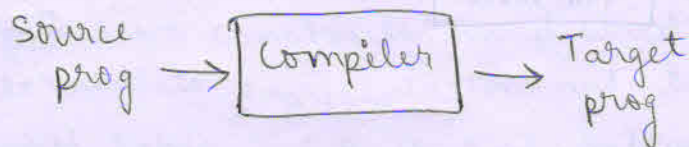# Compiler Design - 10CS63

## UNIT 1 : Introduction

Translator - Any program that converts a high level language program to Machine (Low Language) code.

Compiler - Program that reads code in one language i,e source code and translates it into another language i,e target language is a compiler.

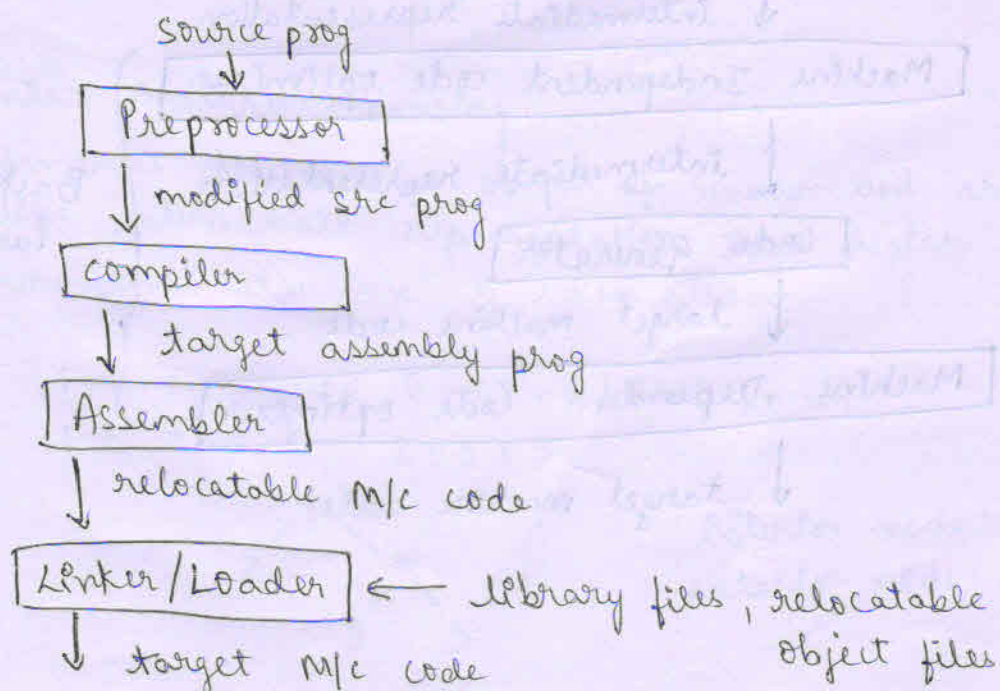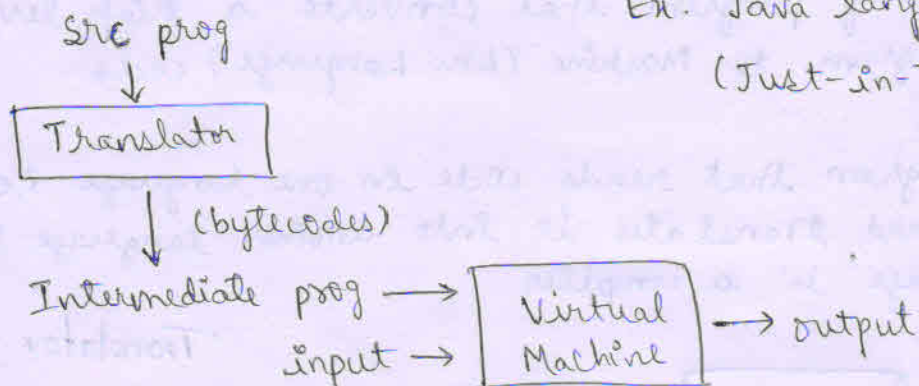Source prog → [Compiler] → Target prog

Translator

Interpretter - A kind of language processor which does not produce target program as a translation, but directly execute the operations specified in source program, on inputs supplied by the user.

