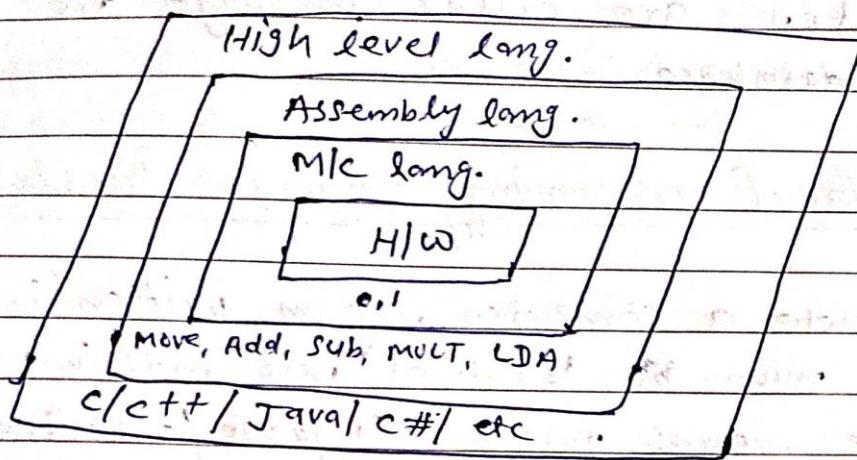


Unit-I

* Why we study of compiler [importance of compiler Design] ⇒



⇒ Machine Language ⇒ This language contains only 0's & 1's instructions code. It is not required any type of translator. But it's very difficult for programmer to write programming code in this lang.

⇒ Assembly language ⇒ It is similar to machine language {ADD, SUB, MOVE, MULT, LDA} but it's different to machine language (it's only combination of 0's & 1's). It's used new symbols like Add, Move, Sub, MULT, LDA, etc. It's required Assembler for decode the programming code into machine code.

It's difficult language for programmer but less difficult as machine language. It's required memorisation for complex & large program for more instruction.

⇒ High level language ⇒ Today we use all this language in all s/w. means all s/w made in High level language.

It is user friendly language, easily user understandable. High level language required translator because translator translate the High level language into low level language (Mic language).

Example:- c/c++/ Java/ C# etc.

⇒ 4GL (fourth Generation language) ⇒ 4GL are closer to Human language. It is

easy to learn the language like:- oracle, vctt, SQL etc. are called language are used to access database.

* Translator [programming language translator] ⇒

To execute a computer program written in High level language must be translated into machine understandable language means machine language / machine code.

" A translator is a program that takes as a input program written in one of the programming language (source language / source code) & produces as a output a program in another language (object / machine language / Target code).

If the source language is High level language such as FORTRAN, C, C++, JAVA, C# etc & the object language is low level language such as assembly language / machine language.

⇒ Types of translators ⇒ There are mainly three types of programming language translators :-

- Assembler
- Compiler
- Interpreter

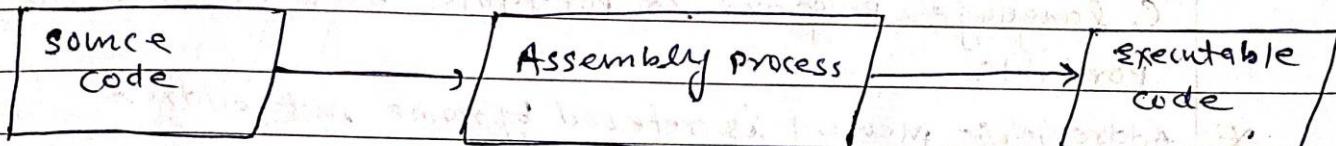
* Assembler ⇒ Assembler is a translator but it translate high level language into Assembly language.

Assembly language means which is basically mnemonics code like GO, HALT, JUMP, NOT, MOVE, LDA, ADD, SUB MULT etc code.

It translate machine language into hexa decimal digit code.

Assembler Assembler is program (sw) हेरा ए ट्रान्सलेटर machine dependent हेरा ए

Assembler converts assembly language into patterns of bits that the computer processor can use to perform its basic operations. This bit pattern is basically machine language.



Assembly lang.

Assembly process

executable code

Assembly lang.

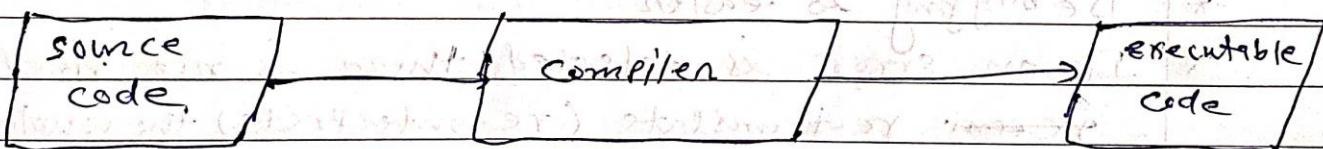
bit pattern or
Machine lang.

Example:- NASM, MASM (it's Microsoft assembler).

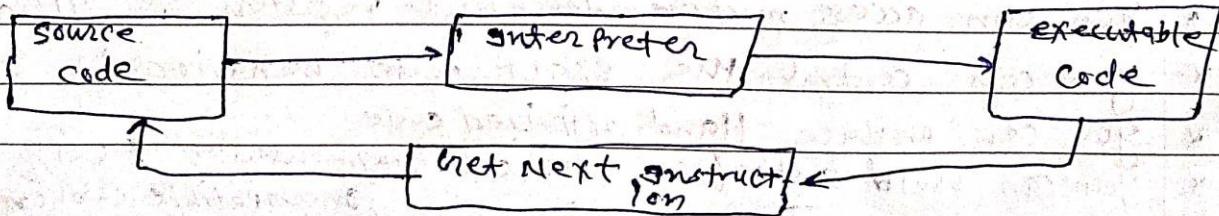
* Compiler \Rightarrow Compiler is a program that translates a high level language into machine code.

