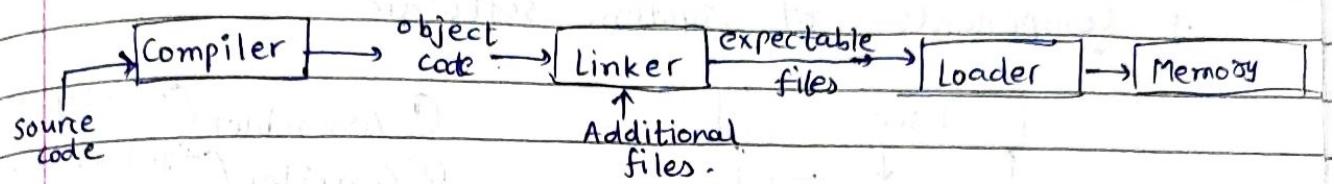


# Language Translator

(Jeffrey D. Ullman, Ravi Sethi,  
Alfred V. Aho)  
GBS Publication)

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_



- Q. Meanings = High level language (User Oriented languages)  
= Assembly Language, Machine Level language,  
Translator, Linker, Loader

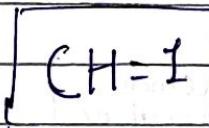
\* Real time application of Language Translator. (2 marks)

↳ Healthcare

Assembler → Assembly to low  
compiler → compiles whole

- Computer → S/w + H/w

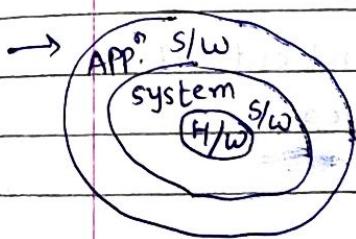
interpreter → compile line by  
line



## → Language Translation overview

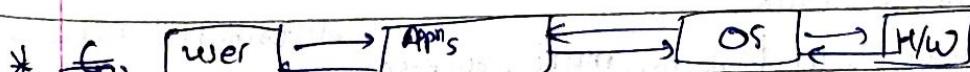
### \* need and component of system

- Application software = browser, word, bank managing system
- System Software = OS, compiler, Assembler, interpreter.

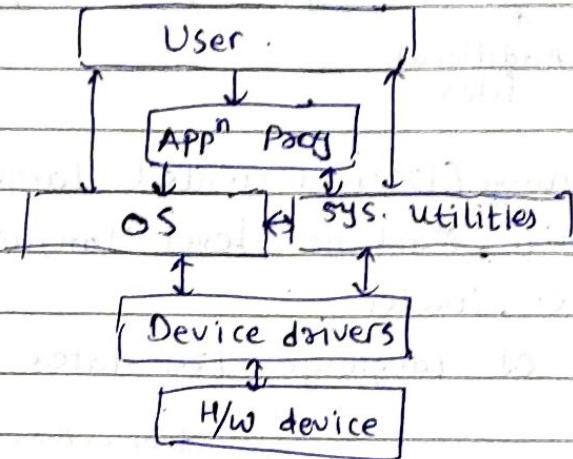


Q. \* (2 marks) System software is barrier b/w  
H/w & App<sup>n</sup> S/w.

- Also sys. s/w will ~~convert~~ convert high level lang. to machine lang.
- Also helpful to store database
- Manage f<sup>n</sup> storing data, retrieving files & <sup>Scheduling</sup> task
- To make effective execution of each prog.
- To make effective utilization of all the resources



## \* Components of system software



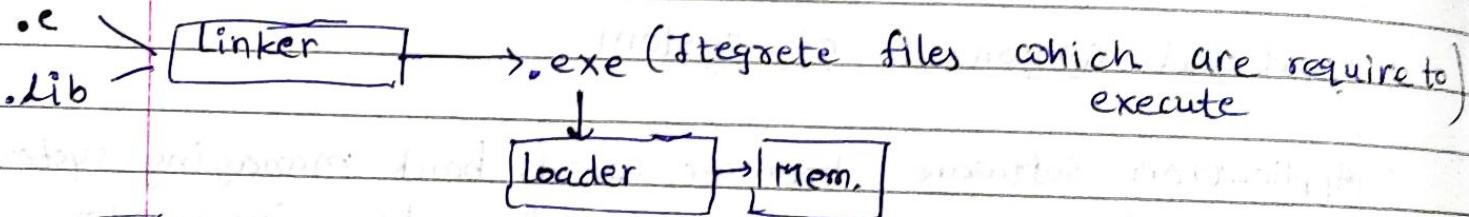
- ① Assembler (Assembly → Machine)
- ② Compiler (High level → Machine)
- ③ Interpreter (Line by line high level machine language)
- ④ Editor (edit particular file as per requirements)
- ⑤ Linker
- ⑥ Loader
- ⑦ Debugger
- ⑧ Macro
- ⑨ OS
- ⑩ device driver

→ Assembly lang → Assembler → Machine lang

→ High level → Compiler → Binaryx Machine

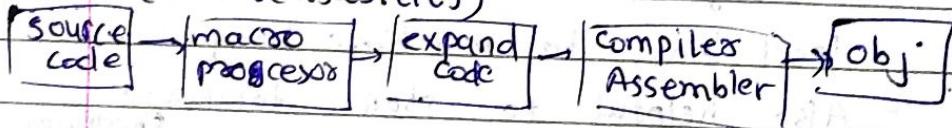
→ prog. statement → Interpreter → Machine lang → statement execution

→ Editor Types = Line editor, screen, word processor, structure, etc.



→ Debugger is a computer prog. used to find out the errors which are also called as bugs in source prog.

→ Macro (Code reusability)



→ OS is the system Prog. that enable user to communicate with computer H/W

→ It has other prog. to run & control them.

e.g. Linux, Windows, Unix.

