

# Compiler Design

PAGE NO.	
DATE	/ /

Compiler :- Executes whole code in one time

↳ C, C++, ... → Languages for software development

Interpreter :- line by line execution

↳ Small Code Use → Web based Language

→ Compiler execute whole program at a time.

↳ Large Code Use

Ex. C, C++, Java, Pascal, COBOL, BASIC,

↳ Programming Language

\* Interpreter example :- Ruby, Python, ASP.NET, Perl

↓  
Scripting language

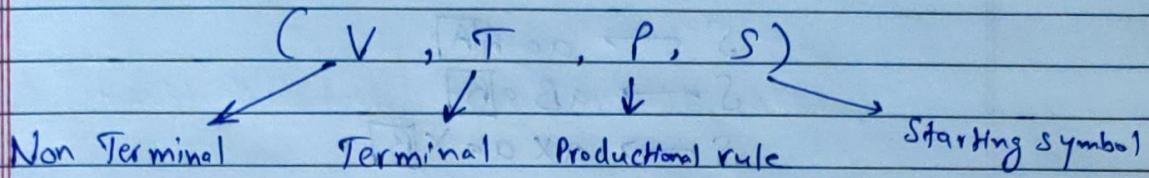
COBOL → Common Business oriented Language

BASIC → Beginner All Purpose symbolic Instruction Code.

\* Markup Language → HTML, XML → No use  
they just help

in ~~Interpre~~  
with scripting language

\* Grammar :- 4 Tuple



① Terminal :- Small letter  
 $(a-z), (0-9)$

② Non Terminal :-  
 $(A-Z)$

③ Production Rule :-  $S \rightarrow abc$   
 ↓  
 P.R.

④ Starting Symbol :-  $S \rightarrow abc$   
 ↓  
 S.S.

# Regular Grammar (R.G.) → Finite length.  
 Least Powerful

↳ 4 Tuple  $(V, T, P, S)$

↳ Two form ↳ left linear Regular Grammar \*\*  
 ↳ Right linear Regular Grammar \*\*

① left linear Regular Grammar :-

$S \rightarrow Aaa$  ↳ Non Terminal on left  
 $S \rightarrow Bba$   
 $S \rightarrow Xaaa$

(2)

## Right linear Regular Grammar I -

$$S \rightarrow aa A[A]$$

$$S \rightarrow aB a[B]$$

$$S \rightarrow ax a = X[B]$$

★

$$S \rightarrow [A]aa[A]$$

$$S \rightarrow [BB]$$

$$S \rightarrow [Z]XAB[B]$$

$$S \rightarrow [A]aXa[B]$$

left, Right both

Q.

$$\begin{array}{l} S \rightarrow aa \\ S \rightarrow \epsilon \\ S \rightarrow aaaa \end{array} \quad \begin{array}{c} N.T. \\ X \\ X \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} R.G. \quad \begin{array}{l} X \\ X \end{array}$$

★

If it is R.G., then it will be C.F.G., C.R.S.G.,  
Unrestricted grammar

Q. ①

$$\begin{array}{l} S \rightarrow aAA \\ S \rightarrow aab \\ A \rightarrow a \end{array} \quad \begin{array}{l} \text{Right linear grammar} \\ \text{R.G. } \checkmark \end{array}$$

ii

$$\begin{array}{l} S \rightarrow aaAa \\ S \rightarrow A[A] \\ S \rightarrow B \\ S \rightarrow a \\ B \rightarrow b \end{array} \quad \begin{array}{l} \text{Both linear grammar} \\ \text{Regular Grammar} \end{array}$$

(i)  $S \rightarrow aAa$   
 $S \rightarrow aBa$

R.G. X

(ii) (i)  $A \rightarrow aaAa$   $\boxed{A}$  } R.G. ✓ (Right Linear)  
 $B \rightarrow bb$

(ii)  $B \rightarrow b$   
 $A \rightarrow aAa$   
 $A \rightarrow aa$

} R.G. X

(iii)  $A \rightarrow aAA$   
 $B \rightarrow BB\epsilon$   
 $\downarrow$   
Empty

} Both Linear  
} R.G. ✓

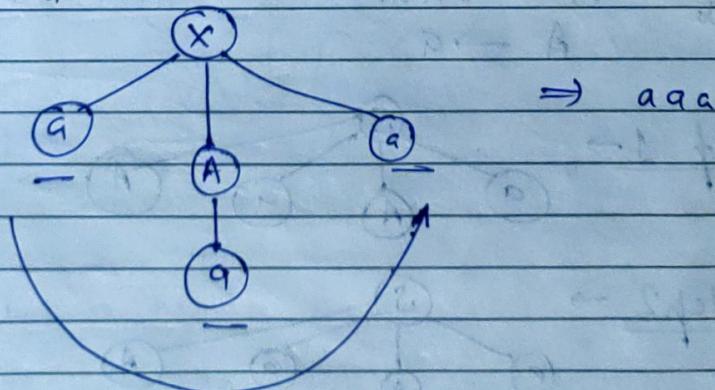
(i) Derivation Tree

Starting Node  $\rightarrow$  Root Node

Vertex  $\rightarrow$  Non Terminal

leaf Node (leaves)  $\rightarrow$  Terminal

$X \rightarrow aAa$   
 $A \rightarrow a$

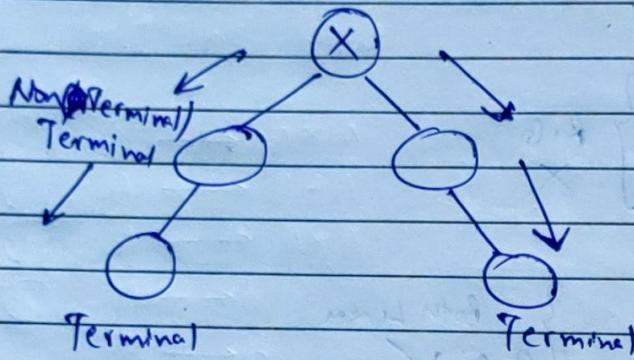


(i)

Top-down derivation tree :-

i) Start with Root Node.

ii) Goes down toward Leaf Node (leaves)



(ii)

Bottom-up derivation tree :-

i) Start with leaf Node (leaves Node)

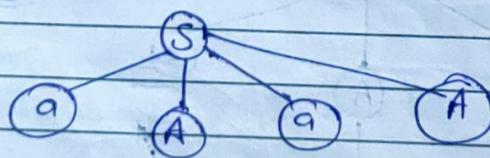
ii) Towards Root Node.

