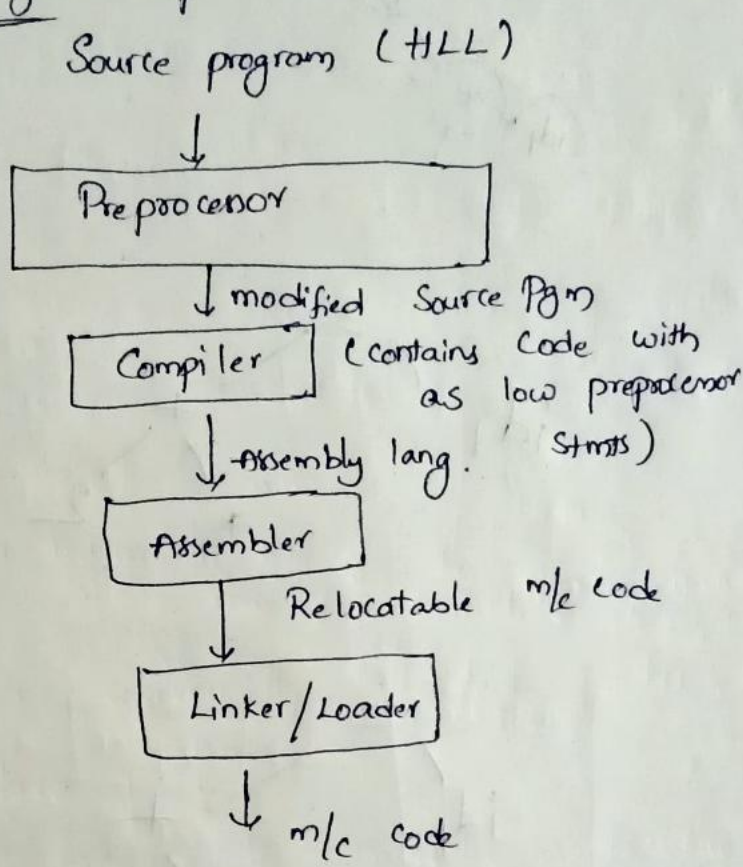




low level lang. eg: Assembly code (os, is).

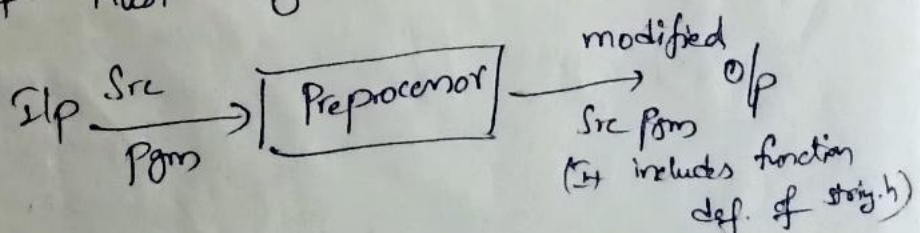
Lang. processor system steps:-



Preprocessor:- begins The stmts which are processed before compiler are known as

Preprocessor stmts. functions.  
 eg: `printf()`, `clrscr()`, `strlen()`  
       |                  |                  |  
       stdio.h          conio.h          string.h

\* It must begins with # functions.



① file inclusion: A file is included

```
#include <file.h>
#include <stdio.h>
```

② MACRO - Define symbolic constant  
 #define PI 3.14      Variable/constant



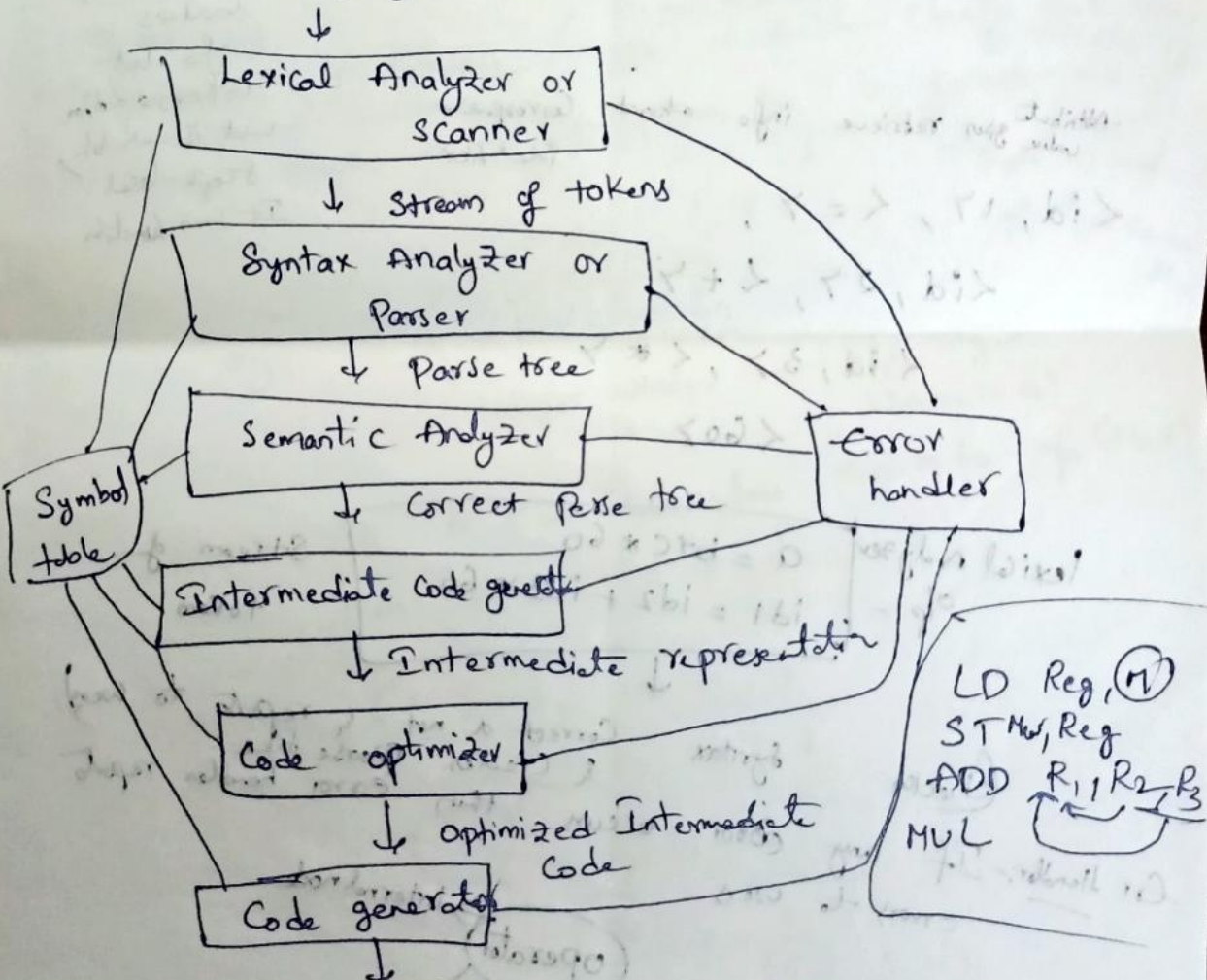
# Relocatable m/c Code :-

In OS, Paging & Segmentation Concepts,

- ① Logical - The Pgm generated addr.
  - ② Physical - Addr. available in memory (main)
  - ③ offset -  $\{ \text{logical} + \text{offset} \}$  in physical addr.
- relocatable addr.

## Phases of a Compiler :-

Source program



Lexeme → A sequence of characters.

$a = b + c * 60$

Token → A keyword or identifier or operator or variable or constant or value

for, while } C, C++, Java, C#, Constant

| <u>lexeme</u> | <u>Token</u>        |
|---------------|---------------------|
| a             | identifier          |
| =             | Assignment operator |
| b             | Identifier          |
| +             | Addition "          |
| c             | Identifier          |
| *             | Multiplication "    |
| 60            | Constant            |

$\langle \text{tokenname, attribute} \rangle$   
value

↓  
info stored in  
Symbol  
Table

Each Token is represented in a pair.

Attribute value gives retrieve info. about Corresponding identifier

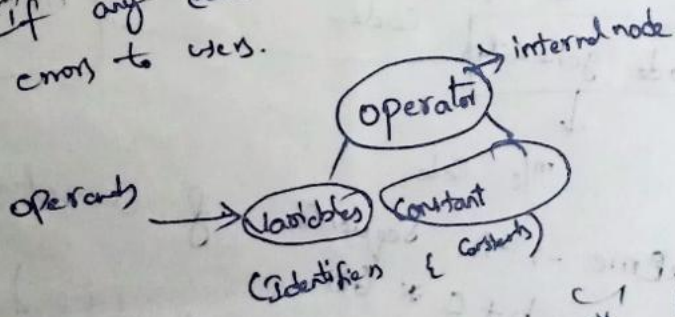
DS which contains info about Tokens & identifier what is variable, scope-label, It uses identifier

- $\langle \text{id, 1} \rangle, \langle = \rangle,$
- $\langle \text{id, 2} \rangle, \langle + \rangle,$
- $\langle \text{id, 3} \rangle, \langle * \rangle$
- $\langle 60 \rangle$

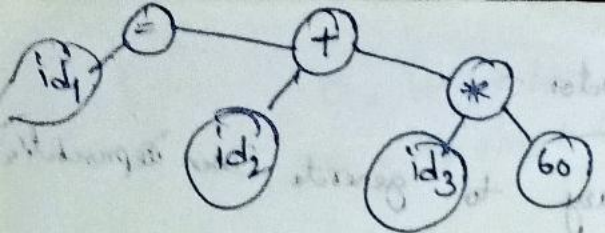
Lexical Analyzer of:-  
 $a = b + c * 60$   
 $\text{id1} = \text{id2} + \text{id3} * 60$

Stream of tokens

Err Handler:- If any error occurs, user. checks Syntax error occurs. Correct or not & reports to user & creates parse tree. Error handler reports



out of + & \* is \* highest priority.



Operator have ↑ Priority

Parse tree:

Syntax Analyzer: It checks whether or not.

tokens are following legal syntax only

If not legal syntax only then error handler.  
 If follows rules, then generates parse tree

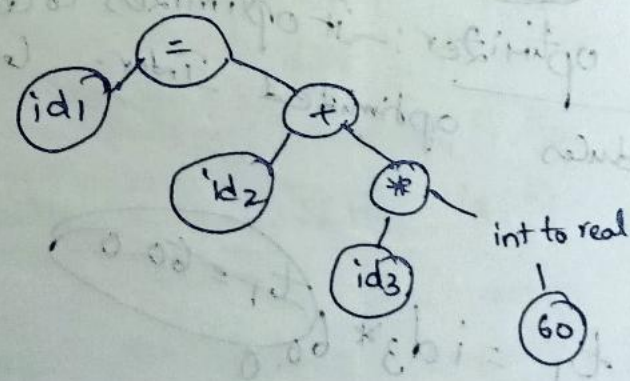
(Meaning)

Semantic Analyzer

(Type checker) checks data types of variables if correct or not. there is need of type conversion

↳ checks meaning of parse tree

eg:-  $a = b + c * 60$   
 integer variables — o/p integer. (Stores in a)  
 If b, c — float values then \* 60 — o/p (float)  
 then a is float. So.



# Intermediate code generator:

It is easy to generate Inter representation

Many ways (notations)

## Three Address Code:-

①

② temp. Variable

③ Some instructions must contain fewer than 3 operands & addr

= Operator

The instr. should contain 3 addr

A three address code assign instr should have at most right operator max. of 3

② Compiler must have a temp. var. to store

Results.

## Inter Representation:-

$$t_1 = \text{intoreal}(60); \quad \text{zero} \quad \text{" operator$$

$$t_2 = id_3 * t_1 \quad \text{one} \quad \text{"$$

$$t_3 = id_2 + t_2 \quad \text{zero} \quad \text{"$$

~~$$t_4 = id_1$$~~

Code optimizer: - It optimizes code & produces optimized Inter. Code.

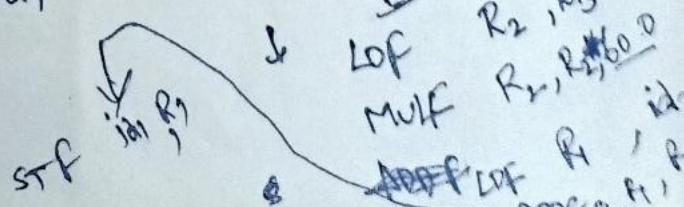
~~$$t_1 = id_3 * 60.0$$~~

~~$$t_2 = id_2 + t_1$$~~

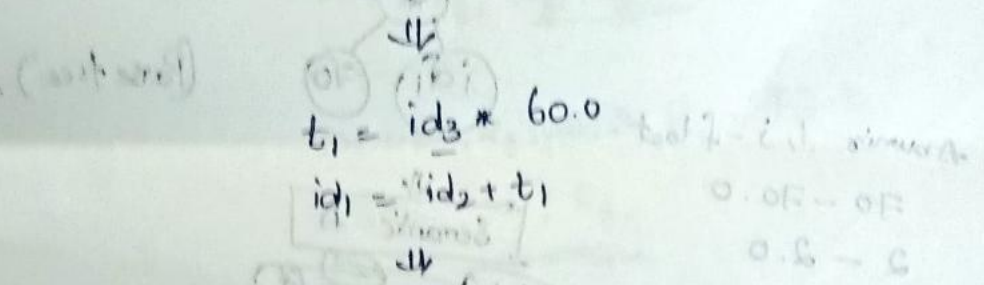
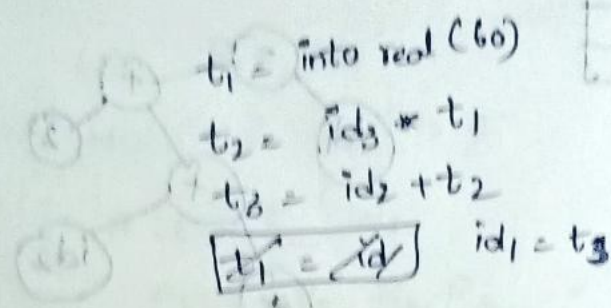
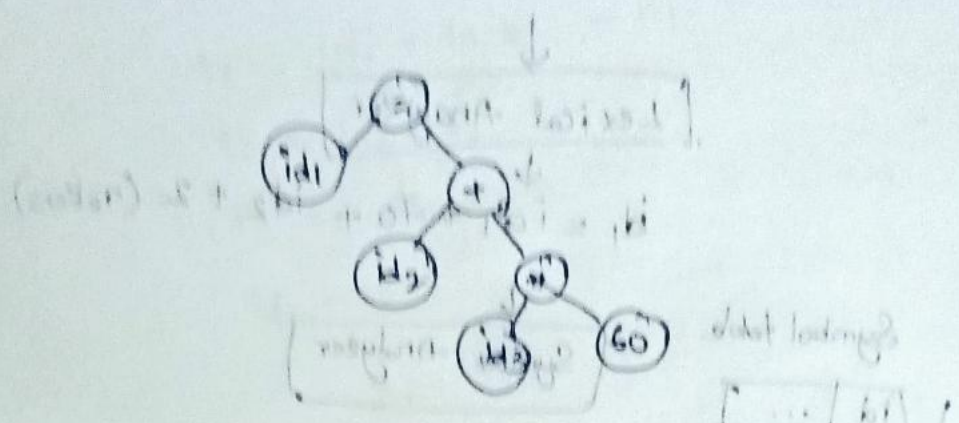
$$id_1 = id_2 + t_1$$

