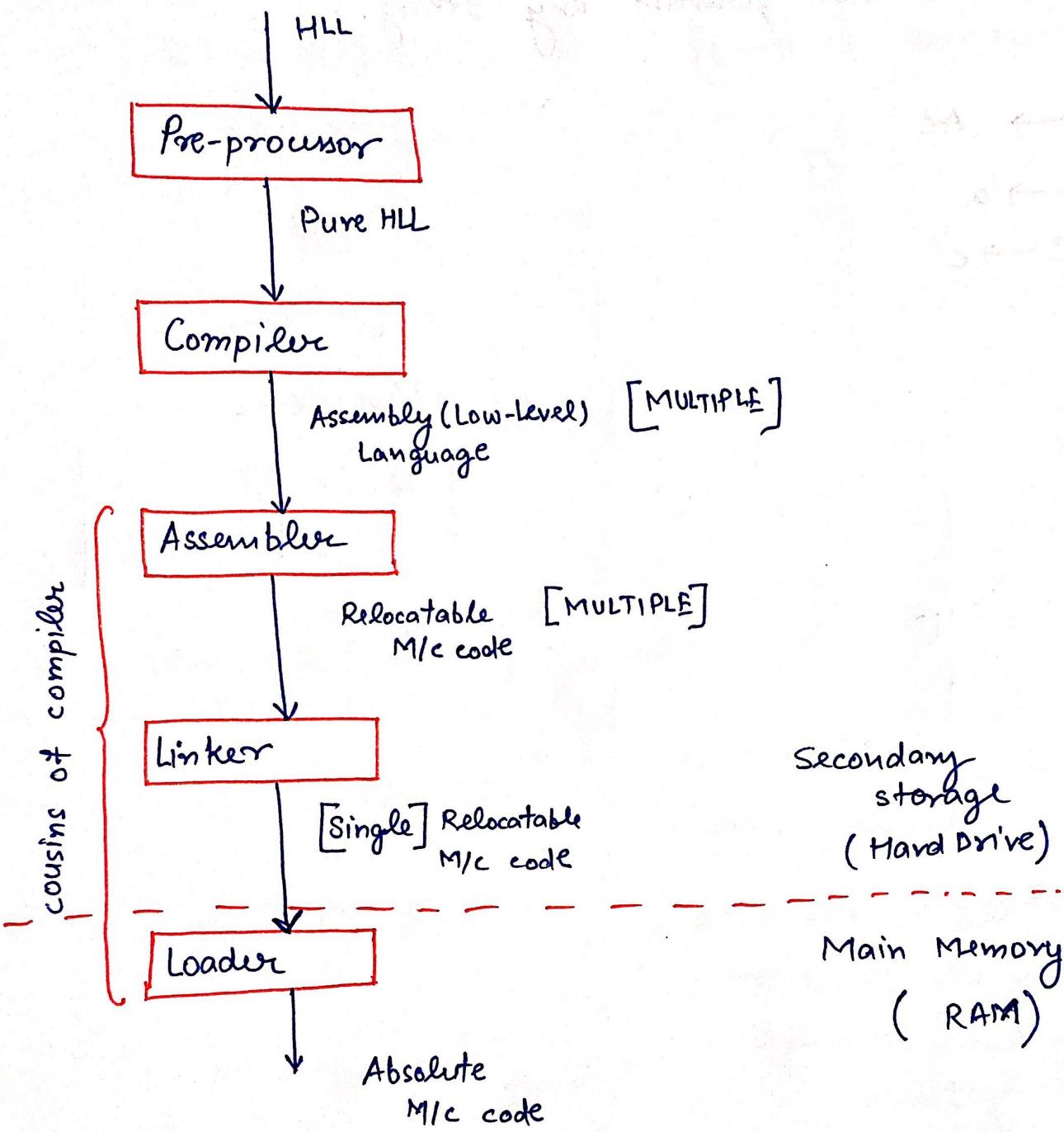


Introduction to compiler



(i) Pre-processor:- Convert HLL to Pure HLL i.e substitute header file with there code.

include < > → ;
define < > → main();
[main()
{ → ;
} ; → ;

(ii) Compiler :- Converts Pure HLL to Assembly (Low-level) Language.

main()

```
{  
    x = a + b;  
}
```



main:

```
MOV a, R0;  
MOV b, R1;  
ADD R0, R1;  
MOV R1, x;
```

(iii) Assembler :- Converts Assembly language to

relocatable M/c code. Relocatable means independence of type of memory.