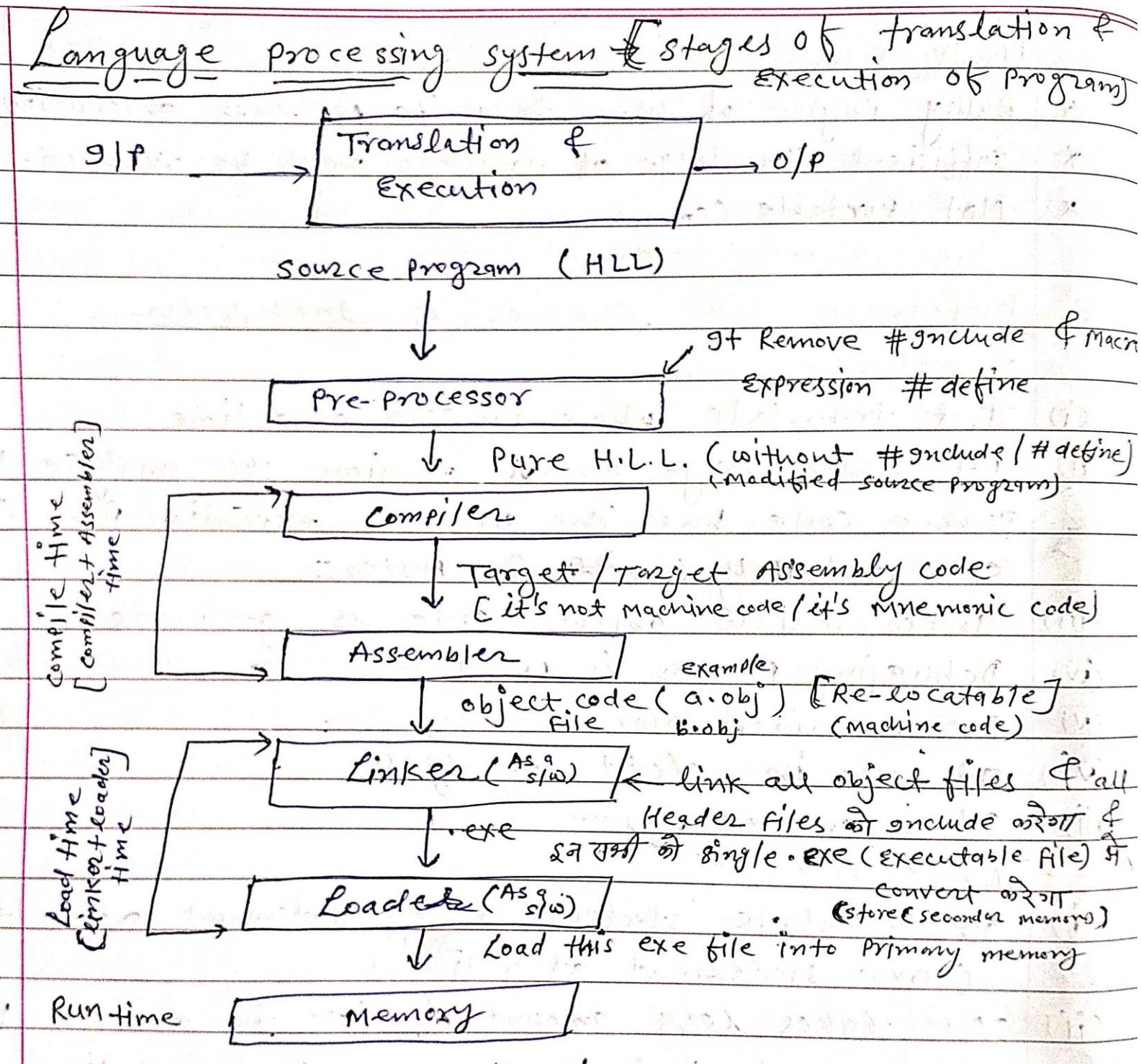
Example:-

The Turbo Pascal Compiler translates a program written in a high level language into machine code that can be run on a personal computer (PC). Example:- C/C++/FORTRAN



* Interpreter \Rightarrow An Interpreter is another type of translator which converts high level language into machine code. Main difference b/w Compiler & Interpreter is that an interpreter translates one line at a time & then executes it. No object code is produced. So program has to be interpreted each time it is to be run. Example:- LISP/Python/BASIC



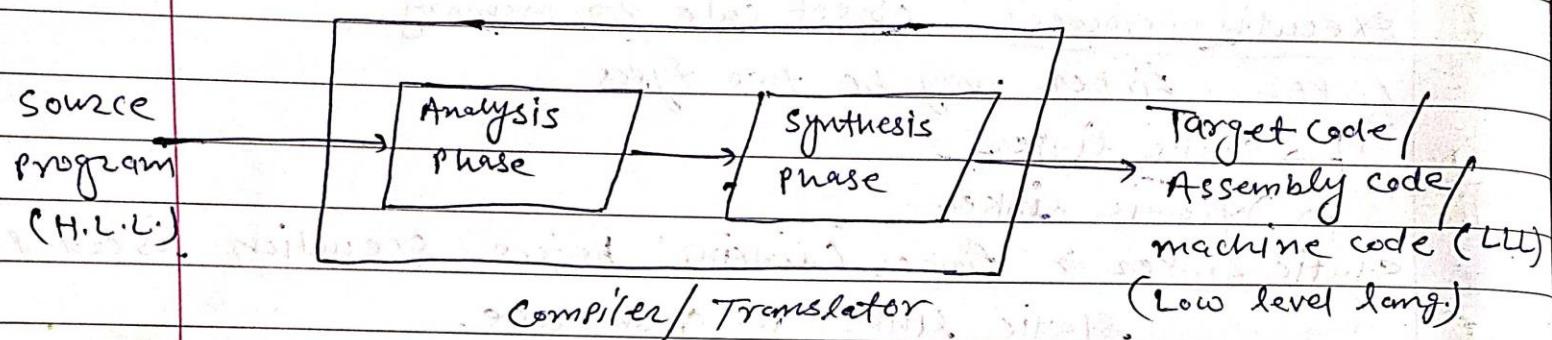


~~X~~ Compiler stages / stages in translation / compiler phases

Compilers are mainly two phases

→ Analysis phase

→ Synthesis phase



* Analysis Phase ⇒ This phase is used for source program recognition phase.
It consists of three phases :-

- Lexical analysis
- Syntax analysis
- Semantic analysis

* Synthesis - Phase ⇒ This phase is used for executable object program recognition phase.