(i)

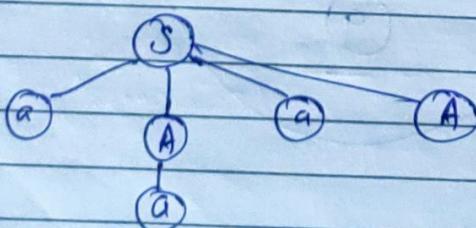
left linear derivation tree :-

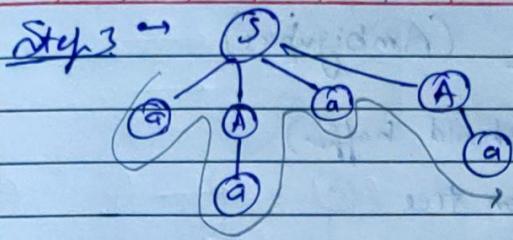
$$\begin{array}{l} \text{Root} \\ \text{Node} \end{array} \leftarrow S \rightarrow aAaA \\ A \rightarrow a$$

Step 1 →



Step 2 →





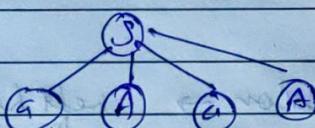
$$S \rightarrow aaaa$$

(ii) Right linear derivation tree :- (Right to left)

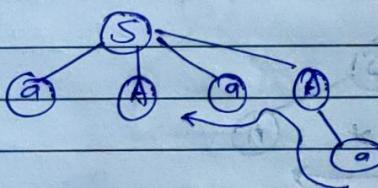
$$S \rightarrow aAaA$$

~~$A \rightarrow a$~~

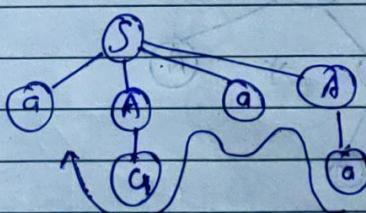
Step 1 :-



Step 2 :-



Step 3 :-



$$S \rightarrow aaaa$$

#

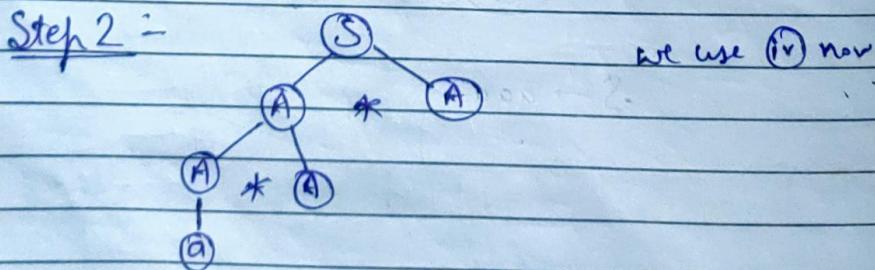
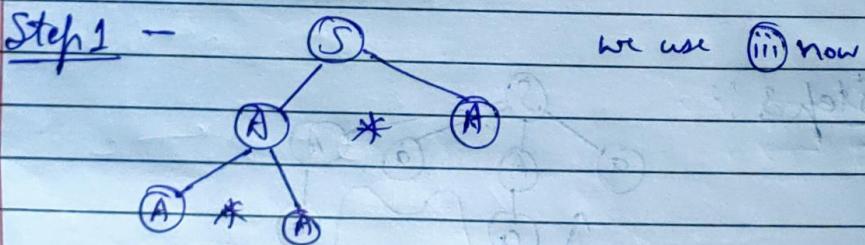
## Ambiguity (Grammar) (Ambigibus)

(Any one should happen)

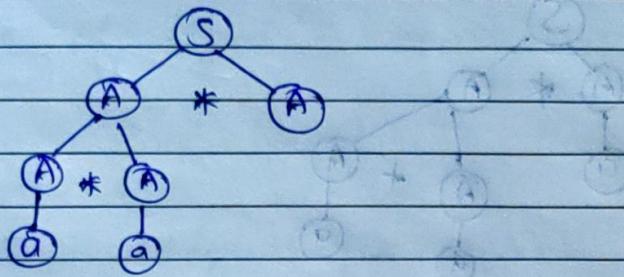
- (i) left linear derivation tree
- (ii) Right " " " "
- (iii) left and right " " " "

- Q. i)  $S \rightarrow A * A$   
 ii)  $A \rightarrow A + A$   
 iii)  $A \rightarrow A * A$   
 iv)  $A \rightarrow a$   
 v)  $B \rightarrow b$

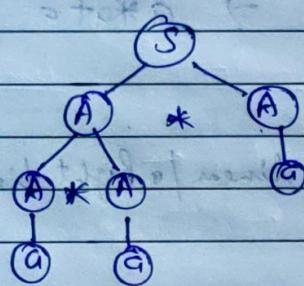
Take any expression  $\rightarrow$  left linear  $\rightarrow$   
 we take (i)



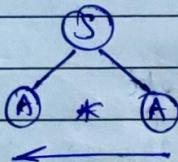
Step 3 :-



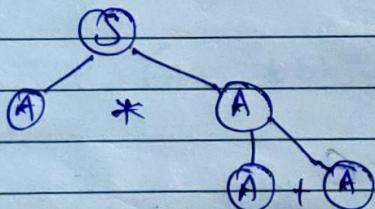
Step 4 :-



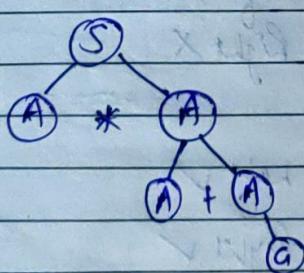
Right Linear :-



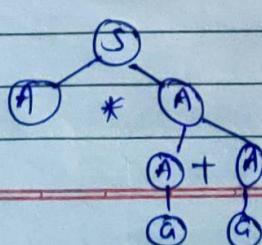
Step 1 :-



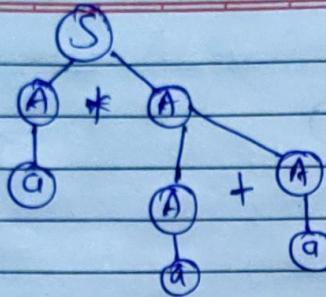
Step 2 :-



Step 3 :-



Step 4:-



$\Rightarrow a * a + a$

Q. Find out left linear / Right Linear Regular Grammar.

$$① \quad S \rightarrow a = A B$$

$$S \rightarrow a A B$$

Right

$$A \rightarrow a$$

$$B \rightarrow b$$

$$② \quad S \rightarrow A a a \text{ left}$$

$$S \rightarrow B b$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$③ \quad S \rightarrow E + E$$

left X

right X

$$④ \quad S \rightarrow A a a A$$

$$S \rightarrow B A$$

$$A \rightarrow a$$

$$B \rightarrow b$$

left ✓

right ✓

#

## Power of Sigma :-

$$\Sigma = \{a, b\}$$

i)  $\Sigma^0 \rightarrow$  only Empty string  
 $\Sigma^0 \rightarrow$  Empty

ii)  $\Sigma^1 \rightarrow$  Set of all string of length 1

$$\Sigma^1 \rightarrow (a, b)$$

iii)  $\Sigma^2 \rightarrow$  Set of all string of length 2

$$\Sigma^2 \rightarrow (aa, ab, ba, bb)$$

iv)  $\Sigma^3 \rightarrow$  Set of all string of length 3

$$\Sigma^3 \rightarrow (aaa, bbb, abb, baa, bab, aba, aab, bba)$$

v)  $\Sigma^4 \rightarrow$  Set of all string of length 4

$$\Sigma^4 \rightarrow (aaaa, \dots)$$

vi)  $\Sigma^n \rightarrow$  Set of all string of length n

# ① Kleene Closure (\*)

$$\Sigma = (a, b)$$

$$\Sigma^* = (\text{Empty}, a, b, aa, ab, \dots)$$

② Kleene Plus (+)

$$\Sigma = (a, b)$$

$$\boxed{\Sigma^+ = \Sigma^* - (\text{Empty})}$$

$$\Sigma^+ = (a, b, aa, bb, \dots)$$

# Regular Language :- (L)

Language accepted by Regular grammar is called Regular Language.

$$\Sigma = (a, b)$$

$$a^* = (\text{Empty}, a, aa, aaa, \dots)$$

$$a^+ = (a, aa, aaa, \dots)$$

$$\Sigma = (a, b)$$

$$L^1 = (a, b)$$

$$L^2 = (aa, bb, ab, ba)$$

$$L^3 = (aaa, aab, abb, bbb, \dots)$$

$$L^4 = (aaaa, aaab, aabb, abbb, \dots)$$

$$\Sigma = (a, b, c, d)$$

$$L' = (a, b, c, d)$$

$$L^2 = (ab, ac, ad, ba, bd, \dots)$$

$$L^3 = (abc, bcd, \dots)$$

### \* Five Properties of Regular Language :-

(1) Union - if we have 2 Regular language  $L_1$  and  $L_2$ . Then Union of  $L_1, UL_2$  are Regular Language.