Source prog →
Input → [Interpretter] → output
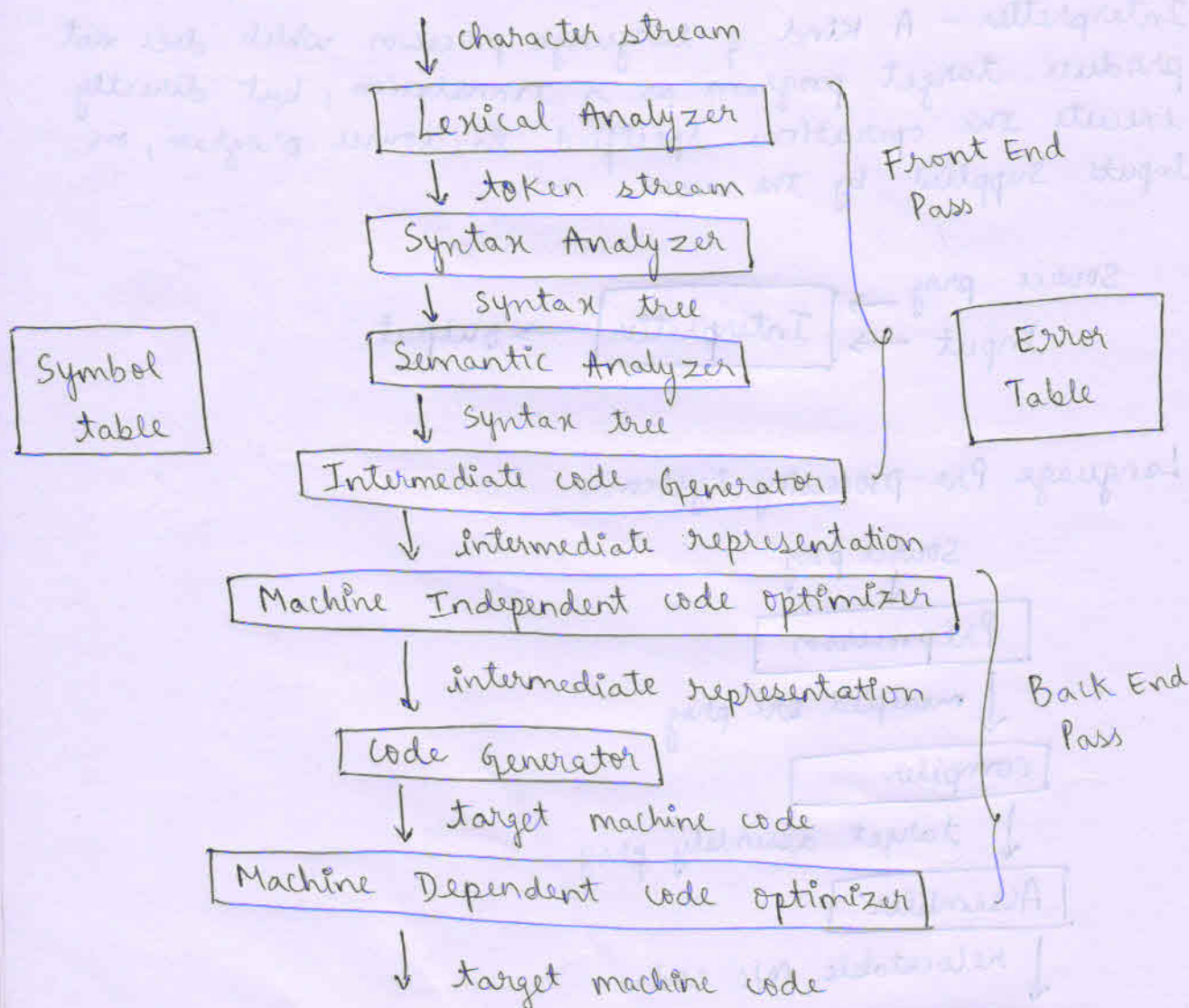
Language Pre-processing system:

Source prog
↓
[Preprocessor]
↓ modified src prog
[Compiler]
↓ target assembly prog
[Assembler]
↓ relocatable m/c code
[Linker/Loader] ← library files, relocatable object files
↓ target m/c code

○ A Hybrid compiler :

Src prog
↓
Translator
↓ (bytecodes)
Intermediate prog ⟶
input ⟶  Virtual Machine ⟶ output

Ex: Java lang processor (Just-in-Time)

Structure of a compiler :

↓ character stream
Lexical Analyzer
↓ token stream
Syntax Analyzer
↓ syntax tree
Semantic Analyzer
↓ syntax tree
Intermediate code Generator
↓ intermediate representation
Machine Independent code optimizer
↓ intermediate representation
Code Generator
↓ target machine code
Machine Dependent code optimizer
↓ target machine code

Symbol table

Error Table

Front End Pass

Back End Pass

→ 2 main parts:

① Analysis — breaks up source prog into constituent pieces & imposes grammatical structure on them. Based on the structure it creates intermediate representation of source prog. Collects information about prog, stores it in the "Symbol Table". (Front End of compiler)

② Synthesis — constructs the desired target prog from the intermediate representation and information in the symbol table. (Back End of compiler)

→ 7 phases:

① Lexical Analysis — Scanning
- On reading character stream of src prog, it groups them into meaningful sequences called "Lexemes".
- for each lexeme, Analyser produces as output a token of form:

  < token-name , attribute value >
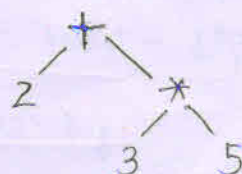
  abstract symbol ⌐                    ⌐→ points to entry in the
  used in parser ←                        symbol table for this token

② Syntax Analysis — Parsing
- parser uses tokens i,e output of scanner and creates a tree like intermediate representation that depicts the "Grammatical Structure" of token stream

  Ex: For Grammar  E → E+E | E*E | num
      For Input    2+3*5



  interior node: operators
  exterior node: arguments

③ Semantic Analysis

- uses syntax tree and information in symbol table to check source prog for semantic (meaning) consistency with lang. definition.

- It gathers type information and saves it in either syntax tree or symbol table for use in ICG.

- Type checking - compiler checks whether each operator has the matching operands

- coercions - Lang specification may permit some type conversion

④ Intermediate Code Generation (ICG)

- The Intermediary code during processing may be in the form of syntax tree or reduced form of source code.

- properties :
  → should be easy to produce
  → should be easy to translate into target M/C.

⑤ Code Optimization (M/C independent)
- to improve intermediate code to get better target code
→ Better in terms of : faster, shorter, less power consuming code
- Instead of using int to float operation, replace integer by its floating-point value <u>directly</u>

⑥ Code Generation
- input from intermediate representation maps to target lang.
- If target lang is M/c code - the instructions are translated into sequences of M/c instruction to perform same task.
- Judicious assignment of registers to hold variables is done.

→ Compiler construction Tools :

① Parser Generators – automatically produce syntax analyzers from a grammatical description of a prog lang.

② Scanner Generators – produce lexical analyzers from a regular expression description of tokens of lang.

③ Syntax directed translation engines – produce collections of routines for walking a parse tree and generate ICG.