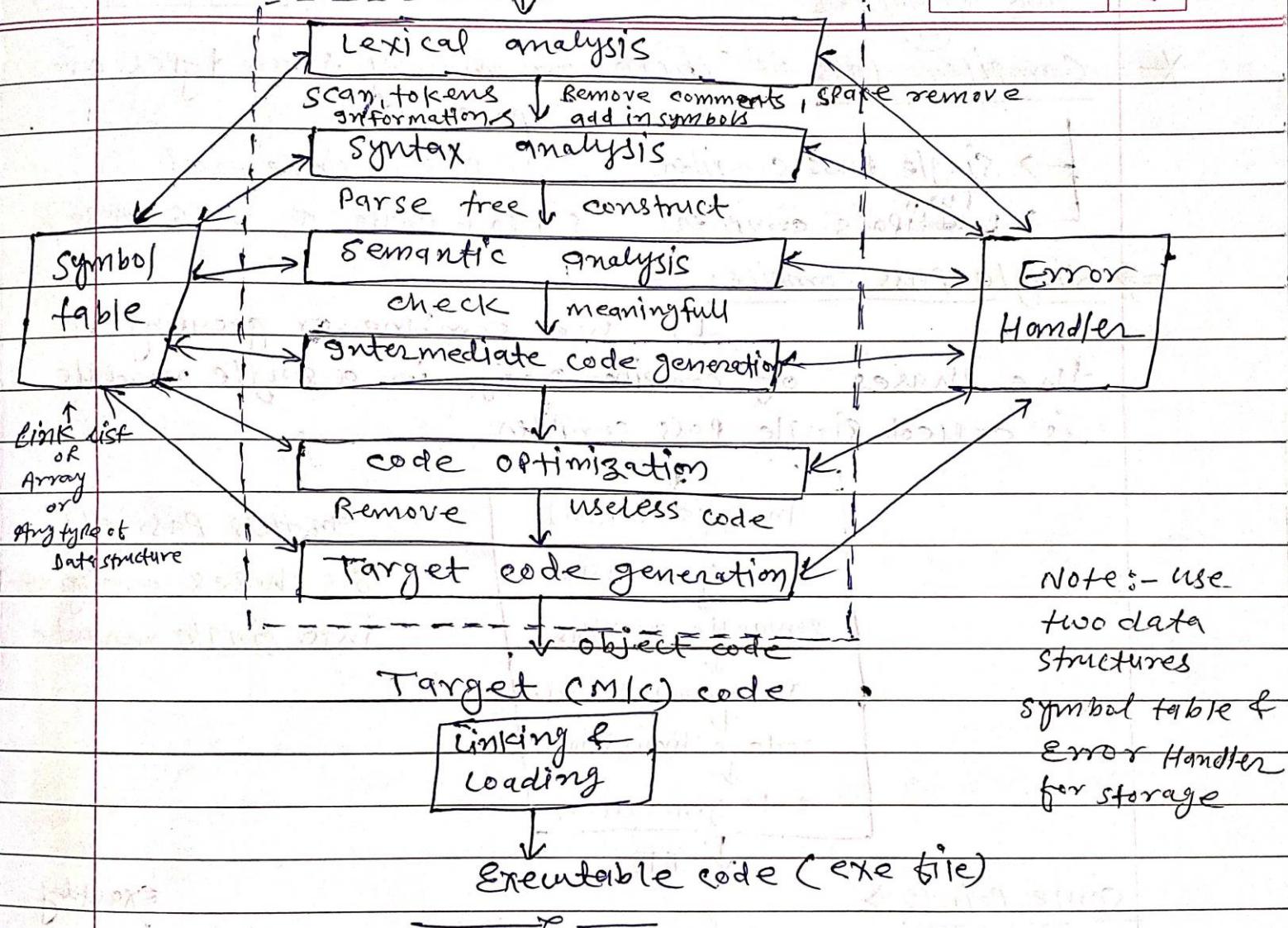
It consists of three phases :-

- Intermediate code generation
- code optimization
- Target code generation

Source program (code) [Pure H.L.L.]

Page No.:
Date:

youva



* Analysis Phase → This phase is used for source program recognition phase. It consists of three phases :-

- Lexical analysis
- Syntax analysis
- Semantic analysis

* Synthesis Phase → This phase is used for executable object program recognition phase.

It consists of three phases :-

- Intermediate code generation
- Code optimization
- Target code generation

① Lexical analysis → The initial phase of any translator/compiler is a lexical analysis phase. Its main task is scanning & generating tokens & removing comments & variable information stored in symbol table. remove spaces

Example:- int a, b;

↑ ↑ ↑
 Keyword Variables Delimiters

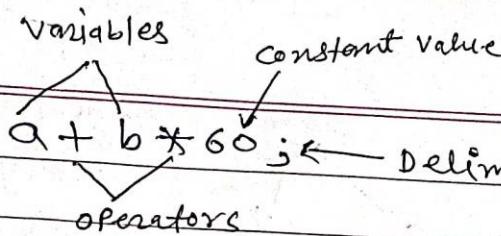
1	a	int
2	b	int

Remove comments like

/* --- */ //

(OR)

Example:-



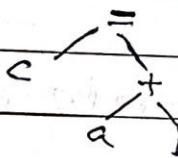
Note:- INT लिखने पर error generate होता है क्योंकि C lang. is not allow this type of INT, C allows only int.

(2) Syntax Analysis \Rightarrow This is the second phase of compiler
 It construct parse tree of tokens
 (lexical analysis tokens)

Parser construct parse tree

Example:- $c = a + b$

Tokens are generate so parse tree construct



Ex:- $a + b$
 lexical not generate error because tokens are generate But syntax analysis generate error because not construct parse tree

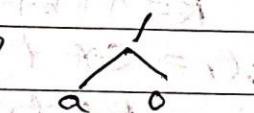
Example:-

while (1) {
 printf ("Hello");
 }
 This program lexically are correct because tokens are generate but syntax error are generate because syntax mistake while (1) & ~~printf ("Hello");~~

(3) Semantic analysis \Rightarrow It is the central phase of compiler. It check the meaning full statement & maintenance symbol table, error detection & correction. It check meaning full tree

Example:- $9/0$

Parse tree \Rightarrow

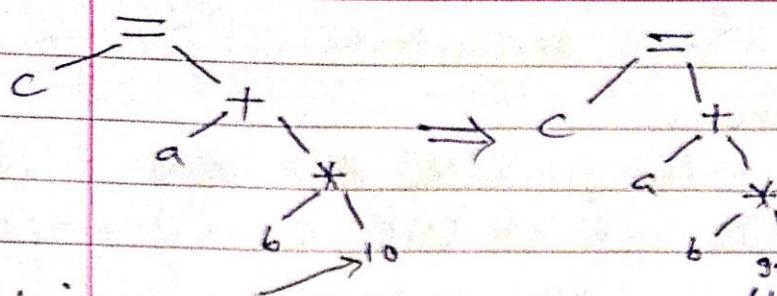


Lexically are correct because it generates token, syntactically is correct because it generate parse tree but it is not correct semantically because it's not meaning full value of $9/0$ [first of number can't divide zero, C language it allow only integer]

Example:- $float a, b, c;$

$c = a + b * 10;$

1	a	float
2	b	float
3	c	float



logically correct because tokens are generate & syntactically is correct because syntax are construct parser tree
But semantically is not correct because

It is an integer value so generate error so 10 is convert into float value 10.0

Note:- Both of these syntax logically, syntactically & semantically are correct at first & 2nd but 3rd phase of error detection chance exist on 3rd stage.

(4) Intermediate code generation \Rightarrow It generate the intermediate code.

It divide the complex statement into smallest sub statement using TAC (three address code) technique [Maximum three address code]

Example:- $R = a + b + c * d$

$$t_1 = c * d, \quad t_2 = b + t_1, \quad t_3 = a + t_2, \quad R = t_3$$

(5) Code optimization (it's optional phase) \Rightarrow

It remove the unusable code which not affected on actual result. It's optional phase

Example:- In above phase example

$$\begin{aligned} R &= a + b + c * d \\ t_1 &= c * d \\ t_2 &= b + t_1 \\ t_3 &= a + t_2, \quad R = t_3 \end{aligned} \quad \left\{ \begin{array}{l} \text{optimize} \\ \hline \end{array} \right\} \quad \begin{aligned} t_1 &= c * d \\ t_2 &= b + t_1 \\ R &= a + t_2 \end{aligned}$$

(6) Target code generation \Rightarrow It is the last phase of any compiler (translator)

It's generate target code (mnemonic code) M/C code

Example:- MOVE C, R1

Pass of compiler

Date:

youva

* Compiler pass \Rightarrow There are mainly two types of compiler pass :-

→ Single Pass compiler { All pass in single pass }

→ Two Pass compiler { 1 to 4 phase & 5 & 6 phase }.

→ Single pass compiler \Rightarrow

If we combine or grouping all the phases of compiler design in a single module

Example :- ~~float~~ $a, b, c;$

$$a = b + c * 60$$

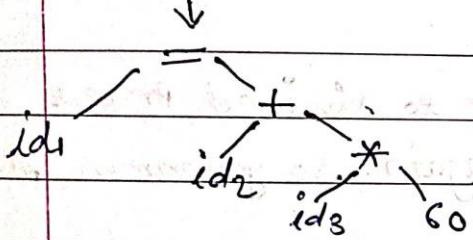
① Lexical analysis

$$id_1 = id_2 + id_3 * 60$$

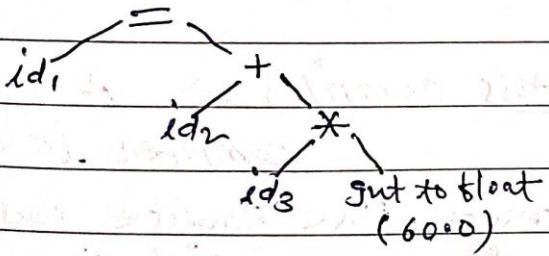
$id_1, op_1, id_2, op_2, id_3, op_3, \text{constant}$