Ex:  $L_1 = (a, b)$

$L_2 = (a, c)$

$L_1 \cup L_2 = (a, b, c)$

(2) Complement - if we have 2 Regular language  $L_1$  and  $L_2$ . Then complement of  $L_1$  and  $L_2$  are Regular Language.

$L_1 = (a, b)$

$L_2 = (a, c)$

$L_1' = (b, a)$

$L_2' = (c, a)$

(3) Concatenation - if we have 2 regular language  $L_1$  and  $L_2$ . Then concatenation of  $L_1, L_2$  are regular language.

$L_1 = (a, b)$

$L_2 = (a, c)$

$L_1 \cdot L_2 = (aa, ac, ba, bc)$

(4)

Kleene Closure (\*) :- if we have 2 Regular language  $L_1$  and  $L_2$ . Then Kleene closure of  $L_1^*$  and  $L_2^*$  are Regular language.

$$L_1 = (a, b)$$

$$L_2 = (a, c)$$

$$L_1^* = (\text{Empty}, a, b, aa, bb, ab, ba, \dots)$$

$$L_2^* = (\text{Empty}, a, c, aa, ac, cc, ca, \dots)$$

(5)

Intersection ( $\cap$ ) :- if we have 2 Regular language  $L_1$  and  $L_2$ . Then intersection of  $L_1 \cap L_2$  are regular language.

$$L_1 = (a, b)$$

$$L_2 = (a, c)$$

$$L_1 \cap L_2 = (a)$$



Precedence  $\rightarrow$

(1) Kleene Closure (\*) Highest Priority

(2) Second Concatenation

(3) Union (lowest)

## Normal Form

$\nwarrow$  CNF (Chomsky Normal Form)

$\searrow$  GNF (Greibach Normal form)

### 1) CNF :-

Conditions -

- (i) Non-Terminal Represent to empty symbol.

$$A \rightarrow \epsilon$$

- (ii) Non-Terminal Represent to Terminal.

$$A \rightarrow b, A \rightarrow c$$

- (iii) Non-Terminal represent atleast to 2 Non Terminal  
( $V \rightarrow V_1V_2$ )

$$A \rightarrow BC, A \rightarrow AB, A \rightarrow BD$$

### 2) GNF :-

Conditions -

- (i) Non-Terminal Represent to empty symbol.

$$A \rightarrow \epsilon$$

- (ii) Non-Terminal Represent to Terminal

$$A \rightarrow b, A \rightarrow c$$

(iii)

Non-Terminal Represent Terminal followed by Non Terminals.  $(a\alpha^*)$

$$A \rightarrow aB, A \rightarrow aBC, A \rightarrow aBCD, A \rightarrow aABCD$$



Only one terminal followed by terminals

(i)

$$A \rightarrow d \checkmark$$

$$A \rightarrow \epsilon \checkmark$$

$$A \rightarrow Aa X$$

$$A \rightarrow a \checkmark$$

CNF X

GNF X

Regular Grammar  
(VTPS)  $\checkmark$

(ii)

$$A \rightarrow a \checkmark$$

$$A \rightarrow \epsilon \checkmark$$

$$A \rightarrow BC \checkmark$$

CNF  $\checkmark$ 

GNF X

(iii)

$$A \rightarrow a \checkmark$$

$$A \rightarrow b \checkmark$$

$$A \rightarrow aaA X$$

$$A \rightarrow \emptyset \checkmark$$

CNF X

GNF X

(iv)

$$A \rightarrow a \checkmark$$

$$A \rightarrow \epsilon \checkmark$$

$$A \rightarrow aBCDEF GHIJ \checkmark$$

CNF X

GNF  $\checkmark$ 

(v)

$$A \rightarrow a \checkmark$$

$$A \rightarrow \epsilon \checkmark$$

$$A \rightarrow BC \checkmark$$

CNF  $\checkmark$ 

GNF X

(vi)

$$A \rightarrow S \quad \checkmark \quad CNR X$$

$$A \rightarrow b \quad \checkmark \quad G.N.P. \quad \checkmark$$

$$A \rightarrow aBC \quad \checkmark$$

## # CFL (Context free Language) :-

### Properties :-

- (i) Union :- if we have 2 Context free language  $L_1$  and  $L_2$  then union of ~~Context free~~  $L_1 \cup L_2$  are CFL.
- (ii) Complement :- if we have 2 Context free language  $L_1$  and  $L_2$  then complement of  $L_1$  and  $L_2$  are ~~CFL~~ NOT CFL.
- (iii) Concatenation :- if we have 2 CFLs  $L_1$  and  $L_2$  then concatenation of  $L_1 \cdot L_2$  are CFL.
- (iv) Kleene Closure (\*) :- if we have 2 CFLs  $L_1$  and  $L_2$  then Kleene Closure  $L_1^*$  and  $L_2^*$  are CFL.
- (v) Intersection :- if we have 2 CFLs  $L_1$  and  $L_2$  then intersection of  $L_1 \cap L_2$  are NOT CFL.

(Parse Tree)

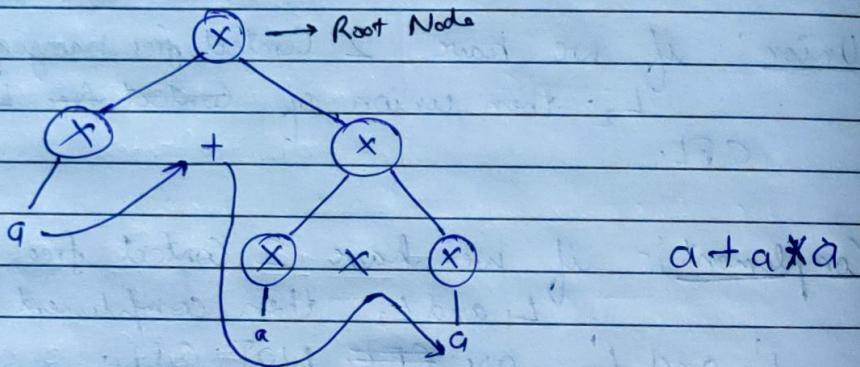
- Q. Find out Left Most Derivation tree  
 Final string "a + a \* a"