④ code generator generators – produce CG from collection of rules for translating each operation of Intermediate lang into M/c lang for a target M/c.

⑤ Data flow analysis engines – facilitate gathering of data about how values are transmitted from one part of prog to every other part.

⑥ compiler construction toolkits – provide integrated set of routines for constructing various compiler phases.

Application of compiler Technology :

① Implementation of high level prog. lang using modern OOPS concept like,
→ Data Abstraction
→ Inheritance properties

② optimization for computer architectures
→ Parallelism
(i) at instruction level – multiple operations executed together
(ii) at preprocessor level – different threads run seperately.

→ Memory Hierarchy
Building very <u>Large</u> or <u>Fast</u> storage, but not both

③ Design New computer Architectures

→ RISC - reduces complex memory addressing, support data structure access, procedure invocation ....

→ Specialized Architectures -
   Data flow M/C, vector M/C, VLIW & SIMD M/C.

④ Program Translations

(i) Binary Translation - Increases s/w availability

(ii) Hardware Synthesis - Verilog, VHDL - reduces time & effort

(iii) Database Query Interpretter - SQL queries effective retrieval

(iv) compiled simulation - model run, to validate design.

(v) Reduce redundancy in code

⑤ Software Productivity Tools

(i) Type checking - to catch program inconsistency

(ii) Bounds checking - Lang. provides range checking like for the buffer overflow, security, optimize range check, sophisticated analyses, error detection tools.
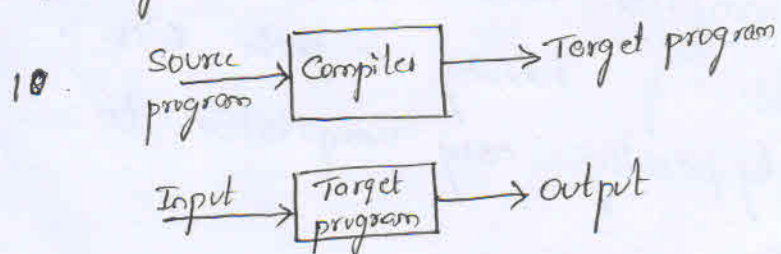
(iii) Memory Management Tools - (Garbage collection)
   ○ Automatic memory management tracks all memory related errors - leaks...

1. Write the difference between Compiler and Interpreter

| Compiler | Interpreter |
| --- | --- |
| 1. Compiler translates the entire program in one go and then executes it | 1. Interpreter first converts high level language into an intermediate code and then executes it line by line. The intermediate code is executed by another program |
| 2. It produces efficient object code therefore programs runs faster | 2. No intermediate object code is generated |
| 3. Error reporting is time consuming (displayed after entire pgm is checked) | 3. Errors are displayed for every instruction interpreted if any (Error reporting is immediate) |
| 4. Conditional control statements are executed faster | 4. Conditional control statements executed slower |
| 5. Memory requirement is more ∵ single object code is generated | 5. Memory requirement is less |
| 6. Program need not be compiled every time | 6. Everytime high level program is converted into lower level pgm |
| 7. Difficult to use | 7. Easy to use for beginners |
| 8. Translate once and then run the result (stand-alone code, faster execⁿ) | 8. read - check - execute loop → slower, not stand-alone |
| 9. Eg :- c, c++ | 9. Eg :- python, prolog |

10.

Source program → [Compiler] → Target program

Input → [Target program] → output

10. Source pgm Input → [Interpreter] → output

→ Examples showing detail phases of compiler :

① position = initial + rate * 60

$\downarrow$

| Lexical Analysis |

$\langle id,1 \rangle \langle = \rangle \langle id,2 \rangle \langle + \rangle \langle id,3 \rangle \langle * \rangle \langle 60 \rangle$
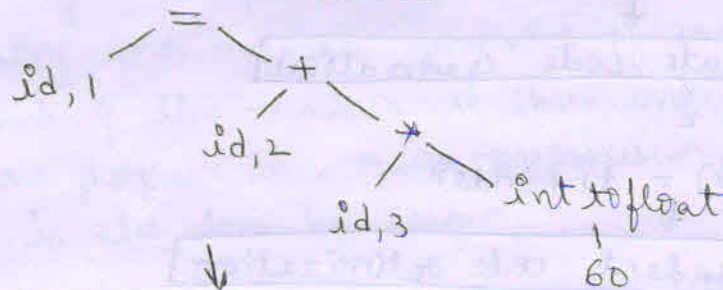
$\downarrow$

| Syntax Analysis |

```
        =
      /   \
   id,1    +
         /   \
      id,2    *
             /  \
          id,3   60
```

Syntax Tree

| | |
|---|---|
| 1 | position |
| 2 | initial |
| 3 | rate |

$\downarrow$

| Semantic Analysis |

```
        =
      /   \
   id,1    +
         /   \
      id,2    *
             /  \
          id,3   inttofloat
                    |
                    60
```

$\downarrow$

| Intermediate code Generation |

$t_1 = \text{int tofloat}(60)$
$t_2 = id_3 * t_1$
$t_3 = id_2 + t_2$
$id_1 = t_3$

$\downarrow$

| M/c independent code optimization |

$t_1 = id_3 * 60.0$
$id_1 = id_2 + t_1$

$\longrightarrow$ | Code Generation |

LDF   R1, id3
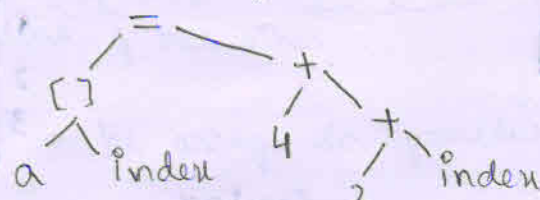MULF  R1, R1, #60.0
LDF   R2, id2
ADDF  R2, R2, R1
STR   id1, R2