MOV a, R0



opcode	src	dest
1101	10001000	10001111

Binary format

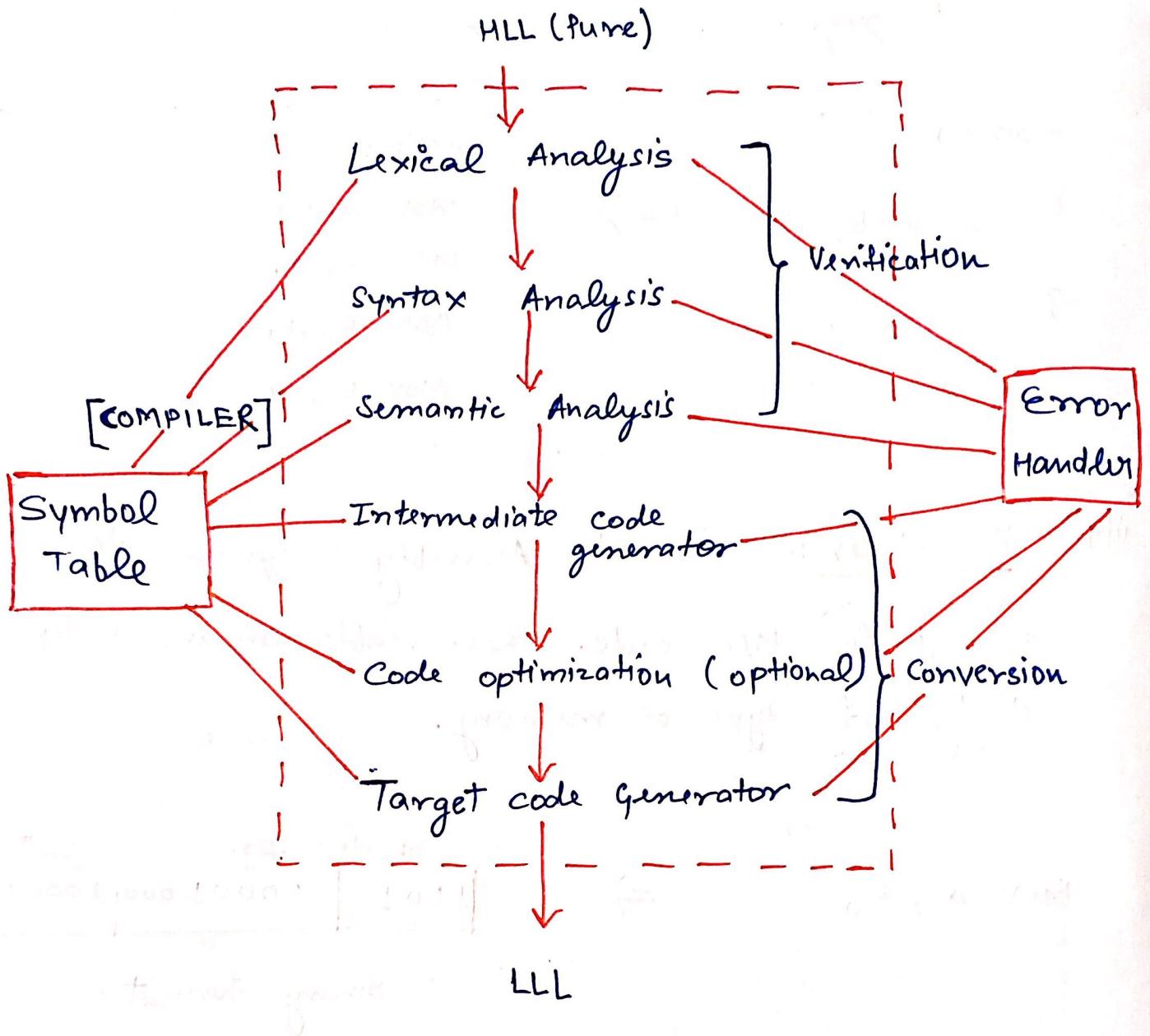
(iv) Linker :- It links multiple M/c code file into a single M/c code (relocatable)

