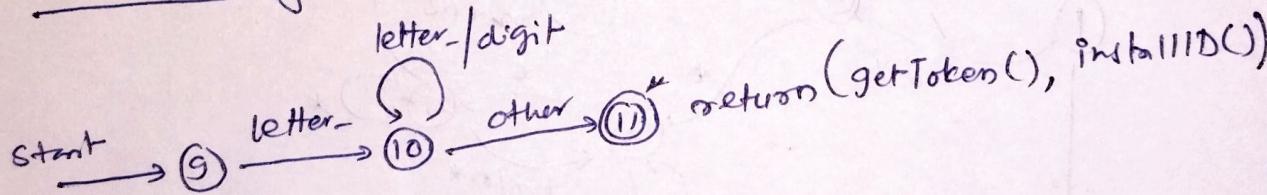
relop → < | > | <= | >= | = | <>

ws → (blank | tab | new | line)*

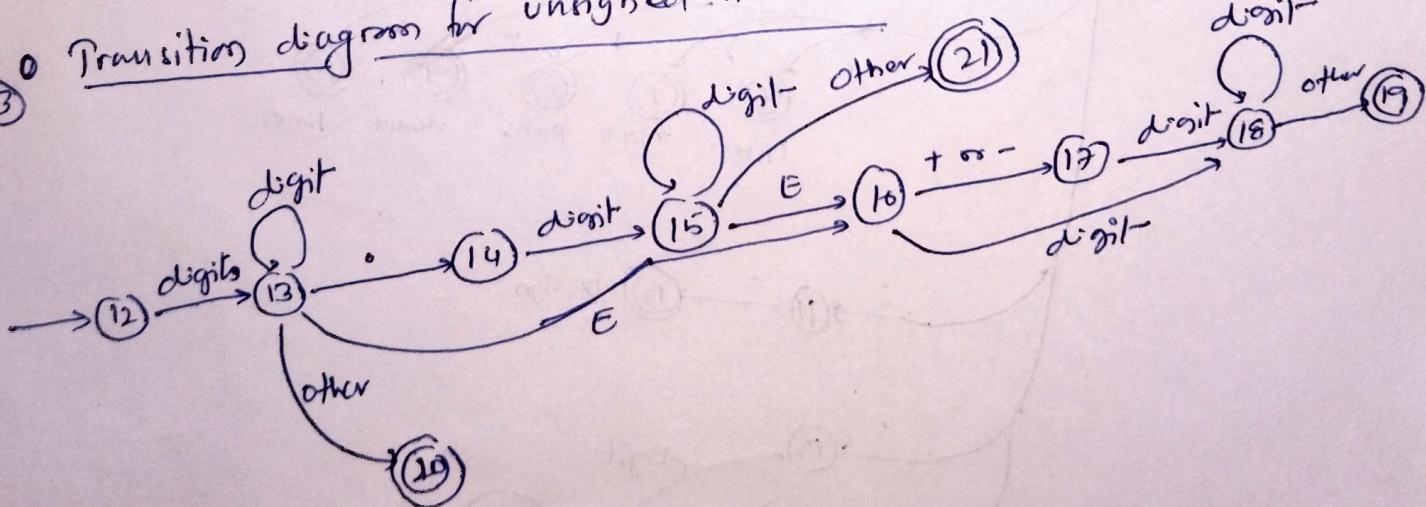
① ② Transition diagram from relop →



② ③ Transition diagram for id's & keyword →



③ ④ Transition diagram for unsigned numbers →



Q Tokens for a calculator language →

Let us consider the following set of tokens:-

assign → :=