②      $a[index] = 4 + 2 + index$

↓

| Lexical Analysis |

$\langle id, 1 \rangle \langle [ \rangle \langle id, 2 \rangle \langle ] \rangle \langle = \rangle \langle 4 \rangle \langle + \rangle \langle 2 \rangle \langle + \rangle \langle id, 2 \rangle$

↓

| Syntax Analysis |



|   |       |   |
|---|-------|---|
| 1 | a     |   |
| 2 | index |   |

| Semantic Analysis |



↓

| Intermediate code Generation |

$t_1 = 4 + 2$

$a[index] = t_1 + index$

↓

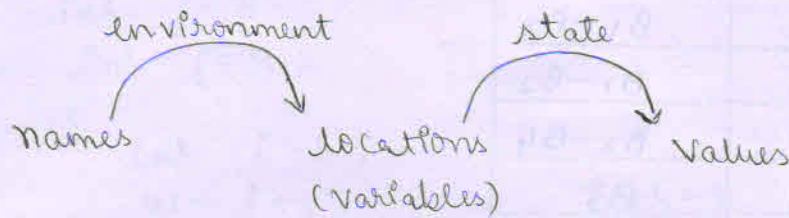| MIC Independent code optimization |

$a[index] = 6 + index$

↓

| code Generation |

```
mov  index, R0      // R0 = index
mov  &a , R1        // R1 = starting address of array a
add  R0, R1         // R1 = R0+R1
mov  #6 , R2        // R2 = 6
add  R1, R2         // R2 = R1+R2 = R1+6
mov  R2, &R1        // store R2 in &R1 i,e & a's value
```

→ Environments and States:

environment       state

names     locations     values
       (variables)

- Environment is mapping from names to locations in the store
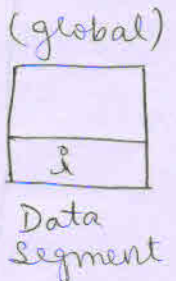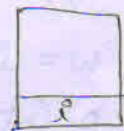- State is mapping from locations in store to their values.

Dynamic Mapping Exceptions:

(i) Static Binding of Names to Locations – global variable
declaration – location in store once for all.

Ex:      int i;        // global i              (global)

       void fun (...) {
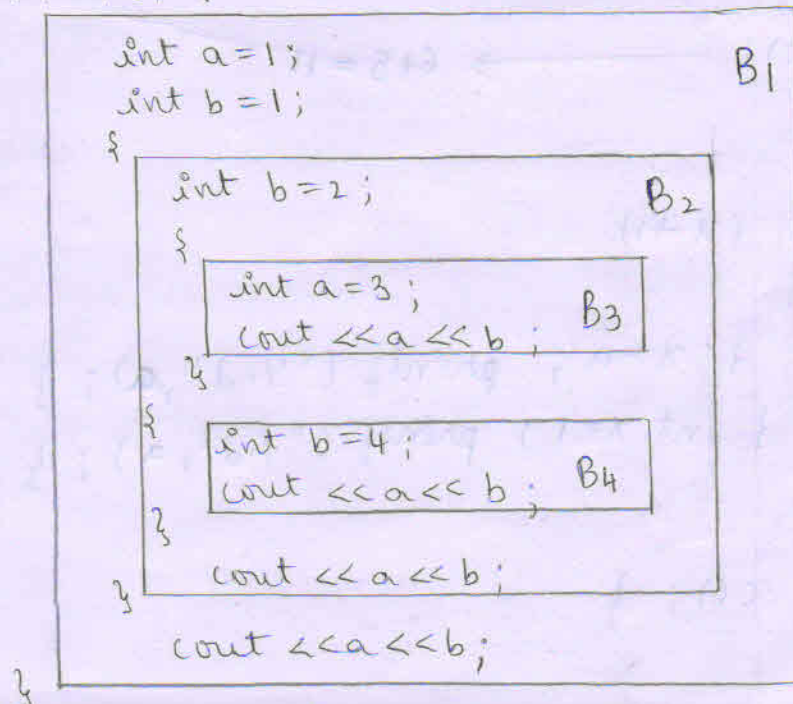          int i;       // local i
       }

Data Segment

(ii) Static Binding of Locations to Values – declared constants

Ex:     # define ARRAYSIZE 1000     // static bind

Static Scope and Block Structure:

```
(1) main () {
        int a = 1;                                      B1
        int b = 1;
        {
            int b = 2;                          B2
            {
                int a = 3;
                cout << a << b;         B3
            }
            {
                int b = 4;
                cout << a << b;     B4
            }
            cout << a << b;
        }
        cout << a << b;
    }
```

→

| Declaration | Scope |
|---|---|
| int a=1 ; | B1-B3 |
| int b=1 ; | B1-B2 |
| int b=2 ; | B2-B4 |
| int a=3 ; | B3 |
| int b=4 ; | B4 |

② main ()
  {
  int w,x,y,z;
  int i=4 ; int j=5;

```
    {
        int j=7; i=6;
        w = i+j ;
        printf (w);
    }
```
→ 6+7 =13

  x = i+j ;
  printf (x);  ─────────────→ 6+5 = 11

```
    {
        int i=8 ;
        y = i+j;
        printf (y);
    }
```
→ 8+5 = 13

  z = i+j ;
  printf (z);  ─────────→ 6+5 = 11
  }

③  # define  a  (x+1)
   int x= 2;
   void b()   { x=a; printf ("%d",a); }
   void c()   { int x=1; printf ("%d",a) ; }
   void main ()
   { b() ;   c(); }