(v) Loader :-

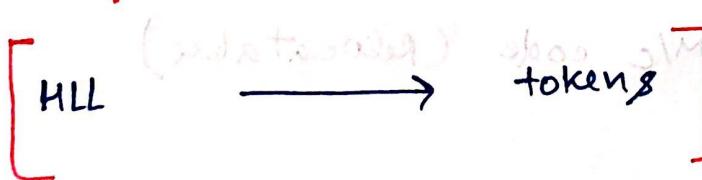
It executable code LOAD L into Main

memory (RAM) from secondary memory

Phases of compilation :-



(i) Lexical Analysis :-



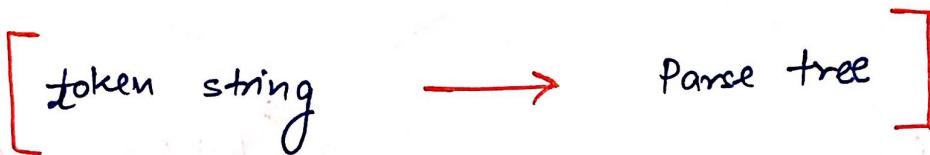
$x = a + b;$

$id = id + id ;$
tokens string

x	id
a	id
b	id

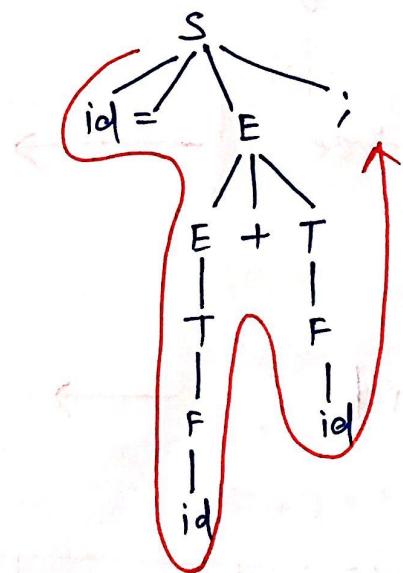
=	op
+	op

(ii) Syntax Analysis :- It verify whether the string of token is syntactically correct or not. It is done by Parser which use a CFG (rules) to generate Parse tree.



If parse tree is deriving the same token string then I/p is syntactically correct, otherwise it raises Syntax Error.

$$\begin{array}{l}
 \text{P: } S \rightarrow \text{id} = E; \\
 E \rightarrow E + T \mid T \\
 T \rightarrow T * F \mid F \\
 F \rightarrow \text{id}
 \end{array}$$



Since it derives same I/p string
so it is correct

(iii) Semantic Analysis :- It modifies the parse tree and check whether it is meaningful ie it checks attributes of identifiers.

Eg:

$$id = id + id$$

$$\checkmark a = b + c ;$$

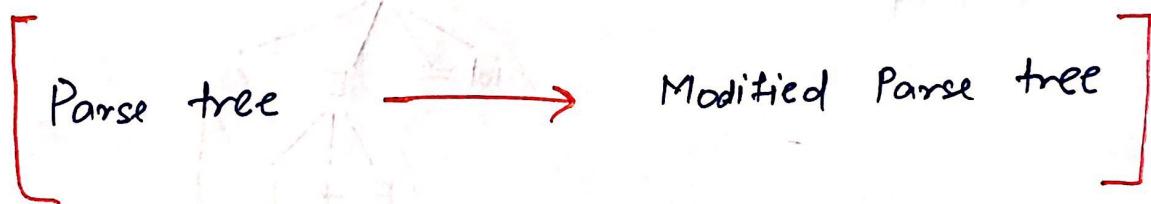
No error

$$2 = b + c ; \times$$

result can't be assigned to a constant

Error

As, we can see both statements are syntactically correct but $2 = b + c ;$ is not semantically correct i.e. meaningless.

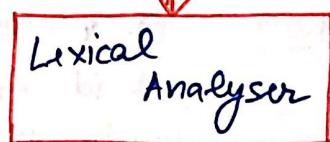


$$id = id + id \longrightarrow id_{\cdot}^{type} = id_1^{type} + id_2^{type}$$

s.no	id/type	V.name	V.type
1	id	x	int
2	id	a	int
3	id	b	int

Eg:-

$$n = a + b * 60$$



x — identifier — (id,1)

= — assignment op

a — identifier — (id,2)

+ — add op

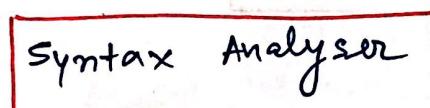
b — identifier — (id,3)