- LD Reg, M
- ST M, Reg
- ADD R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>
- MUL

- LOF R<sub>2</sub>, id<sub>3</sub>
- MULF R<sub>2</sub>, R<sub>2</sub>, 60.0
- ADD R<sub>1</sub>, R<sub>2</sub>, id<sub>2</sub>



$a = b * c * 60$   
 $id_1 = id_2 + id_3 * 60$



```

LDF R2, id3
MULF R2, R2, 60.0
LDF R1, id2
ADF R1, R1, R2
STF id1, R1

```

Q. Consider the following fragment of 'C' code:

```

int i, j;
i = i * 70 + j + 2;

```

write the o/p of all phases of compiler for the above code.

$$i = i * 70 + j + 2$$

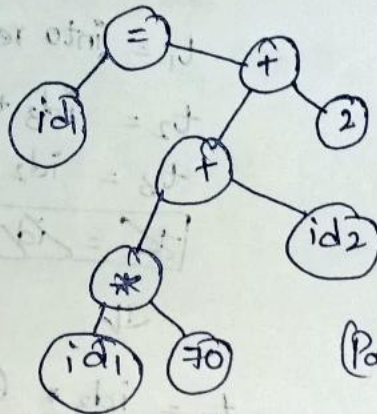
↓  
**Lexical Analyzer**

id<sub>1</sub> = id<sub>1</sub> \* 70 + id<sub>2</sub> + 2 (Tokens)

Symbol table

|   |    |     |
|---|----|-----|
| 1 | id | ... |
| 2 | id | ... |

↓  
**Syntax Analyzer**



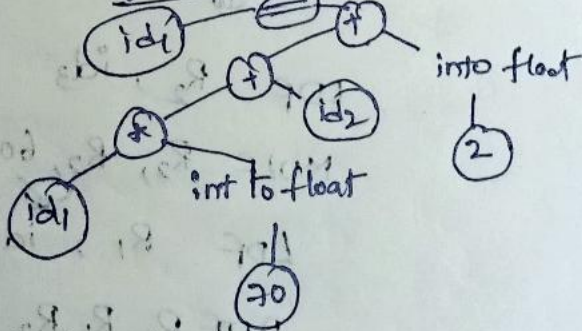
(Parse tree).

Assume  $i, j$  - float.

70 - 70.0

2 - 2.0

↓  
**Semantic A**



↓  
**Intermediate code generator.**

$$t_1 = \text{into float}(70)$$

$$t_2 = \text{id}_1 * t_1$$

$$t_3 = \text{id}_2 + t_2$$

$$t_4 = \text{into float}(2)$$

$$t_5 = t_3 + t_4$$

$$\text{id}_1 = t_5$$

- ① max. one operand
- ② Temp. results
- ③ store instr. < 3 ops

Code optimizer

↓

$$t_1 = id_1 * 70.0 - R_1$$

$$t_2 = id_2 + t_1$$

$$t_3 = t_2 + 2.0$$

$$id_1 = t_3$$

↓

LDF R<sub>1</sub>, id<sub>1</sub>

MULF R<sub>1</sub>, R<sub>1</sub>, #70.0

~~LDF R<sub>2</sub>, id<sub>2</sub>~~

~~ADDF R<sub>1</sub>, R<sub>1</sub>, R<sub>2</sub>~~ ✓

~~LDF R<sub>3</sub>~~

~~ADDF R<sub>1</sub>, #2.0, R<sub>1</sub>~~, #2.0

STF id<sub>1</sub>, R<sub>1</sub>

## (v) Intermediate Code:-

The code may be kept as an array of text strings, a linked list of strings, or a temp text file based on intermediate code type.

## (vi) Temporary files:-

To store intermediate code we use these temp files.

## Issues in Compiler Structure:-

i) Analysis & Synthesis:-  
There are 2 phases of compilation.

- (a) Analysis: It breaks up the source program into constituent pieces & create an intermediate representation of the src prog.
- (b) Synthesis: It constructs the desired target program from the intermediate representation.  
(major impacts on it are reliability, efficiency, usefulness & maintainability)

## ii) Front & Backend:-

F.t: The operations depends only on the src lang.  
B.t: The ops generating the target code.

Scanner, Parser & semantic analyzer  
front end Code generator back end Code optimization  
in both intermediate Code is the medium.

### (iii) Passes:-

It is nothing but the repetitions for  
Problem the entire Src Prog several times before  
generating final code.

for each we end up representation  
Some data in diff. format.

If a compiler is following one pass  
it is less efficient. when : Compiled to  
more than one pass Compilers.

## Major DS in a Compiler:-

Symbol Table :- It stores info about program identifier that will be across all phases.

\* The first 3 phases helps in filling up the Symbol Table.

\* In lexical analyzer phase, the first phase interact with symbol table & symbol table is created by the lexical analysis.

\* A symbol table is a data structure that contains a record for each identifier with fields for the attributes of the identifier.

name of variable

\* Symbol table has facilities to manipulate (add/delete) an element.

eg:-  $x = a + b * 20$

| Symbol | Category   | Type  | Attribute         | Memory location |
|--------|------------|-------|-------------------|-----------------|
| x      | identifier | float | id1<br>assignment | 1000            |
| =      | open       |       |                   |                 |
| a      | id         | float | id2<br>add        | 1004            |
| +      | op         |       |                   |                 |
| b      | id         | float | id3<br>multiply   | 1008            |
| *      | op         |       |                   |                 |
| 20     | Constant   | int   | id4               |                 |
|        | Constant   | int   | product           | 1012            |

Scanner/lexical analyzer, Parser & Semantic analyzer helps to fill above table.

ii) Literal table :-

\* It is also a ds used in Compiler if it is a part of symbol table.

\* It maintains the details of Constant & Strings

| Literal | Category | attribute | memory locs. |
|---------|----------|-----------|--------------|
| 60      | Const    | int       | 1000         |

\* It reduces the size of the Pgm in memory by allowing reuse of Constants & Strings.

Parse tree :

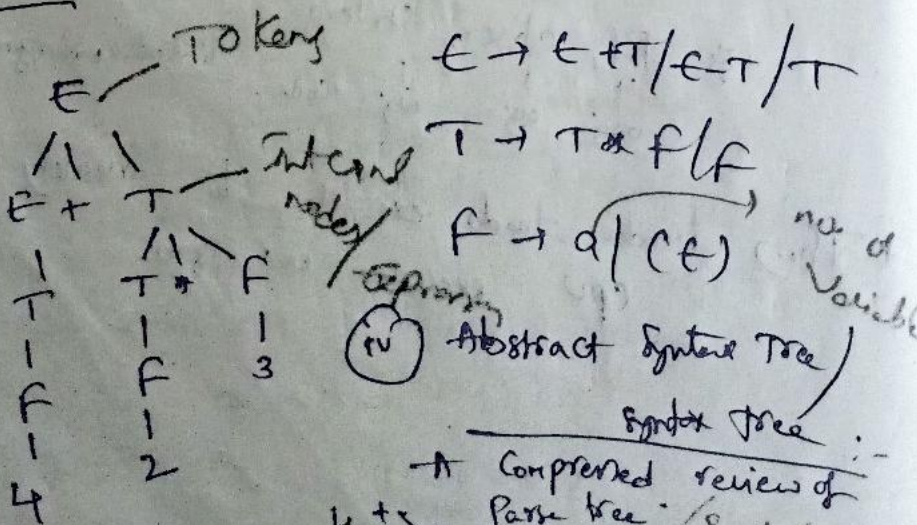
\* It is dynamically allocated pointer based structure.

\* It describes the syntactic structure of 'lp'.

\* In Parse tree, the terminal nodes represent the tokens & internal nodes represent the expressions, (bottom up).

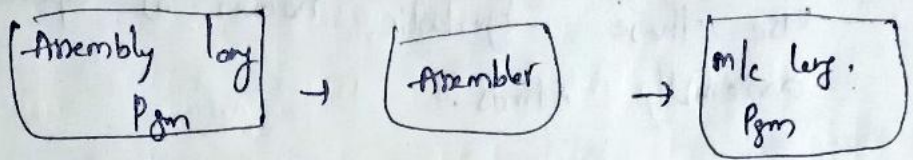
Expression :-  $4 + 2 * 3$

Grammar



\* Compressed review of Parse tree. (Syntactic analysis)

Assembler:- It is a system software which translates, Pgm written in assembly lang into m/c lang.



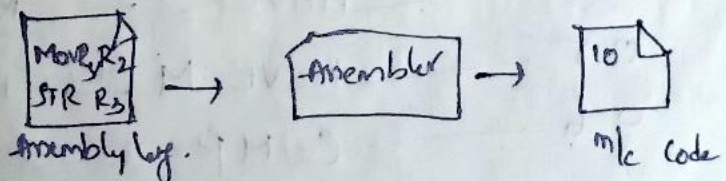
Assembly lang:- It is a kind of low level Pgmng lang, which uses symbolic codes or mnemonics as instructions.

\* Some examples of mnemonics include ADD, SUB, LDA, & STA that stand for addition, subtraction, load accumulator, store accumulator respectively.

Applications of Assembly lang:-

\* It is used for H/w manipulation, access to specialized processor instructions, or to address critical performance issues.

\* Typical uses are device drivers (CD, HDD), low-level embedded systems (Keyboard, water tank Indicator) & real-time systems (Computer, notepad).



Elements of Assembly lang:-

① Mnemonic Operation Code: It is used for m/c instructions (also called mnemonic opcodes) are easier to remember & use than numeric operation codes. Eg: ADD, SUB, MOV, etc.

② Symbolic operands: - A programmer can associate symbolic names with data or instructions & use these symbolic names as operands in assembly stmts.

eg:- ADD R<sub>1</sub>, R<sub>2</sub>, R<sub>3</sub>

③ Data Declaration:- Data can be declared in a variety of notations including the decimal notation.

eg:- NUM1 03 OR  
NUM1 0011H

| Instruction<br>op code | Assembly mnemonic<br>Instructions. | Remarks.                |
|------------------------|------------------------------------|-------------------------|
| 00                     | STOP                               | Stop execution          |
| 01                     | ADD                                | Addition                |
| 02                     | SUB                                | Subtraction             |
| 03                     | MULT                               | Multiplication          |
| 04                     | MOVER                              | Move Memory to Register |
| 05                     | MOVE M                             | Move Register to memory |
| 06                     | COMP                               | Comparison              |
| 07                     | BC                                 | Branch on Condition     |
| 08                     | DIV                                | Division                |
| 09                     | READ                               | Read memory             |
| 10                     | PRINT                              | Write to memory.        |

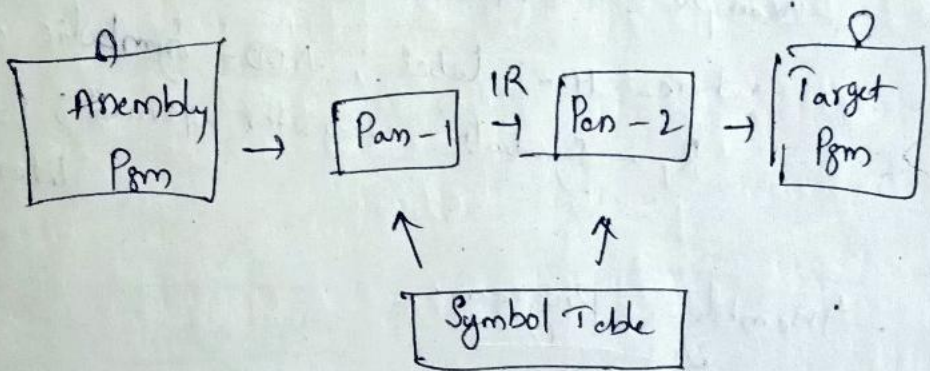
## Assembly process :-

Convert .asm file into .obj file.

Pass 1: Complete scan .asm file. find all labels, instructions & calculating corresponding address.

Pass 2: Convert all the instructions into m/c lang. format.

Symbol Table: stored all the info. of assembly lang. data, variables, instructions, address etc,



Example of Assembly lang. with Assembler Directives

START: This instruction starts the execution of Program from location 200 & label with START provides name for the Program (TO HW is name of Pgm).

MOVER: It moves the content of literal (= '3') into register operand R1.

MOVEM: It moves the content of register into memory operand (x)

MOVER: It again moves the content of literal (= '2') into register operand R2 and its label is specified

LTORG: It assigns address to literal  
(Current LC Value)

DS (Data Space): It assigns a data space of 1  
to symbol X.

END: It finishes the Pgm execution.

[Label] [op code] [operand].

Example: M ADD R<sub>1</sub>, =3'  
 where M - Label; ADD - Symbolic op code;  
 R<sub>1</sub> - Symbolic Register operand; (=3') - literal.

Assembly Programs:-

| Label | op-code             | operand              | LC value<br>(Location Counter) |
|-------|---------------------|----------------------|--------------------------------|
| JOHN  | START               | 200                  |                                |
|       | MOVER               | R <sub>1</sub> , =3' | 200                            |
|       | MOVEM               | R <sub>1</sub> , X   | 201                            |
| Li    | MOVER               | R <sub>2</sub> , =2' | 202                            |
|       | LTORG<br>(Director) |                      | 203                            |
| X     | DS                  | 1                    | 204                            |
|       | END                 |                      | 205                            |

# Design of working of Assembler.

## ① Analysis phase (PASS 1 of Assembler):

- \* To build Symbol Table for synthesis phase to proceed.
- \* Determines address of each symbols called as memory allocation.
- \* Location Counter used to hold addr. of next inst.
- \* Isolated label, mnemonic opcode, operands, constants etc.
- \* Validate meaning & addr. of each stmts.