plus → +

minus → -

times → *

div → /

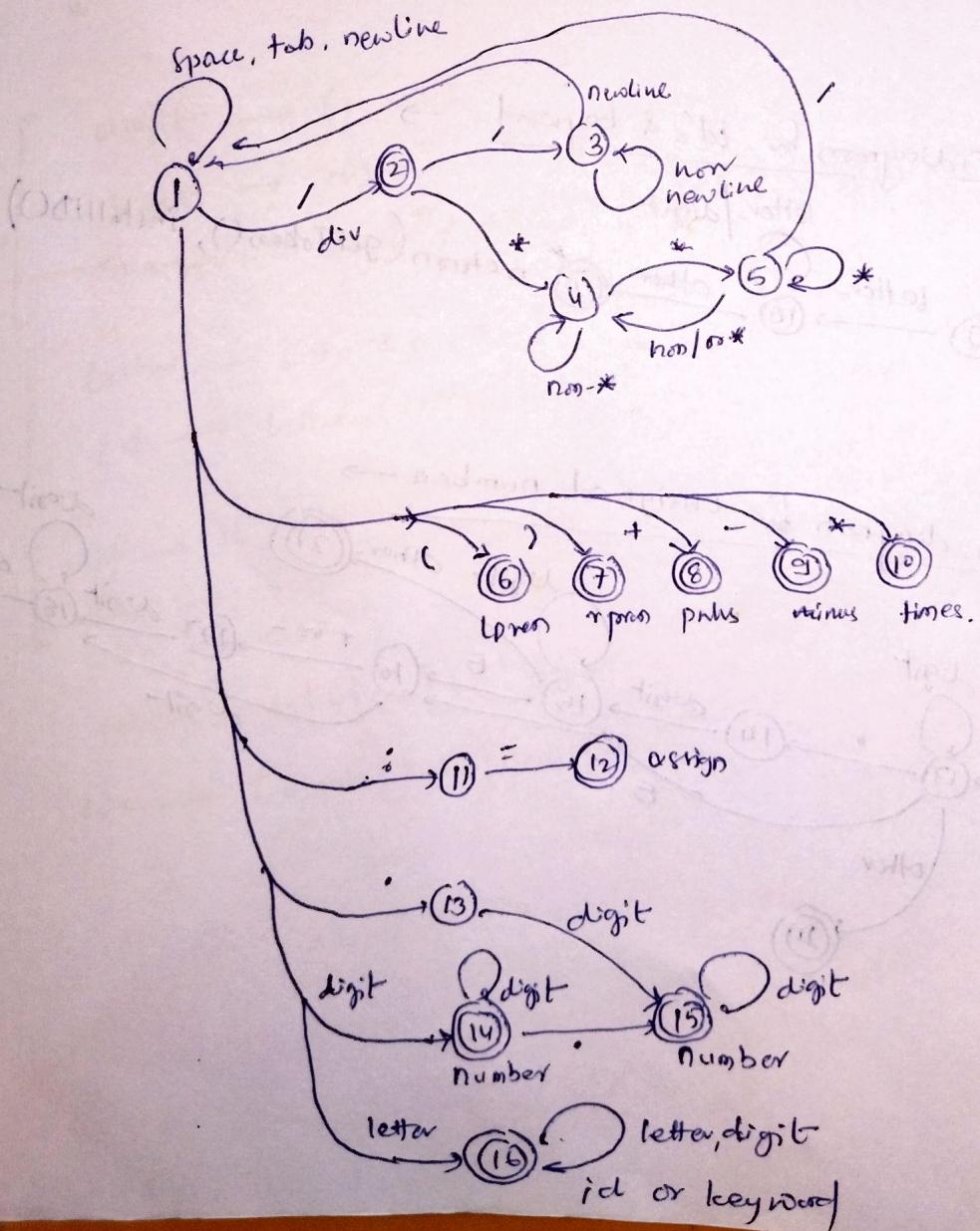
(paren → (

)paren →)

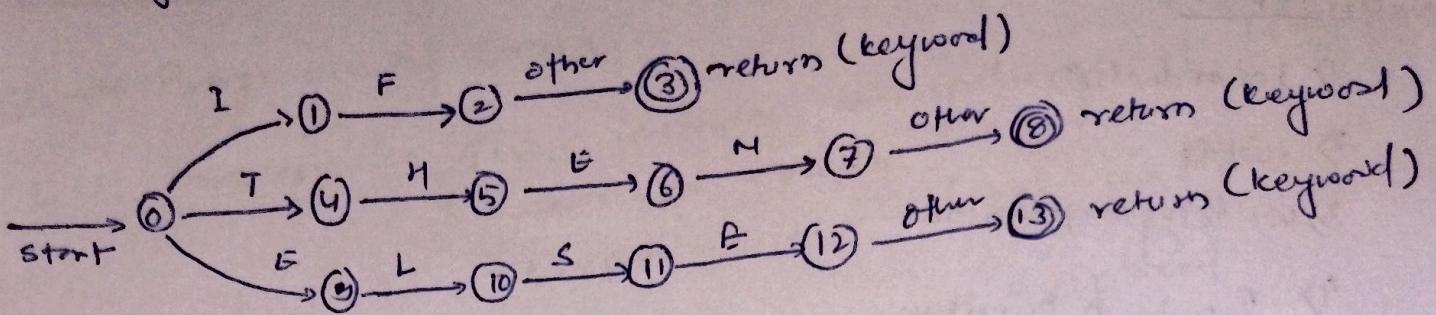
id → letter (letter/digit)*

number → digit digit* | digit* (. digit | digit .) digit*

number → digit digit* | digit* (. digit | digit .) digit*



④ Keyword → IF | THEN | ELSE



⑤ Delimiter / whitespace.

delimiter → space | tab | newline.