* — mul op

60 — int const

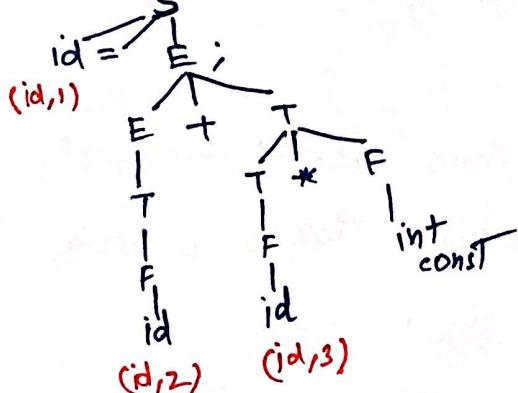


$$(id,1) = (id,2) + (id,3) * 60$$



CFG:

$$\begin{aligned} S &\rightarrow id = E ; \\ E &\rightarrow E + T \mid \text{int const} \\ F &\rightarrow F * F \mid F \\ F &\rightarrow id \end{aligned}$$

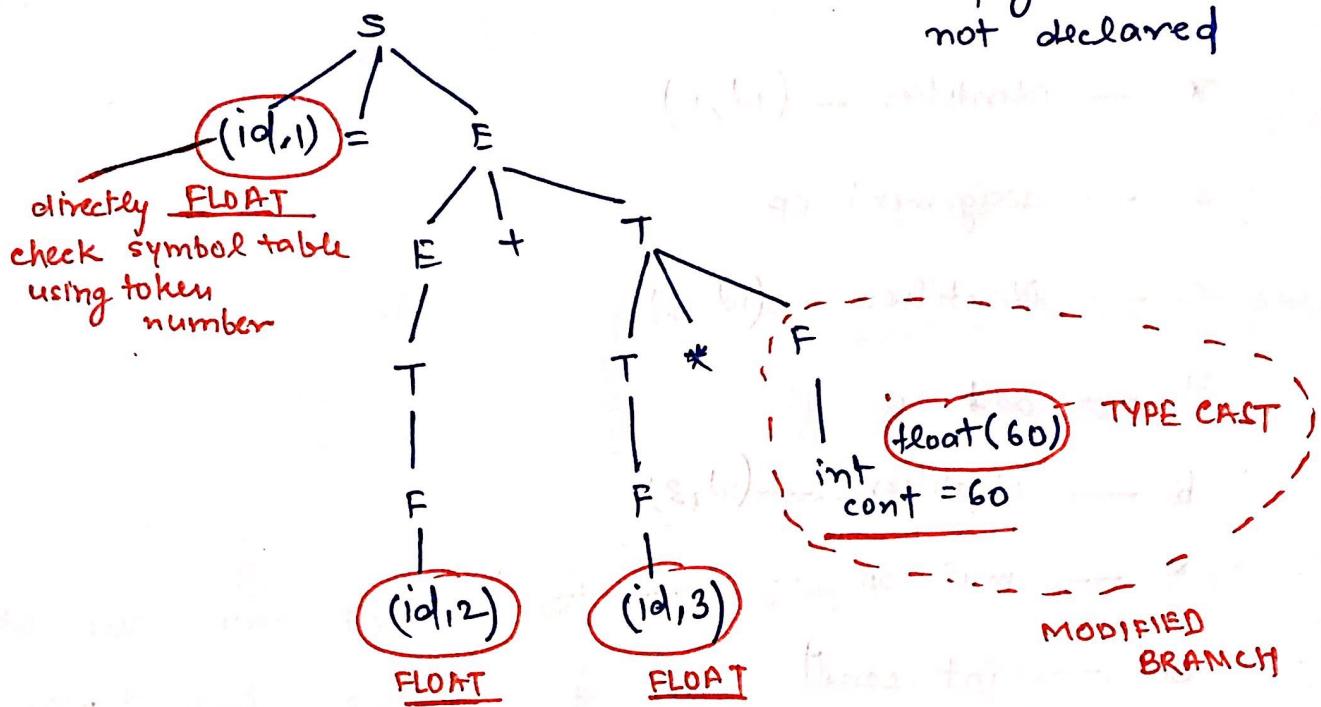


if parse tree possible
then no Error

Semantic Analyser

NOTE:

if .TYPE column
of id token is
empty then var is
not declared



$$(\text{Float} = \text{Float} + \text{Float} * \text{Float})$$

Intermediate code Generator

$$\begin{cases}
 t_1 = b * 60.0 \\
 t_2 = a + t_1 \\
 x = t_2
 \end{cases}$$

Code optimization

$$t_1 = b * 60.0$$

$$n = a + t_1$$



Target Code Generator

MUL 60.0, R₀

Mov b, R₀

ADD R₀, R₁

Mov a, R₁

MOV R₁, n

Short format

NOTE:- Most of the time semantic analyser store into in Symbol Table, but if some undeclared variable comes than Lexical analyser store it with leaving all the fields empty and release a token for future reference.