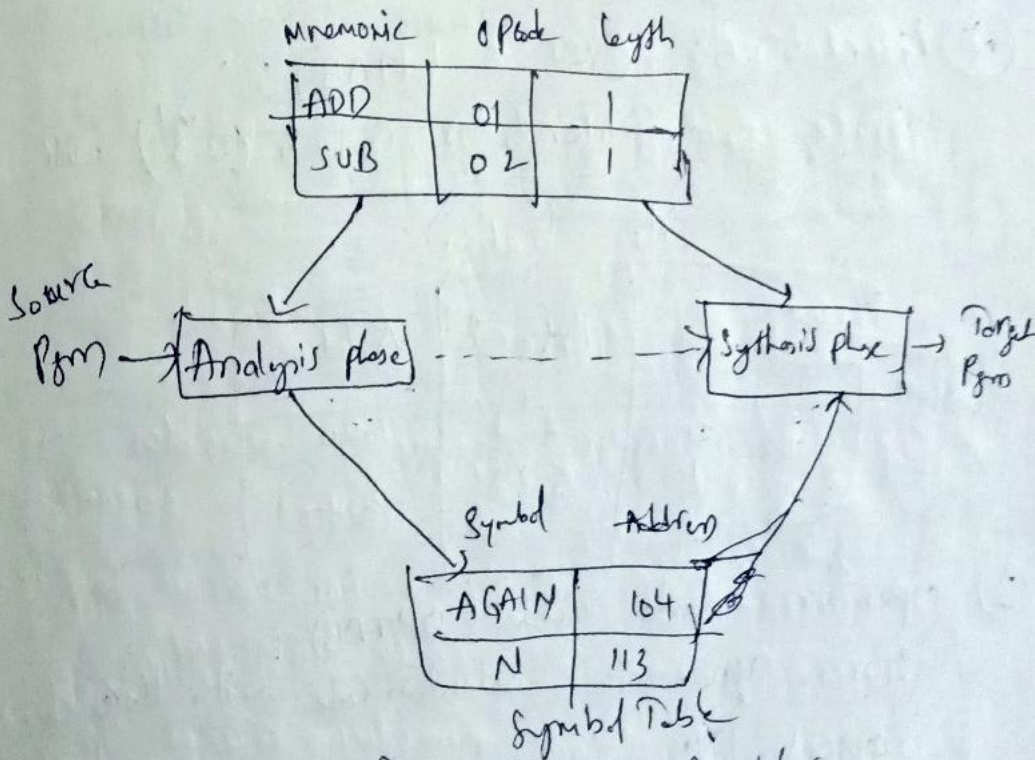


fig: Design of Assembler.

## ② Synthesis phase (PASS 2 of Assembler):

- \* Use DS's generated by analysis phase
- \* To build m/c inst's for every assembly stmts as per mnemonic code & their addr. allocation as per Src Code.
- \* Synthesis is m/c insts.

# Part 1 of Assembly (Analysis phase)

DS in Assembly lang.

① Symbol Table (ST or SYHTAB):

Store value of address, origin to the label

| Label | Addr. |
|-------|-------|
| JOHN  | 200   |
| L1    | 202   |
| x     | 204   |

② Literal Table (LT or LITAB):

Store each literal or constants (= '3') with its location.

Starts with 0.

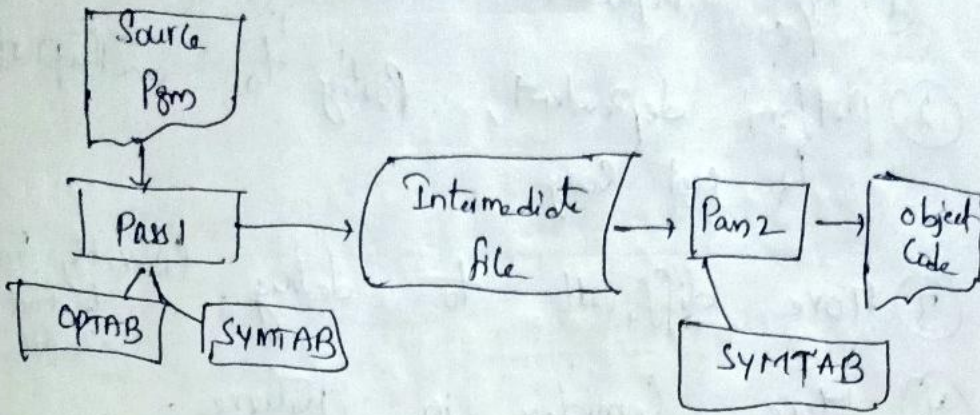
| Index | Literal | Address |
|-------|---------|---------|
| 0     | '3'     | 200     |
| 1     | '2'     | 202     |

③ Operation Code Table (OPTAB):

Stores Mnemonic operation code with their opcodes & length.

| Mnemonic | opcode | length (bytes) |
|----------|--------|----------------|
| MOVW     | 3      | 2 - int        |
| MOVB     | x      | 1 - char       |
| MOVD     | 2      | 2 - int        |

## PASS 2 of Assembler (Synthesis phase)



- \* In the Second pass the inst's are again read & are assembled using the symbol table.
- \* Basically, the assembler goes through the block of Pgm & generates m/c code for that inst.
- \* Then the assembler proceeds to the next inst.
- In this way, the entire m/c code Pgm is created.
- \* Convert mnemonic operators codes, symbolic Table operands with there equivalent m/c code.
- \* Convert data constants to Internal machine representations.
- \* Convert Complete Pgm into obj file.

Adv:-

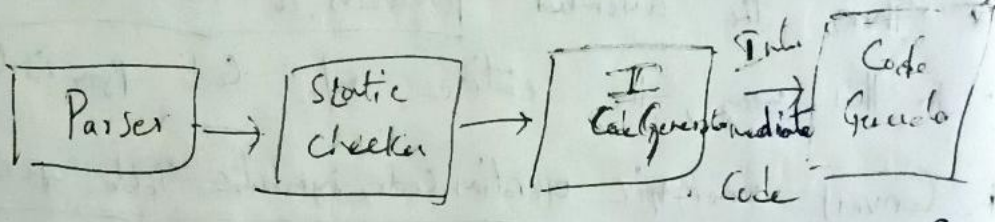
- ① H/w oriented.
- ② Improves readability by using DS.
- ③ Useful in embedded system.
- ④ Less resources, managing size & code.
- ⑤ Access this device of system code easily as compare to high level.

Disadv:-

- ① m/c dependent - (H/w)
- ② platform dependent, Porting to another platform is not easy.
- ③ More difficult to debug. (Directly in contact with H/w)
- ④ More Complex in nature.

Intermediate Code generation:-

It is used to translate the source code into the m/c code. It lies b/w HLL & m/c lang.



\* If the compiler directly translates Src code into m/c code without generating intermediate code then a full native compiler is required for each new m/c.

\* Intermediate code generator receives i/p from its predecessor phase & semantic analyzer phase. It takes i/p in the form of an annotated syntax tree.

\* using intermediate code, the second phase of the compiler is changed according to the target m/c. Intermediate code can be represented in 2 ways.

High level <sup>Intermediate</sup> code

It can be represented as Src Code.

Low level code

It is close to the target m/c. It is used for m/c dependent optimization.

The different forms of Intermediate Code:-

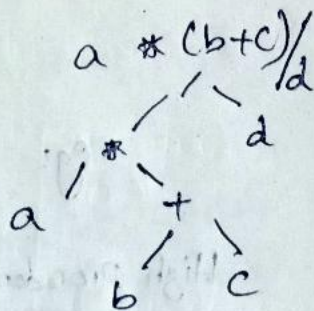
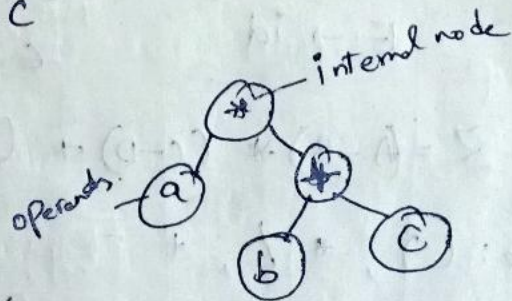
- ① Abstract Syntax tree.
- ② Polish Notation
- ③ Three Address Code.

① Abstract Syntax Tree :-

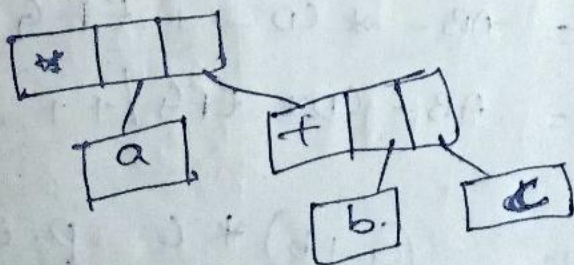
These are compact than a parse tree & can be easily used by compiler.

eg-  $a * b + c$

for arithmetic operators, associativity from left to right



It is represented as



② Polish Notation:- It is of 3 forms.

Infix (a+b)

Prefix (+ab)

Postfix (ab+) Notation:- It is also called as

Suffix or Reverse Polish Notation.

It is a linear representation of syntax tree

eg:-  $x+y \Rightarrow xy+$

Productions

$$E \rightarrow E_1 \text{ op } E_2$$

$$E \rightarrow (E_1)$$

$$E \rightarrow id$$

Semantic Rule (SDG)

$$E \cdot Code = E_1 \cdot Code \parallel E_2 \cdot Code$$

$$E \cdot Code = E_1 \cdot Code$$

$$E \cdot Code = id$$

eg:-  $x = (A-B) * (C-D) + [E + (F/G)]$

High precedence: ( ), \* / , + -

$$x = AB - * CD - + [E + (F/G)]$$

$$= AB - * CD - + EFG / +$$

$$= AB - CD * EFG / + +$$

eg:-  $(a+b) * c$  p.e is ab+ca

$a + (b * c)$  p.e is abc \* +

$(a-b) * (c/d)$  p.e is ab - cd / \*

## Three Address Code:-

- \* It is an intermediate code. It is used by the optimizing compilers.
- \* In 3-address code, the given expression is broken down into several separate instructions. These inst's can easily translate into assembly lang.
- \* Each 3-address code instruction has at most 3 operands. It is a combination of assignment & a binary operator.

In TAC, there is at most one operator on the right side of an instruction.

eg:-  $x + y * 2$   $t_1, t_2$

$$t_1 = y * 2$$

$$t_2 = x + t_1$$

### Quadruple

eg:  $x + y * 2$  (inst)

$$t_1 = y * 2$$

$$t_2 = x + t_1$$

### Triple

operator  
Argument 1  
Argument 2  
Result

### Indirect Triple

Each inst. consists of 4 fields.

$$a = b * -c + b * -c$$

$$t_1 = -c$$

$$t_2 = b * t_1$$

$$t_3 = -c$$

$$t_4 = b * t_3$$

$$t_5 = t_2 + t_4$$

$$t_3 = b * t_1$$

$$t_4 = t_2 + t_3$$

highest priority  
unary  
minus

| Addr | op | arg1           | arg2           | Result         |
|------|----|----------------|----------------|----------------|
| (0)  | -  | c              |                | t <sub>1</sub> |
| (1)  | *  | b              | t <sub>1</sub> | t <sub>2</sub> |
| (2)  | -  | c              |                |                |
| (3)  | *  | b              | t <sub>3</sub> | t <sub>4</sub> |
| (4)  | +  | t <sub>2</sub> | t <sub>4</sub> | t <sub>5</sub> |

Too many temp variables  
 assembly table so that it takes  
 hence we are going for  
 i.e., stored in  
 more time.  
 3-address code

|     | op | Arg1 | Arg2 |
|-----|----|------|------|
| (0) | -  | c    | (0)  |
| (1) | *  | b    | (1)  |
| (2) | -  | c    |      |
| (3) | *  | b    | (2)  |
| (4) | +  | (1)  | (3)  |

With 1 on memory,  
 represent 3-address representation

we can

Indirect Triple. pointer to the triple.

|     |     |
|-----|-----|
| 100 | (0) |
| 101 | (1) |
| 102 | (2) |
| 103 | (3) |
| 104 | (4) |

Pointers ← { 100 } → Triples

① Assignment instructions of the form  $x = y \text{ op } z$ , where  $\text{op}$  is a binary operator.

Eg:-  $x = y + z$ ,  $x = y * z$ ,  $x = y / z$

② Assignments of the form  $x = \text{op } y$  where  $\text{op}$  is a unary operator such as unary -, inc, dec, logical NOT.

Eg:-  $x = -y$ ,  $x = ++y$

③ Copy instructions of the form  $x = y$  where value of  $y$  is assigned to  $x$ .

Eg:-  $y = 10;$   
 $x = y;$

④ Unconditional jump goto L.

Eg:- goto L;

L: stmts to be executed

goto 1000;

5) (a) Conditional jumps of the form If  $x$  relop  $y$  goto  $L$ .

(b) " " " " " " "  $x$  goto  $L_1$   
else goto  $L_2$ .

(relop - relational operator)

(a) if  $x$  relop  $y$  goto  $L$ .

eg:- if  $x > y$ , goto  $L$ .

checks

(b) if  $x$  goto  $L_1$  else goto  $L_2$ .

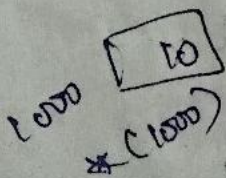
↓  
Contains value 0, 1 (F) (T)

1 —  $L_1$   
0 —  $L_2$ .

6) Functions  
Procedure calls & returns are implemented using following Parameter  $x$   
 $y = \text{Call } P, n$  (Parameter no. of arguments)

return  $y$

7) Address & Pointer assignments of the form  
 $x = \&y$  (Address of  $y$ )  
 $\leftarrow x = \&y$  (The reference op) value at address  
 $*x = y$



8) Indexed Copy Instructions of the form.  
 $x = y[i]$  (indexed).  
 $x[i] = y$

Q: Construct Quadruples, Triples, Indirect Triples for the statements  $(a+b) * (c+d) - (a+bc)$

Sol:-

Three Addr. Code:-  
 $t_1 = a + b$   
 $t_2 = c + d$   
 $t_3 = t_1 * t_2$   
 $t_4 = t_1 + c$   
 $t_5 = t_3 - t_4$

Quadruple:-

|     | op | Arg1  | Arg2  | Result |
|-----|----|-------|-------|--------|
| (0) | +  | a     | b     | $t_1$  |
| (1) | +  | c     | d     | $t_2$  |
| (2) | *  | $t_1$ | $t_2$ | $t_3$  |
| (3) | +  | $t_1$ | c     | $t_4$  |
| (4) | -  | $t_3$ | $t_4$ | $t_5$  |

Triple:-

|     | op | Arg1 | Arg2 |
|-----|----|------|------|
| (0) | +  | a    | b    |
| (1) | +  | c    | d    |
| (2) | *  | (0)  | (1)  |
| (3) | +  | (0)  | c    |
| (4) | -  | (2)  | (3)  |

| Pointers | Triples |
|----------|---------|
| 100      | (0)     |
| 101      | (1)     |
| 102      | (2)     |
| 103      | (3)     |
| 104      | (4)     |

## Literal & Symbol Tables:-

Literals and Constants:

Literals: Literal is an operand with syntax

= '<value>'

\* Its location can't be specified in assembly Pgm. & its not known to assembly lang. Pgm. So its value can't change during execution of Pgm.

eg:-

ADD AREG, =5' this means

ADD AREG, FIVE

-----  
-----  
-----  
FIVE DC(5)

Constant:

\* Its value can be changed during exe.

\* Eg. ADD AREG, 5

Is translated into an instruction with 2 Operands. AREG and the value 5 as an immediate operand.