Total 7 tokens generates

② Syntax Analysis



③ Semantic analysis



④ Intermediate code generator

$$t_1 = 60.0, t_2 = id_3 * t_1,$$

$$t_3 = id_2 + t_2$$

$$id_1 = t_3$$

⑤ Code optimization

$$t_1 = id_3 * \text{int to float}(60.0)$$

$$id_1 = id_2 + t_1$$

⑥ Target code generator

MOVE id₃, R₂

MULT 60.0, R₂

MOVE id₂, R₁

ADD R₂, R₁

MOVE R₁, id₁

\rightarrow We can not optimize very well due to limited process

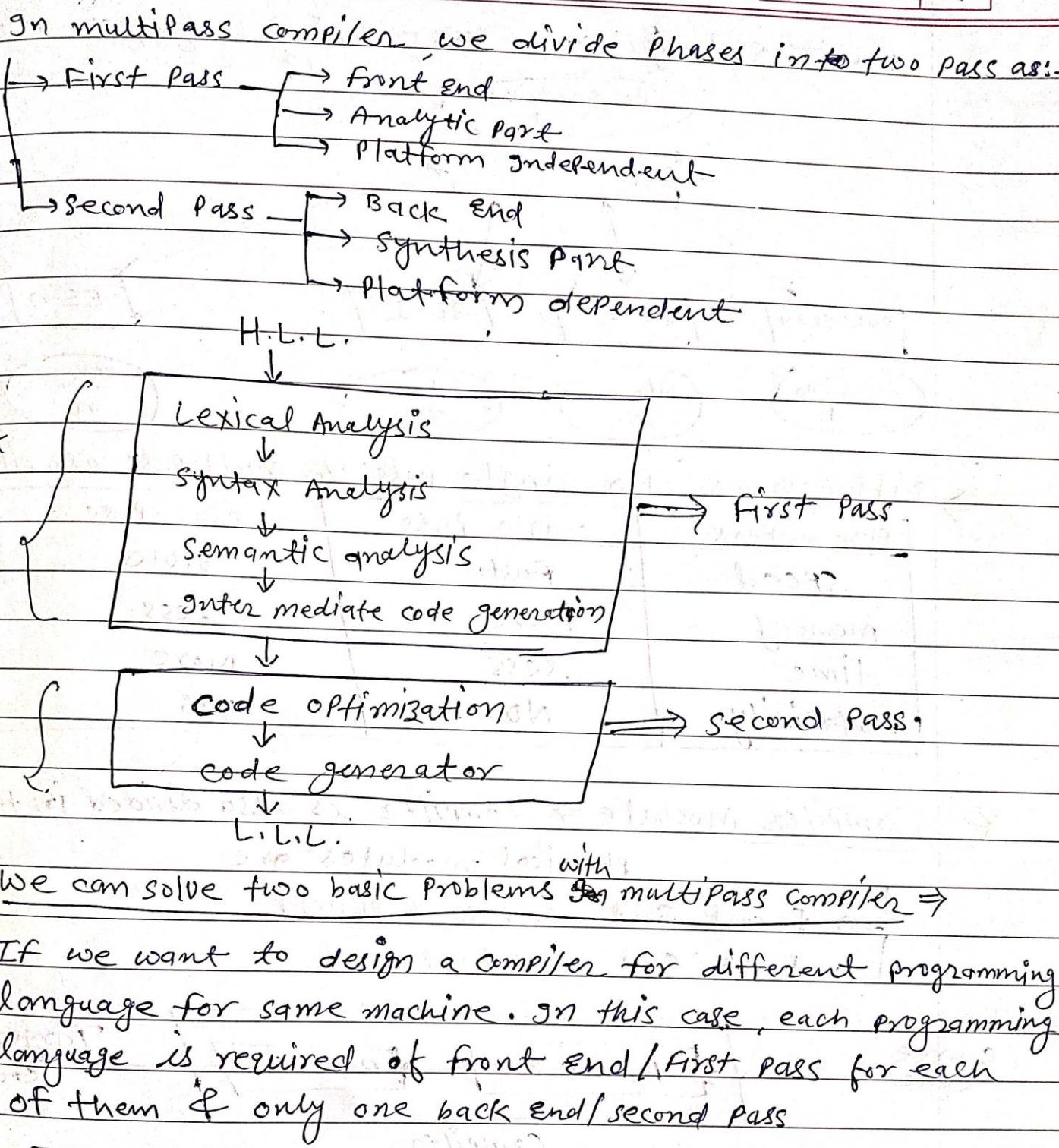
\rightarrow We can not backup & process. It again so grammar should be limited or simplified.

\rightarrow Command interpreters such as bash|sh|tcsh can be consider

→ → →

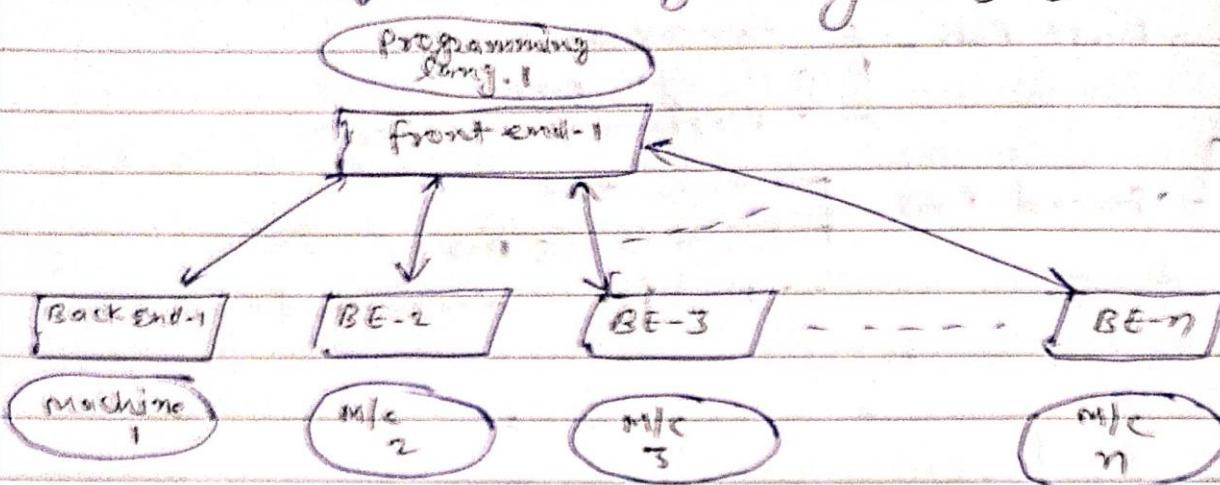
\Rightarrow Two pass compiler / multipass compiler \Rightarrow A two pass compiler is a type

of compiler that ~~passes~~ processes the source code or abstract syntax tree of a program multiple time



ALL → Assembly level lang
 3AC = Three Address Code
 FE = Front end , BE → back end

front end for same programming languages as :-

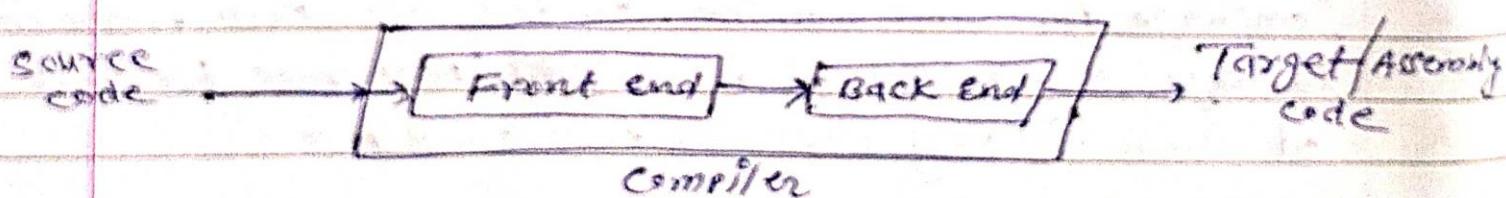


→ Differences b/w single pass Vs multipass compiler →

Parameters	single Pass	multipass	portability
speed	Fast	slow	HLL ↓ M/c(n)
memory	more	less	FE
time	less	more	3AC
portability	NO	yes	M/c ALL