Types - Single user, Multi tasking, multiprocessor, distributed, network

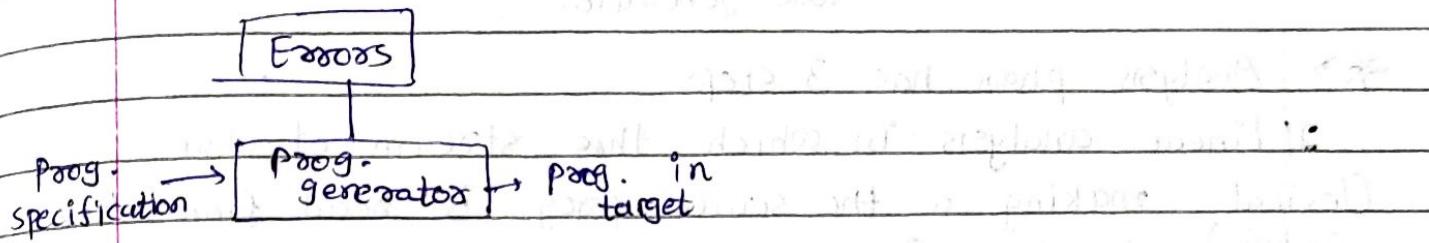
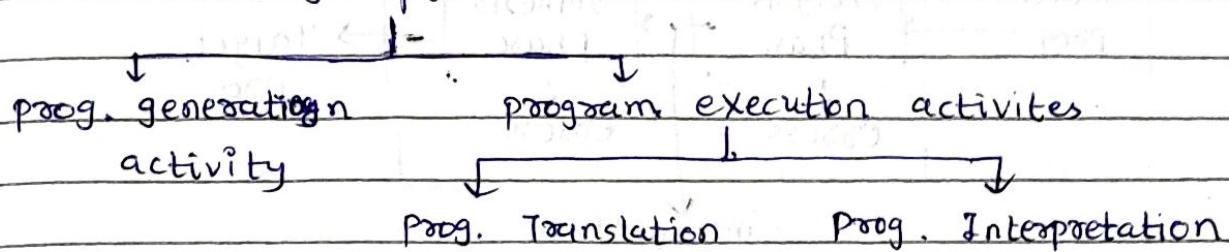
→ device driver → It is a system Prog. used to control no. of devices which are attached to the computer, it tells us that how device will work or certain commands which are generated by the user

e.g. mouse driver, KB, WiFi, Bluetooth

## \* Language processing Activities.

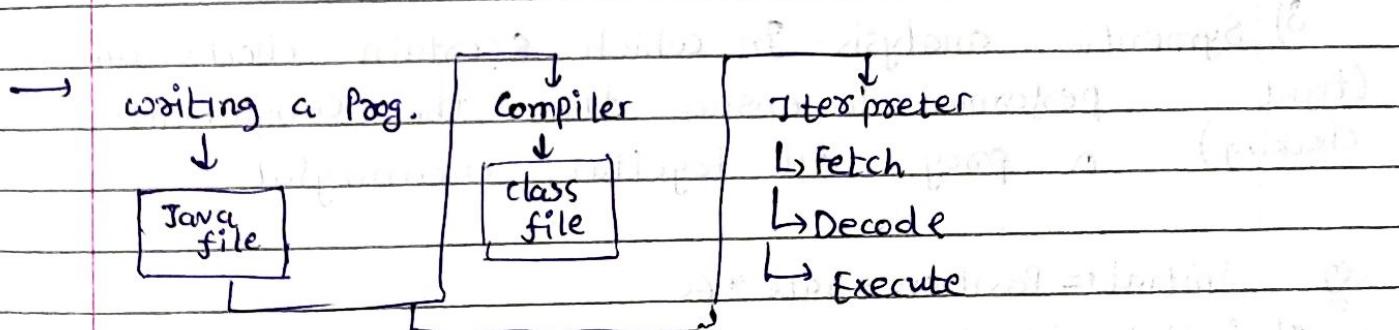
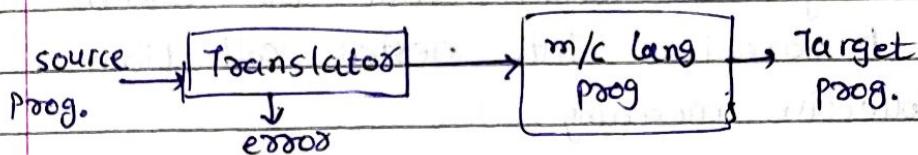
D.1 Language Processor = Convert source code into machine code.

D.1 Language Processing = Activity which is carried out by the language processor



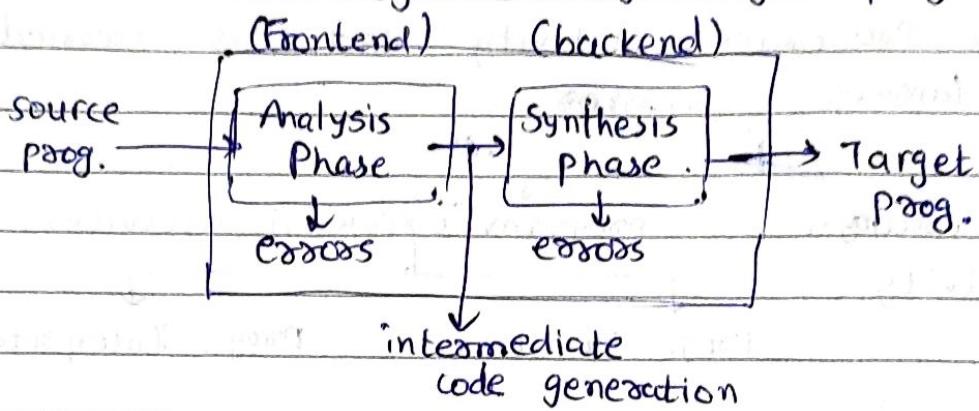
9/12/22

## \* Program Translation



# \* Fundamental of Language processing.

Lang: Processing = Analysis of source Prog. + synthesis of target prog.



→ Analysis phase has 3 steps

1) Linear analysis - in which this stream of char (lexical analysis) making a the source prog. is read from left to right & grouped into tokens that are sequence of char having a collective meaning

2) Hierarchical analysis - In which char or tokens are grouped hierarchical into nested collections with collective meaning.