Consider the following Assembly lang. Program:

Source program

or

Variant I

Variant II

```

START 100
READ N
LOOP MOVER AREG, N
      SUB  AREG, (=1)
      BC  GT, LOOP
      STOP
A DS 1
END
    
```

Symbol Table

| Symbol  | Address | Length |
|---------|---------|--------|
| 1) N    | -       | 01     |
| 2) LOOP | 101     | 1      |
| 3) A    | 105     | 1      |

Literal Table

| Literal | Address |
|---------|---------|
| =1      | 106     |

Pool Table

| Literal No. |
|-------------|
| #1          |

| Instruction opcode | Assembly Mnemonic |
|--------------------|-------------------|
| 00                 | STOP              |
| 01                 | READ              |
| 02                 | MOVER             |
| 03                 | SUB               |
| 04                 | BC                |

| Declarative | Start |
|-------------|-------|
| DS          | 01    |

| ASSEMBLER | DIRECTIVE |
|-----------|-----------|
| START     | 01        |
| END       | 02        |

write intermediate code form using Variant I or II.

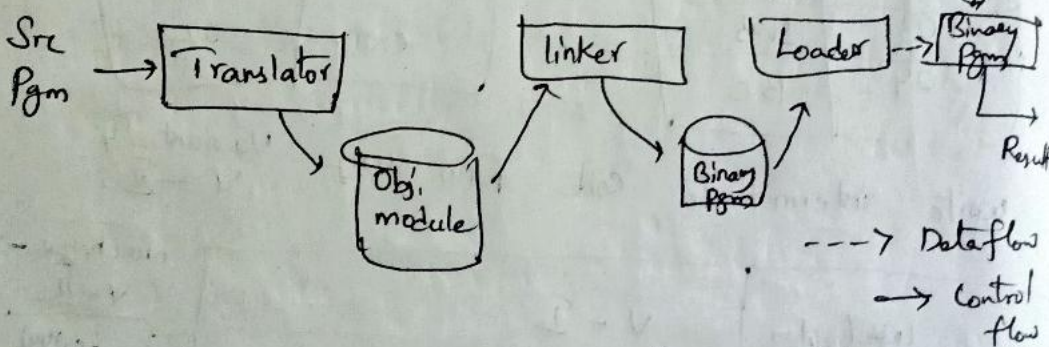
| Addr Location Counter | V - I                | No change → Reg. code size symbol code |
|-----------------------|----------------------|--|
| 100                   | (AD, 01) (C, 100)    | (AD, 01) (C, 100)                      |
| 101                   | (IS, 01) (S, 01)     | (IS, 01) N                             |
| 102                   | (IS, 02) (1) (S, 01) | (IS, 02) AREG, N                       |
| 103                   | (IS, 03) (1) (L, 01) | (IS, 03) AREG, (L, 01)                 |
| 104                   | (IS, 04) (4) (S, 02) | (IS, 04) GT, LOOP                      |
| 105                   | (IS, 00)             | (IS, 00)                               |
| 106                   | (DL, 01) (C, 1)      | (DL, 01) (C, 1)                        |
| 107                   | (AD, 02)             | (AD, 02)                               |



# Relocation & linking concepts

Execution of a Pgm written in a language L involves following steps:

- \* Translation of the Pgm by translator of language L
- \* Linking of the Pgm with other Pgm needed for its execution by a linker.
- \* Relocation of the Pgm to execute from the specific memory area allocated to it by a linker.
- \* Loading of the Pgm in the memory for the purpose of execution by a loader.



Translated, linked and load time addresses:

Translation Time (or translated) address: Address

assigned by translator.

Linked Address: Address assigned by the linker.

Load Time (at load) Address: Address assigned by

Loader.

Translated Origin: Addr. of the origin assumed by translator which is specified by the Programmer in an ORIGIN stmt.

Linked Origin: Addr. of the origin assigned by the linker while producing a binary Pgm.

Load time (or load) Address: - Addr. of the origin assigned by the loader while loading the Pgm for execution.

Topics:-

a) Program Relocation  
• Performing "

b) Linking

- EXTERN & ENTRY stmt
- Resolving external references
- Binary processes.

c) Object Module

Pgm Relocation:-

Addr. Sensitive Pgm:-

- ① - AA - Set of Absolute Addresses, may be inst. addr. or data addr.
- ② -  $AA \neq \phi$ ; means inst. or data occupy memory words with specific addr. Addr. sensitive
- ③ Such a Pgm is called

Pgm which

① Addr.

② Addr.

Sensitive instruction: an instr. which contains addr.  $a_i \in AA$ .

Constant: a data word which contains an address  $a_i \in AA$ .

Important: -Addr. Sensitive Pgm 'P' can execute  
Correctly only if the start addr. of  
the memory area allocated to it is  
the same as its translated origin.

$\boxed{\text{Start Addr.} = \text{Translated Addr.}}$

Thus, to execute correctly from any  
memory area, addr. used in each address  
sensitive instruction of P must be  
'corrected'.

Def :- Pgm Relocation is the process of  
modifying addr. used in the address  
sensitive instructions such that the Pgm can  
execute correctly from the designated area  
of memory.

→ Linker performs relocation if  
Linked origin  $\neq$  Translated origin

→ Loader performs relocation if  
Load origin  $\neq$  Linked origin

→ In general: Linker always performs  
relocation, whereas some  
loaders do not.

Absolute loaders do not perform relocation, thus

load origin = linked origin.

Thus, load origin & linked origin are used interchangeably.

Performing Relocation:-

\* Relocation factor (RF) of P is defined

as

$$[ \text{relocation factor } P = l_{\text{-origin } P} - t_{\text{-origin } P} ]$$

where, RF is the address

\* Symbol is operand

$$[ t_{\text{symb}} = t_{\text{-origin } P} + d_{\text{symb}} ]$$

where,  $d_{\text{symb}} = \text{offset}$

$t_{\text{symb}} = \text{translation time addr.}$

$l_{\text{symb}} = \text{link time addr.}$

$t_{\text{-origin } P} = \text{translation origin}$

$$[ l_{\text{symb}} = l_{\text{-origin } P} + d_{\text{symb}} ]$$

where,  $l_{\text{-origin } P} = \text{linked origin,}$

$$* l_{\text{symb}} = t_{\text{-origin } P} + \text{relocation factor } P$$

$$\begin{aligned} & \rightarrow t_{\text{symb}} \\ & = t_{\text{symb}} + \text{relocation factor } P \end{aligned}$$

IRR: Inst. Requiring Relocation.

It's set of insts that perform relocation in Pgm P.

Steps:

- \* Calculate Relocation factor (RF)
- \* Add it to Translation Time Addr (etc) for every inst. which is member of IRR.

eg:-  $RF = 900 - 500 = 400$   
IRR contains translation addresses 540 & 538.

(inst: read A)

Addr. is changed to  $540 + 400 = 940$  &

$538 + 400 = 938$ .

Linking:-

AP is an App Pgm Consisting of a set of Pgm units  $SP = \{P_i\}$   
A Pgm unit  $P_i$  interacts with another Pgm unit  $P_j$  using insts & addresses of  $P_j$ .

For this it must have.

Public Def.: A symbol <sup>pub-symb</sup> defined in Pgm unit which may be referenced in other Pgm.

external Reference: A reference to a symbol <sup>ext-symb</sup> which is not defined in Pgm unit containing the reference.

Who will handle these 2 things?

EXTRN & ENTRY Stmts:-

\* The ENTRY Stmt lists the public def. of Program unit.

- it lists those symbols defined in Pgm unit which may be referenced in other Pgm units.

\* The EXTRN Stmt lists the symbols to which external references are made in the Pgm unit.

Ex: - (1) TOTAL is ENTRY Stmt. (public def.)  
(2) MAX & ALPHA are EXTRN Stmts. (external reference)

\* Assembler don't know addr. of EXTRN symb & so it puts 0's addr. field of inst. wherever these symb are found.  
\* what if we don't refer these var with EXTRN stmts?

\* Assembler gives errors.

\* This requirement arises for resolving external references.

## Resolving External Reference :-

\* Before AP executes, every external reference should be bound to correct link time address.

\* who will do this binding?

— Linker.

\* Linking: It is the process of binding an external reference to correct link time addr.

\* External reference is said to be unresolved until linking is

performed & resolved once linking is completed.

## Binary Programs.

Is a m/c lang. Pgm comprising set of Program units  $SP$  such that for all  $P_i \in SP$ ,

①  $P_i$  relocated at link origin

② linking is performed for each external reference in  $P_i$ .

### Linker Command:-

- linker  $\langle$  link origin  $\rangle$   $\langle$  object module name  $\rangle$
- $[$ ,  $\langle$  execution start address  $\rangle$
- linker converts obj. module into → the set of Pgm units ( $SP$ ) into → a binary Pgm.
- If link addr = load addr., loader simply loads binary Pgm into appropriate memory area for execution.

### Object Module:-

- OM contains all the info. necessary to relocate & link the Pgm with other Pgms.
- OM consist of four components.

- ① Header: Contains translated origin, size & execution start addr. of P.
- ② Program: Contains n/c lang. Pgm  
Corresponding P.
- ③ Relocation Table: (RELOC TAB) Describing IRRP (Instructions Requiring Relocation).  
Each entry contains field: Translated Address' - of address sensitive instruction.
- ④ Linking Table: (LINK TAB) contains info. concerning public def. & external reference. Has three fields:  
Symbol: symbolic name  
Type: PD/TEXT i.e. public def. or external reference.  
Translated Addr: If PD, addr. of first memory word allocated. If TEXT, addr. of memory word containing reference of symbol.

## Example:

→ for previously given assembly program,  
object module contains following info:

① Translated Origin: 500, size: 42,

Exe. start. Addr = 500

② Machine Language instructions

③ Relocation Table: 

|     |
|-----|
| 500 |
| 538 |

④ Linking Table:

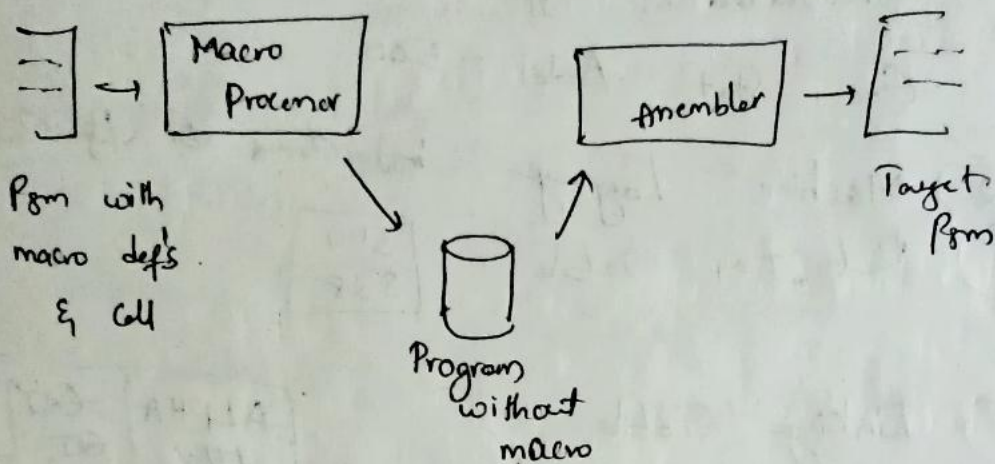
- Loop is not here as it is not PD.  
← Public Def

|       |     |     |
|-------|-----|-----|
| ALPHA | EXT | 518 |
| MAX   | EXT | 519 |
| A     | PD  | 520 |

# Macro Processors & Loaders

## 2. Unit

### Macro instructions & features:-



A schematic of a macro preprocessor.  
eg:- Macro: PRINT\_MSG(msg) expands to printf("%s", msg);  
Why Macro? Code: PRINT\_MSG("hello") becomes printf("%s", "hello")

### In C pgmng : Concept of function.

- \* function is block of stmts related to Such task which we want to execute repeatedly in program.
- \* The function is defined once & can be repeatedly call whenever necessary.
- \* Suppose function call 100 times, compiler transfers control call to definition & definition to call 100 times.
- \* It waste lots of time.

Step 2

```

#include <stdio.h>
void function_name() {
    _____
}
int main() {
    _____
    function_name();
    _____
}

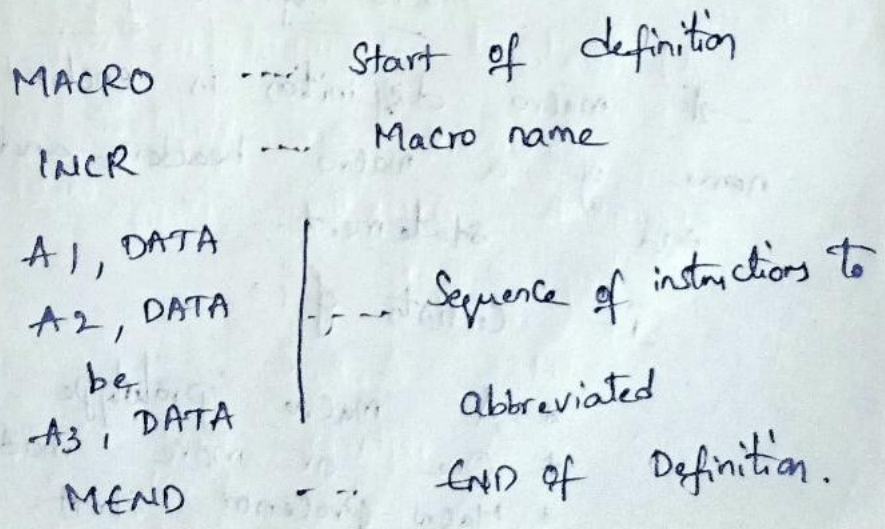
```

Fig: working of functions

Macro :-

It is sequence of instructions assign by name which can be used at any location in the program.

- \* Macro instructions are single line abbreviations for sequential instructions.
- \* Using a Macro, programmer can define a single "instruction" to represent block of code.



## Macro Processor:

Macro instructions are an extension of the basic assembly lang. that can simply debuggig & pgm modification during translation.

\* To process macro instructions, most assemblers use pre-processors known as Macro processors.

### Adv:-

- 1) Simplify & reduce the amount of repetitive coding.
- 2) Reduce the possibility of errors caused by repetitive coding.
- 3) Make an assembly Pgm more readable.

Disadv:- The macro is the size of Pgm (low level HLL)

The reason is, the pre-processor will replace all the macros in the Pgm by its real definition prior to the compilation process of the program.

### Macro Definition & Call

A macro definition is enclosed b/w the name of a macro header and a macro end statement.

It consists of:-

- \* A macro prototype stmt
- \* One or more model stmt.
- \* Macro Processor stmt.