* Compiler Module → compiler is also divided in two physical modules are

- front end [1 to 4 phases]
- Back end [5 & 6 phases]

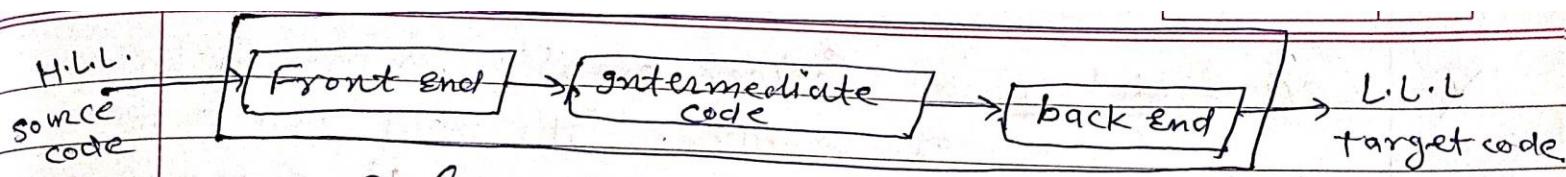


→ Front end → All the activities which is Hardware independent is called front end.

First five activities (1 to 5 phases) are front end.

→ Back end → All the activities which is H/w dependent is called back end. Last (6 No. phase) is activity is back end.

All different company intel, IBM, wipro, microtakce etc has same kind of front end but they have different back end.



अगर एक language के intermediate code same हो तो
 उसके back end भी same हो जायेगा और front end अलग-अलग
 तरह (types) के हो सकते हैं इसके compiler construction पर
 effect पड़ेगा [compiler construction और efforts कम हो जाते हैं]

* Tokens \Rightarrow Tokens are terminal symbols of the source language.

Examples:- Identifiers, numbers, keywords, punctuations, symbols etc

Identifiers:- { Total, count, sum, main ... }

$$\text{index} = 2 * \text{count} + 17;$$

Constants:- { 10, 25, 50, 5, ... }

Numbers:- { 0, 1, 2, 3, ... }

Keywords:- { int, while, if, for, ... }

Punctuations:- { !, , ;, :, ?, - }

Operators:- { +, -, *, /, % }

Symbols:- { (,), [,], @, \$, *, ... }

Lexemes	tokens
index	identifier
=	operator
2, 17	constant
*	operator
count	identifier
+	operator
;	punctuations or delimiters

\Rightarrow Comments, space, tabs, blank, newline these are not tokens

* Pattern \Rightarrow Pattern is a rule describing all these lexemes that can represent a particular token in a source language.

Identifier rules:- Start with alphabet or _, followed by any alphanumeric character.

* Lexemes:- Lexemes are matched against the pattern

A specific instance of a token

{Count = count + temp;}

Tokens:- id, operator, punctuation (;

lexeme ;

* Advantage / Disadvantage of compiler ⇒

Advantage :- * fast in execution

* NO Bulky program (means simple program)

* NO Hardware knowledge required (register, capacitor)

* Portable (means s/w run on different machine)

C language program is portable but compiler is not
Portable

* Addressing :- Memory is referred by name not address

* Object executable code produced by a compiler can be distributed / executable without having (without compiler)

* Object code can be used without re-compilation.

Disadvantages :-

* Debugging a program is ~~not~~ much harder

* When an error is found then whole program has to be re-compiled. * more memory required.

* Advantage / Disadvantage of interpreter ⇒

Advantages :-

* Good for coding error in program

* Debugging is easier

* If an error is detected there is no need to ~~re-compile~~ re-translate (re-interpret) the whole program.

* It uses less memory because few lines (statements) have to be loaded into the memory if no object code generated.

Disadvantages :-

* Slow :- speed is the biggest disadvantage

* No object code is produced so translation has to be done every time when ~~ever~~ ever program is run.

* When program is run then interpreter must be present.

Advantage / Disadvantage of Assembler ⇒

Advantages :-

* You can access machine-dependent registers & I/O devices

* You can control the exact code behavior in critical sections

* You can produce Hand-optimized code

* You can build interfaces b/w code fragmentation & using low-level incompatible conventions

Disadvantages:-

- * Bulky program of user to write for each instruction
- * Sufficient knowledge of Hardware must be required.
- * Not Portable (small)

Difference b/w Compiler Vs interpreter compilers:-

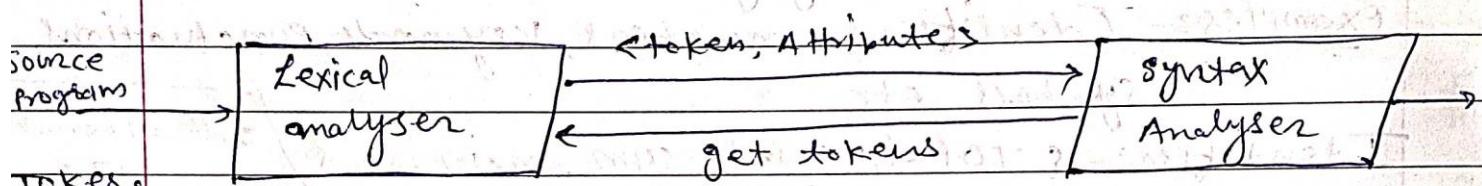
- (i) It translates whole program at a time
- (ii) It takes large amount of time to analyse the source code but the overall execution time is comparatively (अपेक्षित) faster.
- (iii) Intermediate object code is generated
- (iv) Debugging process is hard.
- (v) It requires more memory.
- (vi) It is used in C/C++ languages.
- (vii) It is complex program.

Interpreter:-

- (i) It translates statement by statement at a time (one statement at a time)
- (ii) It takes less amount of time to analyse the source code but the overall execution time is slow.
- (iii) No intermediate object code is generated
- (iv) Debugging process is easy
- (v) It requires less memory.
- (vi) It is used in Python/Ruby languages.
- (vii) It is easy program.

* Lexical analyser \Rightarrow Lexical analyser is the first phase of compiler. Its main task is to read (scan) the input symbols (characters) & produce as output a sequence of tokens that the parser uses for syntax analysis.

Some times, Lexical analysers are divided into a cascade of two phases :- The first phase is called "Scanning" & the second phase is called "Lexical analysis".