eg:-

float x, a, b;

By semantic
analyser

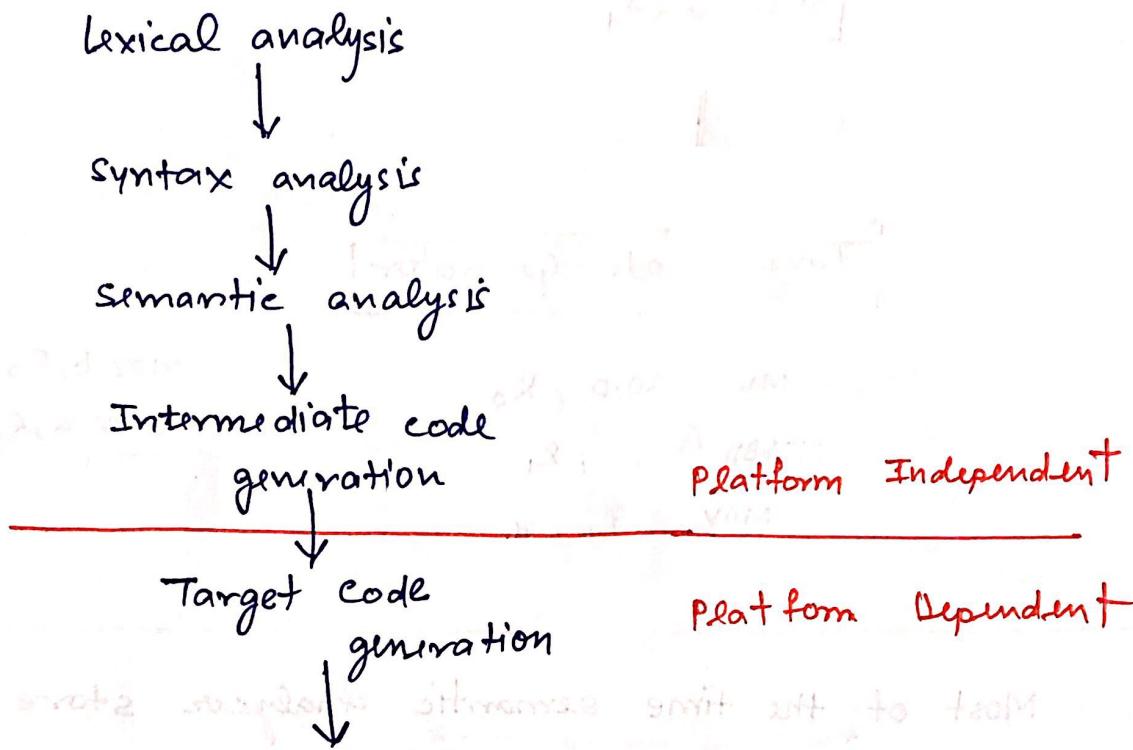


s.no	name	type
1.	x	float
2.	a	float
3.	b	float

NOTE:-

Symbol Table only contains user defined variable, func, objects, interfaces.

- After compilation, O/P is M/C dependent code i.e. it can only run on the machine in which compilation is done.



O/P of ICG : Platform Independent
as it is not written in machine language

O/P of TCG : Platform dependent
Code

- So compiler O/P depends upon whether TCG is included or excluded.