3
2

④
```
int  w, x, y, z ;
int  i = 3 ;
int  j = 4 ;
{
    int  i = 5 ;
    w = i + j ;                    5 + 4 = 9
}
    x = i + j ;                    3 + 4 = 7
    {
        int  j = 6 ;
        i = 7 ;
        y = i + j ;                7 + 6 = 13
    }
    z = i + j ;                    7 + 4 = 11
```

10. What is printed by the following C code.

a)
```
#define a (x+1)

int x=2;
void b() { x=a; printf("%d\n", a); }        → 3
void c() { int a=1; printf("%d\n", a); }     → 2
void main() { b(); c(); }
```

b)
```
#define a (x+1)

int x=2;
void b() { x=a; printf("%d\n", a); }      → 3
void c() { printf("%d\n", a); }           → 4
void main() { b(); c(); }
```
∴ redesignment for the same variable x
x=3, a=3+1=4

c)
```
#define a (x+1)

int x=2;
void b() { int a=1; printf("%d\n", a); }    → 2
void c() { printf("%d\n", a); }             → 3
void main() { b(); c(); }
```

d)
```
#define a (x+1)

int x=2;
void b() { int x=a; printf("%d\n", a); }    → 4  ∴( x=2+1=3
void c() { printf("%d\n", a); }             → 4      again a=3+1=4)
void main() { b(); c(); }
```

→ Parameter Passing Mechanisms :

(1) Actual parameters – parameters used in call of procedure

(2) Formal parameters – parameters used in procedure definition

① Call by value

- The actual parameter is evaluated (if an expression) or copied (if a variable) ; The value is placed in the location belonging to corresponding formal parameter of called procedure.

- It has all computations involving formal parameter done by called procedure is <u>local</u> to that procedure.

② Call by reference

- The address of actual parameter is passed to the callee as the value of corresponding formal parameter

- Uses of formal parameter in code of callee are implemented by following this pointer to location indicated by caller.

- changes to formal parameter ⟹ Appear as changes in actual param

- If actual parameter is expression, it is evaluated before the call and it's value stored in a location of its own.

- changes to formal parameter change value in this location, But – No effect on data of caller.

③ Call by name

- used in early prog – Algol 60.

- it requires callee execute as if actual parameter were substituted literally for formal parameter in the code of the callee as if formal parameter were macro standing for the actual parameter.

→ Examples :

① call by value

```
int add ( int a, int b)
{
    return (a+b);
}
main ()
{
    :
    c = add (10,20)
    :
}
```

② call by reference

```
int  add ( int *a, int *b)
{
    return (a+b);
}
main ()
{
    :
    int p=10;
    int q=20;
    c = add ( &p, &q);
    :
}
```

③ call by Name - Aliasing

```
int add (int a , int b)
{
    return (a+b);
}
main ()
{
    int p=10;
    int q=20;
    c= add ( &p, &q);
    :
}
```

- Aliasing :
→ Interesting consequence of call by reference parameter passing where references to objects are passed by value.
- It is possible that two formal parameters refer to the same location → Such variables are ALIAS to one another.
- Though they may be <u>distinct</u> formal parameters, they may be Alias of one other.

<u>Ex:</u>   Let a be array in procedure p

P
{

‖ q(x,y)   call
   q(a,a) ;

}

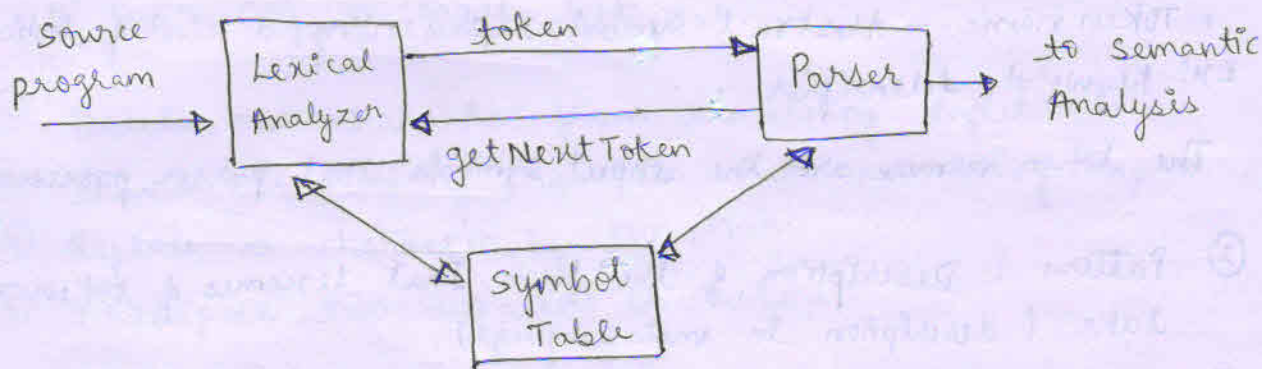array names are references to location ⟹ Alias

$x[i] = y[i]$

# Chapter-1 Introduction

Questions

1. Define Compilers?
2. Differentiate b/w compilers & Interpreter?
3. Explain the long processor System?
4. Describe the analysis-Synthesis model of the Compiler or Explain in detail the Various phases of Compiler with an example?
5. Explain in detail the Various phases of Compilation for the i/p string

   a. $P = i + n * 60$
   
   b. $x = a * b + a * b$
   
   c. $a = (b+c) * (b+c) * 2$
   
   d. $a[index] = 4 + 2 + index$

6. Why is it necessary to group phases of Compiler
7. What Is the purpose of Compiler Const$^n$ tool. Describe the different Compiler Construction tool we used?
8. Analyse the s/w productivity toal and explain
9. Explain the different parameter passing technique with an example?