$$X \rightarrow X + X \quad - \textcircled{1}$$

$$X \rightarrow X$$

$$X \rightarrow X \times X \quad - \textcircled{3}$$

$$X \rightarrow a$$



Right

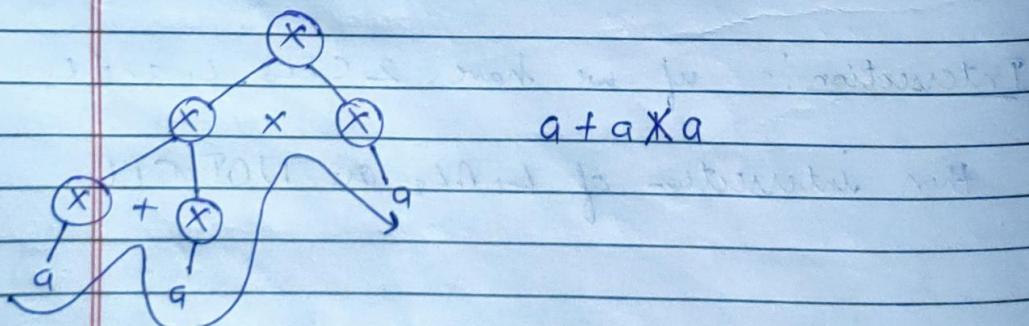
- Q. ~~Left~~ Most Derivation tree → <sup>(Parse)</sup><sub>tree</sub>

$$X \rightarrow X + X \quad - \textcircled{3}$$

$$X \rightarrow X$$

$$X \rightarrow X \times X \quad - \textcircled{1}$$

$$X \rightarrow a$$



\* If more than one parse tree are possible, like in this example left & right derivation tree are

different), therefore this is Ambiguous Grammar.

\* If the structure of left and Right derivation tree are same then grammar is Unambiguous grammar.

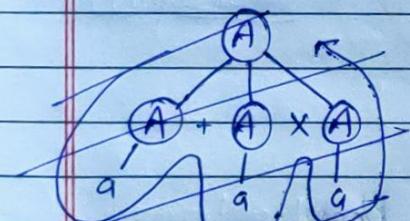
Q. Find out left most derivation tree and right.

$$A \rightarrow A + A \cdot X A \Rightarrow A \rightarrow A + (AXA)$$

$$A \rightarrow AXA + A$$

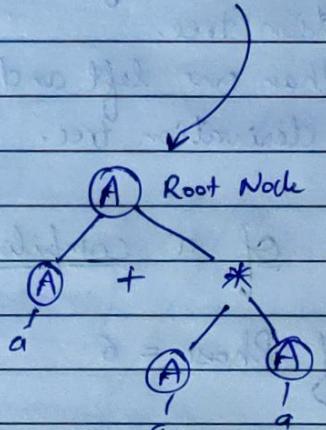
$$A \rightarrow a$$

left most  $\rightarrow$



$a + a \cdot X a$

Sahi tha yeh bhi



diff. structure & diff. strings  $\rightarrow$  Ambiguous

\* For Unambiguous -

(i) ek se jyada left ban jaayen

(ii) " " " right " "

(iii) ek se jyada dono ban jaayen

### Ambiguous

- (i) More than one left and right derivation tree.
- (ii) Different structure of Derivation tree.
- (iii)
  - (1) More than one left derivation tree.
  - (2) More than one right derivation tree.
  - (3) More than one left and right derivation tree.

### Unambiguous

- (i) Only one left and one right derivation tree.
- (ii) Same structure of derivation tree.

$A \rightarrow A + A \leftarrow A$   
 $A + A \rightarrow A \leftarrow A$   
 $\vdots \leftarrow A$

## # Phases of a compiler -

\* No. of Phase = 6

\* First 3 Phase called Analysis Phase.

\* Last 3 Phase called Synthesis Phase.

### Analysis Phase -

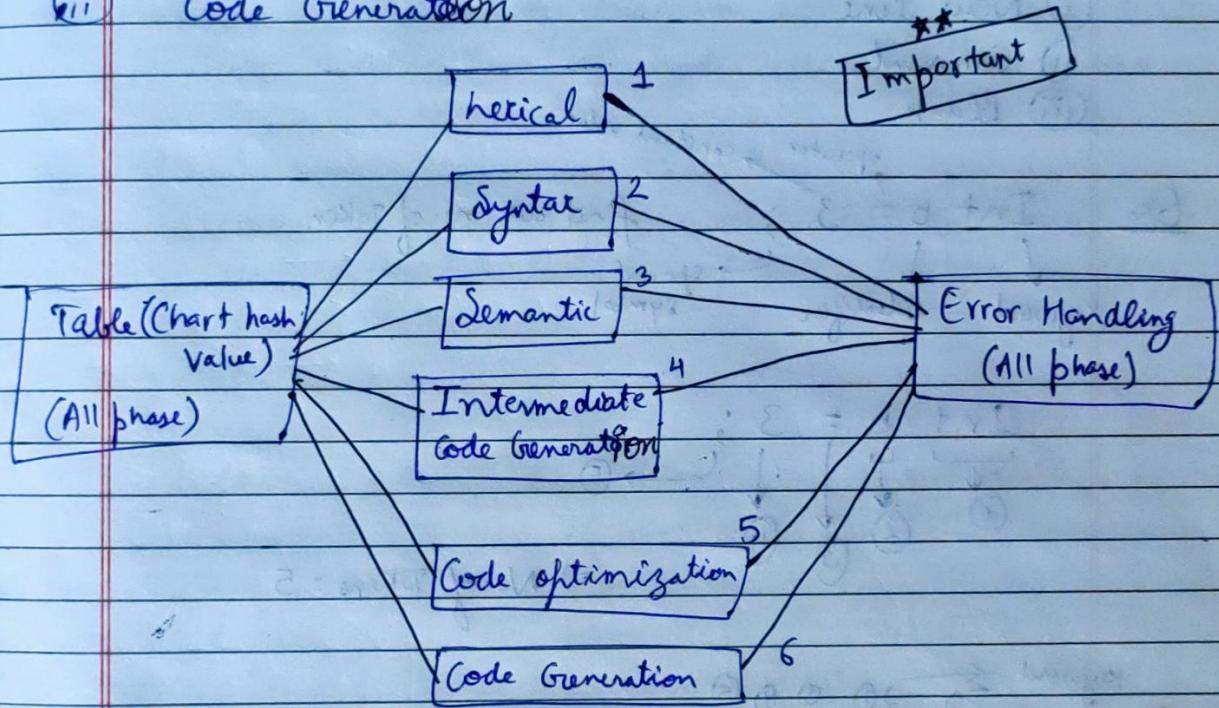
- (i) Lexical Phase
- (ii) Syntax Phase
- (iii) Semantic Phase

### Synthesis Phase -

- (i) Intermediate Code Generation

## ii) Code optimization

## iii) Code Generation



★ Token :- Sequence of character which represent unit of information in the code.