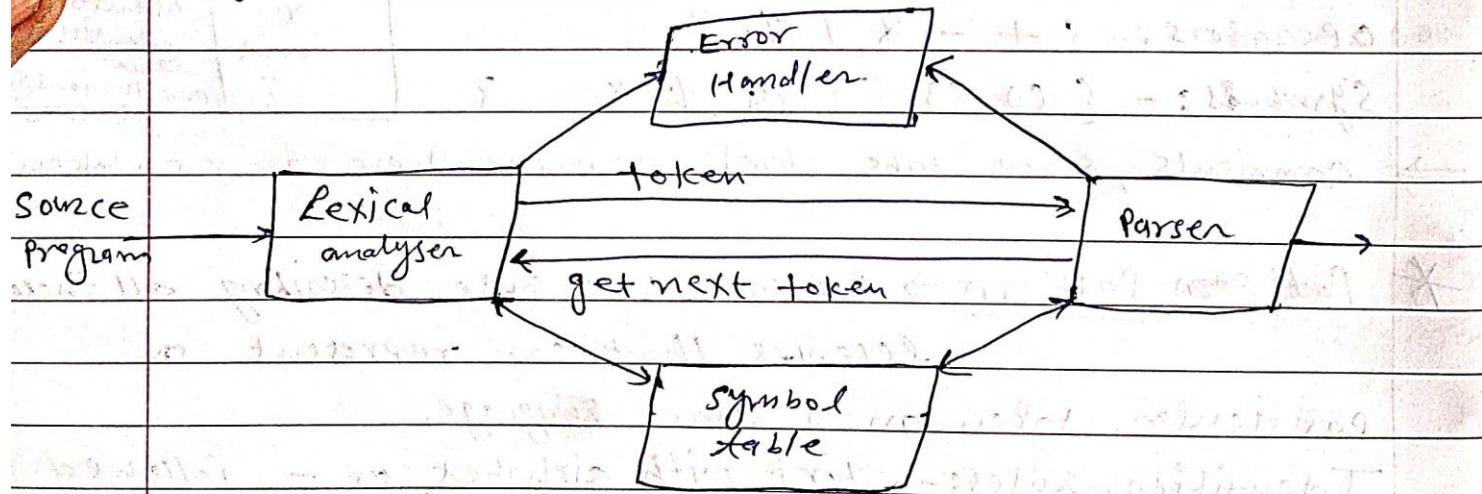


operators :- + / - * / / mod / %

keywords :- If | while | do | then | for |

Letters :- a | b | c | ... | z | A | B | C | ... | Z

Digits :- 0 | 1 | 2 | 3 | ... | 9

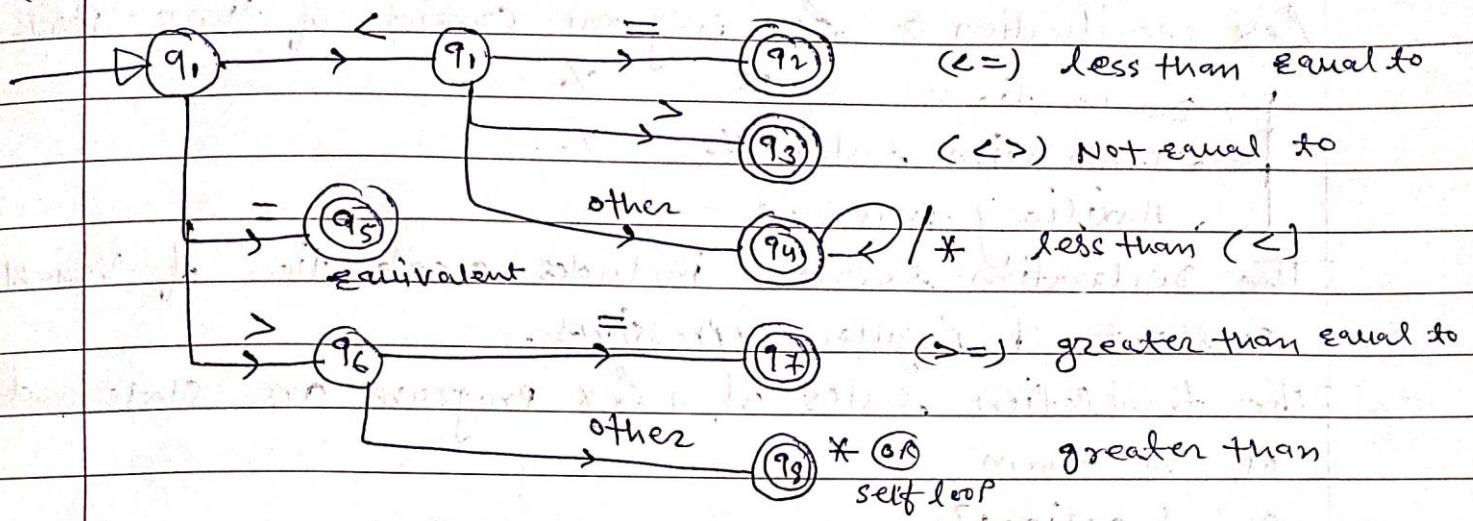


\rightarrow functions of Lexical analyser \Rightarrow

- Remove white space \Rightarrow White space of blank's tabs & new lines etc. appear b/w tokens they are detected & dropped by Lexical analyser.
- Remove comments :- like // or /* --- */ (comment) में को से किया होगा वो सिर्फ remove होगा।
- Handling of constants :- int, white, for, if, 10, 20, 5, 100 etc
- Handling of operators :- + / - * / / % / mod
- Recognition of identifiers & keywords :- Count, sum, main If, else, while, out, for

→ Recognition of tokens by Lexical analyser ⇒

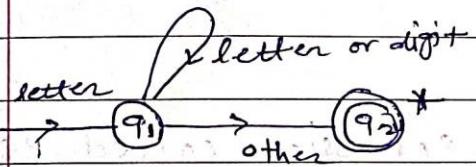
(i) Recognition of relation operators ⇒



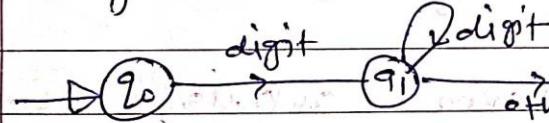
Operators relations :-

<, >, <=, >=, <>, =

(ii) Recognition of identifiers & keywords :-

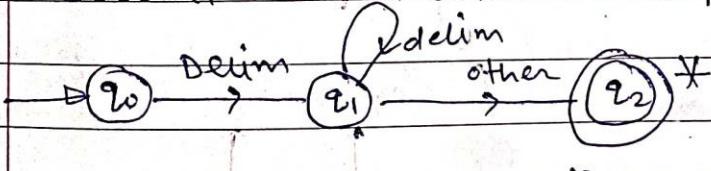


(iii) Recognition of number constants ⇒



(iv) Recognition of white space ⇒

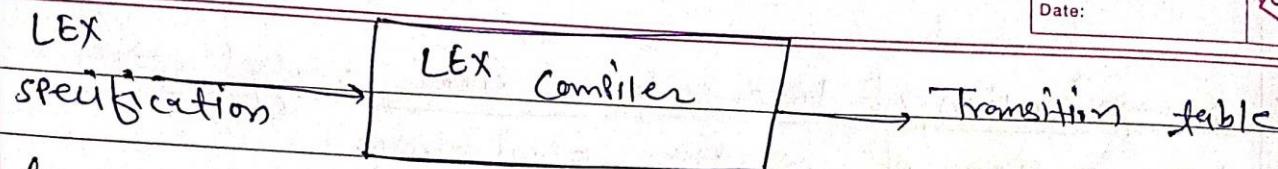
delimiter or Delim = blank / tab / newline



* Lex (LEX) Concept ⇒ Lex is a compiler writing tools which facilitated writing of lexical analyser & hence a compiler.

There are several tools that have been built for constructing a lexical analyser.

LEX has been widely used to specify lexical analysers for variety of languages. This tool is referred to LEX compiler & input specification is called Lex language.



Lex specification \Rightarrow Lex program consists of three parts:-

- Declaration % %.
- Translation rules %. %.
- Auxiliary procedures

The declaration section includes declaration of variables, constant & regular expressions.