# Macro prototype start :-

## Syntax:

<macro name> [ <formal Parameters> ]

eg:-

- Macro definition (logic of macro)

- MACRO

MNAME &par1, &par2, &par3

MOVER &par1, &par2

ADD &par1, &par3

MOVEM &par1, &par2

MEND

Macro call

MNAME REG1, A, B

| Formal par | Actual par |
|------------|------------|
| &par1      | REG1       |
| &par2      | A          |
| &par3      | B          |

Expanded Code

MOVER REG1, A

ADD REG1, B

MOVEM REG1, A

Macro expansion (Macro processor):-

Replacement of macro call by corresponding sequence of instructions is called as Macro expansion.

| Sequence     | Source  | Expanded Source |
|--------------|---------|-----------------|
| Macro        |         |                 |
| incr         |         |                 |
| A            | 1, data |                 |
| A            | 2, data |                 |
| A            | 3, data |                 |
| MEND         |         |                 |
| :            |         | A 1, Data       |
| :            |         | A 2, Data       |
| INCR         |         | A 3, Data       |
| :            |         |                 |
| :            |         |                 |
| INCR         |         | A 1, Data       |
| :            |         | A 2, Data       |
| :            |         | A 3, Data       |
| data DC F'5' |         | data DC F'5'    |
| :            |         |                 |
| END          |         |                 |

Macro Arguments:-

A) Positional argument (Actual Par)

Argument are matched with dummy arguments according to order in which they appear.

INCR A, B, C

|   |          |        |       |      |
|---|----------|--------|-------|------|
| A | replaces | first  | dummy | arg. |
| B | "        | second | "     | "    |
| C | "        | third  | "     | "    |

9) Keyword arguments - (found for)

\* This allows reference to dummy args by name as well as by position.

eg:-  
INCR 4 arg1 = A, 4 arg3 = C, 4 arg2 = B

eg:  
INCR 4 arg1 = 4 arg2 = A, 4 arg = C

Nested Macro:-

In a macro, Model statement can constitute a call on any other macro. These are called as nested macro.

eg:-

Source Code

Def. of Macro ADDI

```
MACRO
ADDI 4 arg
L 1, 4 arg
A 1, =F'10
ST 1, 4 arg
MEND
```

Def. of Macro ADDS

```
MACRO
ADDS 4 arg1, 4 arg2,
4 arg3
ADDI 4 arg1
ADDI 4 arg2
ADDI 4 arg3
MEND
```

| Source Code | Expanded Code (Level 1) | G.C (L2) |
|-------------|-------------------------|----------|
|-------------|-------------------------|----------|

```

:
MACRO
ADDI 4arg
    L 1, 4arg
    A 1, =f'10'
    ST 1, 4arg

```

```

MEND
MACRO
-ADDS 4arg1, 4arg2, 4arg3
    ADDI 4arg1
    ADDI 4arg2
    ADDI 4arg3
MEND

```

Expansion of ADDS

```

:
:
-ADDS data1, data2, data3 -> { ADDI data 1
:                               { ADDI data 2
:                               { ADDI data 3
:

```

```

data1 DEF'5'
data2 DEF'6'
data3 DEF'7'
END

```

```

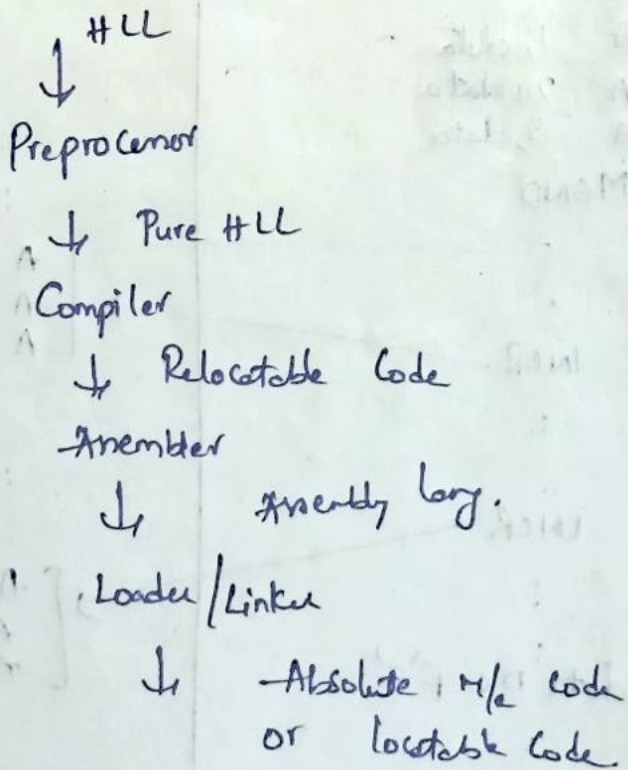
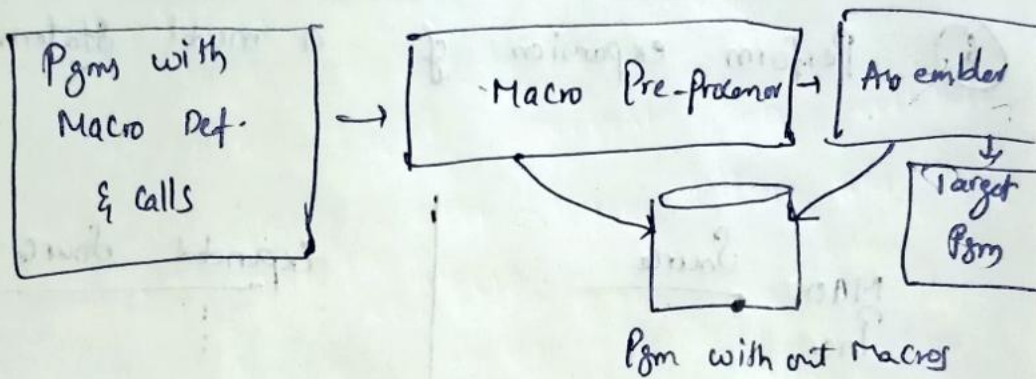
:
:
:
data DEF'5'
data DEF'6'
data DEF'7'
END.

```

## Macro pre processor:-

It takes assemble lang. Pgm which contains macros as i/p & translate into assembly lang. Pgm without macros.

The o/p of macro preprocessor can be given to assembler to convert it into target lang. Pgm.



# Design working of Macro Preprocessor.

## Design overview:

- ① Identify macro calls in the Pgm.
- ② Determine the values of formal Parameters.
- ③ Maintain the values of expansion time variables declared in a macro.
- ④ Organize expansion time control flow.
- ⑤ Determine the values of Sequencing Symbol.
- ⑥ Perform expansion of a macro statement.

Macro Source

```
A 1, data  
A 2, data  
A 3, data  
MEND
```

INCR

INCR

data D C F's'

END

Expanded Source

```
A 1, Data  
A 2, Data  
A 3, Data
```

```
A 1, Data  
A 2, Data  
A 3, Data
```

Data structures :-

Macro Name Table (MNT)  
Parameter Name Table (PNTAB)  
EV Name Table (EVNTAB)  
SS " " (SSNTAB)  
Keyword Parameter Default Table (KPDTAB)  
Macro Definition Table (MDT)  
Actual Parameter Table (APTAB)  
EV Table (EVTAB)

## UNIT –III

### Scanning, Parsing, and Compilers

**Language grammars and ambiguity, Lexical analysis – regular expressions, token generation, Syntax analysis – parsing techniques (top-down, bottom-up), Semantic analysis and intermediate code generation, Code optimization techniques – constant folding, dead code elimination**

It focuses on the **front-end phases of a compiler**, which involve **analyzing the structure of a program** written in a high-level language and converting it into an intermediate form that can be optimized and translated into machine code.

It mainly covers the **process of reading and understanding source code** — from recognizing words and symbols (scanning/lexical analysis) to understanding the grammatical structure (parsing/syntax analysis) and producing intermediate code for further compilation.

#### Language Grammars and Ambiguity

Defines the formal structure (syntax) of programming languages using **context-free grammars (CFGs)**.

Discusses **ambiguous grammars**, where a statement can be parsed in more than one way, and how to resolve ambiguity through grammar refinement or precedence rules.

#### · Lexical Analysis (Scanning)

The **first phase** of compilation.

Converts the **stream of characters** from the source code into **tokens** using **regular expressions** and **finite automata**.

Handles tasks such as removing comments and whitespace, recognizing identifiers, keywords, literals, and operators.

#### Token Generation

Defines how tokens are created and categorized (e.g., IDENTIFIER, NUMBER, OPERATOR).

Uses **regular expressions** to specify token patterns.

Implements **longest-match rules** and **priority handling** to manage overlapping patterns.

## Syntax Analysis (Parsing)

The **second phase** of compilation. Uses **parsing techniques** (top-down and bottom-up) to analyze the sequence of tokens and check whether they follow the language grammar.

Builds **parse trees or syntax trees** representing program structure.

## Semantic Analysis & Intermediate Code Generation

Checks for **semantic correctness** (e.g., type compatibility, variable declarations).

Produces an **intermediate code representation (IR)** — a bridge between high-level source code and low-level machine code.

## Code Optimization Techniques

Improves efficiency without changing the meaning of the program.

Examples: **constant folding, dead code elimination, strength reduction.**

## What is a Language Grammar?

A **grammar** defines the **syntactic structure** (rules) of a programming language — i.e., how valid statements and expressions are formed.

It is a formal description of how tokens (like identifiers, numbers, and operators) can be combined to form valid programs.

## Components of a Grammar (Context-Free Grammar - CFG)

A **Context-Free Grammar (CFG)** is defined as a 4-tuple:

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

Where:

**V (Variables / Non-terminals):** Symbols that represent groups or categories of syntactic structures (e.g., E for Expression).

**T (Terminals):** Actual symbols or tokens of the language (e.g., id, +, \*, (, )).

**P (Productions):** Set of rules that describe how terminals and non-terminals can be combined.

**S (Start Symbol):** A special non-terminal that represents the whole program (e.g., E).

**Example of a Grammar:** Let's define a grammar for arithmetic expressions:

$E \rightarrow E + EE \rightarrow E * EE \rightarrow ( E )E \rightarrow id$

Here:

Non-terminal: E (Expression)

Terminals: id, +, \*, (, )

Start symbol: E

This grammar defines expressions like:

id + id

id \* (id + id)

(id + id) \* id

### Derivation and Parse Tree

A **derivation** shows how a string is produced using the grammar rules.

A **parse tree** visually represents the structure of the derivation.

#### Example:

Input: id + id \* id

Possible derivations:

$E \rightarrow E + E \rightarrow id + E \rightarrow id + E * E \rightarrow id + id * id$

The parse tree shows how operators combine operands.

### What is Ambiguity?

A grammar is **ambiguous** if a **single string** can have **more than one parse tree (or derivation)**.

That means the same expression can be understood in more than one way — leading to confusion in meaning or order of execution.

#### Example of Ambiguity

Consider the grammar above:

$E \rightarrow E + EE \rightarrow E * EE \rightarrow id$

For the input:

id + id \* id

There are **two different parse trees** possible:

**(a) + evaluated first**

$$\begin{array}{c} E \\ / \backslash \\ E + E \\ / \ \backslash \\ id \ id \ id \end{array}$$

Meaning:  $(id + id) * id$

**(b) \* evaluated first**

$$\begin{array}{c} E \\ / \backslash \\ E + E \\ / \ \backslash \\ id \ E \ id \\ / \backslash \\ id * id \end{array}$$

Meaning:  $id + (id * id)$

Hence, the grammar is **ambiguous**, because it does not specify operator **precedence** or **associativity**.

**What is Lexical Analysis?**

Lexical Analysis is the **first phase of a compiler**.

It reads the **source code (characters)** and converts it into a sequence of **tokens**.

□ Think of it like a translator that breaks a sentence into individual words before understanding its meaning.

A **token** is the **smallest meaningful unit** in a program.

Examples of tokens:

**Code            Tokens Generated**

`int num = 10;`    `int, num, =, 10, ;`

`a = b + c;`      `a, =, b, +, c, ;`

Each token has a **type** and a **value** (lexeme):

`int` → (KEYWORD, "int")

`num` → (IDENTIFIER, "num")

`10` → (NUMBER, "10")

= → (OPERATOR, "=")

; → (SYMBOL, ";")

## Role of the Lexical Analyzer

The **Lexical Analyzer (Scanner)**: Reads source code character by character.

Groups characters into **tokens** using patterns called **regular expressions**.

Removes spaces, tabs, and comments.

Passes the tokens to the **Parser** (next phase)

## What are Regular Expressions?

A **regular expression (regex)** is a pattern that describes how certain strings (like identifiers, numbers, or keywords) are formed.

In lexical analysis, **regular expressions are used to define the rules for each token type.**

## Common Regular Expression Patterns

| Token Type            | Regular Expression     | Example               |
|-----------------------|------------------------|-----------------------|
| Identifier            | [A-Za-z_][A-Za-z0-9_]* | sum, _total, count1   |
| Integer Constant      | [0-9]+                 | 5, 123, 99            |
| Floating Point Number | [0-9]+\.[0-9]+         | 3.14, 10.0            |
| Keyword               | \if                    | else                  |
| Operator              | \+                     | -                     |
| Whitespace            | [\t\n]+                | (space, tab, newline) |
| String Literal        | `"([^\`\\`])*"`        | \.)*"`                |

---

## 6. Example: Step-by-Step Tokenization

### Input Code:

```
int total = price + tax;
```

### Process:

| Step | Input Characters | Matched Token | Token Type |
|------|------------------|---------------|------------|
| 1    | int              | int           | KEYWORD    |
| 2    | space            | ignored       | —          |
| 3    | total            | total         | IDENTIFIER |
| 4    | space            | ignored       | —          |

| Step | Input Characters | Matched | Token | Token Type |
|------|------------------|---------|-------|------------|
| 5    | =                | =       |       | OPERATOR   |
| 6    | space            | ignored | —     |            |
| 7    | price            | price   |       | IDENTIFIER |
| 8    | space            | ignored | —     |            |
| 9    | +                | +       |       | OPERATOR   |
| 10   | space            | ignored | —     |            |
| 11   | tax              | tax     |       | IDENTIFIER |
| 12   | ;                | ;       |       | SYMBOL     |

### Output Tokens:

(KEYWORD, "int")  
 (IDENTIFIER, "total")  
 (OPERATOR, "=")  
 (IDENTIFIER, "price")  
 (OPERATOR, "+")  
 (IDENTIFIER, "tax")  
 (SYMBOL, ";")

### Important Rule – Longest Match Rule (Maximal Munch)

When two patterns can match the same string,  
 The **longest match** is chosen.

#### Example:

For input ==,

= (operator) and == (relational operator) both match,

but == is **longer**, so it's chosen.

### Keywords vs Identifiers

Both match the same pattern `[A-Za-z_][A-Za-z0-9_]*`,  
 but if the word is a **reserved keyword** (like if, for, int),  
 the lexer labels it as a **keyword** instead of an identifier.