- i) Identifier :- (a, b)
- ii) Keyword → INT, VOID
- iii) Constant value → 5, 6, 0, 1
- iv) Special symbol → , ;
- v) operator → Bitwise, Arithmetic



## Non-Tokens :-

i) New line

ii) Comment,

iii) Blank

Ex:-

Int b = 3 ; find out no. of Token

Int b = 3 ;  
 ↓      ↓      ↓      ↓      ↓  
 ①    ②    ③    ④    ⑤  
 ⇒ No. of Tokens = 5

Ex:-

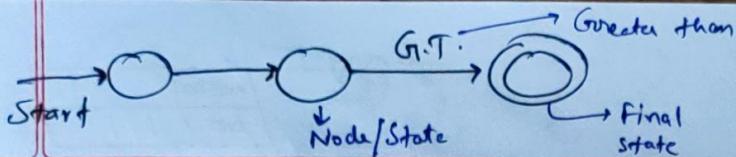
Int main () {

⑥    ⑦    ⑧    ⑨    ⑩  
 printf (" My Name is \_\_\_\_\_ ");  
 ⑪ return 0;  
 }  
 ⑫    ⑬  
 ⑭ → ⑭

1,2,3 ← entry tokens

- which line

which unit ← which

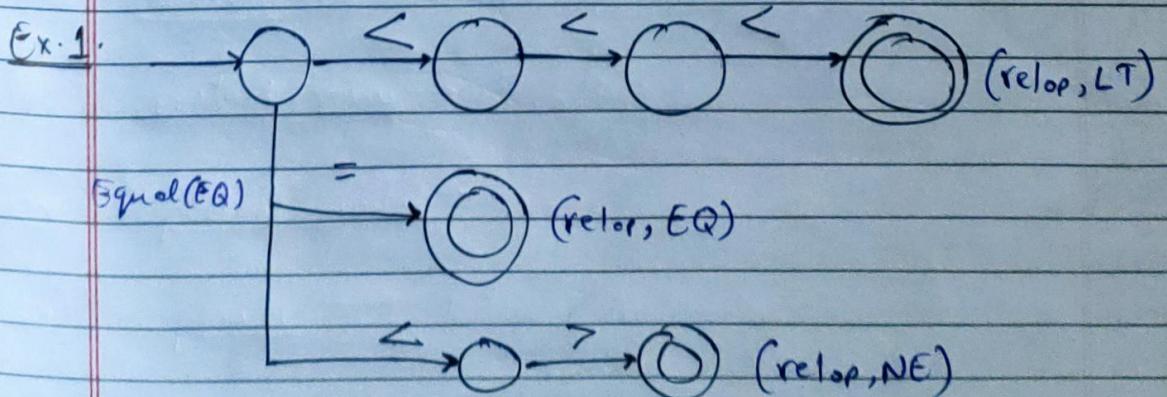


PAGE NO.	
DATE	/ /

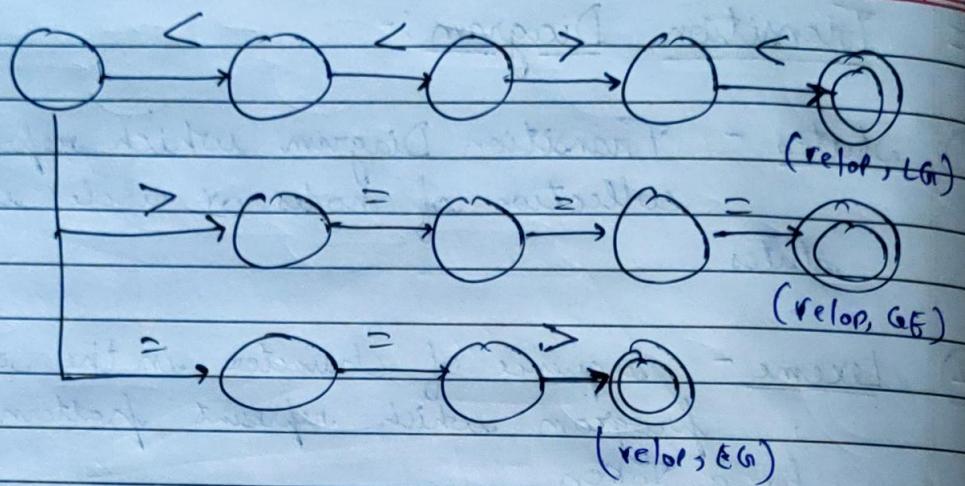
## # Transition Diagram :-

- ① States - Transition Diagram which represent collection of nodes or circle is called states.
- ② Lexeme - Sequence of character in the source program which represent pattern matching for a Token.

lexeme	Token Name	Attribute
if	if	—
Then	Then	—
else	else	—
Any Identifier	Id	—
Any Number	Number(Num)	—
<	relOp	LT
>	relOp	GT
=	relOp	EQ
<=	relOp	LE
>=	relOp	GE
<u>&lt;&gt;</u>	relOp	NE (Not Equal) <u>&lt;&gt;</u>



Ex-2.



① ( $LT$ ,  $GT$ )

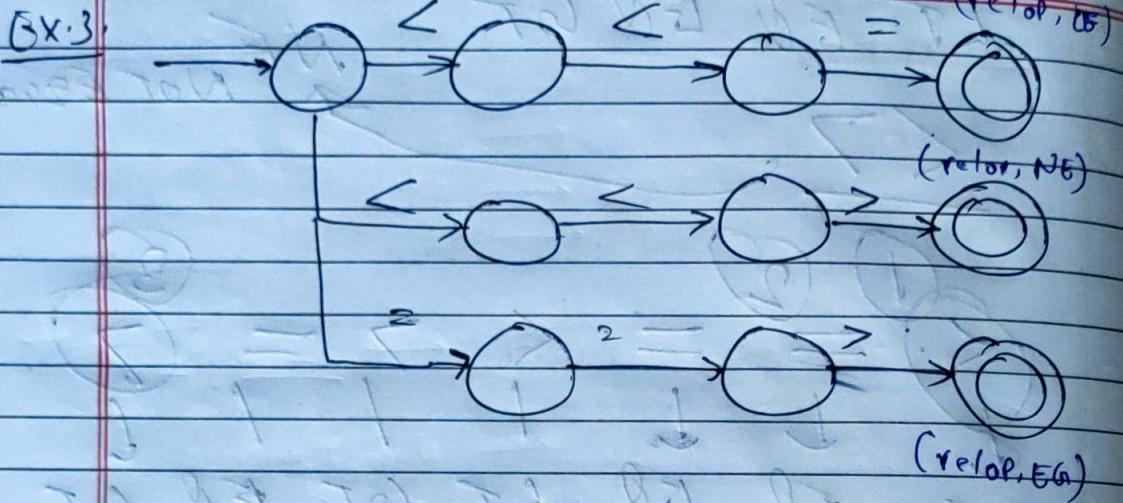
$\downarrow \quad \downarrow$   
~~(LG)~~ (NE)

② ( $GT$ ,  $EQ$ )

$\downarrow \quad \downarrow$   
(GE)

③ ( $EQ$ ,  $GT$ )

(EG)