JAVA compiler :- Excludes TCG

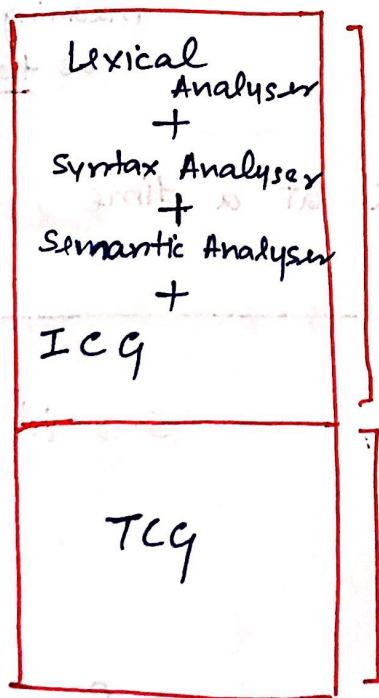
Platform
independent

working on C/C++ compiler :- Include TCG

(Java compiler) work on the
platform dependent

<u>JAVA compiler</u>	<u>C/C++ compiler</u>
- Half compilation	- Full compilation
- Exclude TCG	- Include TCG
- Platform independent	- Platform dependent
- Portable	- Not Portable
- Execution slow (it has to first convert into M/c code)	- Execution fast

STRUCTURING OF COMPILER



Back-end

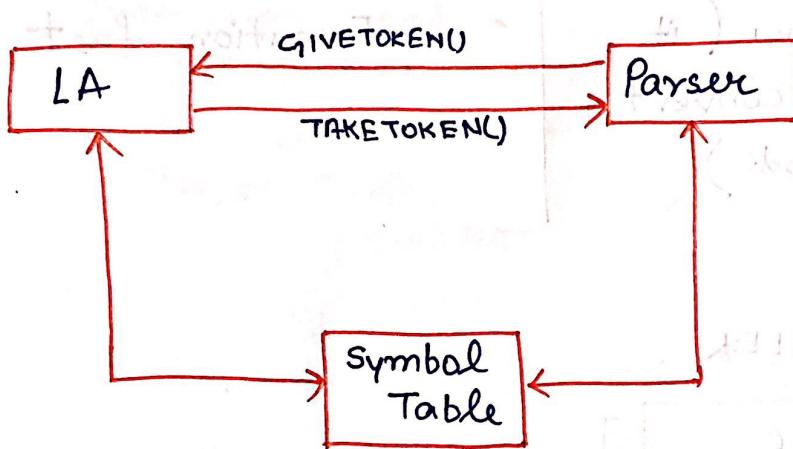
For more & more to add

Single pass compiler :- Total compilation process (phases) at a time (Less time, more space)

Multi passes compiler :- compilation process at time

frames. (More time, Less space)

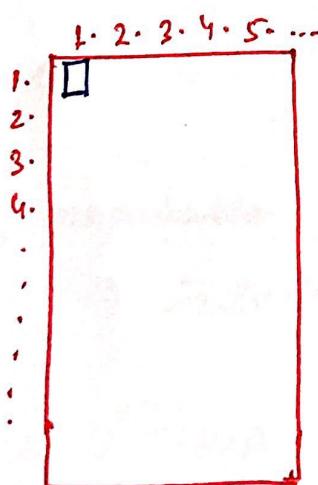
Lexical Analyser :-



- LA is first phase of compiler, also known as scanner

- It divides into meaningful strings known as tokens

- LA and Parser both work at a time, side-by-side,



$$S \rightarrow id \ ? = E \mid AB$$

PST

CD

Parser at S

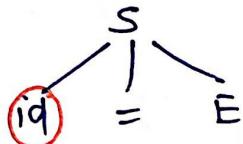
LA is at row=1 & column=1

Type :- $id = E$

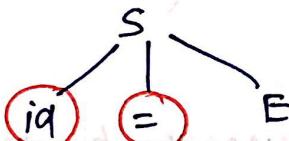
Type :- $Id + E$

- LA send 1st token $(id, 1)$ to parser

- Parser checks production for id alphabet



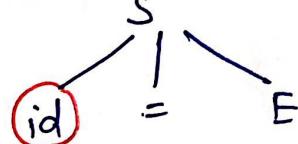
- Then Parser again ask for token from Parser



- LA send 1st token $(id, 1)$

- Parser checks

for id



Now LA sends

'+' token, But
Parser Expect for
'=' so Error
is generated

- Types of tokens :-

(i) Identifiers

(ii) Operators

(iii) Keywords

(iv) Constants

Responsibilities of IA :-

(i) Dividing given program into tokens.

e.g. int ₁²a₃,⁴b₅;

Tokens: <int/1>, <a/2>, </1/3>, <b/4>

</> / 5 >

(ii) It will eliminate whitespace character (tab, new-line, gap)

(iii) It will eliminate comment lines { /* */ }

(iv) It will help in giving error messages by providing (row, column) number.