3) Syntactic analysis - In which certain checks are performed to ensure that the components of a prog. fit together meaningful

e.g. initial := Position + rate \* 60

1) initial → identifier  
:= → assignment symbol

position → identifier  
+ → plus sign(symbol)

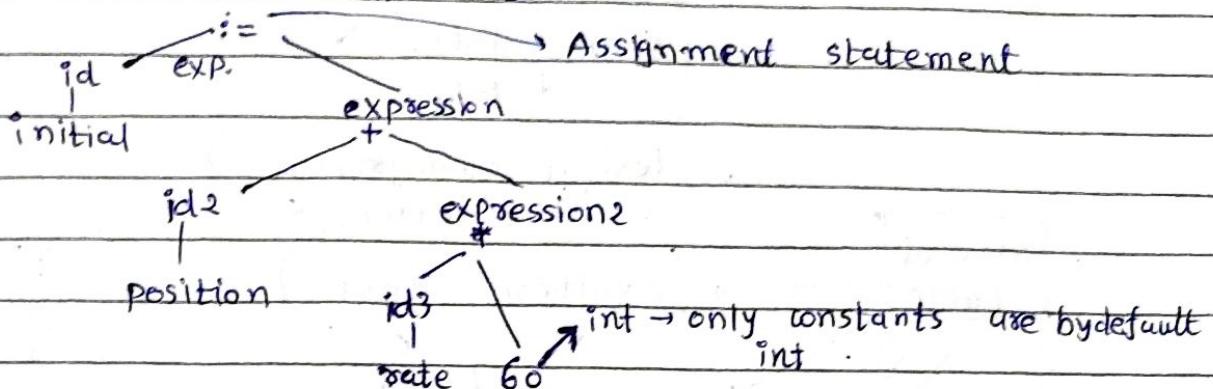
rate → identifier

\* → mul. sgn (Symbol)

60 → int (constant)

2)  $\text{initial} = \text{position} + \text{rate} * 60$

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



3) \* (Intermediate code generation) So only here tight checking

$$\rightarrow \text{temp1} = \text{inttoreal}(60)$$

$$\text{temp2} = \text{rate} * \text{temp1}$$

$$\text{temp3} = \text{position} + \text{temp2}$$

$$\text{initial} = \text{temp3}$$

→ Synthesis has 2 activities

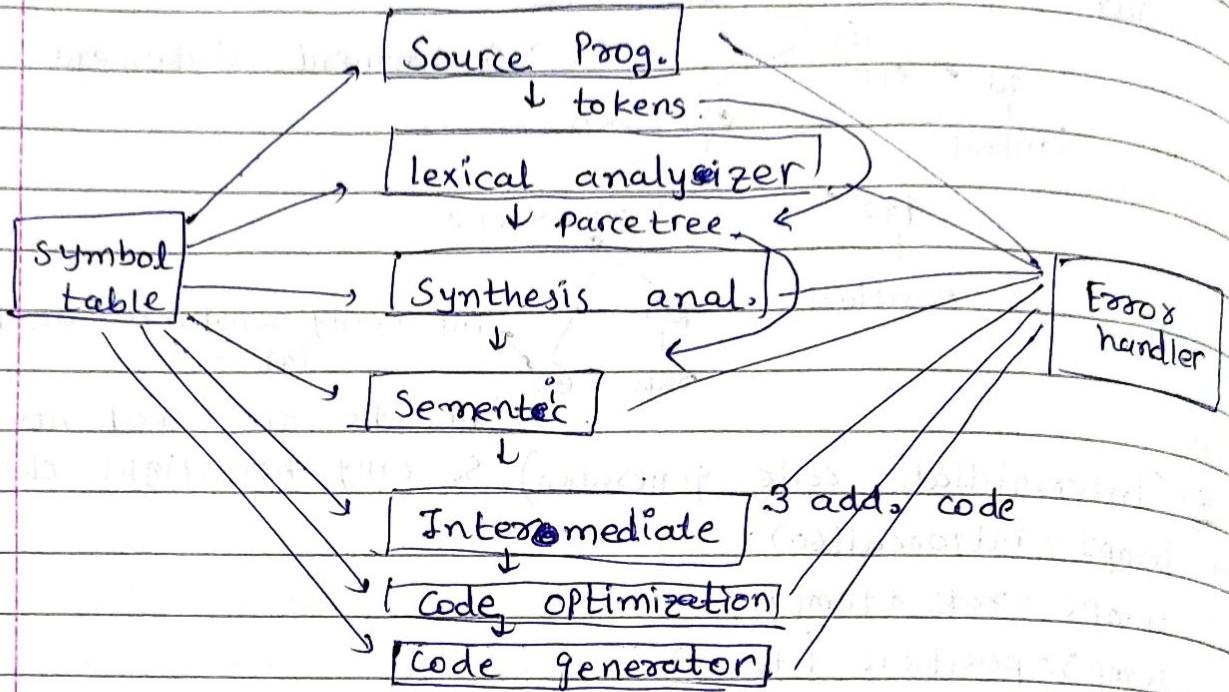
1) Code optimizer , 2) Code generator

(1) optimization Intermediate  $\Rightarrow \text{temp1} = \text{id}_3 * 60.0$   
Code  $\text{id}_1 = \text{id}_2 + \text{temp1}$

(2) generation

MOVF R1, id3
MULF R1, #600
MOVF R2, id2
ADDF R2, R1
MOVF id1, R1.

## \* Phases of Compiler



$$\text{Q. Total} = a/b + c * 60$$

## Token generation

Total → identifier

= → assignment symbol

a → identifier  
/ → divide sign  
(symbol)

b → identifier  
+ → Plus sign  
(symbol)

C → identifier  
\* → mul. sign  
c symbol

60 → int (constant)

### 3) Intermediate code generator.

temp1 = a/b

temp2 = C\*~~60~~ temp4 → temp4 = inttoreal(60)

$$\text{temp3} = \text{temp1} + \text{temp2}$$

Total = temp3

→ Code optimization

$$\text{temp}2 = C * 60.0$$

$$\text{Total} = \text{temp}1 + \text{temp}2.$$

→ Code generation.