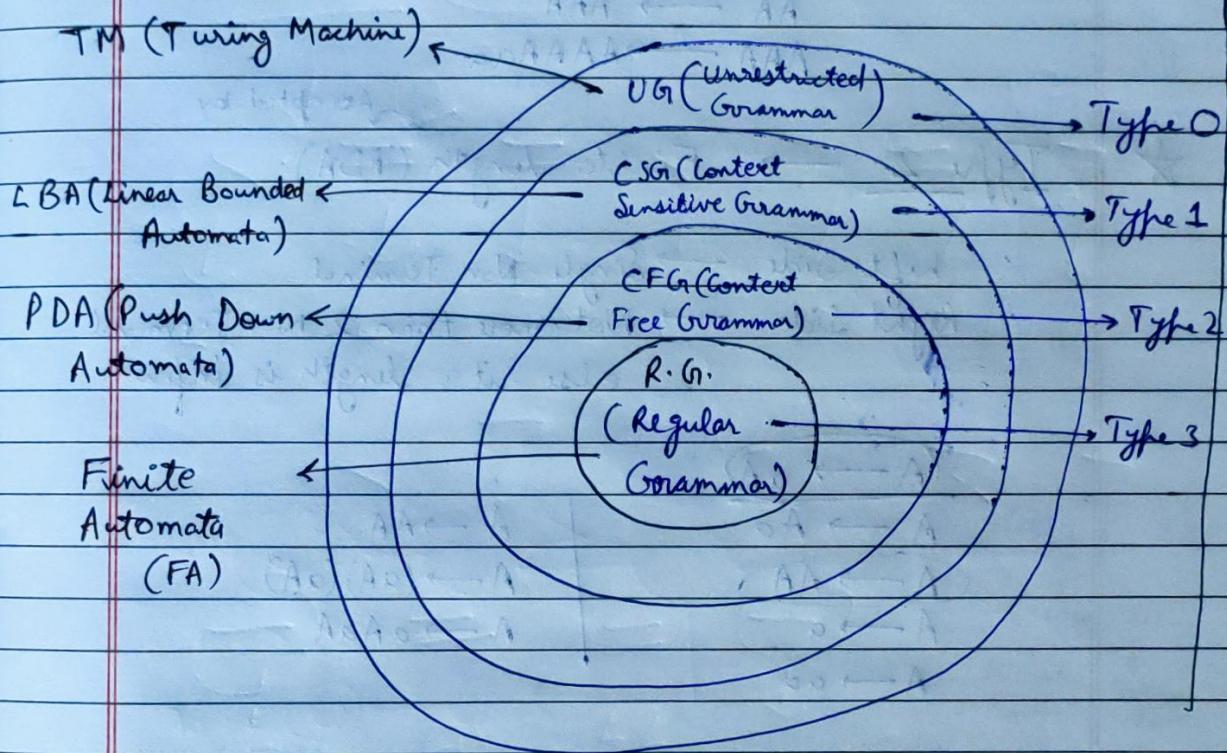


① (LT, EQ)  
 (LE)

② (LT, GT)  
 (EG) ↓  
 N6

③ (EQ, GT)  
 (EG)

More Complex  
More Powerful



\* Type 0 → Infinite length (TM)

Example →

$$\begin{array}{l}
 A \rightarrow \text{AAAA} \\
 AA \rightarrow A \\
 AAA \rightarrow A \\
 AA \rightarrow AA
 \end{array}
 \quad \mid \quad
 \begin{array}{l}
 \text{AAA} \rightarrow A \\
 \text{AA} \rightarrow \text{AaA} \\
 aA \rightarrow a
 \end{array}$$

\* Type 1 → Infinite length (LBA)

left < Right

PAGE NO. \_\_\_\_\_  
DATE \_\_\_\_\_

$A \rightarrow AA$

$A \rightarrow A(AA)$

Example →

$$AA \longrightarrow AAA$$

$$A \longrightarrow AA$$

$$AA \longrightarrow AAA$$

$$AAA \longrightarrow AAAAa$$

$$A \rightarrow \underline{AaA}$$

$$A \rightarrow (Aa) \underline{GA} \rightarrow Aa \underline{aa}$$

Accepted by

★ Type 2 → Finite length (PDA)

left Side → Single Non Terminal

Right Side → Not more than 2 Non Terminal  
else it's length is infinite

$$\begin{array}{l} A \rightarrow aA \\ A \rightarrow \overline{Aa} \\ A \rightarrow \underline{AA} \\ A \rightarrow a \\ A \rightarrow aa \end{array}$$

$$\begin{array}{l} A \rightarrow AA \\ A \rightarrow (aA)(aA) \\ A \rightarrow aAaA \end{array}$$

★ Type 3 → Finite length (FA)

$$A \rightarrow a$$

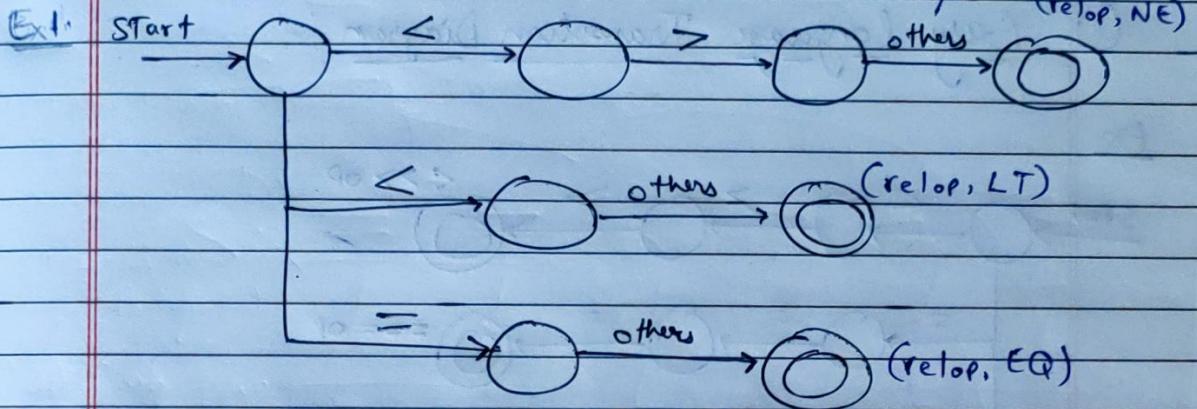
$$A \rightarrow aA$$

$$A \rightarrow \underline{Aa}$$

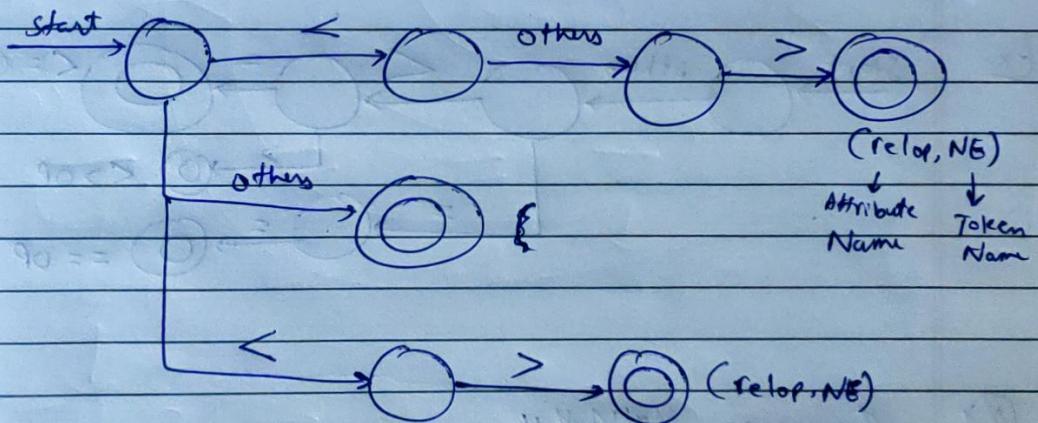
$$A \rightarrow aa$$

Type	Grammar	Language	Automata
Type 3	RG	Regular Language	FA
Type 2	CRG	CFL	PDA
Type 1	CSG	CSL	LBA
Type 0	UG	Recursive Language	TM