Q. Find number of tokens present in the following C program :

int main()

{ /* finding max of a & b */

int a=10, b=20, max;

if (a<b)

 max=b;

else

 max=a;

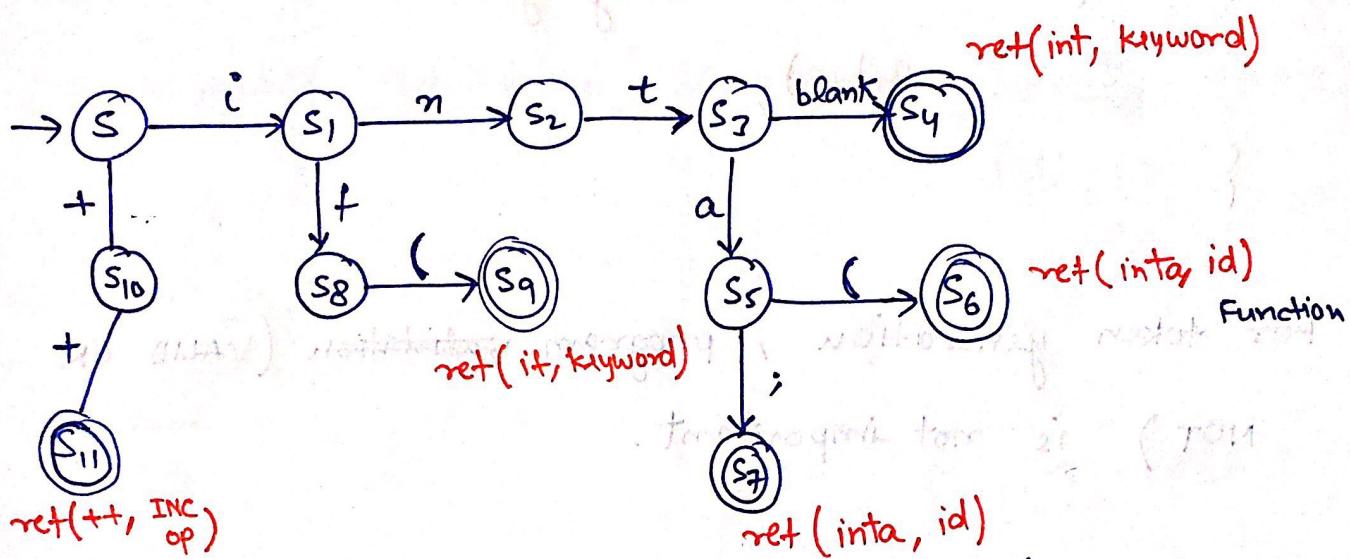
return (max);

(37)

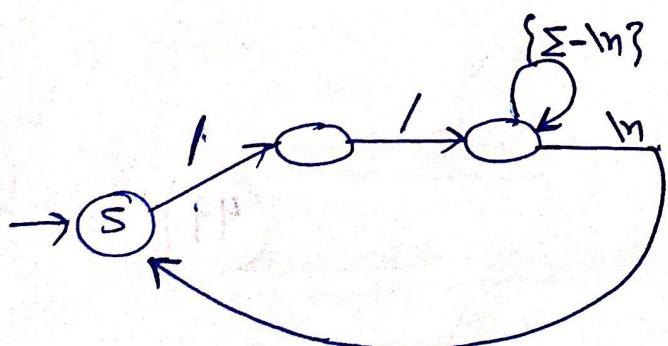
// Token Creation :- LA use FSA (finite state Automata)

to create token.

Abstract DFA :-

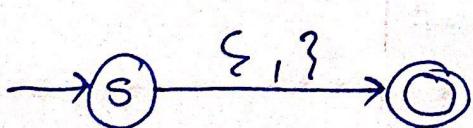


- DFA is provided to LA, which includes finite states.
It includes every possible combination for identifying tokens.