MOVF RI, a

MOVF R2, b

DIVF RJ, R2

MOVF R2, C

MULF R2, #60.0

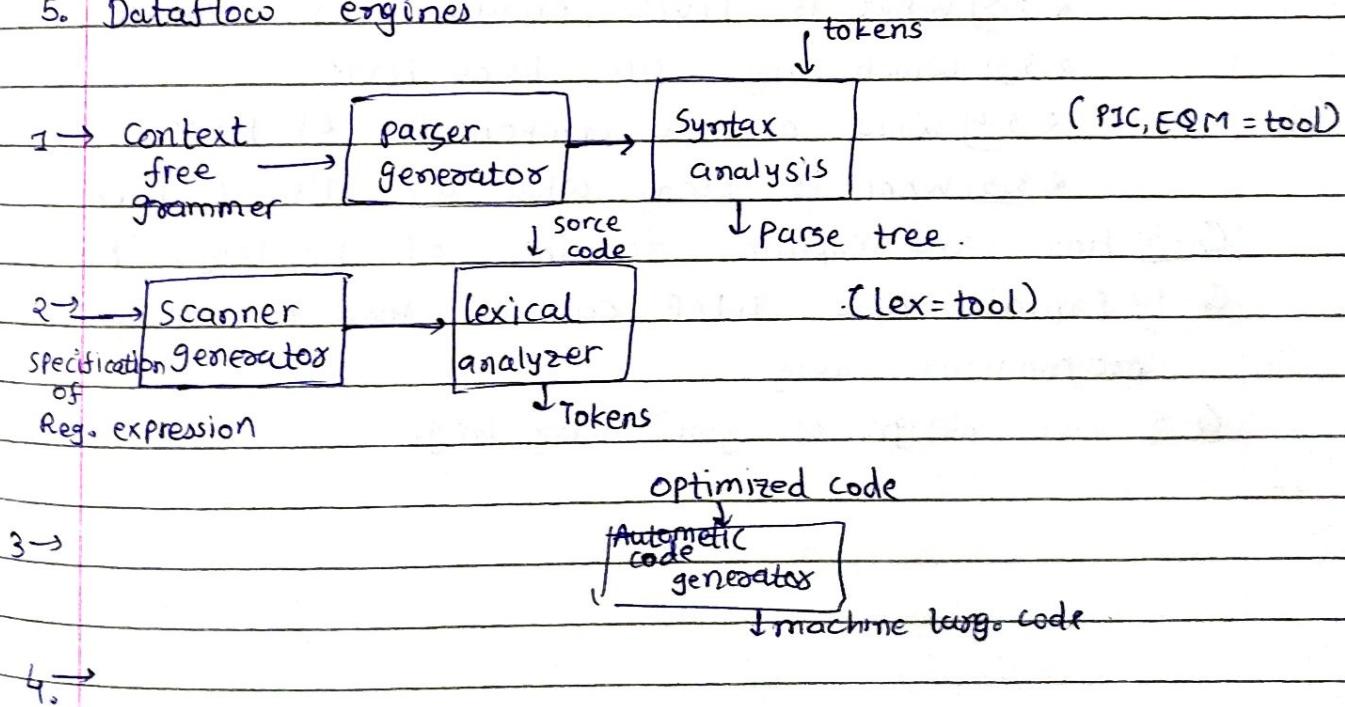
Get total  
Register  
(total)

MOVF id1, RI

## \* Compiler construction tools.

\* Rules.

1. Parser generator
2. Scanner generator
3. Syntax directed translation engines. ( generate intermediate code )
4. Automatic code generators.
5. Dataflow engines



Lab I. Aim = study compiler and various related tools.

Q-1] Write a simple prog. in C/C++ lang. and observe compiling, linking & loading process using GCC compiler (Find difference in size of files generated) (Find which files generated after each step)

Q-2] Study & ans. the following Ques.

a.1] Discuss following

compiler, interpreter, Linkers, loaders

Q-2] GCC

2.2.1] What is GCC, explain it's features

2.2.2] In which lang. is GCC written

2.2.3] What standard optimizations are supported by GCC

2.2.4] Explain following options of GCC

-O, -S, -c, -save, -f, V, temps, -c(static

size<exefilename>

objdump -hrt <objfile>

Q-3] LLVM

2.3.1] Understand & discuss following keyword

OpenCL, OpenCL, OpenMP, clang, LLVM, Bison.

2.3.2] What is LLVM, explain it's feature

2.3.3] Which lang. are Using LLVM

2.3.4] What are the components of LLVM

2.3.5] What is clang, what are it's features

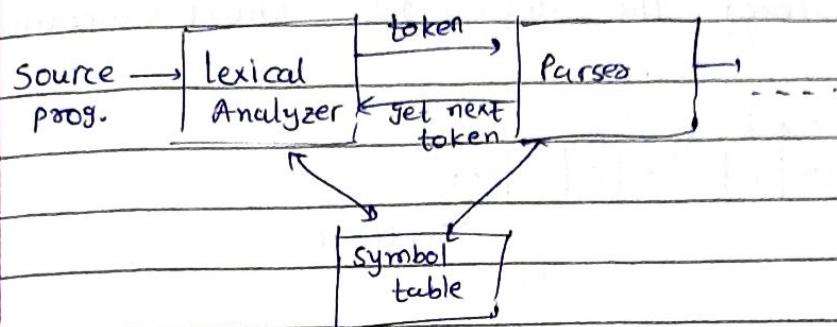
Q-3] Find and explain features of toy lang. "BhaiLang"

Q-4] Explain how JFLAP can be used to design programming lang.

Q-5] Give design of your toy lang.

14/12

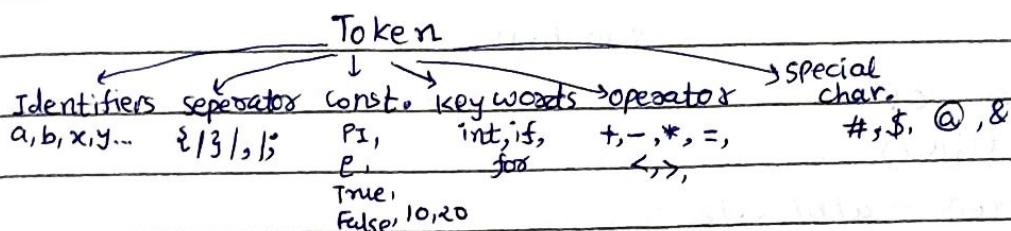
## \* The role of lexical analyzer.



→ How lexical Analyzer works.  
 (will give token)  
 lexeme → Tokens

↓  
 How it works.

↳ int a, b → after each char, it will check for tab, space, new line,  
 so lexeme will stop. (i, n, t, l, )



Q. int main()

= int a=20, b=30;

if(a < b)

return b;

else

return a;

special char =

Identifier = a, b

Separator = {}, (,), ;,

const. = 20, 30

operator = <, ==

key word = int, if, else, main, return

?

→ printf("i=%d, &i=%f", --)

b/w "" is one lexeme

→ It eliminates comments & White space (tab, space, new line)

→ if no relation b/w char, found token and return (a\*b)

→ How it will handle error

1) exceeding length.