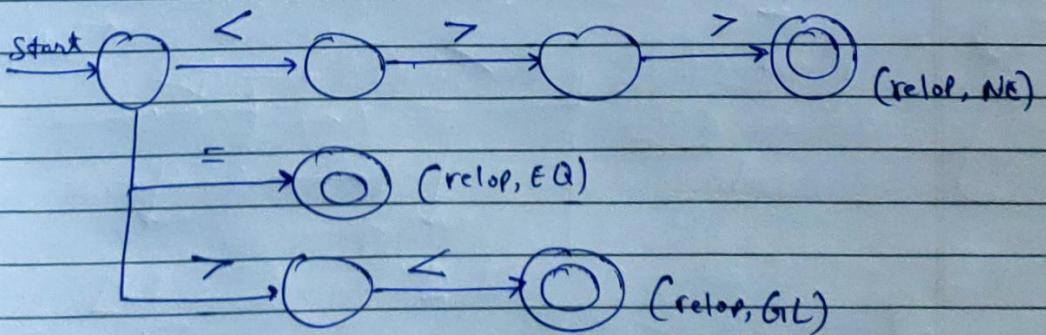
Transition Diagram :- (Leevin Method) Ignore it



Ex.2.



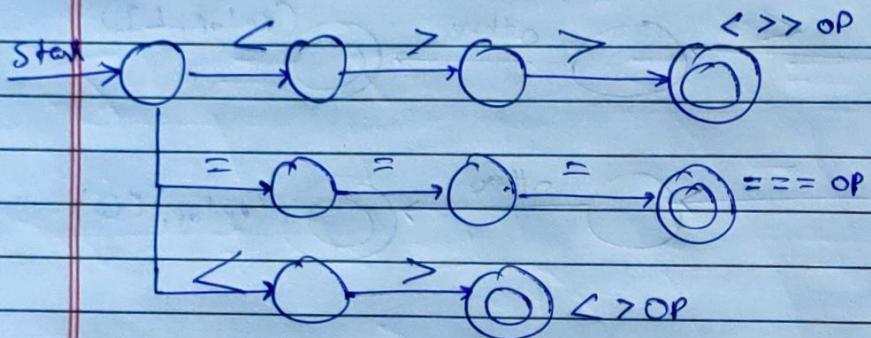
Ex.3.



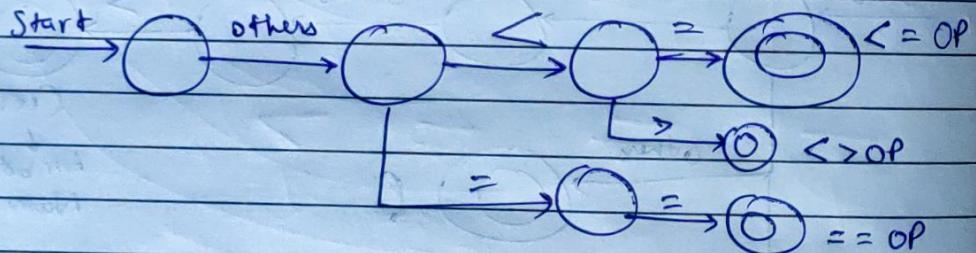
## \* Transition Diagram :-

### ① Easy Language Transition Diagram :-

Ex.



Ex.



②

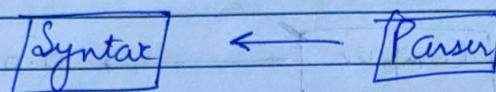
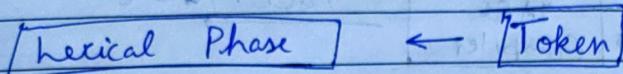
Language Method

## LEX Tools

- ① Lexical Analyzer Generator
- ② Lexical Phase
- ③ Support all operating system.

## YACC Tools

- ① Yet Another Compiler Compiler
- ② Syntax
- ③ Unix operating system
- ④ LALR → Look Ahead Left to Right (Parsing)

Ex.