# Chapter-3 - Lexical Analysis

→ Lexical Analysis :

Interaction between Lexer and Parser :

Source program → Lexical Analyzer ⇄ (token / get Next Token) ⇄ Parser → to semantic Analysis

Lexical Analyzer ⇄ Symbol Table ⇄ Parser

→ Task of Lexer :

① Identification of Lexemes
② Stripping out comments
③ Removing whitespace ( blank, \n, \t)
④ corelating error messages generated by compiler
⑤ Keep track of line numbers to show error
⑥ If source program uses Macro-preprocessor, the expansion of macros is also done by scanner.

→ Lexer - cascade of 2 processes :

① Scanning consists of simple processes that do not require tokenization of input, such as deletion of comments & compaction of consecutive whitespace characters into one.

② Lexical Analysis proper in more complex portion, where scanner produces sequence of tokens as output.

Lexer versus Parser : Seperate phases because :

① simplicity of design - important consideration
   (compiler)
② compiler efficiency improved - use specialized technique for lexical Analysis (Input Buffering)
③ compiler portability is enhanced.

→ Tokens, Patterns, Lexemes:

① Token: A pair consisting of a token name and an optional attribute value.

• Token name — Abstract symbol representing a kind of lexical unit Ex: Keyword, identifier ....

The token names are the input symbols that parser processes.

② Pattern: Description of the form that lexemes of token may take (description in meta language).

Ex:   token name: identifier
   pattern: $[\_a-zA-Z]^+ [a-zA-Z0-9]^*$

③ Lexeme: Sequence of characters in source program that <u>Matches</u> the pattern of a token and is identified by the lexer as an instance of that token.

Ex:   token name: Keyword
   pattern: $[i][f]$
   lexeme: if

| Token | Informal Description | Sample Lexemes |
|---|---|---|
| if | characters i, f | if |
| else | characters e, l, s, e | else |
| comparision | $<, >, <=, >=, ==, b=$ | $<=, <>$ |
| id | letter followed by letter and digits | pi, score |
| number | numeric constants | 3.14, 0, 6.9e8. |
| literal | enclosed within " " | "core dumped" |

→ Lexical Errors : Recovery options

① Panic mode recovery - delete successive characters from remaining input until lexer finds well known token at beginning of input left out.

② Delete one character from remaining input
③ Insert one missing character into remaining input
④ Replace a character by the other
⑤ Transpose two adjacent characters.

Examples :

$$fi\ (\ a<b) \implies if\ (a<b)$$
$$int\ a,\ ; \implies int\ a\ ;\ \ or\ int\ a,b\ ;$$

→ Input Buffering : To speed up reading of src prog.

① Single Buffer | 1-Buffer Technique
We use only one single buffer to store processed character from large no. of characters from source prog.
Main overhead is that if,

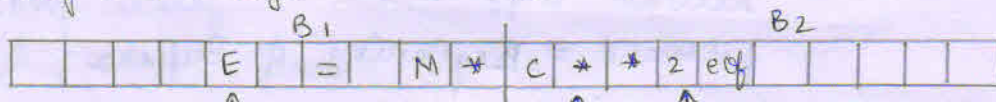| lexeme size > Buffer size |

we loose the lexeme

World
5 Bytes

$$\boxed{W}\ \boxed{0}\ \boxed{r}\ \boxed{l}^{d}$$

4 Bytes

. It reloads data, removes old data.

② 2-Buffer Technique ⟨ without sentinel / with sentinel
We use two buffers that are alternately reloaded,
Each buffer of same size N, N = Size of a disk block. (4096 Byte)
. Using read system call, N characters are read.

B₁                B₂
| | | | | E | = | | M | * | c | * | * | 2 | eof | | | | | |
    ↑lexeme begin    ↑lexeme begin  ↑forward    i/p < Buffsize

→ special char eof marks end of src file and this char is different from any other char of src prog.
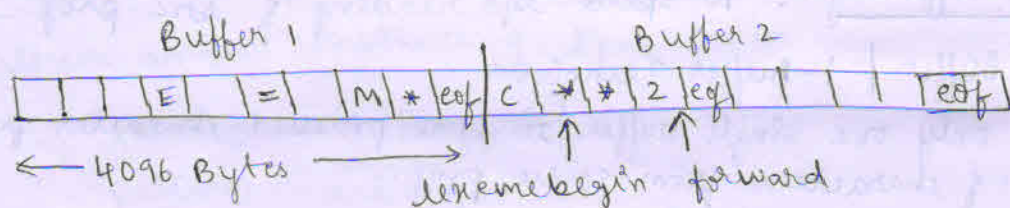
→ Two pointers maintained :

① Lexeme Begin — marks beginning of current lexeme whose extent we are attempting to determine.

② Forward — scans ahead until a pattern match is found When forward reaches end of next lexeme, ** we retract one position back and return token.

- We need 2 checks in 2 Buffer without sentinels :

1) Advancing forward requires whether we reached the end of one of the buffer, if Yes Reload other buffer and make forward point to newly loaded buffer beginning.

2) Before returning token check whether valid or not.

→ Sentinels : (2 Buffer technique with Sentinels)
Using sentinel character at the end which is a special char that is not part of src prog  (usually eof)

Buffer 1                          Buffer 2

| | | | E | | = | | M | * | eof | C | * | * | 2 | eof | | | | | eof |

← 4096 Bytes →      lexemebegin   forward

Here check if reached end of Buffer or not.
Look Ahead is atmost 1 char, make previous char as returned valid token.