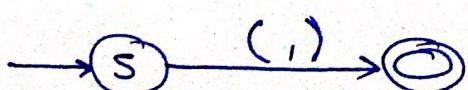


DFA for comment

' Σ ' means any alphabet



DFA for brackets



DFA for parenthesis

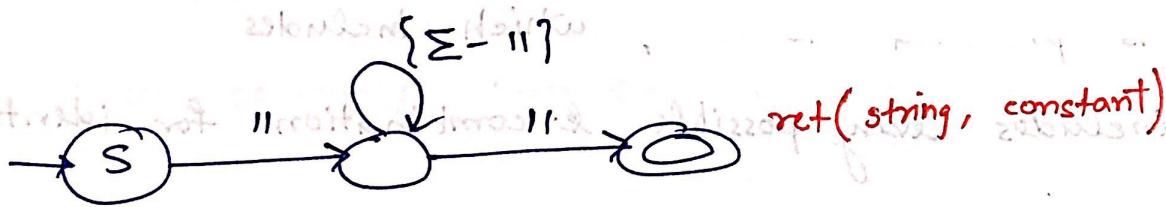
Q. find no. of tokens for the following program;

```
int main()
{
    u = a + b * c;
    int x, a, b;
    printf("you & your family goto hell\n.%d\n.%d",
    a, b, c);
}
```

(32)

For token generation, program (VALID OR NOT) is not important.

whole string is counted as one token.



Q. Find no. of tokens -

```
main()
{
```

u = a + b - c;

y = ++|++|-+ ==|=*=/*|*++++++---==--

* * * + + == * * * | + - + - ;

comment

printf(" %d , %d , %d ", a, *a, **a, ***a, &a);

(49)

3 — math()

story of the time by means of a series of maps.

$$8 - n = a + b - c; \quad \text{составляющая}$$

~~9 - y = ++|++|--|-|=|=|*|/*++----*---***/~~

$$5 - +|-+|-;$$

23 — printf("-1.d -d d", a, *a, **a, &a, **a);

~~bitwise AND op~~ | ~~a~~ id | ~~base~~ * | * | * | a
 bitwise AND op address op = $*(*(*(a)))$

Q. Find no. of tokens:

3 main()

18

$$b \quad u = a + b;$$

(24)

13

y = xab-123 abcdef *** ghi *** / ***

1

$$\underline{n = ab + cd - ghabcd / \ast \ast \ast \ast \ast} / \underline{abcd}$$

1

comment not counted

NOTE:-

Lexical Error :- When LA not able to generate token from given string.

Eg:- comment started but not ended
program is ended.

① main() ② ③

④

/* -----

} Comment
not ended

LA not able to
generate token.

Q. Find number of token:

main()

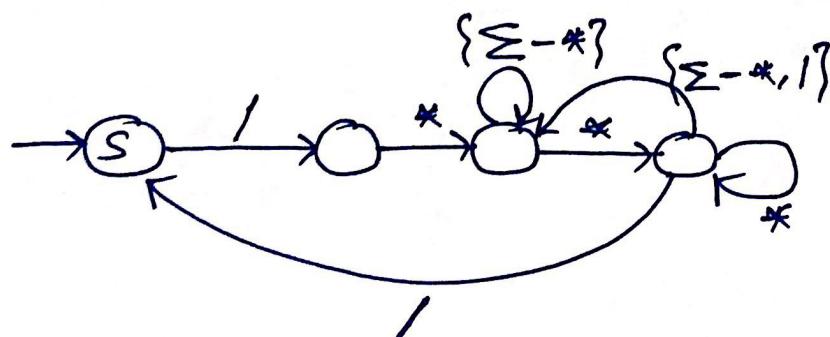
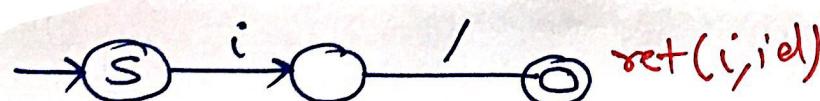
{

int n=10;

i nt n = 10;

i /* comment */ nt n=10;

(22)



Multi-line
comment

Q. Which one of the following string is said to be token without looking next input char in C lang?

X(a) return

X(b) if

X(c) main

X(d) +

✓(e) ;

✓(f),

✓(g) (or) with backslash instead of comma

Q. A lexical analyzer uses the following patterns to recognize three tokens T_1, T_2 , & T_3

$T_1 : a^q (b|c)^* a$

$T_2 : b^q (a|c)^* b$

$T_3 : c^q (b|a)^* c$

I/p: bbaaacbc; which one of the following is the sequence of tokens it outputs? which has longest matching?

(a) $T_1 T_2 T_3$

(b) $T_1 T_1 T_3$

(c) $T_2 T_1 T_3$

(d) $T_3 T_3$