The translation rules of a Lex program are statements of the form

P₁ { action-1 }

P₂ { action-2 }

P₃ { action-3 }

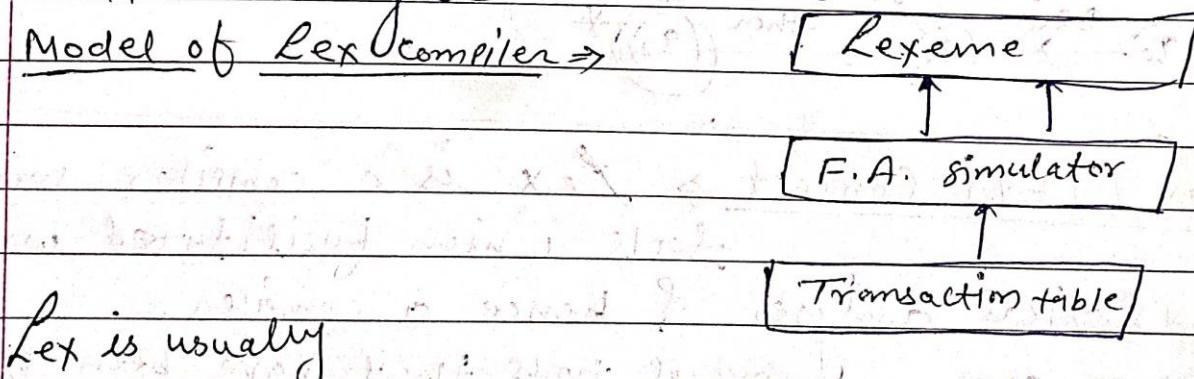
⋮

P_n { action-n }

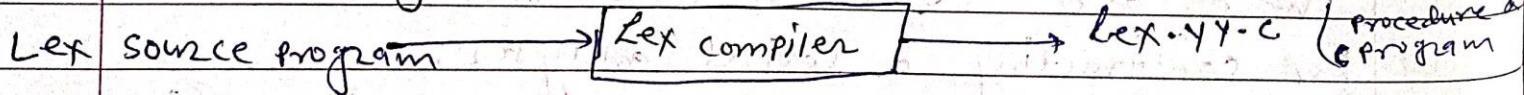
Where each P is a regular expression called pattern & each action is called program that is to be execute of

Third section holds whatever auxiliary procedures are needed by the actions. These procedures can be ~~compiled~~ compiled separately & located with lexical analyser

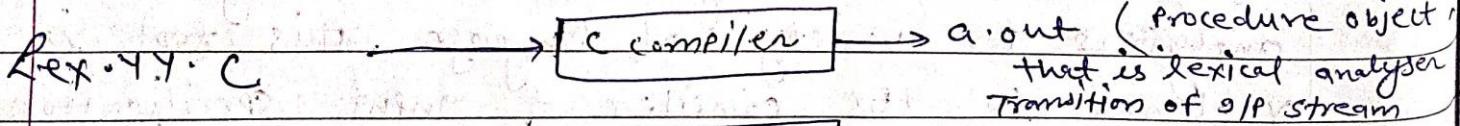
Model of Lex Compiler \Rightarrow



Lex is usually



Lex.yy.c



Input stream

a.out

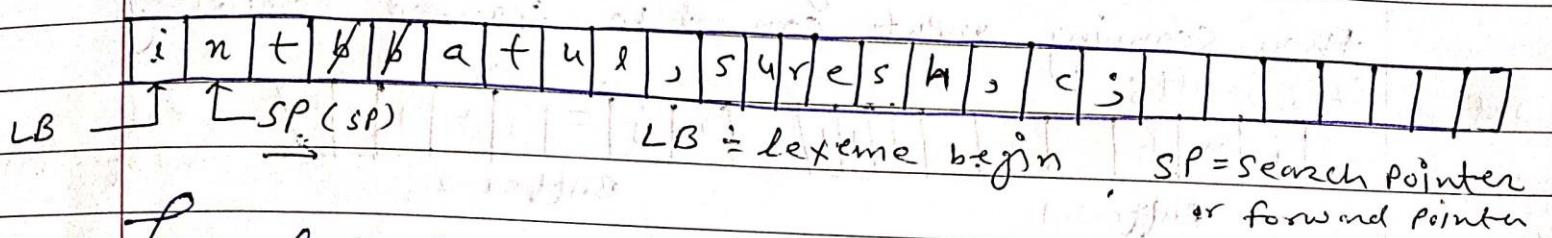
a.out (procedure object that is lexical analyser transition of I/O stream)

→ sequence of tokens generate

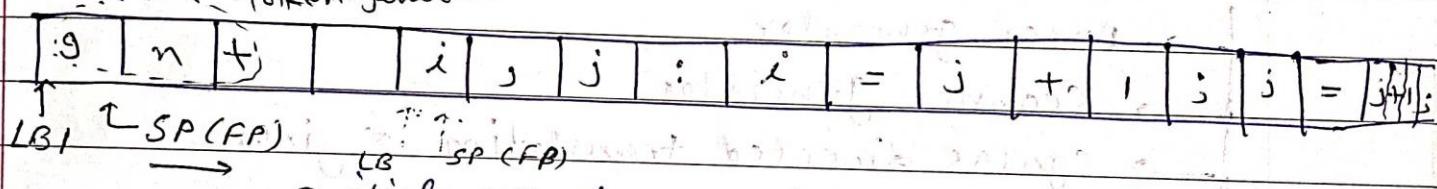
* Input Buffer / g/p Buffering \Rightarrow input buffer is done for increasing efficiency of the compiler.

A two buffer input scheme is useful when look ahead on the input is necessary to identify tokens.

Example:- gnt, atul, suresh, c;

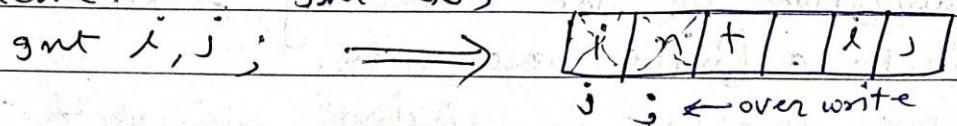


The Lexical analyzer scans the input from left to right one character at a time. It uses two pointers begin pointer (LB) & forward to keep track of the pointer of the input scanned.



Initial both pointers point to the first character of the input symbol (g/p string). Forward pointer moves ahead to search for the end of lexeme. As soon as the blank symbol (blank space) is encountered, it indicates end of lexeme.

\Rightarrow One (single) Buffer scheme \Rightarrow In this scheme, only one buffer is used to store the input string but the problem with this scheme is that if lexeme is very large (long) then if it crosses the buffer boundary, to scan rest of the lexeme the buffer has to be refilled that makes overwriting the first of lexeme.



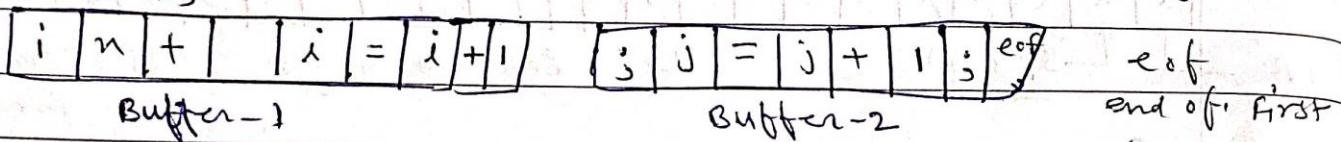
\Rightarrow Two buffer (double buffer) scheme \Rightarrow To overcome the problem of one buffer scheme, in this method two

buffers are used to store the input string.

The first buffer & second buffer are scanned alternately. When end of the current buffer is reached the other buffer is filled.