$$a = b + c \times d$$

↓

LEX

$$Id_1 = Id_2 + Id_3 \times Id_4$$

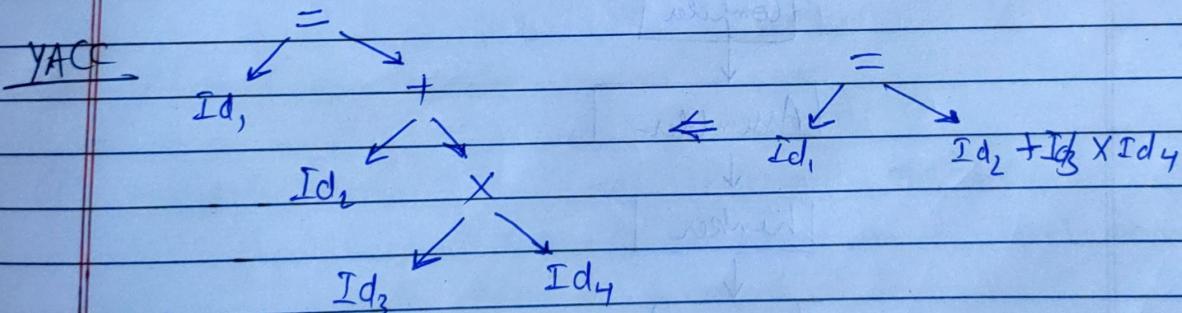
1 2 3 4 5 6 7

Token = 7

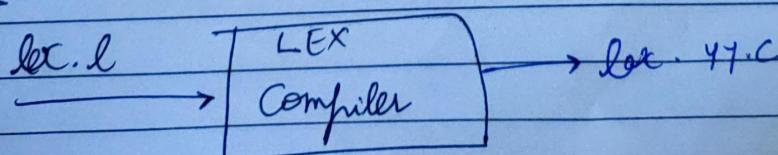
$$a = b t c \times 2$$

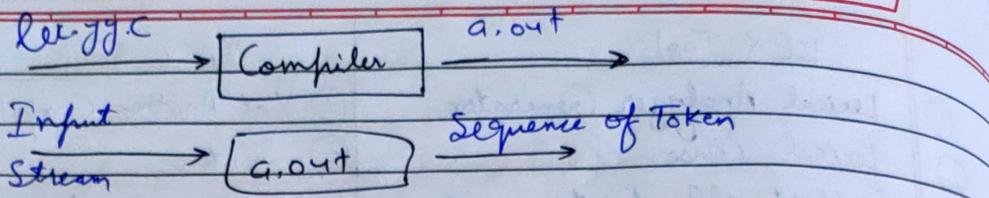
$$Id_1 = Id_2 + Id_3 \times Num$$

1 2 3 4 5 6 7

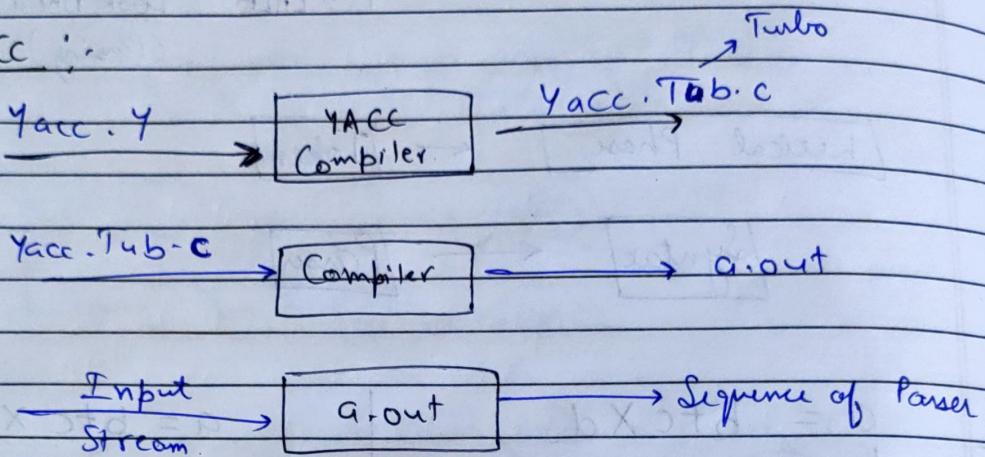


Parse Tree

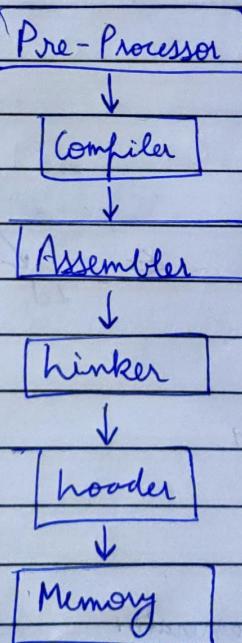
# LEX → Lexical Analyzer Generator



# YACC :



# Relation b/w Compiler , Assembler , linker and loader :



- \* (i) Pre-Processor :- It takes source code as a input and share to next phase.
- \* (ii) Compiler :- i) It convert High Level language to machine language.  
ii) It executes full program (whole) at a time
- \* (iii) Assembly :- i) It is also called low level language.  
ii) It converts Assembly language to machine language.  
iii) Difficult to understand.
- \* (iv) Linker :- i) It take machine code as a input and attached library file with it and share to next phase loader.
- \* (v) Loader :- i) It is part of operating system (OS).  
ii) It loaded whole code and run them.

## # Phases of Compiler:

- (i) Syntactical Phase :- i) It is first phase of compiler.  
ii) It take source code as a Input and generate token as a output.
- (ii) Syntactic Phase :- i) It is second phase of Compiler.  
ii) It take token as a Input and generate parse tree as a output.

(iii)

Semantic :- ① It is third phase of compiler.

② Semantic error

③ Type checking → Static Type checking (Compiletime)

Ex: C, C++, Java



Dynamic Type checking (Run Time)

Ex: Perl, Ruby, Python.

(iv)

Intermediate Code Generator :-

① It is 4<sup>th</sup> phase of Compiler

② It takes machine code as a Input and generate intermediate code as a output.

(v)

Code Generator :-

① It is 5<sup>th</sup> phase of Compiler.

② It takes Intermediate Code as Input and generate code as a output.

(v)

Code optimization :- 5<sup>th</sup> Phase

① It is best phase of compiler.

② It reduce complexity → Space Complexity

→ Time Complexity

## # Types of Compilers:-

### ① Multipass compiler →

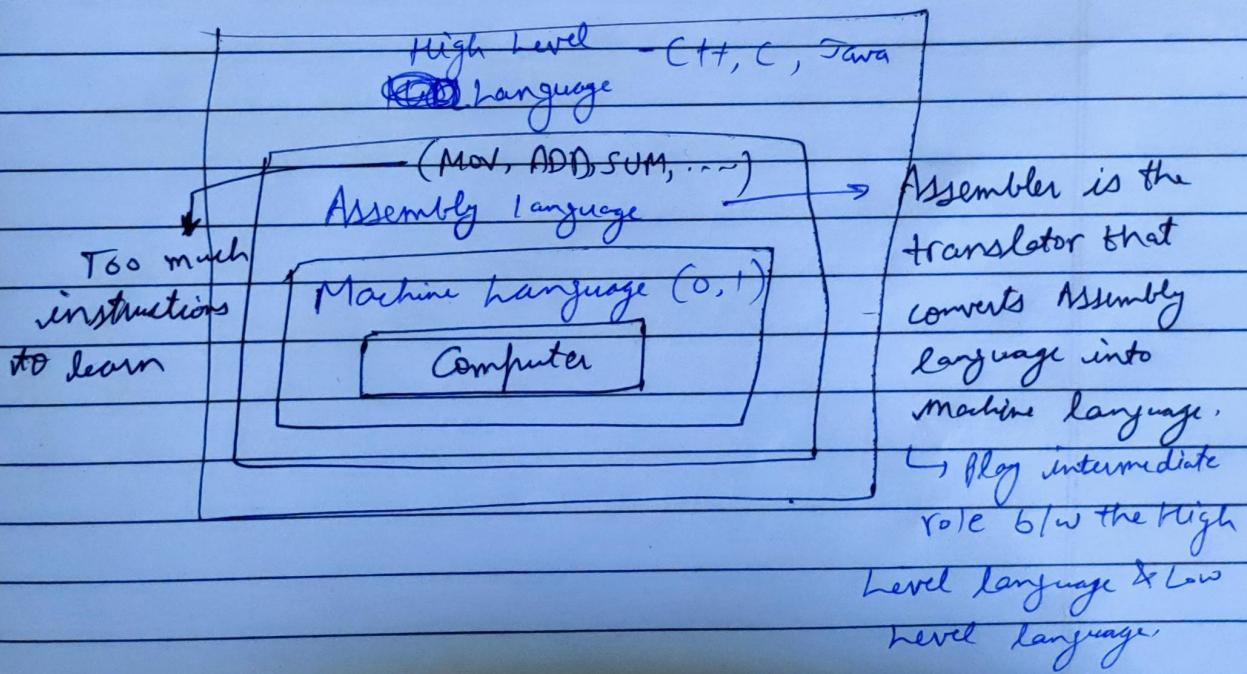
- ① It is also called 2 or 3 compiler.
- ② It is step by step execution.
- ③ Ex. → COBOL, ALGOL, Fortran.

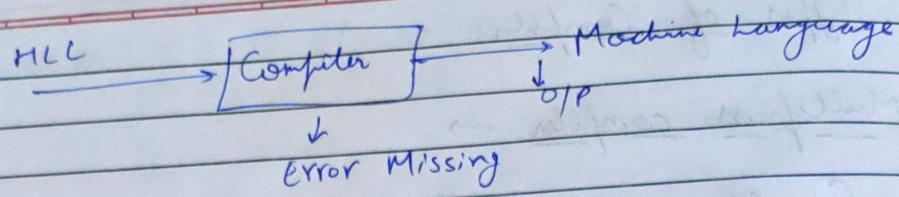
### ② Source to Source Compiler :-

- ① It takes input as high level language and O/P as a high level language.

Ex. ① C# Compiler

- ② C++ Compiler
- ③ Java Compiler
- ④ COBOL Compiler
- ⑤ ALGOL Compiler





Compiler → C, C++, Java

Interpreter → Python, Ruby, MATLAB

- ★ Compiler always creates the target code or source file (in assembly) & forwards to machine language.  
Interpreter doesn't make.

★ Speed → Compiler >> Interpreter