2) Unmatched string. → /\* --- \*/

3) illegal char.

## \* Sentinels:-

- buffers - It stores the data for short amount of time

$$\text{Eg. } E = M * C ** 2 .$$

## \* Recognition of token

- 1- Recognition of identifiers,
  - 2- " " delimiter
  - 3- " " relational operators,
  - 4- " " key words
  - 5- " " numbers

## \* Recog. of Identifiers

- letters  $\rightarrow a/b/\dots/z$
  - A/B/ $\dots/z$
  - digit  $\rightarrow 0/1/\dots/9$
  - $RE \rightarrow id \rightarrow \text{letter}(\text{letter/digit})^*$
  - $* \rightarrow \text{delim} \rightarrow \text{blank / tab / newline}$
  - $RE \rightarrow ws \rightarrow \text{delim}^+$

$\alpha_{ij}$  The tokens and associated attribute values for the statement  $E = M * C * ^2$

E Sid swedes -

= <assignment> or blank

assignment op, assignment op, blank op.

1

\* < mult-op, >

$\gamma c <$

1

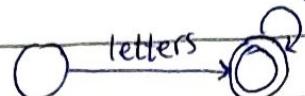
\*\* {exponential\_op, }  
3

2 | 5

2 | < Constant/Num, >

15/12

→ Identifier's FA :-



digits/letters.

start

letters

letters

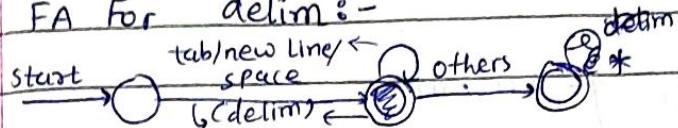
digits  
letters

others

\*

(return gettoken()  
install-i.d.)

→ FA For delim :-



others

delim

\*

\* Recognition of Relational Algebra.

Relop = &lt;/&gt; / &gt;= / &lt;= / == / !=

return (relOp, NE)

&gt; return (relOp, LE)

return (relOp, GE)

Start --&gt; &lt; --&gt; return (relOp, EQ)

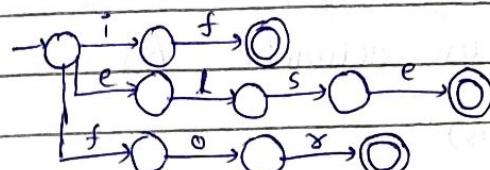
others

return (relOp, LT)

&gt; &gt; --&gt; return (relOp, GT)

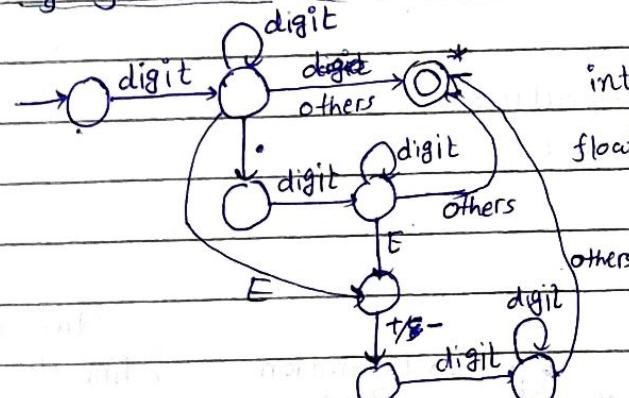
Other --&gt; return (relOp, LT)

\* Recognition of keywords (if, else, for ...)

digit<sup>+</sup> (. digit)\*num → digit<sup>+</sup> (. digit)\*

digit

\* Recognition of Numbers



int = digit other

floating-point = digit . digit other

digit

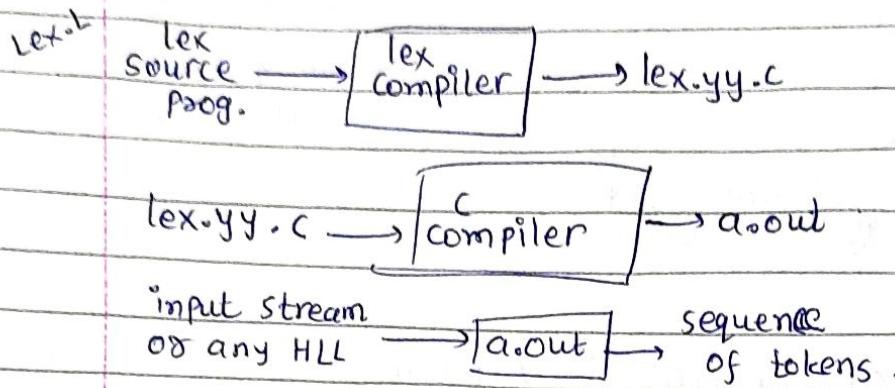
→ num → (digit)<sup>+</sup> (digit)<sup>+</sup> ? if 0 → the not used

1 → used

num → (digit)<sup>+</sup> (. digit)<sup>+</sup> ? (E (+/-))<sup>?</sup> digit<sup>+</sup> ?

16/12

## \* Language For specifying Language Analyzer.



## \* Lex Specifications.

abc	$\wedge a$	$\text{yywrap}() \rightarrow$ it calls by lex for when i/p is over else 0 $\text{yylex}() \rightarrow$ it reads i/p string & generates tokens to the reg. expression
$[a-z]$	$a\$ \rightarrow$ string ends with a	
$[a-z]^*$	$[0-q]^+$	
$[A-Z, a-z]^+$	Pointer to yytext $\rightarrow$ input string	

→ declaration

after declaration we will use delimiters. % %

transition rule → Pattern (actions)  
% %

auxiliary procedures

% %  
# include <stdio.h>  
/\* ----- \*/ ← comment.  
% %

"hi" prints ("How are you?") → Transition Part  
other else like → \* { printf ("Wrong String"); } → Transition Part

% %

(like if-else)

hi → How are you  
will print  
or wrong str  
will be print

→ main()

{

printf("Enter i/p"); } i/p read  
yylex();

}

[h|i|i|EOF]

int yywrap();  
{ return 1; }

} generate token

until

EOF it will return 0, at EOF return 1.