### Example with Regular Expressions

Let's match some tokens:

| Input | Matched Regex                       | Token      |
|-------|-------------------------------------|------------|
| sum1  | <code>[A-Za-z_][A-Za-z0-9_]*</code> | IDENTIFIER |
| 123   | <code>[0-9]+</code>                 | NUMBER     |
| while | while                               | KEYWORD    |
| +     | <code>\+</code>                     | OPERATOR   |

| Input   | Matched Regex | Token  |
|---------|---------------|--------|
| "Hello" | "([\^"])"     | "\."*" |

## What is Token Generation?

After lexical analysis (scanning), the compiler **groups characters** into meaningful units called **tokens**.

This process is known as **token generation**.

Each **token** represents a logical category such as **identifier**, **number**, **operator**, or **keyword**.

## Structure of a Token

A token usually has 2 or 3 parts:

| Part                        | Meaning                                       | Example                               |
|-----------------------------|---|---------------------------------------|
| <b>Token Type</b>           | The category of the token                     | IDENTIFIER, NUMBER, KEYWORD, OPERATOR |
| <b>Lexeme</b>               | Actual string from the source code            | "count", "=", "10"                    |
| <b>Attribute (optional)</b> | Extra info (like value or symbol table index) | value = 10                            |

## Example:

### Input:

```
int sum = a + b;
```

### Generated Tokens:

| Token Type | Lexeme | Attribute (if any) |
|------------|--------|--------------------|
| KEYWORD    | int    | —                  |
| IDENTIFIER | sum    | symbol table index |
| OPERATOR   | =      | —                  |
| IDENTIFIER | a      | symbol table index |
| OPERATOR   | +      | —                  |
| IDENTIFIER | b      | symbol table index |
| SYMBOL     | ;      | —                  |

## How Tokens Are Generated

**Regular expressions** define patterns for each token type.

Example: Identifier  $\rightarrow$  [A-Za-z\_][A-Za-z0-9\_]\*

Number  $\rightarrow [0-9]^+$

Operator  $\rightarrow \{ + | - | = \}$

**Lexical analyzer (scanner)** reads the source code **character by character**.

## Syntax Analysis (Parsing Techniques)

### What is Syntax Analysis?

After token generation, the **parser** checks if the token sequence follows the grammar rules of the language.

This phase is called **Syntax Analysis** or **Parsing**.

□ The parser builds a **parse tree** or **syntax tree** showing the grammatical structure of the source code.

### Example

#### Grammar:

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid id$

#### Input Tokens:

id + id \* id

The parser checks whether these tokens follow the grammar and builds a tree like:

E  
/  
E + T  
| /\  
T T F  
| | |  
F F id  
| |  
id id

When a sequence of characters matches a regex pattern,  
→ the **token is created** and sent to the **parser**.

The analyzer also ignores whitespaces and comments.

### Example – Step-by-Step

#### Input:

total = 100 + tax;

#### Process:

total → matches identifier pattern → (**ID**, "total")

= → matches operator pattern → (**OP**, "=")

100 → matches number pattern → (**NUM**, "100")

+ → (**OP**, "+")

tax → (**ID**, "tax")

; → (**SYMBOL**, ";")

#### Output Token Stream:

(ID, total), (OP, =), (NUM, 100), (OP, +), (ID, tax), (SYMBOL, ;)

### Syntax Analysis (Parsing Techniques)

#### What is Syntax Analysis?

After token generation, the **parser** checks if the token sequence follows the grammar rules of the language.

This phase is called **Syntax Analysis** or **Parsing**.

The parser builds a **parse tree** or **syntax tree** showing the grammatical structure of the source code.

#### Example

##### Grammar:

$E \rightarrow E + T \mid TT \rightarrow T * F \mid FF \rightarrow (E) \mid id$

### Input Tokens:

id + id \* id

- The parser checks whether these tokens follow the grammar and builds a tree like:

```
      E
     /\
    E + T
   | /\
  T T F
  | | |
  F F id
  | |
  id id
```

### Types of Parsing Techniques

There are two major types of parsing:

- **1. Top-Down Parsing**

Starts **from the start symbol (S)** and tries to **derive the input tokens**.

It expands non-terminals into possible productions until the input string is produced.

Works like **prediction** → “Can this rule produce the input?”

### Common Methods:

#### Recursive Descent Parsing

#### LL(1) Parsing

### Example (Recursive Descent):

Grammar:

$E \rightarrow T + E \mid TT \rightarrow id$

Input: id + id

Step-by-step:

Start with E.

Expand  $E \rightarrow T + E$ .

Match  $T \rightarrow id$ .

Input id matched.

Match +.

Expand remaining  $E \rightarrow T$ .

Match  $T \rightarrow \text{id}$ .

□ Accepted — the input is valid.

### Advantages:

Simple to implement by hand.

Easy to understand.

### Disadvantages:

Can fail with **left-recursive** grammars (e.g.,  $E \rightarrow E + T$ ).

Not suitable for complex grammars.

## 2. Bottom-Up Parsing

Starts **from the input tokens** and tries to **reduce** them back to the start symbol.

Works like **construction in reverse**  $\rightarrow$  “Can we reduce this string to S?”

### Common Methods:

#### Shift-Reduce Parsing

#### LR(0), SLR(1), LALR(1) Parsers

### Example (Shift-Reduce):

Grammar:

$E \rightarrow E + T \mid TT \rightarrow \text{id}$

Input: id + id

Steps:

**Shift:** read id  $\rightarrow$  stack: id

**Reduce:** id  $\rightarrow$  T

**Reduce:** T  $\rightarrow$  E

**Shift:** read +

**Shift:** read id

**Reduce:**  $id \rightarrow T$

**Reduce:**  $E + T \rightarrow E$

□ Reduced to start symbol  $E \rightarrow$  Accepted!

### Advantages:

Very powerful, handles most programming language grammars.

Used in modern parser generators like **YACC**, **Bison**.

### Disadvantages:

More complex to understand and implement.

### Top-Down vs Bottom-Up (Comparison Table)

| Feature        | Top-Down Parsing             | Bottom-Up Parsing          |
|----------------|------------------------------|----------------------------|
| Direction      | From start symbol to input   | From input to start symbol |
| Method         | Expands productions          | Reduces productions        |
| Common Types   | Recursive Descent, LL(1)     | Shift-Reduce, LR(1), LALR  |
| Implementation | Simple, hand-written         | Complex, automated         |
| Grammar Type   | Cannot handle left recursion | Can handle left recursion  |
| Use            | Simpler languages            | Real-world languages       |

### Semantic Analysis

After parsing, the compiler knows the structure (syntax) of the program — but now it must check whether the **meaning** of that structure is **correct**. That process is called **semantic analysis**.

**Purpose:** Semantic analysis checks for logical or meaning-related errors that are not caught by syntax rules.

### Main Tasks of Semantic Analysis

| Task                       | Example  | What Happens                               |
|----------------------------|--|--|
| Type Checking              | <code>int a; a = "hello";</code>                                 | Error — assigning string to int            |
| Variable Declaration Check | <code>x = 10;</code>   | Error — x not declared                     |
| Function Parameter Check   | <code>sum(5);</code> but function <code>sum(int a, int b)</code> | Error — missing parameter                  |
| Scope Checking             | Using a variable outside its scope                               | Error — undefined variable                 |
| Array Bounds Check         | <code>arr[10]</code> for array size 5                            | Error — index out of range (runtime check) |

## How It Works

Semantic Analyzer uses a **Symbol Table**, which stores:

- Variable names

- Data types

- Scope (global/local)

- Function info (parameters, return type)

When analyzing, it:

- Looks up each identifier in the symbol table.

- Checks if usage matches declaration.

- Ensures expressions are type-compatible.

## Example

```
int a, b; float c;  
a = 5;  
b = a + 2;  
c = a + b;
```

All assignments are **type-consistent** (int + int = int, which can be stored in float).

If you wrote:

```
a = "text";
```

**Semantic Error** → Type mismatch (cannot assign string to int).

## Intermediate Code Generation

Once syntax and semantics are correct, the compiler converts the source code into an **intermediate form**, called **Intermediate Code (IC)**.

This IC is not machine code, but it's easier to **optimize** and **translate** to any target machine.

**Purpose:** Acts as a bridge between source code and machine code.

- Makes the compiler **machine-independent**.

- Simplifies optimization.

## Common Intermediate Representations

| Type                       | Example Code                  | Description   |
|----------------------------|-------------------------------|---|
| <b>Three-Address (TAC)</b> | $t1 = a + b$<br>$t2 = t1 * c$ | Each instruction has at most 3 addresses (operands) |
| <b>Quadruples</b>          | (op, arg1, arg2, result)      | Explicit representation for each operation          |
| <b>Triples</b>             | (op, arg1, arg2)              | Uses positions instead of names for results         |

### Example

#### Source Code:

$a = b + c * d;$

#### Intermediate Code (Three-Address Form):

$t1 = c * d$   
 $t2 = b + t1$   
 $a = t2$

#### Quadruple Representation:

##### Op Arg1 Arg2 Result

\* c d t1

+ b t1 t2

= t2 — a

### Code Optimization Techniques

#### What is Code Optimization?

After generating intermediate code, the compiler tries to make it **faster** and **more efficient**, without changing what the program does.

Goal: **Better performance**, **less memory usage**, and **smaller code**.

#### Types of Code Optimization

| Type                       | Description  |
|----------------------------|--|
| <b>Machine-Independent</b> | Performed on intermediate code (works for all CPUs)    |
| <b>Machine-Dependent</b>   | Performed after code generation (specific to hardware) |

We'll discuss two common **machine-independent** techniques:

#### (a) Constant Folding

##### Definition:

If an expression involves **only constants**, the compiler can **evaluate it at compile time** instead of runtime.

### **Example:**

```
int a = 5 * 10;
```

### **Constant Folding result:**

```
int a = 50;
```

Another example:

```
int b = (2 + 3) * 4;
```

→ Compiler changes to:

```
int b = 20;
```

It Saves runtime computation — compiler already knows the value.

### **(b) Dead Code Elimination**

#### **Definition:**

Removes code statements that **never affect the program's output** or are **never executed**.

#### **Example 1 – Unreachable Code:**

```
if (false) {  
    printf("This will never run");  
}
```

□ Compiler removes this block — **dead code**.

#### **Example 2 – Unused Variable:**

```
int a = 10;  
a = 20;  
a = 30;printf("%d", a);
```

Only the last assignment matters → previous ones are **dead** and removed.

#### **Example 3 – After return statement:**

```
return x;  
y = x + 1; // Dead code (never executes)
```

Compiler deletes  $y = x + 1$ .

### **Benefits of Optimization**

Reduces program size

Increases execution speed

Decreases memory and CPU usage  
Improves battery life (for mobile devices)

### Quick Comparison Table

| <b>Concept</b>                      | <b>Purpose</b>                              | <b>Example</b>        |
|-------------------------------------|---|-----------------------|
| <b>Semantic Analysis</b>            | Checks meaning and type correctness         | int a; a = "hi";<br>□ |
| <b>Intermediate Code Generation</b> | Converts source to machine-independent code | t1 = a + b            |
| <b>Constant Folding</b>             | Pre-computes constant expressions           | 2 + 3 → 5             |
| <b>Dead Code Elimination</b>        | Removes unnecessary code                    | Code after<br>return  |

## UNIT IV

### Linkers, Debuggers, and Shell Programming

Symbol resolution and relocation, Linking (static vs dynamic), relocation records, Debugging techniques and breakpoints, Unix/Linux shell environment, Shell commands, variables, redirection, pipes, control statements, Shell script functions and script-based automation

#### What is *linking* (big picture)?

A *linker* (link editor) combines compiled object files (.o) and libraries (.a, .so) into a final program (an executable or a shared library).

Its two main jobs are:

**Symbol resolution** — match every *reference* (use of a variable or function) to a *definition* (where it is defined).

**Relocation** — fix up addresses in the code/data so each reference points to the correct final memory address.

Linking can happen in different phases:

**Link-time** (when building the executable — static linking resolves most symbols then)

**Load-time / run-time** (dynamic linking may resolve some symbols when the program starts or on first use)

#### Symbol resolution

A symbol is a name representing code or data: function names, global variables, static/local symbols, etc.

#### Symbol kinds

**Defined symbol:** an object file provides its location (e.g., `int x = 10;` or `void foo(){...}`).

**Undefined symbol:** a file references it but does not define it (e.g., `extern int x;` or a call to `printf`).

**Local vs Global:** local symbols are internal to object file; global symbols are visible for resolution across files.

**Strong vs Weak:** weak symbols allow a linker to prefer a strong definition if one exists (useful for default implementations).

## How resolution proceeds (conceptual algorithm)

The linker collects all object files and libraries given on the command line.

It builds a global table of **defined symbols** and **undefined references**.

For each undefined symbol, the linker searches the defined symbols (in object files and then in libraries).

If a definition is found, the undefined reference is **resolved** to that definition.

If not found → **link error**: “undefined reference to foo”.

If there are multiple definitions of the same global symbol:

Two *strong* definitions → **error** (multiple-definition).

Strong + weak → strong wins.

For static archives (.a), the linker pulls in only those object modules from the archive that define currently unresolved symbols; order matters.

## Quick example

main.c

```
extern int x;int main() { return x; }
```

def.c

```
int x = 42;
```

Commands:

```
gcc -c main.c # -> main.o (has undefined symbol x)
```

```
gcc -c def.c # -> def.o (defines x)
```

```
gcc -o prog main.o def.o # linker resolves x -> def.o
```

Result: prog has symbol x resolved to the data in def.o.

## Notes about archives and order

With libstuff.a, the linker scans archive members and **adds only those .o files that satisfy unresolved symbols**. If symbols appear later on the command line, they may remain unresolved — so ordering matters or use --start-group.

## Symbol visibility & versioning

Compilers/linkers can hide or expose symbols (visibility attributes) and symbol versioning is used in shared libraries to manage API/ABI changes. These affect how symbol resolution behaves at runtime.

## Relocation

### Why relocation?

Object files are generated with internal addresses that are **relative** or placeholders. The final executable or shared object will be loaded at particular memory addresses; the linker/loader must **patch** instructions/data to point to real addresses.

### What a *relocation entry* contains

For each address in the code/data that depends on a symbol, the object file contains a relocation record including:

**offset** (where to patch)

**type** (absolute vs PC-relative, etc.)

**symbol index** (which symbol this relocation refers to)

**addend** (extra constant)

### Two common relocation types (conceptual)

**Absolute relocation:** final value = address\_of(symbol) + addend → used for global data pointers.

**PC-relative (relative) relocation:** final value = address\_of(symbol) + addend – place → used for relative jumps/calls (helps with position independency).

### What the linker does (static linking)

For each relocation entry, the linker computes the final address of the referenced symbol and **writes (patches)** the computed value into the output section at the relocation offset.

It also merges/concatenates sections from multiple .o files into one output section (e.g., .text, .data) and adjusts offsets accordingly.