→ Look Ahead code with sentinel :

```
Switch ( * forward ++ )
{
    case eof : if (forward is at end of Buffer 1) {
                    reload Buffer2;
                    forward = Beginning of Buffer2; }
               else if (forward is at end of Buffer 2) {
                    reload Buffer1;
                    forward = Beginning of Buffer1; }
```

```
else    /* eof within a Buffer marks end of input */
        terminate lexical Analysis
    break;
    cases for other char
}
```

———o ——o ——o ——o ——o ——

(1) Alphabet – finite set of symbols    Ex: $\Sigma = \{0,1\}$
string – finite sequence of symbols from $\Sigma$    Ex: 0101
Language – countable set of strings over $\Sigma$.

(2) Prefix of string – string obtained by removing zero or more
symbols from end of string.
Ex: ban, banana, $\epsilon$ are prefixes of banana.

(3) Suffix of string – string obtained by removing zero or more
symbols from beginning of string.
Ex: nana, banana, $\epsilon$ are suffixes of banana

(4) Substring – string obtained by deleting any prefix and
any suffix from string.
Ex: banana, nan, $\epsilon$ are substrings of banana

(5) Proper prefix – prefixes, which is not $\epsilon$ or equal to string
Ex: ban, banan

(6) Proper suffix – suffix which is not $\epsilon$ or equal to string itself
Ex: anana, na

(7) Proper substring – substring from string which is not $\epsilon$ or the
string itself
Ex: anan, banan, anana

Subsequence – string formed by deleting zero or more not
necessarily consecutive positions of string.
Ex: baan, anaa -- for banana

→ Operations on Languages :

| operation | definition & notation |
|---|---|
| Union of L & M | $L \cup M = \{ s \mid s \text{ is in } L \text{ or } s \text{ is in } M \}$ |
| concatenation of L & M | $LM = \{ st \mid s \text{ is in } L \text{ and } t \text{ is in } M \}$ |
| Kleene closure of L | $L^* = \bigcup_{i=0}^{\infty} L^i$ |
| Positive closure of L | $L^+ = \bigcup_{i=1}^{\infty} L^i$ |

→ Regular Definition :
For some alphabet set $\Sigma$, sequence of regular definition :

$$d_1 \to r_1$$
$$d_2 \to r_2$$
$$\vdots$$
$$d_n \to r_n \qquad \text{where}$$

1) each $d_i$ is new symbol ( not in $\Sigma$ & other $d_i$ )

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

Ex ① C identifiers :

letter $\to$ A | B | --- | Z | a | b | --- | Z | _

digit $\to$ 0 | 1 | ---- | 9 |

id $\to$ letter ( letter | digit )*

Ex ② unsigned numbers :

digit $\to$ 0 | 1 | ---- | 9 |

digits $\to$ digit digit*

optional fraction $\to$ . digits | $\epsilon$

optional exponent $\to$ ( E ( + | - | $\epsilon$ ) digits ) | $\epsilon$

number $\to$ digits optional fraction optional exponent

# Algebraic Laws for Regular Expressions

| LAW | DESCRIPTION |
|---|---|
| 1. $r\|s = s\|r$ | $\|$ is commutative |
| 2. $r\|(s\|t) = (r\|s)\|t$ | $\|$ is associative |
| 3. $r(st) = (rs)t$ | Concatenation is associative |
| 4. $r(s\|t) = rs\|rt$ ; $(s\|t)r = sr\|tr$ | Concatenation distributes over $\|$ |
| 5. $\epsilon r = r\epsilon = r$ | $\epsilon$ is the identity for concatenation |
| 6. $r* = (r\|\epsilon)*$ | $\epsilon$ is guaranteed in a closure |
| 7. $r** = r*$ | $*$ is idempotent |

→ Recognition of Tokens :

stm → if expr then stmt | if expr then stmt else stmt | ∈
expr → term relop term | term
term → id | number
where,

number → $[0-9]^+ (. [0-9]^+)? (E [+-]? [0-9]^+)?$
id → $[\_a-zA-Z] [\_a-zAZ0-9]^+$
if → if
then → then
else → else
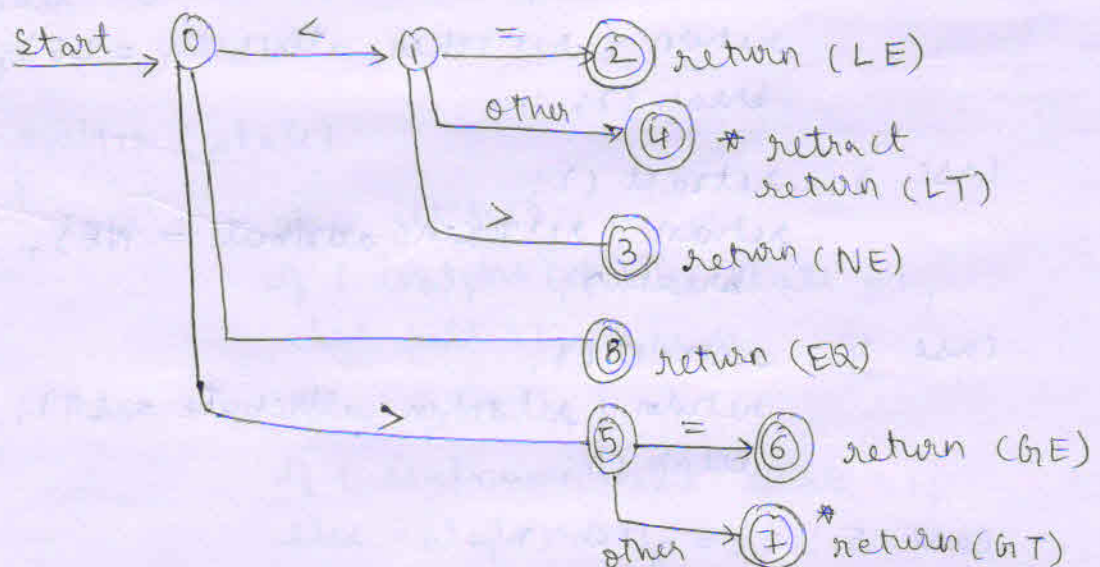relop → < | > | <= | >= | = | <>


white space :    ws → (blank | tab | newline)$^+$


→ Transition diagram :