Q. Find lexemes, tokens & patterns of the following que.

	lexemes	token	Pattern
	keyword	keyword	
int	int	int	int
x		identifier	letter(letter/digit)
=		operator	=,<,>
5		constant	constant
;		separator	separator

Q. Find out lexical errors

→ Int x; → Not lexical error.

⇒ lexical error Rules

1. - exceeding length of identifier → int x=124681025262099985;
2. - Appear of illegal char. → printf("G<sup>th</sup>-I"); \$
3. - Unmatched string → String s = "abc / s = abc" ( <sup>not both ""</sup> )
4. - Identifiers start with letters → int 6a = 25
5. - Replacing a char with an incorrect char → int x = 1\$2
6. - White space b/w char → int a b = 25

→ int x,y; → Not error

→ float m = 5..7 → Error → 5

→ int 1sal; → Error → 4

→ int A = #12; → Error → 5

→ String s = success' → Error → 3

→ #include<stdio.h>

main() {

int x=5, y=6;

char \*a;

a=&z;

✓ x = zabi;

printf("%d", x);

letter/digit

}

**eg-1**

start → ① letter → ② other → ③ return (gettoken() install\_id)

int start=0;

token nexttoken() {

while(1) {

switch(state) {

case 0:

c = nextchar();

if(isletter(c)) state=1;

else state=fail;

break;

case 1:

c = nextchar();

if(isletter(c))

state=1;

else if(isdigit(c)) state=2;

else state=2;

break;

case 2:

retract(1);

install\_id();

return gettoken(c);

// retract is a lookahead pointer | char  
// It is used bcoz delimiter is not a  
part of identifier.

// it is used to excess the buffer & it  
called as attribute val.

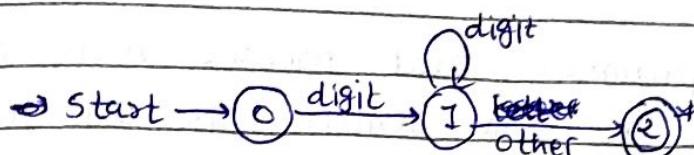
// to obtain a token

}

```

int fail() {
    switch(start) {
        case 0:
            start = 0;
            break;
        default:
            /* compiler error */
    }
    return start;
}

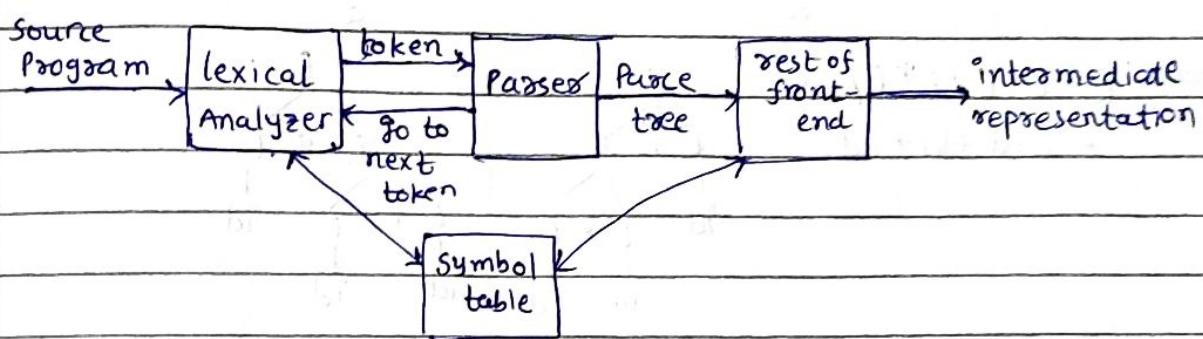
```



21/12

## (2) \* Syntax Analysis

### → Role of the parser



- \* Context free grammar.

- Context free grammar is a four tuples.

$$\rightarrow G = (V, T, P, S)$$

- $V$  = finite set of symbols called as non terminal or variable.  
(It may be in uppercase)  
usually  $S$  is used for starting symbol. Some italic names are also allowed.

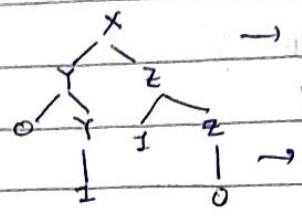
eg.  $\exp \text{ stmt}$

$$S \rightarrow OA$$

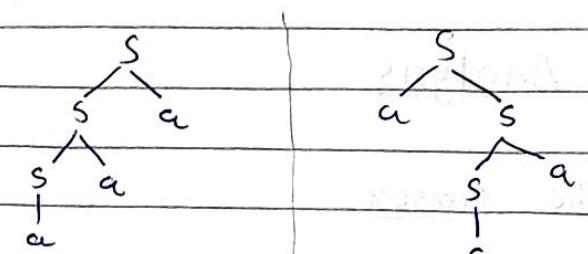
$$A \rightarrow OB$$

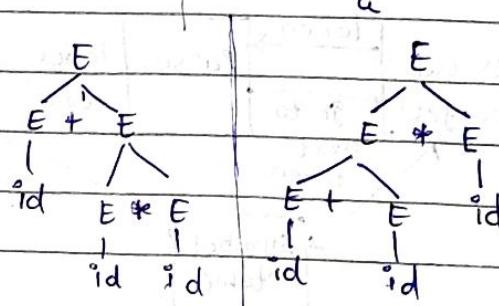
$$B \rightarrow E$$

- $T = \{ \text{set of symbols that are called as terminal symbols} \}$   
 Lowercase letters, operator symbols  $\rightarrow +/-$ , punctuation symbols  
 digits, bold phase string  $\rightarrow \text{id, if}$
- $P = \{ \text{set of productions} \} (A \rightarrow B)$
- $S = \{ \text{it is member of } V, \text{ called as start symbol.} \}$

$\Rightarrow X \rightarrow YZ$        $w: 0110$       
 → No other parse tree is possible for  $w$   
 $Y \rightarrow 0Y/1$   
 $Z \rightarrow 1Z/0$        $\rightarrow$  So it is non-ambiguous.

→ Ambiguity = A grammar that produce more than one parse tree for same sentence is said to be ambiguous grammar.

$\Rightarrow S \rightarrow Sa/aS/a$        $w: aaa$       

$\Rightarrow E \rightarrow E+E/E^*E/\text{id}$        $w: \text{id id}^*$       

Rules:- (to remove ambiguity)

→ \$ and ↑ is right associated, remaining all are left associative.