### What the loader/dynamic linker does (dynamic linking)

Shared libraries may be loaded at addresses determined at load time (ASLR, address conflicts). The *dynamic linker* (ld.so on Linux) will process relocation entries in the shared objects and executable, patching addresses either:

**At load time (eager relocation)** — all relocations fixed when loading,  
or

**Lazily** — function addresses are fixed on first call (via PLT/GOT mechanism).

Position-Independent Code (PIC) reduces the number of relocations needed for code by using relative addressing and indirection through the GOT.

### **Example of relocation: from main.o**

If main.o has an instruction like `call foo` but the object file stores a placeholder, relocation tells the linker: “at offset 0x20 in .text there’s a call instruction that needs to be updated to call symbol foo’s final address + addend”.

## **PIC, GOT and PLT — dynamic linking avoids heavy relocation**

### **Position-Independent Code (PIC)**

Compiled with `-fPIC` on Unix-like systems.

Uses relative addressing so the same code can run regardless of absolute load address (necessary for shared libraries).

Reduces or eliminates relocations inside .text.

### **GOT (Global Offset Table)**

A table of addresses used by PIC code to access global data and function addresses. The code loads addresses from GOT entries rather than embedding absolute addresses.

**PLT (Procedure Linkage Table):** A trampoline area used to call external functions. The PLT entry initially jumps to the dynamic linker resolver. On first call, the resolver finds the real address, writes it into the GOT, and subsequent calls go directly to the function.

Enables **lazy binding** (resolve a function only when it is first called).

## **Static linking vs Dynamic linking**

### **Static Linking**

All used library code is copied into the final executable at link time. The executable contains everything it needs.

link with static libraries (e.g., `libfoo.a`) or use `gcc -static` to prefer static `libc`.

#### **Pros:**

Single self-contained binary — no external dependencies at runtime.

Predictable behavior (no surprise changes when system libraries update).

Sometimes slightly faster startup (no runtime relocation of libraries).

**Cons:**

Large executable size (duplicate code if many programs statically link the same library).

Cannot share memory pages of the library code across processes.

Harder to patch a bug in a library (need to rebuild/redeploy all apps).

Licensing issues (some libraries/OSS licenses restrict static linking).

**Dynamic Linking (Shared Libraries)**

Executable contains references to shared libraries (.so / .dll) and the dynamic loader resolves and loads those libraries at program start (or earlier/later).

compile shared objects (`gcc -fPIC -shared -o libmylib.so mylib.o`) and link with them (`gcc main.o -L. -lmylib`).

**Pros:**

Smaller executables (shared code stays in .so files).

Multiple processes share the same in-memory library pages.

Libraries can be updated independently — a bugfix in a .so benefits all apps that use it.

**Cons:**

Programs can break if incompatible library changes are installed (dependency hell / ABI mismatch). Slight runtime overhead: loader work and potential relocations at load-time, plus possible lazy resolution overhead on first call.

Less deterministic if system libraries change.

**Practical example commands**

Create object files:

```
gcc -c main.c    # main.o
gcc -c lib.c     # lib.o
```

Create a shared library:

```
gcc -fPIC -shared -o libmylib.so lib.o
gcc -o prog main.o -L. -lmylib # links to libmylib.so (dynamic)
```

Create a statically linked program:

```
gcc -static -o prog_static main.o lib.o# or link with static archive libmylib.a
```

## Runtime linking details & tools

**Dynamic loader** (ld.so / ld-linux.so) is responsible for locating shared libs (via rpath, RUNPATH, LD\_LIBRARY\_PATH, standard system paths) and performing relocations.

**dlopen/dlsym** allow explicit runtime loading and symbol lookup.

Tools to inspect/linking artifacts:

ld / gcc (linker driver)

readelf -s / nm (inspect symbol tables)

ldd (shows dynamic dependencies of an executable)

objdump (view relocations and sections)

## Edge cases & important behaviors

**Multiple definitions:** Multiple strong definitions → linker error. Some platforms permit symbol interposition where a symbol in the executable overrides one in a shared library.

**Archives (.a):** Linker only extracts needed object modules, so order of libraries matters.

**Lazy vs Eager resolution:** Lazy reduces startup cost but moves cost to first call; eager resolves everything at startup.

**ASLR & relocations:** Address Space Layout Randomization moves libraries around; PIC + GOT/PLT help make shared libraries relocatable without modifying code pages.

**Stripping symbols:** strip reduces file size but removes symbol info used for debugging

## What is Relocation?

When you compile a program, each .c or .cpp file becomes an **object file** (.o), which contains **machine code**, **symbols**, and **relocation information**.

However, the **final memory addresses** of functions and variables are **not yet known** — because the linker still needs to **combine** multiple object files and libraries into one executable.

□ So, **Relocation** means **adjusting addresses** in the code and data after linking.

## We Need Relocation

Each object file assumes it starts at **address 0**.

When files are combined, they are placed at **different memory addresses**.

Therefore, **addresses in instructions and data** must be **updated** to point to their actual (final) locations.

Example:

| Symbol | Object File | Original Address | Final Address |
|--------|-------------|------------------|---------------|
| main() | main.o      | 0x0000           | 0x1000        |
| sum()  | sum.o       | 0x0000           | 0x2000        |

If main() calls sum(), the call instruction must be **relocated** from address 0x0000 to 0x2000.

## What Are Relocation Records?

A **relocation record** (or relocation entry) is **metadata** stored in an object file that tells the linker:

“At this place (offset) in the code, you need to insert or adjust the address of a symbol.”

It contains **information for the linker** to fix the addresses.

## Contents of a Relocation Record

Each relocation record typically includes:

| Field        | Description  |
|--------------|--|
| Offset       | The location in the section (e.g., .text, .data) where a fix-up is required. |
| Type         | The type of relocation (absolute, PC-relative, etc.).                        |
| Symbol Index | Which symbol this relocation refers to.                                      |
| Addend       | An additional constant value to add.   |

## □ Example

Let's say we have two files:

**main.c**

```
extern int sum(int, int);int main() {
    return sum(2, 3);
}
```

### **sum.c**

```
int sum(int a, int b) {
    return a + b;
}
```

After compiling:

```
gcc -c main.c # main.o
gcc -c sum.c # sum.o
```

main.o has:

A **symbol table** entry for the external function sum.

A **relocation record** like:

Offset: 0x14Type: R\_X86\_64\_PC32Symbol: sumAddend: -4

This tells the linker:

“At offset 0x14, replace the placeholder with the final address of sum.”

### □ **During Linking**

When the linker (ld) combines the files:

It assigns final addresses to all functions and variables.

It goes through each relocation record.

It updates the instructions in .text and data in .data to the correct absolute or relative addresses.

After relocation → the executable (a.out) has **fixed addresses** that the CPU can execute directly.

### **Relocation Records**

| <b>Concept</b>   | <b>Explanation</b>                                |
|------------------|---|
| <b>Purpose</b>   | Adjust symbol addresses after linking.            |
| <b>Stored In</b> | Object files (.o) under sections like .rela.text. |
| <b>Used By</b>   | Linker and Loader.                                |
| <b>Contains</b>  | Offset, type, symbol, addend.                     |

## Concept      Explanation

**Example Tool** readelf -r main.o shows relocation entries.

## Debugging Techniques and Breakpoints

Debugging is the process of **finding and fixing errors (bugs)** in a program so it runs correctly.

A **debugger** is a special tool that allows developers to:

- Run programs step by step,

- Examine variable values,

- Set **breakpoints**,

- Watch memory/registers,

- And control program execution interactively.

## Common Debugging Tools

### Platform      Debugger Tool

Linux / Unix    gdb (GNU Debugger)

Windows      Visual Studio Debugger

macOS        lldb

IDEs          Code::Blocks, Eclipse, VSCode (internally use gdb/lldb)

## Common Debugging Techniques

### 1 Print (or Trace) Debugging

Add printf() or cout statements in code to check values and flow.

```
printf("x = %d\n", x);
```

Simple but manual; not suitable for large programs.

## Using a Debugger

Compile the program with debug info:

```
gcc -g myprog.c -o myprog
```

Run it inside gdb: `gdb ./myprog`

## Breakpoints

A **breakpoint** is a marker that tells the debugger:

“Stop running the program here, before executing this line.”

You can inspect variable values, memory, and flow from that point onward.

**In gdb:**

```
(gdb) break main      # set breakpoint at start of main()
(gdb) run             # run the program
(gdb) print x        # see the value of variable x
(gdb) next           # go to next line
(gdb) continue       # resume execution
(gdb) delete 1       # remove breakpoint number 1
```

### **Step Execution**

**step** → go into functions line by line

**next** → go to next line (skip inside functions)

**continue** → run until next breakpoint

### **Watchpoints**

A **watchpoint** stops the program whenever a **variable's value changes**.

```
(gdb) watch x
```

Useful to detect where a variable becomes incorrect.

### **Backtrace**

Shows the sequence of function calls (stack trace) that led to a crash or current line.

```
(gdb) backtrace
```

### **Core Dumps**

If a program crashes, the OS can produce a **core file** (snapshot of memory).  
You can open it with gdb:

```
gdb ./a.out core
```

and inspect the state of the program at the moment of crash.

### **Memory and Logic Checks**

Tools like: **Valgrind** → finds memory leaks and invalid accesses.

**AddressSanitizer (ASan)** → runtime check for out-of-bounds or use-after-free.

## Unix/Linux Shell Environment — Overview

The **Shell** is a **command-line interface (CLI)** between the user and the **Unix/Linux operating system**.

It lets you **run commands, write scripts, automate tasks, and manage the system**.

Common shells:

**Bash (Bourne Again Shell)** – default on most Linux systems

**Sh, Zsh, Csh, Ksh, Fish**

## Shell Environment

The shell environment consists of all **settings, variables, and configurations** that define how the shell behaves for a user.

It includes:

Environment variables

The working directory

Command history

Aliases and functions

## Example Commands:

```
echo $HOME    # shows user's home directory
echo $PATH    # shows directories
               searched for commands
pwd           # print working directory
whoami       # shows
current user
```

## Key Point:

The shell environment defines the **user's working context** when using the command line or running shell scripts.

## Shell Commands

Shell commands are **instructions** entered in the terminal to perform actions like viewing files, copying data, or running programs.

There are two main types:

**Internal (built-in)** – handled by the shell itself (e.g., cd, echo, exit)

**External** – separate executable programs (e.g., /bin/ls, /usr/bin/grep)

### Common Examples:

| Command        | Purpose                    |
|----------------|----------------------------|
| ls             | List files and directories |
| cd folder      | Change directory           |
| pwd            | Show current directory     |
| cp file1 file2 | Copy files                 |
| mv file1 file2 | Move or rename files       |
| rm file        | Delete file                |
| cat file       | View file contents         |
| man command    | Show manual for a command  |

### Example:

```
ls -l /home
```

→ Lists all files in /home in long format.

Shell Variables:

Variables store **data or values** that can be used later in commands or scripts.

### Types:

**User-defined variables** – created by the user in scripts or terminal.

**Environment variables** – predefined system variables used by the shell.

### Examples:

```
name="Sona"echo "Hello $name" # prints Hello Sona
echo $HOME # predefined environment variable
PATH=/usr/bin:/bin # set pathexport PATH # make it available to child
processes
```

### Key Point:

Variables help store information temporarily and make scripts dynamic.

Redirection:Redirection allows you to **change where input or output comes from or goes to**.

By default:

Input → keyboard

Output → screen

You can **redirect** these to files.

### Operators:

| Operator | Meaning                      | Example                 |
|----------|------------------------------|-------------------------|
| >        | Redirect output (overwrite)  | ls > files.txt          |
| >>       | Redirect output (append)     | echo "Hi" >> notes.txt  |
| <        | Take input from a file       | sort < data.txt         |
| 2>       | Redirect error output        | gcc file.c 2> error.txt |
| &>       | Redirect both output & error | command &> all.txt      |

### Example:

```
cat file1 file2 > combined.txt
```

Combines two files and writes output to combined.txt.

### □ Pipes (|)

A **pipe** connects the **output of one command** directly as the **input to another** — without using intermediate files.

Used to **chain commands** together for efficient data processing.

### Syntax:

```
command1 | command2 | command3
```

### Examples:

```
ls | wc -l
```

□ Counts the number of files in the current directory.

```
ps aux | grep firefox
```

Filters the list of running processes to show only “firefox”.

### Key idea:

Pipes help you combine simple commands to perform complex tasks.

### Control Statements (in Shell Scripting)

Control statements manage the **flow of execution** in shell scripts — similar to loops and conditions in programming languages.

### Types and Examples:

### (a) If-Else Statement

```
if [ $a -gt $b ]then
  echo "a is greater"else
  echo "b is greater"fi
```

### (b) For Loop

```
for i in 1 2 3do
  echo "Number $i"done
```

### (c) While Loop

```
count=1 while [ $count -le 5 ]do
  echo "Count = $count"
  count=$((count + 1))done
```

### (d) Case Statement

```
case $choice in
  1) echo "Start";;
  2) echo "Stop";;
  *) echo "Invalid";;esac
```

### Shell Script Function:-

A **function** is a **block of reusable code** inside a shell script. Functions let you **group commands together** and **call them multiple times**, which makes scripts **shorter, cleaner, and easier to manage**.

### Basic Syntax

```
function_name() {
  # commands
}
# orfunction function_name {
  # commands
}
```

We can **call the function** by writing its name:

```
function_name
```

Functions can **take arguments** like scripts do.

### Using Arguments

```
greet() {
  echo "Hello, $1!" # $1 = first argument
}
```

```
greet "Sona"
```

### **Output:**

Hello, Sona!

\$1, \$2, ... \$n → positional parameters for function arguments

\$@ → all arguments

return → returns an exit status (0–255)

echo → can return text for assignment

### **Example: Reusable Function**

```
# Function to check if a file exists
check_file() {
  if [ -f "$1" ]; then
    echo "File $1 exists."
  else
    echo "File $1 does not exist."
  fi
}
```

```
check_file file.txt
check_file data.csv
```

Advantages:

Avoid repeating the same code multiple times

Makes scripts **organized and readable**

Supports **modular scripting**

### **Script-Based Automation**

Automation means using **scripts** to **perform tasks automatically** without human intervention.

Shell scripts can automate:

File management (backup, delete, move)

System monitoring (disk usage, CPU load)

Software installation or updates

Scheduled jobs using **cron**

## Use Automation

Saves time on repetitive tasks

Reduces human error

Ensures tasks run consistently

Useful for **system administration** and **DevOps**

## Basic Automation Example

### Backup Script

```
#!/bin/bash# Backup /home/user/Documents to /home/user/Backup
backup() {
  src="/home/user/Documents"
  dest="/home/user/Backup"
  timestamp=$(date +%Y%m%d_%H%M%S)
  mkdir -p "$dest"
  cp -r "$src" "$dest/Documents_${timestamp}"
  echo "Backup completed at $timestamp"
}
```

backup # call the function

Can be scheduled using **cron**:

```
crontab -e# Add line to run every day at 2 AM
0 2 * * * /home/user/scripts/backup.sh
```

**Output:** Automatically creates timestamped backups every day.

### Tips for Script Automation

Always test scripts manually before scheduling.