① For relational operator, regular definition is
relop → < | > | <= | >= | = | <>

```
Code:
    state = 0;
    TOKEN getrelop ()
    {
        TOKEN retTOKEN = new (relop);
        while (1)
        {
            switch (state)
            {
                case 0:  c = newchar();   or   c = getch();
                if ( c == '<' ) state = 1;
                else if ( c == '>') state = 5;
                else if ( c == '=')  state = 8;
                else fail ();   break;

            case 1:  c = getch();
            if ( c == '=')  state = 2;
            else if ( c == '>')  state = 3;
            else if ( c == '...')  state = 4;  // other
            else fail ();  break;

            case 2:  retract ();
                    return ( retTOKEN. attribute = LE);
                    break ();

            case 3:  retract ();
                    return ( retTOKEN. attribute = NE);
                    break ();

            case 4:  retract ();
                    return ( retTOKEN. attribute = LT);
                    break ();

            case 5:  c = getch ();
                if ( c == '=')  state = 6;
                else if ( c == '...')  state = 7;  // other
                else fail ();  break ();
```

```
case 6 :  retract ();
          return ( retToken . attribute = GE);
          break ;

case 7 :  retract ();
          return ( retToken . attribute = GT);
          break ;

case 8 :  retract ();
          return ( retToken . attribute = EQ);
          break ;
        }
     }
  }
```

② for identifier

```
letter → [a-z A-Z _]
digit → [0-9]
id → letter ( letter | digit ) *
```



```
state = 0;
for (;;)
{
  switch (state)
  {
    case 0 :  ch = getch ();
              if ( isalpha (ch) )   state = 1;
              else fail ();   break;

    case 1 :  ch = getch ();
              if ( isalnum (ch) )   state = 1;
              else   state = 2;
              break;

    case 2 :  retract ();
              Install ID ();
              return ( retToken );
              break;
            }
         }
      }
```
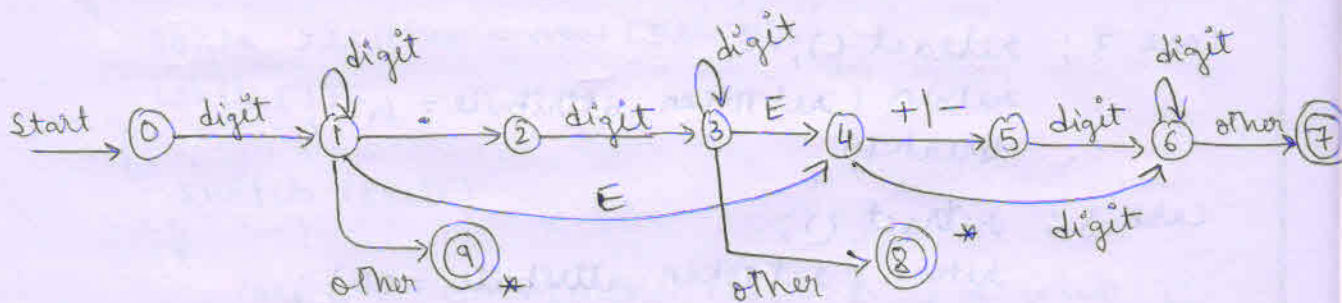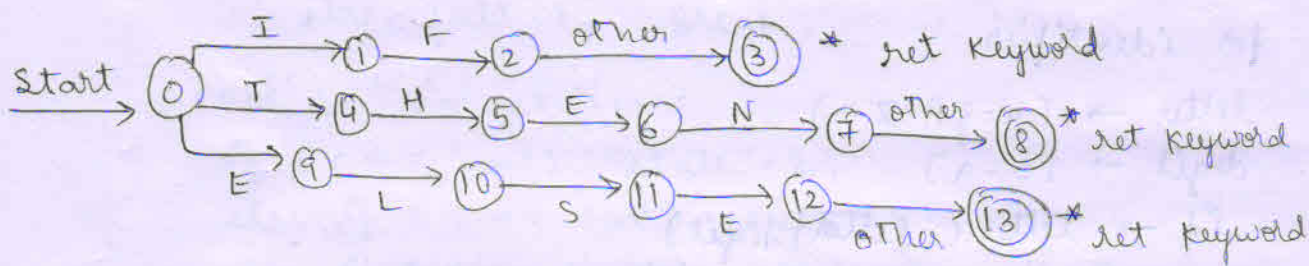
③ unsigned number

digit → [0-9]



④ Keywords

Ex:      Keyword → IF | THEN | ELSE



⑤ Delimiter / Whitespace

delim → space | tab | newline

# chapter-2   Lexical Analysis

Questions

1. Explain Lexical Analysis in detail with block diagram

2. Explain the reason for separating analysis phase of compiler for lexical Analysis and Syntax Analysis

3. What do you mean by lexical errors? How do we recover them

4. Define the terms token, pattern, lexeme with an example.

5. Why 2-buffer technique is used in LA? with an algorithm for lookahead code with Sentinel.

6. Give the formal definitions for operations on languages with notations

7. list the algebraic laws for Regular Expression.

8. Define the term prefix, Suffix, substring, proper prefix, proper suffix, proper substring, subsequence with an example.

9. write regular definition for Identifiers, unsigned numbers, keywords, relational operators and whitespace

10. Draw the transition diagram for relop, identifiers, unsigned number, keywords and white spaces.