→ Priority: \$ < # < @

→ < and < not.

→ - < + < \* < / < ↑ < \$  
 + < \* < ()\*

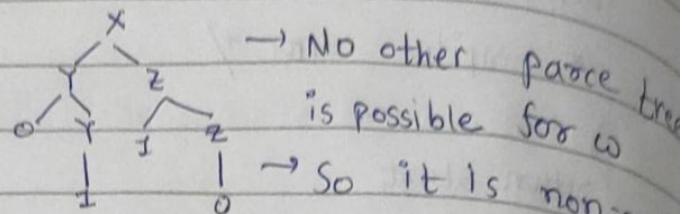
- $T = \text{set of symbol that are called as terminal}$   
(lowercase letters, operator symbols  $\rightarrow +/-$ , punctuation symbols, digits, bold phase string,  $\rightarrow \text{id, if}$ )
- $P = \text{set of productions } (\alpha \rightarrow \beta)$
- $S = \text{It is member of } V, \text{ called as start symbol.}$

$$\Rightarrow X \rightarrow YZ$$

$$Y \rightarrow 0Y/1$$

$$Z \rightarrow 1Z/0$$

w: 0110

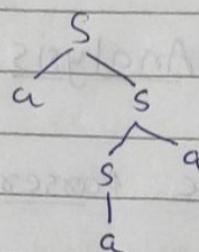
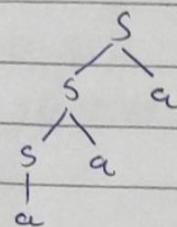


→ No other parse tree is possible for w  
→ So it is non-ambiguous

→ Ambiguity: A grammar that produce more than one parse tree for same sentence is said to be ambiguous grammar.

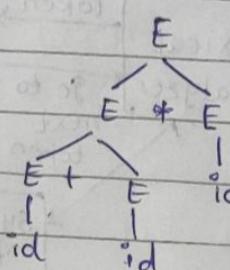
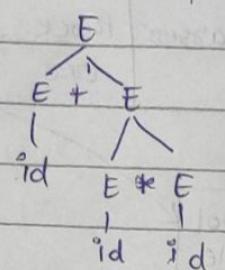
$$\Rightarrow S \rightarrow Sa/as/a$$

w: aaa



$$\Rightarrow E \rightarrow E+E/E^*E/\text{id}$$

w: id + id \* id



Rules:- (to remove ambiguity)

→ \$ and ↑ is right associated, remaining all are left associative

→ Priority: \$ < # < @

→ < and < not.

→ - < + < \* < / < ↑ < \$  
+ < \* < ()\*

eg.  $E \rightarrow E+E / E^*E / id \rightarrow E \rightarrow E+T / T$   
 $T \rightarrow ET^*F / F$   
 $F \rightarrow id$

$E \rightarrow E^*E / E+E / id \rightarrow E \rightarrow E+F / F$  (+ has low priority so, start with  $E+E$ )  
 $F \rightarrow F^*T / T$   
 $T \rightarrow id$

eg.  $E \rightarrow E\$E / id \rightarrow E \rightarrow E\$T / T \xrightarrow{\text{so }} E \rightarrow T\$E / T$  (\$ is right associative so T will go to left)  
 $T \rightarrow id$

eg.  $A \rightarrow A\$A / A\#A / A@A / d \rightarrow A \rightarrow T\$A / T$   
 $T \rightarrow T\#P / P$   
 $P \rightarrow P@N / N$   
 $N \rightarrow d$

22-12

①  $bExp \rightarrow bExp \text{ or } bExp / bExp'$  and  $bExp / \text{NOT } bExp / \text{True} / \text{False}$

②  $R \rightarrow R+R / RR / R^*/a/b/c \checkmark$  Ans  $\rightarrow R \rightarrow R+M/M$   
 $E \rightarrow E+T/T \times M \rightarrow MS/S$   
 $T \rightarrow TF/F \times S \rightarrow S^*/a/b/c$   
 $F \rightarrow F^*/a/b/c \times E \rightarrow$   
 $T \rightarrow$

Ans J  $\rightarrow bExp \rightarrow bExp \text{ or } T/T$

$T \rightarrow T \text{ and } F/F$

$F \rightarrow \text{not } M/M$

$M \rightarrow \text{True} / \text{False}$

## \* Left Recursion

→ Algorithm

input will take some grammar G, then o/p equivalent grammar with no recursion.

→ Method

If we have left recursive pair of production

$A \rightarrow A\alpha / B$ , where B doesn't begin with A

then we can eliminate left recursion by

replacing this pair of production with

$$A \rightarrow BA'$$

$$A' \rightarrow \alpha A' / \epsilon$$

e.g.

$$A \rightarrow A\alpha / B/\gamma$$

$$\Rightarrow A \rightarrow BA' / \gamma A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

e.g.

There are 2 type of left recursion

1. direct left recursion

2. indirect left recursion

I. e.g.

$$A \rightarrow ABd / Aa/a$$

$$B \rightarrow Bcb$$

→ No ambiguity

Ans:

$$A \rightarrow \cancel{Bcb} aA'$$

$$A' \rightarrow aA' / BcbA' / \epsilon$$

$$B \rightarrow bB'$$

$$B' \rightarrow \epsilon B' / \epsilon$$

23/12

e.g.  $A \rightarrow Aab/\alpha$

$$B \rightarrow Bcb/y$$

$$C \rightarrow \epsilon c / \epsilon$$

→ No ambiguity.

Ans:  $A \rightarrow xA'$

$$A' \rightarrow aBA' / \epsilon$$

$$B \rightarrow yB'$$

$$B' \rightarrow cbB' / \epsilon$$

$$C \rightarrow cc^*/\epsilon$$

$$\left. \begin{array}{l} C \rightarrow c \\ C \rightarrow cc' \end{array} \right\} C \rightarrow cc'/\epsilon$$

\*

Indirect Left recursion

$$S \rightarrow Aa/b$$

$$A \rightarrow Ac / Sd / \epsilon$$

→

$$A \rightarrow Ac / Aad / bd / \epsilon$$

$$A \rightarrow bdA' / \epsilon$$

$$A' \rightarrow CA' / adA' / \epsilon$$

Algo. for removing.

will

\*

Left factory - I/p: Take grammar g.