Use **functions** for modular tasks (like backup, logging).

Use **variables** for paths, filenames, and parameters.

Add **logging** for tracking success or failure:

```
echo "$(date) Backup completed" >> backup.log
```

## UNIT V

### Unix/Linux System Programming

**Introduction to system-level programming in C, File I/O system calls (open, read, write, close), Process creation using fork(), exec(), wait(), Inter-process communication (pipes, FIFO), Signal handling and POSIX threads (pthread\_create, pthread\_join), Case studies: background processes, daemon creation, mini shell**

#### Introduction to System-Level Programming in C

System-level programming refers to writing programs that interact closely with the **operating system** and **hardware resources**. Unlike high-level application programming, system programming allows you to:

- Access files, directories, and devices directly.

- Control processes and memory.

- Perform network communications.

- Use system calls provided by the OS.

In **Unix/Linux**, system programming is mostly done in **C**, because C provides:

- Low-level memory access.

- Direct mapping to system calls.

- Efficient execution.

#### Key characteristics:

- Programs interact with **kernel services** through **system calls**.

- Examples: File management, process control, signal handling, inter-process communication.

#### File I/O in Unix/Linux

File input/output (I/O) in Unix/Linux can be done in two ways:

- Standard I/O** – Using functions like fopen(), fread(), fprintf().

  - High-level, buffered.

Easier for general programming.

**System-level I/O** – Using **system calls** like `open()`, `read()`, `write()`, `close()`.

Low-level, unbuffered.

Directly interacts with the kernel.

Provides more control and efficiency.

We will focus on **system-level file I/O**.

### **Common File I/O System Calls in Unix/Linux**

#### **(a) `open()`**

Used to **open a file** and obtain a **file descriptor**.

#### **Syntax:**

```
#include <fcntl.h>
int open(const char *pathname, int flags, mode_t mode);
```

**pathname** – Path to the file.

**flags** – Access mode (read/write/append) and other options.

**mode** – File permission (used when creating a new file).

#### **Common flags:**

| <b>Flag</b>           | <b>Description</b>                |
|-----------------------|-----------------------------------|
| <code>O_RDONLY</code> | Open file for reading only        |
| <code>O_WRONLY</code> | Open file for writing only        |
| <code>O_RDWR</code>   | Open file for reading and writing |
| <code>O_CREAT</code>  | Create file if it doesn't exist   |
| <code>O_TRUNC</code>  | Truncate file to zero length      |
| <code>O_APPEND</code> | Append data to end of file        |

#### **Example:**

```
int fd = open("file.txt", O_RDWR | O_CREAT, 0644);
if (fd == -1) {
    perror("open");
    return 1;
}
```

Returns a **file descriptor** (integer  $\geq 0$ ) if successful.

Returns -1 on error.

### **(b) read()**

Reads data from a file descriptor into a buffer.

#### **Syntax:**

```
#include <unistd.h>
ssize_t read(int fd, void *buf, size_t count);
```

**fd** – File descriptor obtained from open().

**buf** – Buffer to store data.

**count** – Number of bytes to read.

#### **Example:**

```
char buffer[100];
ssize_t bytesRead = read(fd, buffer, sizeof(buffer));
if (bytesRead == -1) {
    perror("read");
}
```

Returns the number of bytes actually read.

Returns 0 if end-of-file (EOF) is reached.

Returns -1 on error.

### **(c) write()**

Writes data from a buffer to a file descriptor.

#### **Syntax:**

```
#include <unistd.h>
ssize_t write(int fd, const void *buf, size_t count);
```

**fd** – File descriptor.

**buf** – Buffer containing data to write.

**count** – Number of bytes to write.

#### **Example:**

```
char *msg = "Hello, Unix!\n";
ssize_t bytesWritten = write(fd, msg, strlen(msg));
```

```
if (bytesWritten == -1) {
    perror("write");
}
```

Returns the number of bytes written.

Returns -1 on error.

#### **(d) close()**

Closes an open file descriptor.

#### **Syntax:**

```
#include <unistd.h>int close(int fd);
```

**fd** – File descriptor to close.

Frees kernel resources associated with the file.

Always close files after use to avoid leaks.

#### **Example:**

```
if (close(fd) == -1) {
    perror("close");
}
```

#### **Complete Example of File I/O in C**

```
#include <stdio.h>
#include <fcntl.h>
#include <unistd.h>
#include <string.h>
int main() {
    int fd = open("example.txt", O_CREAT | O_WRONLY, 0644);
    if (fd == -1) {
        perror("open");
        return 1;
    }

    char *msg = "Hello, Unix/Linux System Programming!\n";
    if (write(fd, msg, strlen(msg)) == -1) {
        perror("write");
        close(fd);
        return 1;
    }

    close(fd);
}
```

```

// Reading the file
fd = open("example.txt", O_RDONLY);
if (fd == -1) {
    perror("open");
    return 1;
}

char buffer[100];
ssize_t bytesRead = read(fd, buffer, sizeof(buffer)-1);
if (bytesRead == -1) {
    perror("read");
    close(fd);
    return 1;
}

buffer[bytesRead] = '\0'; // Null terminate
printf("File content:\n%s", buffer);

close(fd);
return 0;
}

```

### **Explanation:**

Open a file example.txt for writing (O\_CREAT | O\_WRONLY).

Write a message to the file.

Close the file.

Reopen the file for reading (O\_RDONLY).

Read the contents into a buffer.

Print the content.

Close the file.

### **Process Creation in Unix/Linux**

In Unix/Linux, a **process** is an instance of a program in execution. **Creating processes** is fundamental in system programming. The main system calls involved are:

**fork()** – Create a new process.

**exec()** – Replace the current process image with a new program.

**wait() / waitpid()** – Synchronize parent and child processes.

#### **(a) fork() – Creating a new process**

**fork()** is used to create a **child process**. The child is a copy of the parent, sharing almost everything initially, except:

**Process ID (PID)** is unique.

**Parent PID (PPID)** points to the parent.

File descriptors are shared but maintain separate offsets.

**Syntax:**

```
#include <unistd.h>
pid_t fork(void);
```

**Return value:**

Returns 0 in the **child process**.

Returns the **child PID** in the **parent process**.

Returns -1 on **failure**.

**Example:**

```
#include <stdio.h>
#include <unistd.h>
int main() {
    pid_t pid = fork();

    if (pid < 0) {
        perror("fork failed");
    } else if (pid == 0) {
        printf("Child process, PID = %d\n", getpid());
    } else {
        printf("Parent process, PID = %d, Child PID = %d\n", getpid(), pid);
    }

    return 0;
}
```

**Explanation:**

fork() duplicates the process.

Both parent and child execute the code **after fork()**.

PID helps distinguish parent from child.

**(b) exec() – Replacing a process image**

The **exec()** family of functions replaces the **current process image** with a **new program**.

There are several versions: `execl()`, `execp()`, `execv()`, etc.

### Syntax (example with `execl()`):

```
#include <unistd.h>
int execl(const char *path, const char *arg0, ..., NULL);
```

### Example:

```
#include <stdio.h>
#include <unistd.h>
int main() {
    pid_t pid = fork();

    if (pid == 0) {
        // Child executes a new program
        execl("/bin/ls", "ls", "-l", NULL);
        perror("execl failed"); // Only executes if execl fails
    } else {
        wait(NULL); // Parent waits for child
        printf("Parent process finished\n");
    }

    return 0;
}
```

### Explanation:

Child process runs `ls -l`.

Parent waits until child finishes.

### (c) `wait()` – Synchronizing parent and child

`wait()` makes the **parent process wait** until any child terminates.

### Syntax:

```
#include <sys/wait.h>
pid_t wait(int *status);
```

Returns **child PID** on success.

Returns -1 if no children exist.

`status` contains termination information.

### Example:

```

int status;
pid_t child_pid = wait(&status);
if (WIFEXITED(status)) {
    printf("Child exited with code %d\n", WEXITSTATUS(status));
}

```

## Inter-Process Communication (IPC)

Processes can communicate using **shared memory, signals, pipes, FIFOs, and sockets**. We'll focus on **pipes** and **FIFO**.

### (a) Pipes – Anonymous Communication

Pipes are **unidirectional communication channels** between **related processes** (parent-child).

Data written by **write()** can be read by **read()**.

Created using `pipe()`.

#### Syntax:

```

#include <unistd.h>
int pipe(int fd[2]);

    fd[0] → read end

    fd[1] → write end

```

#### Example:

```

#include <stdio.h>
#include <unistd.h>
#include <string.h>
int main() {
    int fd[2];
    pipe(fd); // Create pipe
    pid_t pid = fork();

    if (pid == 0) {
        // Child writes to pipe
        close(fd[0]); // Close read end
        char msg[] = "Hello from child";
        write(fd[1], msg, strlen(msg)+1);
        close(fd[1]);
    } else {
        // Parent reads from pipe
        close(fd[1]); // Close write end
        char buffer[100];
        read(fd[0], buffer, sizeof(buffer));
    }
}

```

```

    printf("Parent received: %s\n", buffer);
    close(fd[0]);
}

return 0;
}

```

### Explanation:

Pipe allows **one-way communication** from child → parent.

Always close unused ends to avoid deadlocks.

### (b) FIFO – Named Pipes

FIFO (First In First Out) is a **named pipe**, can communicate between **unrelated processes**.

Created using `mkfifo()` or `mknod()`.

Accessed like a regular file using `open()`, `read()`, `write()`, `close()`.

### Syntax:

```

#include <sys/types.h>
#include <sys/stat.h>
int mkfifo(const char *pathname, mode_t mode);

```

### Example:

```

#include <stdio.h>
#include <fcntl.h>
#include <unistd.h>
#include <sys/stat.h>
int main() {
    char *fifo = "/tmp/myfifo";
    mkfifo(fifo, 0666); // Create FIFO

    pid_t pid = fork();
    if (pid == 0) {
        // Child writes
        int fd = open(fifo, O_WRONLY);
        write(fd, "Message through FIFO", 20);
        close(fd);
    } else {
        // Parent reads
        int fd = open(fifo, O_RDONLY);
        char buffer[100];
        read(fd, buffer, sizeof(buffer));
        printf("Parent received: %s\n", buffer);
    }
}

```

```
    close(fd);
}

unlink(fifo); // Delete FIFO
return 0;
}
```

### **Explanation:**

FIFO persists in the filesystem.

Can be used by **unrelated processes** for communication.

unlink() removes the FIFO.

### **Signal Handling in Unix/Linux**

**Signals** are **asynchronous notifications** sent to a process to notify it of events like:

Interrupt from keyboard (Ctrl+C)

Illegal memory access

Child termination

Timer expiry

### **Common signals:**

| <b>Signal</b> | <b>Description</b>              |
|---------------|---------------------------------|
| SIGINT        | Interrupt (Ctrl+C)              |
| SIGKILL       | Kill process (cannot be caught) |
| SIGTERM       | Termination request             |
| SIGCHLD       | Child process terminated        |
| SIGALRM       | Timer alarm                     |

### **Handling signals**

We can **catch signals** using signal() or sigaction().

### **Syntax:**

```
#include <signal.h>
void signal(int sig, void (*handler)(int));
```

### **Example – Catching Ctrl+C:**

```
#include <stdio.h>
```

```

#include <signal.h>
#include <unistd.h>
void handle_sigint(int sig) {
    printf("\nCaught signal %d. Exiting safely.\n", sig);
    _exit(0);
}
int main() {
    signal(SIGINT, handle_sigint);
    while (1) {
        printf("Running...\n");
        sleep(1);
    }
    return 0;
}

```

### Explanation:

When the user presses Ctrl+C, handle\_sigint() executes.

Provides safe termination or cleanup.

### POSIX Threads (pthreads)

**Threads** are lightweight processes that share the **same memory space**. POSIX threads allow **multithreading in C**.

### Key Functions:

**pthread\_create()** – Create a new thread

```

#include <pthread.h>
int pthread_create(pthread_t *thread, const pthread_attr_t *attr,
    void *(*start_routine)(void*), void *arg);

```

thread → thread ID

start\_routine → function executed by thread

arg → argument to thread function

Returns 0 on success

**pthread\_join()** – Wait for a thread to finish

```

#include <pthread.h>
int pthread_join(pthread_t thread, void **retval);

```

Blocks the calling thread until the target thread terminates.

### Example – Thread creation and join

```

#include <stdio.h>
#include <pthread.h>
void* print_message(void* arg) {
    char* msg = (char*) arg;
    printf("%s\n", msg);
    return NULL;
}
int main() {
    pthread_t tid;
    pthread_create(&tid, NULL, print_message, "Hello from thread");
    pthread_join(tid, NULL);
    printf("Thread finished execution\n");
    return 0;
}

```

### **Explanation:**

Thread prints a message.

Main thread waits for it to complete.

### **Case Studies**

#### **(a) Background Processes**

Background processes run **without blocking the terminal**.

Achieved by appending & in shell or using fork() + setsid() in C.

### **Example:**

```

pid_t pid = fork();
if (pid == 0) {
    // Child process running in background
    printf("Child running in background\n");
    sleep(10);
    _exit(0);
} else {
    printf("Parent continues without waiting\n");
}

```

Parent doesn't wait; child runs independently.

#### **(b) Daemon Process Creation**

Daemon: **Background service process** with no controlling terminal.

Steps to create a daemon:

fork() and exit parent.

setsid() → create a new session.

Optional second fork() to prevent reacquisition of terminal.

Close standard file descriptors (stdin, stdout, stderr).

Run background tasks.

### Example Skeleton:

```
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>
#include <fcntl.h>
#include <sys/stat.h>
int main() {
    pid_t pid = fork();
    if (pid > 0) exit(0); // Parent exits

    setsid(); // New session
    chdir("/"); // Change working directory

    // Close standard file descriptors
    close(STDIN_FILENO);
    close(STDOUT_FILENO);
    close(STDERR_FILENO);

    // Daemon main loop
    while (1) {
        // Background tasks
        sleep(10);
    }
    return 0;
}
```

### (c) Mini Shell Implementation

A **mini shell** allows basic command execution. Core steps:

Display prompt.

Read input (fgets()).

Parse command and arguments.

fork() → create child process.

execvp() → execute command in child.

wait() → parent waits for child.

### Example:

```
#include <stdio.h>
#include <unistd.h>
#include <stdlib.h>
#include <string.h>
#include <sys/wait.h>
int main() {
    char cmd[100];
    while (1) {
        printf("mini-shell> ");
        fgets(cmd, sizeof(cmd), stdin);
        cmd[strcspn(cmd, "\n")] = 0; // Remove newline

        if (strcmp(cmd, "exit") == 0) break;

        pid_t pid = fork();
        if (pid == 0) {
            char *args[] = {cmd, NULL};
            execvp(args[0], args);
            perror("exec failed");
            exit(1);
        } else {
            wait(NULL);
        }
    }
    return 0;
}
```

Supports simple commands like ls, pwd.

Extensible for pipes, redirection, and background processes.