The only problem with this method is that if length of the lexeme is longer than length of the buffer then scanning input cannot be scanned completely.

int i=i+1; j=j+1;



* Construction tools for compiler \Rightarrow There are different types of tools used in compiler constructions are following as

- Parser generator
- Scanner generator
- Syntax directed translation engines
- Automatic code generators
- Data flow engines

* Boot strapping ^(OR) Bootstrap compiler \Rightarrow

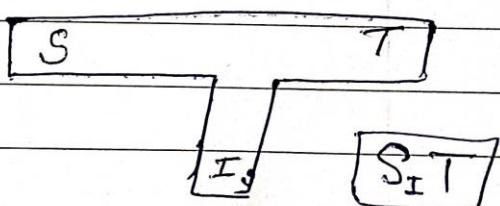
Compiler is a complex enough program that we would like to write it in a friendlier language than assembly language.

Bootstrapping is used to create compilers or to move them from one machine to another by modifying the backend. It is represented by an "T" diagram.

Bootstrapping is widely used in the compilation development. It is used to produce a self-hosting compiler. Self-hosting compiler means it is a type of compiler that can compile its own source code.

A compiler can be characterized by three languages:-

- Source language (S)
- Target language (T)
- Implementation lang. (I)

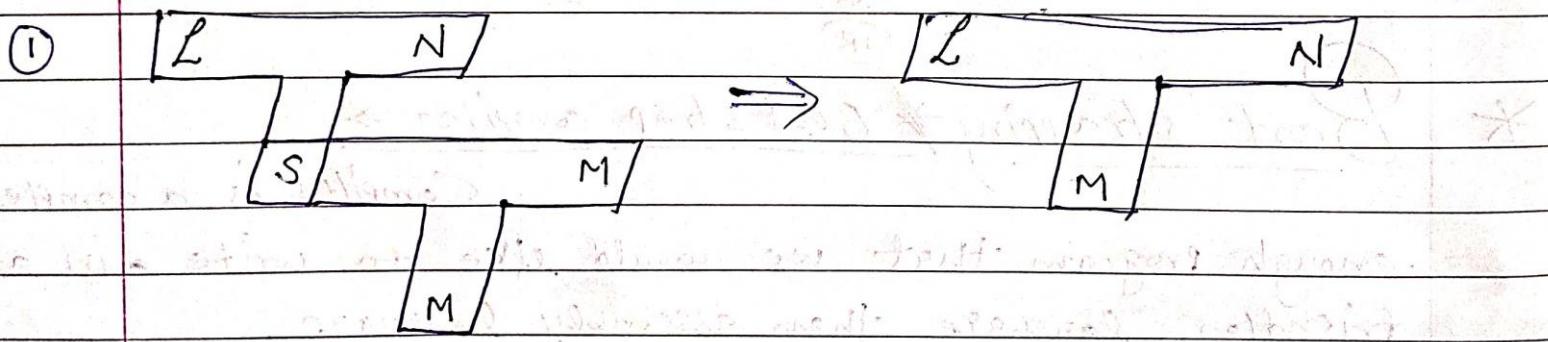


Bootstrapping is the process by which simple language is used to translate more complicated programs, which in turn may handle an even more complicated program soon.

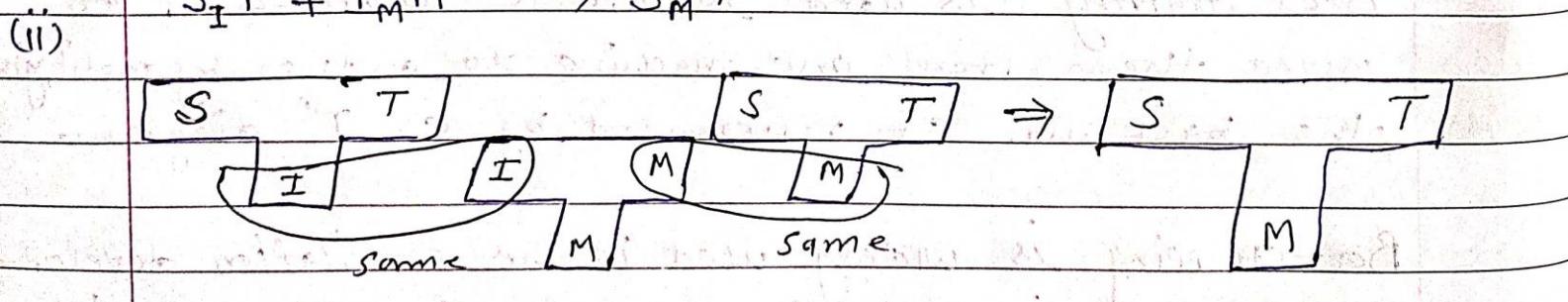
* Cross compiler \Rightarrow A compiler may run on one machine & produce target code for another machine such a compiler called cross compiler.

Example \Rightarrow suppose, we write a cross compiler for a new language L in implementation language S to generate code for machine N that we create S_N . If existing compiler for S runs on machine M & generates code for M it characterises by S_M . If L_N is run through S_M then we get a compiler L_{MN} .

$$L_N + S_M = L_{MN}$$



$$S_I T + T_M M \Rightarrow S_M T$$



- * Symbol table \Rightarrow symbol table is an important data structure created & maintained by compiler in order to keep track of semantics of variable. It stores information.
- \rightarrow It is built in lexical & syntax analysis phases.
- \rightarrow The information collected is collected by analysis phase & it used by synthesis phase.
- \rightarrow It is used by various phases of compiler to achieve compilation time efficiency.
- \rightarrow It is used by various phases of compiler:-
- * Lexical analysis:- creates new table, token entries in this table
- * Syntax analysis:- Add informations regarding attribute type, scope, dimensions, line of reference, use etc. in table.
- * Semantic analysis:- uses variables information in the table to check for semantics (type checking)
- * Intermediate code generation:- Refers symbol table for knowing how much & what type of run time is allocated & table helps in adding temporary variable information.
- * Code optimization:- uses information present in table for machine dependent optimization.
- * Target code generation:- Generates code by using address information of identifiers present in the table.

\Rightarrow Symbol table entries:- Each entry in symbol table is associated with attributes that support compiler in different phases.

Items stored in symbol table:-

Variable name & constants, procedure & function names, literal constants & strings, compiler generated temporaries, labels in source languages.

Information used by compiler from symbol table:-

Data type & name, declaring procedures, offset in storage, If structure or record then pointer to structure table, parameters, number & type etc.

Operation of symbol table:- According to data structure like; list, linked, arrays, hash table, binary search tree

Error Handling in compiler Design \Rightarrow The tasks of the error handling process are to detect each error, report it to the user & then make some recover strategy & implement them to handle error. During this whole process, An error is the blank entries in the symbol table.

-Types of error:- There are mainly two types (source) of error:- Runtime errors

compile time error

\rightarrow Runtime error:- Run time error is an error which takes place during the execution of a program. types:- Lack of sufficient memory, run an application, logical errors occurs when Executed code does not produce the expected result.

\rightarrow Compile time error:- compile time errors rises at compile time before execution of the program.

Compile time errors:-

- \rightarrow Lexical:- misspellings of identifiers
- \rightarrow Syntactical:- missing semi colon or unbalance parentheses
- \rightarrow Semantical :- incompatible value, assignment or type mismatch
- \rightarrow Logical:- code not reachable, no loop, b/w operand & operators

Error recovery methods:-

Panic mode recovery :- $(1+2)+3$

Phase level recovery:-

Error Productions

Global corrections