O/p: Equivalent non-left factory grammar

method: For each non-terminal A find the prefix  $\alpha$

Common two or more of its ~~one~~ alternatives. Replace one of the A production

$\rightsquigarrow A \rightarrow \alpha B_1 / \alpha B_2 / \alpha B_3 / \dots / \alpha B_m / y$

- where  $y$  represents all alternatives that do not begin with  $\alpha$

$A \rightarrow \alpha A' / y$

$A' \rightarrow B_1 / B_2 / B_3 / \dots / B_n$

- Repeatedly apply, that transformation until those two alternatives for a common prefix

eg.  $S \rightarrow iEtS / iEtSe / a$

$E \rightarrow b$

Ans  $\rightarrow S \rightarrow iEtS' / a$

$S' \rightarrow \epsilon / es$

$E \rightarrow b$

eg.  $A \rightarrow xByA / xByA_2A / a$

$B \rightarrow b$

Ans  $\rightarrow A \rightarrow xByA_2A / a$

$A_2 \rightarrow \epsilon / zA$

$B \rightarrow b$

eg.  $A \rightarrow aAb / aA' a$

Ans  $\rightarrow A \rightarrow aA'$

~~$a \rightarrow aE$~~

~~$a \rightarrow aE$~~

$A' \rightarrow Ab / A' E$

$A' \rightarrow AA'' / \epsilon$

$A'' \rightarrow b / \epsilon$

eg.  $A \rightarrow ad / a / ab / abc / b$

Ans  $\rightarrow A \rightarrow aA'$

$A' \rightarrow d / \epsilon / b / bc$

$A' \rightarrow bB' / d / \epsilon$

$B' \rightarrow E / c$

eg. Ambiguity  $\rightarrow$  left recursion  $\rightarrow$  left factory

$E \rightarrow E + E / E^* E / id$

$E \rightarrow E + T / T$

$T \rightarrow T^* F / F$

$F \rightarrow id$

$E \rightarrow * TE'$

$E' \rightarrow * TE' / \epsilon$

$T \rightarrow PT'$

$T' \rightarrow * PT' / \epsilon$

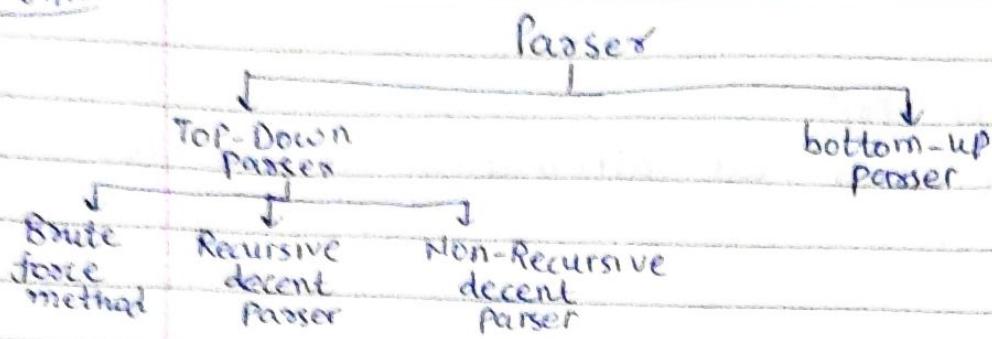
$F \rightarrow id$

$A \rightarrow A\alpha / \beta$

$A \rightarrow \beta A'$

$A' \rightarrow \alpha A' / \epsilon$

27/11

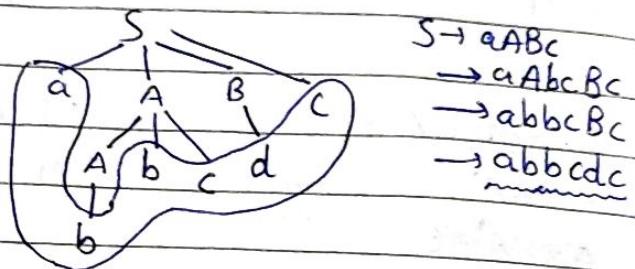


### \* Top-Down Parser.

Generates a parse tree for the given input stream with the help of grammar productions by expanding the non-terminals.

- In this, if we start from start symbol & ends on the terminal. It uses left most derivation.

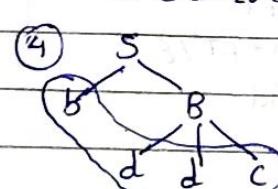
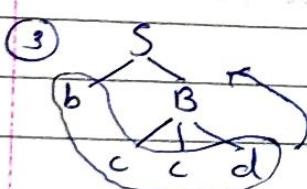
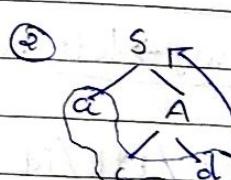
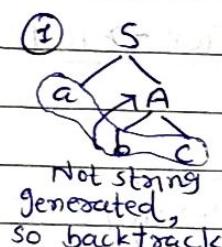
e.g.  $S \rightarrow aABC$   
 $A \rightarrow Abc/b$   
 $B \rightarrow d$   
 string = abbcde



- at every point we have to design decide what is the next production we should use

### \* Brute force Parser.

e.g.  $S \rightarrow aA/bB$   
 $A \rightarrow bc/cd$   
 $B \rightarrow ccd/ddc$   
 string = bdddc



drawback  $\rightarrow$  Time complexity  $2^n$

## \* Recursive decent Parser

eg.  $S \rightarrow PAy$



in Efficient

eg.  $E \rightarrow iE'$

$E' \rightarrow +^i E'/\epsilon$

(1) EC)

(2) E'c)

(3) match(token t)

→ eg.  $i + i \$$

→ EC){

if(lookahead == 'i') {

match('i');

; E'();

if(lookahead == '\$') {

point("success");

}

else

return;

}

eg.  $E \rightarrow b$

→ Ambiguous

$S \rightarrow iEts / iEtSeS / a$

$S \rightarrow iEts \rightarrow iEt^i Et^i SeS$

→ E'c){

if(lookahead == '+') {

match('+');

; E'();

if(lookahead == 'i') {

match('i');

; E'();

else return;

}

eg.  $S \rightarrow ABC$

$A \rightarrow oAi/\epsilon$

$B \rightarrow lBi/\epsilon$

$C \rightarrow lCo/\epsilon$

SC(), AC(), BC(), CC(), match(token)